



DIGITAL ACCESS TO  
SCHOLARSHIP AT HARVARD  
DASH.HARVARD.EDU

HARVARD  
LIBRARY



# Essays in Decision Theory

## Citation

Ridout, Sarah. 2021. Essays in Decision Theory. Doctoral dissertation, Harvard University Graduate School of Arts and Sciences.

## Link

<https://nrs.harvard.edu/URN-3:HUL.INSTREPOS:37368182>

## Terms of use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material (LAA), as set forth at

<https://harvardwiki.atlassian.net/wiki/external/NGY5NDE4ZjgzNTc5NDQzMGIzZWZhMGFIOWI2M2EwYTg>

## Accessibility

<https://accessibility.huit.harvard.edu/digital-accessibility-policy>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#)

HARVARD UNIVERSITY  
Graduate School of Arts and Sciences




DISSERTATION ACCEPTANCE CERTIFICATE


The undersigned, appointed by the  
Department of Economics  
have examined a dissertation entitled  
"Essays in Decision Theory"

presented by Sarah Ridout

candidate for the degree of Doctor of Philosophy and hereby  
certify that it is worthy of acceptance.

Signature   
Typed name: Prof. Tomasz Strzalecki

Signature   
Typed name: Prof. Matthew Rabin

Signature   
Typed name: Prof. Jerry Green

Signature   
Typed name: Prof. Shengwu Li

Signature \_\_\_\_\_  
Typed name: Prof.

Date: March 15, 2021

# Essays in Decision Theory

A DISSERTATION PRESENTED  
BY  
SARAH RIDOUT  
TO  
THE DEPARTMENT OF ECONOMICS

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
DOCTOR OF PHILOSOPHY  
IN THE SUBJECT OF  
ECONOMICS

HARVARD UNIVERSITY  
CAMBRIDGE, MASSACHUSETTS  
MARCH 2021

©2021 – SARAH RIDOUT  
ALL RIGHTS RESERVED.

## Essays in Decision Theory

### ABSTRACT

This dissertation comprises three independent chapters that span topics in decision theory, including ethical decision-making, learning from others, and anticipatory emotions.

In Chapter 1, I model decision-making constrained by ethics, duty, law, or other virtues or principles. In addition to a true preference, the decision maker (DM) is characterized by a set of preferences that he considers justifiable. In each choice setting, the DM maximizes his true preference over the subset of alternatives that maximize at least one of the justifiable preferences. The justification model unites a broad class of empirical work on distributional preferences, discrimination, corruption, philanthropy, and other domains. I provide simple axiomatic characterizations of several variants of the justification model as well as practical tools for identifying true preferences and justifications from choice behavior. I show that identification is partial in general, but full identification can be achieved by moving to between-subject data and imposing some additional structure on true preferences and justifications. Moving to between-subject data also eliminates the consistency motives that may arise in within-subject experiments. I extend the between-subject justification model to information choice and relate its predictions to the “moral wiggle room” literature.

In Chapter 2, I model a decision maker who cares about the relative frequencies with which different items are selected by peers. The decision maker has a Bayesian prior on these relative frequencies, which he updates as he collects more data. The model provides a unified framework for interpreting recent empirical evidence on peer effects in varied domains. I provide an axiomatic characterization of the model and show how different specializations of the representation generate

different responses to peer information. I then extend the model to distinguish peer effects caused by inference about the best alternative from peer effects caused by concerns about relative status.

In chapter 3, I model a decision maker who fantasizes or worries about future risks by reweighing probabilities in an optimistic or pessimistic manner. The degree of optimism or pessimism toward a given lottery may depend on the lottery itself, its context, and the period in which it takes place. The representation resembles a standard recursive expected utility model, but with distorted probabilities instead of physical probabilities. The model is well suited to applications because it (a) has a flexible domain (infinite tree of continuous and/or discrete distributions on a bounded or unbounded interval), (b) admits simple comparative statics, and (c) can generate tractable parametric distorted distributions.

# Contents

TITLE PAGE	i
COPYRIGHT	ii
ABSTRACT	iii
TABLE OF CONTENTS	v
DEDICATION	viii
ACKNOWLEDGEMENTS	ix
<b>I A MODEL OF JUSTIFICATION</b>	<b>I</b>
1.1 Introduction . . . . .	1
1.2 Empirical motivation . . . . .	10
1.3 Primary preference observable . . . . .	14
1.4 Primary preference unobservable . . . . .	31
1.5 Random justification model . . . . .	40

1.6	Related models . . . . .	56
1.7	Conclusion . . . . .	61
<b>2</b>	<b>A DYNAMIC MODEL OF PEER EFFECTS</b>	<b>63</b>
2.1	Introduction . . . . .	63
2.2	Related Literature . . . . .	66
2.3	Model and Characterization . . . . .	70
2.4	Types of peer effects . . . . .	79
2.5	Interpreting peer effects: pure learning vs. social image . . . . .	85
2.6	Conclusion . . . . .	91
<b>3</b>	<b>FANTASIES AND WORRIES AS DISTORTED PROBABILITIES</b>	<b>92</b>
3.1	Introduction . . . . .	92
3.2	Model . . . . .	96
3.3	Parametric Special Cases . . . . .	103
3.4	Comparative statics . . . . .	107
3.5	Axioms . . . . .	111
3.6	Extensions . . . . .	116
3.7	Future Work . . . . .	120
<b>APPENDIX A APPENDIX TO CHAPTER 1</b>		<b>122</b>
A.1	Proofs of results in text . . . . .	122
A.2	Model variants . . . . .	200
<b>APPENDIX B APPENDIX TO CHAPTER 2</b>		<b>211</b>
B.1	Proofs of results in text . . . . .	211



APPENDIX C	APPENDIX TO CHAPTER 3	230
C.1	Domain of choice . . . . .	230
C.2	Probability Distortions . . . . .	233
C.3	Proofs of results in text . . . . .	234
REFERENCES		259

TO MY PARENTS, JENNIFER ANN MULHOLLAN AND ROBERT MASTIN RIDOUT, JR.

# Acknowledgments

I am afraid that my gratitude to the following individuals outpaces my expressive abilities. Within economics, I am particularly grateful to...

...Tomasz Strzalecki, for a great many things, including: an engaging introduction to decision theory, several years of patient guidance on research and other aspects of the profession, and much-needed encouragement during the job market. Tomasz set an excellent example as a researcher, teacher, and advisor. I will strive to follow that example in the years to come.

...Matthew Rabin, for pushing me to explore new directions and new ideas, for reminding me to consider reality as much as tractability and elegance, and for making me laugh a great deal.

...Jerry Green, for sharing his immense knowledge of economic theory, and for encouraging me to pursue research in that area.

...Shengwu Li, for volunteering exceptionally helpful and practical advice, and for building a community of theorists in the department.

...Ian Martin, for offering invaluable help, advice, and encouragement during the PhD application process.

...James Forder, for setting me on the path to becoming an academic economist, and for teaching

me to write in a precise and professional manner.

In my personal life, I am most grateful to...

...my husband, Deepsagar Lambor, for supporting me without question, for tolerating my occasionally prickly disposition, and for filling the past three years with unexpected happiness.

...my mother, Jennifer Mulhollan, for prioritizing my education above almost everything else, and for modeling grit and determination at every turn.

...my father, Robert Ridout, for showing me the joy of learning for learning's sake, and for supporting me equally through both successes and failures.

...my best friends, Hanna Fanous and Katherine Lin, for fifteen years of sisterhood, and for never letting me take myself too seriously.

...my pitbull-basset mix, Kaizer, for making me take breaks at regular intervals, for fighting off the mice in my apartment, and for making me smile through the toughest parts of the PhD.

# 1

## A Model of Justification

### 1.1 INTRODUCTION

Consequential decisions are often subject to the demands of principle.<sup>1</sup> At the same time, complex situations provide some freedom to interpret a general principle or to decide which principle to

---

<sup>1</sup>For helpful comments and suggestions, I am indebted to Christine Exley, Ed Glaeser, Ben Golub, Jerry Green, Yoram Halevy, Peter Klibanoff, David Laibson, Shengwu Li, Yusufcan Masatlioglu, Efe Ok, Pietro Ortoleva, Fernando Payró Chew, Matthew Rabin, Tomasz Strzalecki, and participants at an online decision theory conference in July 2020.

prioritize. For instance:

- A politician is charged with pursuing justice, but has some freedom in deciding which notion of justice to apply.
- A judge is charged with applying the law, but has some freedom in interpreting legal language and resolving conflicting precedents.
- A hiring manager is charged with choosing the most qualified candidate, but has some freedom in determining which qualifications matter most.

To capture these situations and others, this paper presents a two-tier model of preference maximization in which a “primary” preference resolves conflicts among “justifiable” preferences. The set of justifiable preferences captures the DM’s notion of acceptable behavior in the domain at hand, while the primary preference captures the DM’s own inclinations. On any choice set, the DM maximizes his primary preference over those alternatives that maximize at least one of the justifiable preferences. He thereby pursues his own objectives without doing anything obviously objectionable.

For instance:

- The politician selects among the policies that could be implemented by a disinterested seeker of justice.
- The judge selects among the rulings that could be made by an impartial instrument of the law.
- The hiring manager selects among the candidates who could be ranked highest by an unbiased appraiser of talent.

The reader may wonder how an analyst limited to choice data could possibly disentangle the DM’s underlying inclinations from his notion of acceptable behavior. Disentangling the two is indeed impossible if the DM never appeals to more than one justification. In that case, behavior is consistent

with standard preference maximization, so a more general model is not required. The justification model becomes useful when the DM finds it optimal to appeal to different justifications in different situations, creating telltale inconsistencies with preference maximization.

The expansive “moral wiggle-room” literature provides evidence of these inconsistencies. We review this literature in Section 1.2, but provide one instructive example here. [Exley \(2016\)](#) offered subjects binary decisions involving (possibly random) payments to themselves and/or donations to a charity. She found that subjects’ treatment of risk shifted in response to their selfish interests, resulting in cyclic choice patterns. For instance, subjects confronted with pairs of outcomes in

$a =$  experimenter pays \$2.50 to DM

$b =$  50% chance experimenter donates \$10 to charity

$d =$  experimenter donates \$4 to charity

might choose

$$a = c(\{a, b\}) \quad b = c(\{b, d\}) \quad d = c(\{a, d\}).$$

This pattern is inconsistent with preference maximization, but it can be generated by a selfish primary preference and a set of generous justifications with different risk attitudes. Intuitively, the DM chooses  $a$  over  $b$ , but not  $a$  over  $d$ , because choosing  $a$  over  $b$  can be attributed to risk aversion rather than selfishness. To see this in more detail, suppose the primary preference  $\succsim$  is represented by

$$\mathbb{E} [3(\$ \text{ to DM}) + (\$ \text{ to charity})],$$

while the set of justifiable preferences  $\mathcal{M}$  is represented by

$$\{\mathbb{E} [1.5(\$ \text{ to DM})^x + (\$ \text{ to charity})^x] : x \in [0.5, 2]\}.$$

The primary preference has  $a \succ b \succ d$ . The DM can justify choosing  $a$  over  $b$  because  $a \succ_m b$  for sufficiently risk-averse  $\succ_m \in \mathcal{M}$ . He can justify choosing  $b$  over  $d$  because  $b \succ_m d$  for the risk-loving  $\succ_m \in \mathcal{M}$ . However, he cannot choose  $a$  over  $d$  because  $d \succ_m a$  for all  $\succ_m \in \mathcal{M}$ .

The justification model is intended as a general framework that can unite a wide variety of choice settings. To focus on the conflict between primary preferences and justifications, the model abstracts away from other interesting features of justifying behavior. First, it does not offer a theory of justice, rationality, or any other virtue or principle. While the analyst can learn about the DM's notion of acceptable behavior ex post, the model does not impose ex ante restrictions on that notion. An extension of the model allows for restrictions on the justifiable preferences, but leaves the analyst to determine what those restrictions should be.

Second, the model is agnostic about the DM's motivations for limiting himself to justifiable alternatives. One interpretation is that the DM is pretending to be a better person than he actually is. By limiting himself to alternatives that good people could select, he conceals his true inclinations from anyone observing his choice. The other interpretation is that the DM wants to remain within the bounds imposed by principle. He is not pretending to be a particularly virtuous person, but simply refraining from circumscribed actions. The reader may object that the latter interpretation is more plausible than the former, as a DM who appeals to different justifications in different situations reveals that his primary preference is not fully aligned with any of them. The only way to conceal one's primary preference, the objection continues, is to maximize a single justification across all decisions. This objection is most forceful when the DM must make several closely connected decisions in rapid succession before the same audience. If the connection between decisions is obscured, the time interval is long, or the set of observers varies, the DM may not attempt to maintain consistency across decisions. In any case, this objection applies only to versions of the justification model that rely on within-subject data. The random justification model of Section 1.5 uses between-subject data, so it accommodates both interpretations equally well.



### 1.1.1 OVERVIEW OF RESULTS

The main results of this paper fall into two categories. First, the paper provides full behavioral characterizations of several variants of the justification model. With the exception of continuity conditions, which appear only in Section 1.3.3, all the axioms used in these characterizations are simple enough for practical application. For instance, they can be used to test whether apparently altruistic choices are consistent with a selfish primary preference, whether apparently irrational choices are consistent with *any* primary preference, and whether choices that fluctuate with the decision environment are consistent with a stable primary preference. Second, the paper provides tools for identifying both primary and justifiable preferences from choice data. For instance, the identification results can be used to disentangle the DM's true level of generosity from the minimal level he thinks acceptable, or to determine the range of allocation rules the decision-maker considers fair. In most cases, the identification results are byproducts of the behavioral characterizations, so the two groups of results are complementary. While the identification procedures are not always guaranteed to deliver full identification, they deliver all of the information contained in the data. Like the axioms, they rely on simple patterns of choice and remain useful on limited datasets.

The Justification model with Observable primary preference (JO) is introduced in Section 1.3. As the name suggests, this version of the model takes the primary preference as well as the choice correspondence as primitive. There are both expository and practical reasons for this assumption. On the expository side, the results for the Justification model with Unobservable primary preference (JU) build neatly on the JO results. On the practical side, the analyst may wish to test a particular candidate for the primary preference. Some choice settings may suggest a natural candidate (e.g. monetary self-interest). Alternatively, the analyst may be able to elicit a candidate by conducting a separate treatment in which subjects face less pressure to justify their decisions (e.g. because decisions are anonymous or implemented by someone else). Some experiments lend support to this idea.

For instance, [Hamman et al. \(2010\)](#) found that subjects made blatantly selfish decisions when those decisions were carried out by an intermediary.

Theorem 1 is the representation result for JO. Given the primary preference, it provides necessary and sufficient conditions for an individual DM's behavior to be consistent with the justification model. The key axiom is Irrelevance of Unjustifiable Alternatives (IUA). Say that it is “unjustifiable” to choose alternative  $a$  from set  $A$  if  $a$  is not selected from  $A$ , but the primary preference likes  $a$  at least as much as everything that is selected. IUA says that  $a$  is irrelevant for choice on any superset of  $A$  if it is unjustifiable to choose  $a$  from  $A$ . The proof of Theorem 1 shows that a preference belongs to the maximal set of justifiable preferences if and only if it does not sanction any unjustifiable choices: the preference does not rank  $a$  above the rest of  $A$  if it is unjustifiable to choose  $a$  from  $A$ . Therefore, the unjustifiable choices are fully informative about the DM's notion of acceptable behavior.

The remainder of Section 1.3 characterizes three extensions of JO. All three maintain the primary preference as a primitive; this is not dropped until Section 1.4. Section 1.3.2 allows the analyst to restrict the universe of possible justifications, ruling out preferences she considers obviously unjustifiable. The analyst's preconceptions are captured by an (observable) asymmetric and transitive relation on the domain that all justifiable preferences are required to respect. In addition to stochastic or Pareto dominance, an appropriately chosen relation can capture natural ethical requirements such as impartiality or nondiscrimination. Proposition 1 shows that a simple strengthening of IUA is necessary and sufficient for a representation in which all justifications respect the desired restrictions.

Section 1.3.3 extends JO to require continuity of the justifications. Continuity is desirable because it ensures the existence of utility representations, but it raises significant technical difficulties. The proof of Theorem 2, the representation theorem for the continuous case, addresses these difficulties with a continuous version of the Szpilrajn Extension Theorem from [Herden and Pallack](#)

(2002).

The final extension of JO, in Section 1.3.4, takes a simpler route to utility representations. It restricts the domain to the set of lotteries on a finite prize space, and requires both primary and justifiable preferences to take an expected-utility form. Theorem 3 provides a behavioral characterization for the EU extension of JO. A nice feature of the EU extension is that behavior on an arbitrary choice set  $\mathcal{A}$  is entirely pinned down by behavior on binary subsets of the convex hull of  $\mathcal{A}$ . This makes the EU extension particularly tractable, so it can be used as a building block for other models. In fact, the random justification model in Section 1.5 is built up from the EU extension of JO.

Section 1.4 dispenses with the primary preference as a primitive and recovers it as a component of the representation. It provides two complementary characterizations for the Justification model when the primary preference is Unobservable (JU). These characterizations demonstrate that the justification model imposes substantial restrictions on behavior even if the primary preference is a “free parameter.” The first characterization, Proposition 2, is a simple corollary to Theorem 1. As a byproduct, Proposition 2 provides a simple procedure for constructing the full set of primary preferences consistent with choice data. The second characterization, Theorem 4, is structurally similar to Theorem 1. Both results say that the justification model is characterized by irrelevance of unjustifiable alternatives, and both proofs show that the maximal set of justifications is precisely the set of preferences that do not sanction any unjustifiable choices. The only difference is the procedure for classifying some choices as “unjustifiable.” Theorem 1 relies on the primary preference, while Theorem 4 relies only on choice patterns. The required choice patterns are simple and easy to spot. By picking out these patterns, the analyst extracts all of the information contained in the data about the DM’s notion of acceptable behavior.

Section 1.4.3 extends JU to account for evidence that justifying behavior is sensitive to features of the decision environment, such as anonymity (Charness and Gneezy, 2008; Franzen and Pointner, 2012) or whether anyone else stands to be disappointed by one’s decision (Dana et al., 2006). It is

plausible that changes in behavior across environments are driven by variations in the standard for acceptable behavior rather than shifting inclinations. Proposition 4, which builds on Theorem 4, allows the analyst to test this conjecture. It provides necessary and sufficient conditions for a pair of choice functions to be consistent with the same primary preference, but nested sets of justifications. A special case of Proposition 4, spelled out in Corollary 4, ties neatly back to Theorem 1. If  $(c_1, c_2)$  satisfies the conditions in Proposition 4, and  $c_1$  is consistent with preference maximization, then  $c_2$  satisfies IUA conditional on the preference maximized by  $c_1$ . Thus, it does not matter whether the analyst uses private choice (for example) to elicit the primary preference and subsequently applies Theorem 1 to public choice, or whether she applies Proposition 4 to public and private choice together.

Section 1.5 is motivated by the concern that some DMs may limit themselves to justifications consistent with their previous choices. As mentioned above, DMs who wish to conceal their primary preferences have a strong incentive to maintain consistency in justifications, since appeals to incompatible justifications produce the choice patterns that drive identification. Section 1.5 obviates consistency concerns by studying a large population of DMs, each of whom makes only one choice (as in a between-subject experiment). The primitive is a stochastic choice function, and the representation is a pair of distributions: one over primary preferences, and the other over sets of justifications. As in Section 1.3.4, the domain is a set of lotteries, and both primary and justifiable preferences have an EU form. The model permits arbitrary heterogeneity in primary preferences within the EU framework, but requires the sets of justifications to exist on a spectrum from strict to permissive. Intuitively, DMs within the same population are required to have similar principles, but may differ in their commitment to those principles.

Just as choice in the deterministic justification model violates the Weak Axiom of Revealed Preference (WARP), choice in the Random Justification model (RJ) violates Regularity, the stochastic analogue of WARP. Regularity says that the probability of selecting any given item must weakly de-

cline when more items are added to the choice set. Proposition 5 shows that Regularity violations are ubiquitous in RJ. If it is sometimes unjustifiable to choose  $q$  over  $p$ , then it is always possible to find nested menus with  $p$  and  $q$  in their intersection such that the probability of choosing  $p$  is larger on the larger menu.

The main result for RJ is Theorem 5, which establishes uniqueness of the random justification representation. This result is important because it shows that an analyst with enough data can separately identify the distribution of primary preferences and the distribution of sets of justifications, even though she cannot tell whether any *particular* decision in her data set was driven by inclination or principle. The proof of Theorem 5 provides an explicit identification procedure. It shows that the analyst can fully identify the justification distribution using choice sets of manageable size (no more than four elements).

In both the deterministic and random justification models, the DM's only degree of freedom comes from disagreement between justifications. A range of experiments, reviewed in Section 1.2, suggest that richer choice settings can provide additional degrees of freedom. Information avoidance is a major area of study in this literature, so Section 1.5.2 extends RJ to a simple information choice setting. Propositions 6 and 7 together show that RJ can predict information avoidance, but only if the standards that govern information acquisition are weaker than the standards that govern choice between final outcomes. If the same standards govern both types of choices, the DM will either feel unable to avoid information, or find it unprofitable to do so. The simplest way to generate information avoidance is to assume that information choice does not need to be justified at all. This assumption is adopted in the rest of Section 1.5.2, but relaxed in Appendix A.2.4.

Although Proposition 7 predicts information avoidance, it does not imply that information provision is entirely futile. Some (but not all) DMs with less-than-virtuous primary preferences will voluntarily acquire information that induces them to behave more virtuously. This is not because they feel obligated to become informed, but because information reduces the risk of sacrificing for

nothing. For instance, a selfish DM might research a donation opportunity to make sure that his money is well spent when he feels compelled to donate. More generally, Proposition 7 shows that information can be used to bring behavior into alignment with principle, even if information choice does not have to be justified and remaining ignorant is always an option. The qualifier “judiciously chosen” is important, though. The final result of the paper, Proposition 8, shows that injudiciously chosen information can worsen behavior by allowing DMs with less-than-virtuous primary preferences to exploit disagreement between justifications. For instance, a biased hiring manager might treat disfavored candidates’ resumes as a source of justifications for rejecting them.

Section 1.6 relates the justification model to existing theoretical work. Formally, the justification model is a generalization of the models of willpower in [Masatlioglu et al. \(2020\)](#) and dynamic choice in [Strotz \(1955\)](#), as formalized by [Gul and Pesendorfer \(2005\)](#). It is a special case of the models of limited attention in [Masatlioglu et al. \(2012\)](#) and of rationalization in [Cherepanov et al. \(2013b\)](#). The latter deserves special attention because it is motivated by some of the same evidence as the justification model, and the two models are interpreted along similar lines. The key difference is that justifications are preferences, while rationales are completely unstructured binary relations. Thus, justifiers can be viewed as pooling with more virtuous types, while rationalizers cannot. Moreover, the additional structure of the justification model leads to strictly stronger identification than the rationalization model on any dataset that violates WARP.

## 1.2 EMPIRICAL MOTIVATION

Although the formal models introduced in this paper can be applied to non-moral domains, the vast majority of relevant empirical research covers moral decision-making, particularly tradeoffs between oneself and others. This work can broadly be divided into two strands. The first strand shows that subjects appeal to different principles, or pretend to have different tastes, when their

personal interests change. The second strand provides evidence that subjects attempt to evade the constraints imposed by principle, e.g. by avoiding information about the effects of their actions.

Loewenstein et al. (1993) is a classic reference in the first strand. Each subject was assigned the role of plaintiff or defendant in a tort case involving a motorcycle accident. After reviewing the case material, each subject reported two values: the amount she expected the judge to award to the plaintiff, and the amount she considered fair. On average, subjects assigned to the role of plaintiff expected the judge to award \$39,000, and considered \$37,000 to be fair. Subjects assigned to the role of defendant expected \$24,000, and considered \$19,000 to be fair. Loewenstein et al. (1993) interpreted the gap between plaintiff and defendant as evidence of “self-serving assessments of fairness.”<sup>2</sup>

Rodriguez-Lara and Moreno-Garrido (2012) studied a similar tension between disinterested allocation rules and personal gain. In an initial phase of the experiment, subjects earned money by completing a multiple-choice test. After the test was complete, each subject’s correct answers were converted into money at a random “wage,” which was fixed within-subject but could differ between-subject. Finally, subjects were paired up, and one member of each pair (the “dictator”) allocated the money earned by the pair. Rodriguez-Lara and Moreno-Garrido (2012) found that dictators’ choices were well explained by self-serving choice between three impartial allocation principles: egalitarian, accuracy-based, and earnings-based. Dictators who had low wages and accuracy tended to favor an egalitarian division, while dictators who had high wages and accuracy tended to favor division on the basis of earnings. Dictators who had low wages but high accuracy favored division on the basis of accuracy alone (correcting for the unequal “wages” assigned by the experimenter). We return to Rodriguez-Lara and Moreno-Garrido (2012) in Example 1 of Section 1.3.1, and to the notion of impartiality at the end of Section 1.3.2.

---

<sup>2</sup>Self-serving assessments had real effects on behavior. When plaintiffs and defendants were paired up to negotiate a settlement, pairs with more divergent assessments were more likely to reach an impasse.

Norton et al. (2004) studied hiring decisions in the presence of gender bias. Subjects were asked to choose between a male and a female candidate for a traditionally male role. In one treatment, the male had more education and the female had more experience; in the other, these attributes were flipped. Subjects chose the male candidate a majority of the time in both cases. When asked to explain their decisions, very few subjects mentioned gender, instead citing the attribute (education or experience) in which their preferred candidate was superior. The proportion of subjects who said that they considered education more important than experience dropped from 50% when the male was more educated to less than 25% when the female was more educated. We return to Norton et al. (2004) at the end of Section 1.5.2.

Gneezy et al. (2019) took a different approach to showing that principles shift with self-interest. They conjectured that subjects who had already thought about or expressed their principles would find it harder to reshape them when their interests changed. To test this, they placed subjects in the role of investment advisors. Each advisor chose an asset to recommend to another subject, who was not informed about any of the assets. The key feature of the experiment was a small commission for recommending one of the assets. In line with the self-deception hypothesis, Gneezy et al. (2019) found that the bribe mattered if and only if two conditions were met. First, advisors had to learn about the bribe before they had a chance to decide which asset was best. Second, there had to be some ambiguity in the task: advisors did not accept a bribe to recommend a strictly dominated asset. Gneezy et al. (2020) found similar results for a different task, in which subjects were told to select the funniest joke from a set of jokes written by others.

Haisley and Weber (2010) tested the same idea in a different setup. Subjects faced a choice between an equitable allocation that gave moderate payments to the subject and a stranger, and an inequitable allocation that gave a large prize to the subject and a chance of a small prize to the stranger. As in Gneezy et al. (2019), ambiguity increased the scope for selfishness. Subjects were much more likely to choose selfishly when the chance of the small prize was uncertain rather than fixed at 50%.



Also as in [Gneezy et al. \(2019\)](#), subjects were constrained by evaluations they had already made. Some subjects faced a choice between a risky lottery and ambiguous prospect earlier in the experiment. These subjects did not choose more selfishly in the ambiguous case. [Haisley and Weber \(2010\)](#) concluded that these subjects were unable to inflate the value of the ambiguous prospect because the vast majority of them had already expressed ambiguity aversion. We return to the consistency motives demonstrated in [Haisley and Weber \(2010\)](#) and [Gneezy et al. \(2019\)](#) at the beginning of Section 1.5, which motivates the random justification model.

We now proceed to the second strand of the literature, which shows that subjects take advantage of opportunities to evade the constraints imposed by principle. Most papers in this literature study information avoidance. The classic example is [Dana et al. \(2007\)](#), in which about 40% of subjects blatantly avoided learning how their choice would affect another subject's payoffs. More recent papers extend these findings to a wide variety of choice settings. [Kajackaite \(2015\)](#) found that one-third of subjects avoided learning how their actions would affect contributions to a negatively perceived lobbying organization. [Serra-Garcia and Szech \(2019\)](#) found that a majority of subjects chose not to learn about an opportunity to donate to a charity, and that almost half paid to avoid information. [Woolley and Risen \(2018\)](#) found that a majority of subjects preferred not to learn the calorie content of a tempting dessert, or the monetary bonus for a boring task. [Ehrich and Irwin \(2005\)](#), [d'Adda et al. \(2018\)](#), and [Freddi \(2017\)](#) found evidence of information avoidance about the ethical attributes of consumer goods, the environmental impacts of air conditioning, and the refugee crisis, respectively. On the other hand, [Fong and Oberholzer-Gee \(2011\)](#) found that about a third of subjects facing a self-other allocation decision paid a substantial fee to learn whether the other subject was "deserving." We return to [Fong and Oberholzer-Gee \(2011\)](#), and to information choice more broadly, in Section 1.5.2.

A few papers study other ways in which subjects avoid feeling compelled to behave morally. [Dana et al. \(2006\)](#) found that about a third of subjects paid a fee to exit a dictator game. From a purely

financial point of view, exiting is a dominated choice. It makes sense only for dictators who would feel compelled to give more than they would like if they remained in the game. Exiting is attractive to these dictators because the recipient does not observe the exit decision, so the sense of obligation governing this decision is presumably weaker than the sense of obligation governing the game itself. In a related field experiment, [Andreoni et al. \(2017\)](#) found that subjects went out of their way to avoid being asked to donate.

[Hamman et al. \(2010\)](#) studied delegation as a way to avoid responsibility. Each subject had the opportunity to hire an agent to make a self-other allocation on the subject’s behalf. Subjects typically favored agents who had a record of sharing very little, or who announced the intention to share very little. This drove down sharing relative to a control condition in which subjects did not have the opportunity to delegate. We return to the effects of the choice setting on standards for acceptable behavior in Section 1.4.3.

### 1.3 PRIMARY PREFERENCE OBSERVABLE

This section introduces the Justification model with Observable primary preference (JO). Formally, JO does not place any restrictions on the domain of choice,  $\mathcal{A}$ . That said, JO is primarily intended for choice settings with ties to ethics, virtue or law. JO delivers interesting predictions on domains in which the “right” choice is not always obvious, but some choices are outright unacceptable.

JO has two primitives. The first is the primary preference  $\succsim$ , which is a complete and transitive relation on  $\mathcal{A}$ . Intuitively, it is what the DM would choose in the absence of any need to justify his decision. Formally (as the representation theorem will show), it breaks ties between justifications. The paper does not rest on the assumption of observable  $\succsim$ , which is dropped in Section 1.4. The second primitive is a choice correspondence  $c$  that maps each non-empty finite set of alternatives to a non-empty subset. To formalize this, let  $\mathcal{F}(\mathcal{A})$  be the set of non-empty finite subsets of  $\mathcal{A}$ . Then we

have  $c : \mathcal{F}(\mathcal{A}) \rightrightarrows \mathcal{F}(\mathcal{A})$  such that  $c(A) \subseteq A$  for all  $A \in \mathcal{F}(\mathcal{A})$ .

Definition 1 presents the JO representation, which is a set  $\mathcal{M}$  of complete, transitive and anti-symmetric orders on  $\mathcal{A}$ .<sup>3</sup> Intuitively, the elements of  $\mathcal{M}$  are the set of preferences that the DM considers justifiable. These preferences are not required to be continuous, so they are not guaranteed to have utility representations. However, readers who prefer to think in terms of utility functions will not lose anything by doing so. A continuous extension is presented in Section 1.3.3.

**Definition 1** (JO Representation). *A JO representation for  $(\succsim, c)$  is a nonempty set  $\mathcal{M}$  of strict preferences such that, for all  $A \in \mathcal{F}(\mathcal{A})$ ,*

$$c(A) = \arg \max (\mathcal{M}(A), \succsim)$$

$$\text{where } \mathcal{M}(A) = \bigcup_{\succsim_m \in \mathcal{M}} \arg \max (A, \succsim_m).$$

### 1.3.1 CHARACTERIZATION

JO is fully characterized by a pair of axioms. The first one, Optimization, says that the DM is truly indifferent between all the items he actually selects. Although this is a standard assumption, it is possible to imagine morally-motivated DMs who violate it. For instance, a DM might feel that it is acceptable to select a selfish alternative as long as he also selects an unselfish one. This behavior is not captured by the justification model. In any case, Optimization has no bite when  $c$  is a choice function rather than a choice correspondence.

**Axiom 1** (Optimization). *For any  $A \in \mathcal{F}(\mathcal{A})$ , for any  $a, b \in c(A)$ :  $a \sim b$ .*

The second axiom, Irrelevance of Unjustifiable Alternatives (IUA), is the heart of the model. Consider a menu  $A$  and an item  $a \in A$ . If the primary preference likes  $a$  at least as much as everything that is selected from  $A$ , but  $a$  is not itself selected, then it must be unjustifiable to choose  $a$

---

<sup>3</sup>There is no loss of generality in taking the justifications to be antisymmetric.

from  $A$ . (Otherwise, the DM would have chosen it.) IUA says that  $a$  is irrelevant for choice on any superset of  $A$ .

**Axiom 2** (Irrelevance of Unjustifiable Alternatives (IUA)). *For any  $a \in \mathcal{A}$  and  $A \in \mathcal{F}(\mathcal{A})$  such that  $a \in A$ : if  $a \succ c(A)$  and  $a \notin c(A)$ , then for all  $B \supseteq A$ ,  $c(B \setminus \{a\}) = c(B)$ .*

Example 1, loosely based on [Rodriguez-Lara and Moreno-Garrido \(2012\)](#), clarifies the restrictions that IUA imposes on choice data.

**Example 1.** *Suppose the DM must allocate \$12 that he and two other subjects ( $A$  and  $B$ ) earned in an earlier phase of the experiment. The DM earned \$4, while  $A$  and  $B$  earned \$6 and \$2 respectively. Let*

$$a = (5 \text{ to self}, 5 \text{ to } A, 2 \text{ to } B)$$

$$b = (4 \text{ to self}, 6 \text{ to } A, 2 \text{ to } B)$$

$$d = (4 \text{ to self}, 4 \text{ to } A, 4 \text{ to } B).$$

*Suppose that the DM's primary preference is given by  $a \succ b \succ d$ : he believes in rewarding other people's good performance, but above all wants to keep more for himself. IUA says that the DM cannot flip from choosing the performance-rewarding option  $b$  when all three options are present to choosing the egalitarian option  $d$  when the selfish option  $a$  is removed. Intuitively, the DM's notion of fairness cannot change when the attractive but unjustifiable option  $a$  is made unavailable.*

To see why IUA is necessary for JO, suppose that it is unjustifiable to choose  $a$  from  $A$ . Toward a contradiction, fix some  $B \supseteq A$ , and suppose that  $c(B) \neq c(B \setminus \{a\})$ . Since the primary preference is fixed, the set of justifiable alternatives must be changing. Specifically, some alternative  $b \in B$  must be unjustifiable when  $a$  is present, but justifiable when  $a$  is absent. That is, there must be some justifiable preference that prefers  $a$  to  $b$ , but  $b$  to everything else. Since there is no justifiable preference that prefers  $a$  to everything else in  $B$ , we have the desired contradiction.

Theorem 1 is the representation theorem for JO.

**Theorem 1.**  $(\succsim, c)$  has a JO representation if and only if it satisfies IUA and Optimization.

The proof of Theorem 1 proceeds in two parts. The first part defines  $\mathcal{M}$  and establishes that it is not too big: it does not justify any choice that the DM would like to make, but does not. The second part establishes that  $\mathcal{M}$  is big enough: it justifies every choice that the DM actually makes.

The notion of “exclusion from below” is central to the proof. Intuitively, a menu  $A$  excludes an item  $b$  from below if the DM likes  $b$  better than everything in  $A$ , but does not choose  $b$  over  $A$ . A preference “respects exclusion from below” if it does not rank any item above a set that excludes that item from below.

**Definition 2** (Exclusion from below).  $A \in \mathcal{F}(\mathcal{A})$  excludes  $b \notin A$  from below if  $b \succsim A$  and  $b \notin c(A \cup \{b\})$ . A preference  $\succsim_m$  respects exclusion from below if for all  $A, b$  such that  $A$  excludes  $b$  from below, there exists  $a \in A$  such that  $a \succ_m b$ .

We define  $\mathcal{M}$  to be the set of preferences on  $\mathcal{A}$  that respect exclusion from below. To see why  $\mathcal{M}$  is not too big, suppose that the DM would like to choose  $b$  from  $A \cup \{b\}$ , but does not:  $b \succsim c(A \cup \{b\})$  and  $b \notin c(A \cup \{b\})$ . We need to show that  $b \notin \mathcal{M}(A \cup \{b\})$ . This is not entirely obvious because  $A$  might not exclude  $b$  from below: there might be items in  $A$  that are strictly better than  $b$  according to the primary preference. Fortunately, IUA says that any such item can be removed without changing choice:

$$c(A \cup \{b\}) = c(\{a \in A : b \succ a\} \cup \{b\}).$$

Since  $b$  is not chosen over  $A$ , it is not chosen over  $\{a \in A : b \succ a\}$ . Conclude that  $\{a \in A : b \succ a\}$  excludes  $b$  from below, so no preference in  $\mathcal{M}$  ranks  $b$  above everything in  $\{a \in A : b \succ a\}$ . We have  $b \notin \mathcal{M}(A \cup \{b\})$  as desired.

It remains to show that  $\mathcal{M}$  is big enough: it contains a justification for every choice the DM makes. This part is more involved. Fix an item  $b$  and a menu  $A$  such that  $b \in c(\{b\} \cup A)$ . For simplicity, assume that  $A = \{a\}$ ; this affects the argument very little. We need to find a preference  $\succ_m$  that respects exclusion from below (so belongs to  $\mathcal{M}$ ) and has  $b \succ_m a$ . We will construct an appropriate  $\succ_m$  by carefully extending the “exclusion from below” relation.

Exclusion from below satisfies three convenient properties. It is irreflexive, meaning no menu containing  $x$  excludes  $x$  from below. It is proper, meaning  $\emptyset$  does not exclude any item. Finally, it is transitive: if  $X$  excludes  $x$  from below,  $Y$  excludes  $y$  from below, and  $x \in Y$ , then  $X \cup Y \setminus \{x, y\}$  excludes  $y$  from below. It turns out that a relation with these properties can be extended in a neat way, captured in the following Lemma. If  $R \subset \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  is irreflexive, proper and transitive, and if  $\neg(y R x)$ , then the transitive closure of  $R \cup (\{x\}, y)$  is irreflexive and proper too.

The Lemma is used to show that exclusion from below can be extended to an irreflexive, proper and transitive relation  $R$  that has  $\{b\} R a$  and has  $\{x\} R y$  or  $\{y\} R x$  for every distinct  $x, y \in \mathcal{A}$ . This is done in two steps. First, we take  $R_0$  to be the transitive closure of the union of  $(\{b\}, a)$  and exclusion from below. Since it cannot be that  $\{a\}$  excludes  $b$  from below, the Lemma implies that  $R_0$  is irreflexive and proper. Second, we follow the proof of the Szpilrajn Extension Theorem (SET) to extend  $R_0$  to  $R$ . SET says that every irreflexive and transitive relation on a set  $X$  can be extended to an irreflexive and transitive relation that contains  $(x, y)$  or  $(y, x)$  for every distinct  $x, y \in X$ . Since  $R_0$  relates menus to items, not items to items, and since such relations require a non-standard notion of transitivity, SET is not directly applicable. However, the arguments used to prove SET are easily adapted to deliver the desired  $R$ .

We use  $R$  to define  $\succ_m$  in the natural way: for any distinct  $x, y$ ,  $x \succ_m y$  if and only if  $\{x\} R y$ . It is not difficult to see that  $\succ_m$  is complete, antisymmetric and transitive, and satisfies  $b \succ_m a$ . To see why it respects exclusion from below, take any  $X, y$  such that  $X$  excludes  $y$  from below. Since  $R$  extends exclusion from below, we have  $X R y$ . Since  $R$  is transitive and proper, it cannot be that

$\{y\} R x$  for all  $x \in X$ . Conclude that  $\{x\} R y$  for some  $x \in X$ , so  $x \succ_m y$  for some  $x \in X$ . We have found  $\succ_m \in \mathcal{M}$  that justifies choosing  $b$  over  $a$ . With minor modifications for  $|A| > 1$ , we can apply the same argument to find an appropriate justification for each decision the DM makes. Thus,  $\mathcal{M}$  is indeed a justification representation for  $(\succ, c)$ .

### 1.3.2 EXTENSION: RESTRICTING ALLOWABLE CONSTRAINTS

JO abstracts away from the details of the choice setting to focus on the conflict between primary preferences and justifiable preferences. This allows the model to unite a range of disparate choice settings at a high level. However, particular applications may demand more structure. If the domain is a set of lotteries, a preference that violates first-order stochastic dominance cannot reasonably be considered justifiable. If the domain is a set of payments to a group of experimental subjects, a preference that violates Pareto dominance is presumably not justifiable. Some more interesting, but more involved, examples are deferred to the end of this section.

This section modifies Theorem 1 to exclude obviously unjustifiable preferences from the representation  $\mathcal{M}$ . The setting at hand determines which preferences count as “obviously unjustifiable.” Specifically, the domain of choice is endowed with an (observable) asymmetric and transitive relation that embodies the basic requirements of rationality and/or morality on that domain. To remind the reader of the FOSD and Pareto examples, we refer to this relation as “dominance” and denote it  $\succ_D$ .

We will construct a JO representation in which all justifiable preferences are strictly monotone in  $\succ_D$ .

**Definition 3** (Strict  $D$ -monotonicity). *A relation  $\succ_R$  on  $\mathcal{A}$  is strictly  $D$ -monotone if, for any  $a, b \in \mathcal{A}$ :  $a \succ_D b$  implies  $a \succ_R b$ .*

**Definition 4** (Monotone JO representation). *A JO representation  $\mathcal{M}$  is monotone if each  $\succ_m \in \mathcal{M}$  is strictly D-monotone.*

Unsurprisingly, the key axiom for the monotone representation is a strengthening of IUA. It says that dominated items, as well as unjustifiable items, are irrelevant. To formalize this, let  $S(B)$  be the members of  $B$  that are dominated or unjustifiable in  $B$ :

$$S(B) := \bigcup_{b \in B} \{b' \in B : b \succ_D b'\} \cup \bigcup_{B' \in \mathcal{F}(B)} \{b' \in B' : \text{it is unjustifiable to choose } b \text{ from } B'\}.$$

The key axiom for the monotone representation, Irrelevance of Submaximal Alternatives (ISA), says that choice is unchanged when any subset of  $S(B)$  is removed.

**Axiom 3** (Irrelevance of Submaximal Alternatives (ISA)). *For any  $B \in \mathcal{F}(\mathcal{A})$ , for any  $A \subseteq S(B)$ :  $c(B) = c(B \setminus A)$ .*

The reader may wonder why ISA allows removal of several items, while IUA only allows removal of one item. Intuitively, this is because exclusion from below (defined in Section 1.3.1) satisfies a nice transitivity property, which allows us to remove unjustifiable items sequentially rather than all at once. This transitivity property doesn't hold once we introduce dominance, so we aren't always able to remove submaximal items sequentially. We have to explicitly allow removing multiple items at once.

Proposition 1 says that replacing IUA with ISA delivers a JO representation in which all the justifiable preferences respect dominance. The proof is along the same lines as that of Theorem 1, but the construction must now keep track of dominance as well as exclusion from below.

**Proposition 1.**  *$(\succ, c)$  has a monotone JO representation if and only if  $c$  satisfies ISA and Optimization.*



As promised, Example 2 shows how notions of disinterestedness or impartiality can be captured by a dominance relation.

**Example 2.**

1. As in *Gneezy et al. (2019)*, let the DM be an investment advisor, and  $\mathcal{A}$  be a set of investments. Each investment is characterized by a distribution over payoffs  $p$  and a real number  $b$ . The real number is the bribe the DM will receive if he recommends that investment to his client. The following dominance relation is a natural way to formalize disinterestedness:  $(p, b) \succ_D (p', b')$  if  $p \succ_{\text{FOSD}} p'$ .

2. As in *Rodriguez-Lara and Moreno-Garrido (2012)*, let  $\mathcal{A}$  be a set of allocations to  $n$  people. Subjects are indexed from least deserving (1) to most deserving ( $n$ ). An allocation is  $p \in \mathbb{R}_+^n$ , where the  $i$ th entry is the payment to the  $i$ th subject. The following dominance relation is a natural way to formalize impartiality:  $p \succ_D q$  if  $p \succ \pi(q)$ , where  $\pi$  is a permutation of  $\{1, \dots, n\}$  such that

$$i \leq j \implies q(\pi(i)) \leq q(\pi(j)).$$

Intuitively,  $p$  dominates  $q$  if it Pareto-dominates a reshuffling of  $q$  that gives larger payments to more deserving subjects.

1.3.3 EXTENSION: CONTINUITY

If the domain  $\mathcal{A}$  is uncountably infinite, utility representations are not guaranteed to exist for the justifiable preferences in the JO representation. This section extends Theorem 1 to require the existence of utility representations. This material is more technical than the preceding. An application-focused reader can safely skip to Section 1.3.4.

For this section alone, we take the domain  $\mathcal{A}$  to be a separable metric space. Let  $\mathcal{Z} = \{z_1, z_2, \dots\}$  denote a countable dense subset of  $\mathcal{A}$ . A continuous JO representation is built from continuous utility functions on  $\mathcal{A}$  rather than preferences on  $\mathcal{A}$ . Let  $C(\mathcal{A}, \mathbb{R})$  denote the set of continuous functions from  $\mathcal{A}$  to  $\mathbb{R}$ .

**Definition 5** (Continuous JO representation).  $(u, \mathcal{M}) \in C(\mathcal{A}, \mathbb{R}) \times 2^{C(\mathcal{A}, \mathbb{R})}$  is a continuous JO representation for  $(\succsim, c)$  if  $u$  represents  $\succsim$  and, for all  $A \in \mathcal{F}(\mathcal{A})$ ,

$$c(A) = \arg \max_{a \in \mathcal{M}(A)} u(a)$$

$$\text{where } \mathcal{M}(A) := \bigcup_{m \in \mathcal{M}} \arg \max_{a \in A} m(a).$$

In addition to continuity, the representation theorem imposes three technical conditions on  $(u, \mathcal{M})$ . All three conditions use the following bits of terminology. For any  $u \in C(\mathcal{A}, \mathbb{R})$  and any  $A, B \in \mathcal{F}(\mathcal{A})$ , say that  $A$  is strictly (weakly) preferred to  $B$  by  $u$  if

$$\max_{a \in A} u(a) > (\geq) \max_{b \in B} u(b).$$

For any  $\mathcal{M} \subseteq C(\mathcal{A}, \mathbb{R})$  and any  $A, B \in \mathcal{F}(\mathcal{A})$ , say that  $A$  is strictly (weakly) preferred to  $B$  by  $\mathcal{M}$  if, for all  $m \in \mathcal{M}$ ,  $A$  is strictly (weakly) preferred to  $B$  by  $m$ .

The three technical conditions are closedness, local non-satiation and recoverability. Closedness is a continuity-like condition for sets of utilities. A *finite* set of utilities is guaranteed to be closed if all its members are continuous. This implication does not hold for infinite sets of utilities, so closedness has to be imposed separately.

**Definition 6** (Closed).  $\mathcal{M} \subseteq C(\mathcal{A}, \mathbb{R})$  is closed if, for all  $B \in \mathcal{F}(\mathcal{A})$ , the set

$$\{a \in \mathcal{A} : B \text{ is strictly preferred to } a \text{ by } \mathcal{M}\}$$

is open.

Like closedness, local non-satiation extends a familiar condition to a set of utilities. A set of utilities is locally non-satiated if, for any item  $a$ , we can find a menu  $Z$  arbitrarily close to  $a$  such that all utilities in the set strictly prefer  $Z$  to  $a$ . Notice that the utilities in the set do not have to agree on *which* item in  $Z$  is better than  $a$ . For technical reasons, we require all elements of  $Z$  to be in the countable dense subset  $\mathcal{Z}$ .

**Definition 7** (Locally non-satiated).  $\mathcal{M} \subseteq C(\mathcal{A}, \mathbb{R})$  is locally non-satiated if, for any  $a \in \mathcal{A}$ , there exists  $Z \in \mathcal{F}(B_\varepsilon(a) \cap \mathcal{Z})$  such that  $Z$  is strictly preferred to  $a$  by  $\mathcal{M}$ .

The final condition, recoverability, is the least familiar. A continuous JO representation is recoverable if the justifiable utilities agree on ranking  $B$  (strictly) above  $a$  only if the primary utility ranks  $a$  (weakly) above  $B$ . Intuitively, the justifiable utilities do not prevent the DM from choosing  $a$  over  $B$  unless the DM would actually like to do so. This restriction is not as strong as it may seem. Recall that a DM who wishes to choose  $a$  over  $B$  only needs one justifiable utility to justify his choice. He doesn't care whether all justifiable utilities rank  $a$  over  $B$ , or only some do. Thus, agreement in the justifiable utilities has no effect *unless* it is in opposition to the primary utility.

This argument suggests that one can impose recoverability without loss of generality. This is true when continuity is not required; the representation constructed in the proof of Theorem 1 is recoverable.<sup>4</sup> When continuity is required, recoverability is restrictive—but only slightly. Any  $(\succsim, c)$  with a “nice” continuous representation has a recoverable quasi-representation that makes the right predictions on all menus without ties.<sup>5</sup> For the interested reader, the formal result is Proposition 22 in Appendix A.2.1.

---

<sup>4</sup>Of course, there may not be utility representations in that case—but recoverability can easily be expressed in terms of preferences rather than utilities.

<sup>5</sup>Specifically, the quasi-representation predicts that  $\{a \in A : a \sim c(A)\}$  is chosen from  $A$ . If no unchosen item is tied with  $c(A)$ , this is the right prediction. Otherwise, it is a superset of the right prediction.

**Definition 8** (Recoverable).  $(u, \mathcal{M}) \in (C(\mathcal{A}, \mathbb{R}), 2^{C(\mathcal{A}, \mathbb{R})})$  is recoverable if, for every  $(B, a) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$ ,

$$B \text{ is strictly preferred to } a \text{ by } \mathcal{M} \implies \{b \in B : u(b) \leq u(a)\} \text{ is strictly preferred to } a \text{ by } \mathcal{M}.$$

The continuous representation is characterized by four axioms. Optimization is already familiar. C-IUA and IUA are very similar, but C-IUA accounts for additional information that can be gleaned about the justifiable preferences when they are required to be continuous. To see how this works, fix any menu  $A$ , and let  $W(A)$  be the set of items in  $\mathcal{A}$  that are excluded from below by a subset of  $A$ . As established in Section 1.3.1, every justifiable preference must strictly prefer  $A$  to every item in  $W(A)$ . Now consider a sequence of items  $a_i \rightarrow a$  and a sequence of menus  $A_i \rightarrow A$  such that  $a_i \in W(A_i)$  for all  $i$ . Since every justifiable preference strictly prefers  $A_i$  to  $a_i$  for all  $i$ , every justifiable preference must weakly prefer  $A$  to  $a$ . (This is the step that uses continuity.) If every justifiable preference strictly prefers  $a$  to  $b$ , then every justifiable preference strictly prefers  $A$  to  $b$ . Thus,  $b$  must be irrelevant for choice when  $A$  is present. C-IUA generalizes this logic.

Formally, let

$$W(A) := \bigcup_{A' \subseteq A} \{a \in \mathcal{A} : a \succsim A' \text{ and } a \notin c(\{a\} \cup A')\}$$

$$\bar{W}(A) := \{a \in \mathcal{A} : \exists A_i \rightarrow A, a_i \rightarrow a \text{ s.t. } \forall i \ a_i \in W(A_i)\}.$$

**Axiom 4** (C-IUA). If  $D \in \mathcal{F}(\bar{W}(B))$  for some  $B \in \mathcal{F}(W(A))$ , or if  $D \in \mathcal{F}(W(B))$  for some  $B \in \mathcal{F}(\bar{W}(A))$ , then  $c(A) = c(A \setminus D)$ .

The final two axioms, Continuity and Improvability, are the behavioral analogues of Closedness and Local Non-Satiation respectively. Both are integral to the construction of  $\mathcal{M}$  in the sufficiency proof.

**Axiom 5** (Continuity).

1.  $\succsim$  is continuous.
2. For all  $A \in \mathcal{F}(\mathcal{A})$ ,  $W(A)$  is open.

**Axiom 6** (Improvability). For any  $a \in \mathcal{A}$  and any  $\varepsilon > 0$ , there is some  $Z \in \mathcal{F}(B_\varepsilon(a) \cap \mathcal{Z})$  such that  $a \in W(Z)$ .

**Theorem 2.** *The following are equivalent:*

1.  $(\succsim, c)$  satisfies Continuity, C-IUA, Improvability and Optimization.
2.  $(\succsim, c)$  has a recoverable continuous JO representation  $(u, \mathcal{M})$  such that  $\mathcal{M}$  is closed and locally non-satiated.

Just as the proof of Theorem 1 defines  $\mathcal{M}$  to be the set of preferences that respect exclusion from below, the proof of Theorem 2 defines  $\mathcal{M}$  to be the set of continuous utility functions that respect exclusion from below. Just as Theorem 1 constructs a preference in  $\mathcal{M}$  to justify each decision the DM makes, Theorem 2 constructs a continuous utility function in  $\mathcal{M}$  to justify each decision the DM makes. The proofs diverge after that. The proof of Theorem 2 appeals to a result in [Herden and Pallack \(2002\)](#) (HP), which provides sufficient conditions for an incomplete binary relation to have a continuous extension. The main part of the proof shows that a variant of the “exclusion from below” relation satisfies the HP conditions. (This step is non-trivial because the HP conditions are not easily checked, and bear no obvious relation to the axioms.) The HP theorem delivers a continuous utility representation for this relation, which implicitly defines a justifiable utility function.

#### 1.3.4 EXTENSION: EXPECTED UTILITY

This section covers the expected-utility version of JO, in which the primary preference and all the justifiable preferences have expected-utility representations. A larger set of axioms is needed to

achieve this additional structure, but the result is a tractable model suited to application. Additionally, the model in this section is the precursor to the random justification model in Section 1.5.

The set of prizes  $Z$  is assumed to be finite, although the domain  $\mathcal{A} := \Delta(Z)$  is not. It is convenient to assume that there is a dominance relation  $\succ_D$  on  $Z$ . As in Section 1.3.2,  $\succ_D$  is asymmetric and transitive, but need not be complete. It is enough to have two payoffs ranked by dominance. For instance,  $Z$  could contain two monetary payments to the DM as well as payments to others, or donations to various charities. A version of the representation theorem is available without this assumption, but the axioms are slightly more cumbersome.<sup>6</sup>

We require the primary utility to be strictly  $D$ -monotone and the justifiable utilities to be weakly  $D$ -monotone. This assumption is motivated by convenience: weak monotonicity of the justifiable utilities falls out of the neatest set of axioms. Note that the DM will never actually *choose* a strictly dominated alternative. If the dominated alternative is justifiable, the dominating alternative will be too, and the DM will choose the latter but not the former.

**Definition 9** (*D-monotonicity*). *A Bernoulli utility  $u : Z \rightarrow \mathbb{R}$  is weakly (strictly)  $D$ -monotone if  $a \succ_D b$  implies  $u(a) \geq (>) u(b)$ .*

We require the set of justifiable utilities to be compact and convex.

**Definition 10** (*Expected-utility JO representation*). *An expected-utility JO representation consists of a strictly  $D$ -monotone Bernoulli utility  $u$  and a compact, convex set  $\mathcal{M}$  of weakly  $D$ -monotone Bernoulli utilities such that  $p \mapsto \mathbb{E}_p u$  represents  $\succsim$ , and*

$$c(\mathcal{A}) = \arg \max_{p \in \mathcal{M}(\mathcal{A})} \mathbb{E}_p u$$

$$\text{where } \mathcal{M}(\mathcal{A}) := \bigcup_{m \in \mathcal{M}} \arg \max_{p \in \mathcal{A}} \mathbb{E}_p m.$$

---

<sup>6</sup>In the absence of dominance, the justifiable set may contain a preference exactly opposite the primary preference. This case is inconvenient and needs to be handled separately.

In the EU model, it is often easier not to work directly with  $\mathcal{M}$ , but with the set

$$B(p) := \{q \in \Delta(Z) : p \succsim q \text{ and } \{q\} = c(\{p, q\})\} \quad (1.1)$$

for some  $p \in \text{int}(\Delta(Z))$ . Intuitively,  $B(p)$  is the set of lotteries that the DM feels compelled to choose over  $p$ , contrary to his primary preference. Formally,  $B(p)$  is like a sufficient statistic: it encodes everything we need to know about the representation. It appears several times in the axioms.

Each of the four axioms in this section has two parts. For all axioms but Convexity, the first part is a condition on the primary preference, and the second part is a condition on choice behavior that ultimately translates into a condition on justifiable preferences. Since the conditions for a given preference to have a FOSD-monotone EU representation are well known, the first part is completely standard.

**Axiom 7** (Independence).

1. For all  $p, q, r \in \Delta(Z)$  and  $\alpha \in (0, 1)$ ,  $p \succsim q$  implies  $\alpha p + (1 - \alpha)r \succsim \alpha q + (1 - \alpha)r$ .
2. For any  $A \in \mathcal{F}(\Delta(Z))$  and  $p \in \Delta(Z)$ ,

$$c(\alpha A + (1 - \alpha)\{p\}) = \alpha c(A) + (1 - \alpha)\{p\}.$$

The second part of Independence says that the DM's preference does not flip when every option he faces is mixed with a fixed lottery in a fixed proportion. The necessity of this axiom for an expected-utility JO representation is obvious. If  $p$  is the best item in  $A$  that is top-ranked by some justifiable preference, then mixing everything in  $A$  (including  $p$ ) with  $q$  will not change this. Formally, this axiom ensures that  $B(p)$  is a convex cone. The justifiable preferences are among the supporting hyperplanes of this cone.

**Axiom 8** (Continuity).

1. For any  $p \in \Delta(Z)$ ,  $\{q \in \Delta(Z) : q \succsim p\}$  and  $\{q \in \Delta(Z) : q \precsim p\}$  are closed.
2. For any  $p \in \Delta(Z)$ ,  $B(p)$  is open in  $\{q \in \Delta(Z) : q \succsim p\}$ .

There is nothing unusual about the continuity axiom. If the second part fails, the set of justifiable preferences constructed in the proof will be slightly too permissive. Some preferences will be indifferent between pairs of items that they should strictly rank.

For the monotonicity condition, we need to extend first-order stochastic dominance (FOSD) to allow for an incomplete ranking over prizes. Notice that the definition provided here reduces to the usual one when  $\succ_D$  is complete.

**Definition 11** (FOSD).  $p >_{FOSD} q$  if  $p = \sum_i \alpha_i \delta_{p_i}$ ,  $q = \sum_i \alpha_i \delta_{q_i}$ ,  $p_i \succ_D q_i$  or  $p_i = q_i$  for all  $i$ , and  $p_i \succ_D q_i$  for some  $i$ .

**Axiom 9** (Monotonicity).

1.  $\succsim$  is strictly FOSD-monotone.
2. For any  $p, q \in \Delta(Z)$  such that  $p >_{FOSD} q$  and any  $A \in \mathcal{F}(\Delta(Z))$  containing  $p$  and  $q$ ,  $c(A) = c(A \setminus \{q\})$ .

The second part of Monotonicity is really Irrelevance of Dominated Alternatives (IDA). It says that any lottery is irrelevant in the presence of a lottery that strictly dominates it. This type of condition will be familiar from Section 1.3.2.

The final axiom, Convexity, is the least familiar. The first part says that it is unjustifiable to choose lottery  $p$  over menu  $A$  only if it is unjustifiable to choose  $p$  over some mixture of lotteries in  $A$ . The second part is a partial converse of the first. If it is unjustifiable to choose  $p$  over some mixture of lotteries in  $A$ , then  $p$  cannot be chosen over  $A$ .



**Axiom 10** (Convexity). *For any  $A \in \mathcal{F}(\Delta(Z))$  and  $p \notin A$ :*

1. *If  $p \succsim c(\{p\} \cup A)$  and  $p \notin c(\{p\} \cup A)$ , then  $co(A) \cap B(p) \neq \emptyset$ .*
2. *If  $co(A) \cap B(p) \neq \emptyset$ , then  $p \notin c(A \cup \{p\})$ .*

The word “mixture” is important. Convexity does not say that it is unjustifiable to choose  $p$  over  $A$  only if it is unjustifiable to choose  $p$  over some member of  $A$ . This is typically not the case. Unless the set of justifiable preferences is a singleton, we can find a lottery  $p$  and menu  $A$  such that it is justifiable to choose  $p$  over each member of  $A$ , but unjustifiable to choose  $p$  over  $A$ . Indeed, this inconsistency between choices in binary menus and choices in larger menus is part of what makes the justification model interesting.

The name Convexity arises because necessity of the first part follows from convexity of  $\mathcal{M}$ . To see what role Convexity plays in the proof of sufficiency, suppose that the DM prefers  $p$  to everything in  $A$ . Convexity implies that the DM will choose  $p$  over  $A$  if and only if  $B(p)$  can be separated from  $A$  by an appropriately chosen hyperplane. This hyperplane will turn out to be the indifference curve of a justifiable preference.

Theorem 3 is the EU version of Theorem 1. The proof is quite different from that of Theorem 1; it uses mostly geometric rather than order-theoretic arguments. First, Independence and Continuity are used to show that  $B(p)$  is always a convex cone, open in  $\{q \in \Delta(Z) : p \succsim q\}$ . Monotonicity is used to establish the relationship between  $B(p)$  and the primary preference.  $B(p)$  is used to identify a candidate set of justifiable preferences. Each candidate preference has an indifference curve that is a supporting hyperplane of  $B(p)$ . Finally, Convexity and the axioms from Theorem 1 are used to show that the candidate set is neither too large nor too small. No candidate preference would allow the DM to justify a choice he would have liked to make, but did not; and for each choice the DM actually makes, some candidate preference justifies that choice.

**Theorem 3.** *Suppose there exist  $y, z \in Z$  such that  $y \succ_D z$ . The following are equivalent:*

1.  $(\succsim, c)$  satisfies IUA, Optimization, Independence, Continuity, Monotonicity and Convexity.
2.  $(\succsim, c)$  has an expected-utility JO representation.

The expected-utility JO model is subject to the usual uniqueness issues for EU representations: both the primary utility and the justifiable utilities can be arbitrarily (and independently) shifted and scaled. This issue can be dismissed by shifting attention from the set of justifiable utilities to the set of preferences they represent.

There is also a more subtle uniqueness issue: some EU preferences that are *not* positive affine transformations of other preferences can be added to or removed from the set of justifiable preferences without changing anything. Intuitively, this is because these preferences are too far from the primary preference to be of much use. Any choice they justify is also justified by some preference that is closer to  $\succsim$ . Corollary 1 says that we can obtain a unique representation (up to shifting and scaling) by dropping all these surplus preferences, retaining only those preferences needed to explain the DM's behavior. At the other extreme, we can include all the preferences not ruled out by the DM's behavior or by dominance.

**Corollary 1.** *If  $(\succsim, c)$  has an expected-utility JO representation, it has a unique maximal and a unique minimal set of justifiable preferences.*

A nice feature of the EU model is the ease of comparing different DMs' standards for acceptable behavior. Since the proof of Theorem 3 uses  $B(p)$  to obtain a set of justifiable preferences, we may expect a connection between the size of  $B(p)$  and the strictness of the DM's standards. There is indeed a connection. If  $B_1(p)$  and  $B_2(p)$  are nested, then we can find representations in which the sets of justifiable utilities are also nested. If  $B_2(p)$  is the smaller set, then DM 2 can appeal to the larger set of justifiable utilities. This is convenient because we can figure out each DM's standards for acceptable behavior just by looking at his choices on a small set of menus: binary menus that

contain a fixed item  $p$ . If DM 1 chooses a weakly inferior item  $q$  over  $p$  whenever DM 2 does, we can assume that DM 1's standards are at least as strict as DM 2's.

**Corollary 2.** *Suppose that  $(\succsim, c_1)$  and  $(\succsim, c_2)$  have expected-utility JO representations. The following are equivalent:*

1.  $B_1(p) \supset B_2(p)$  for some  $p \in \text{int}(\Delta(Z))$ .
2. *The maximal set of justifiable preferences for  $(\succsim, c_1)$  is strictly smaller than the maximal set of justifiable preferences for  $(\succsim, c_2)$ .*

The result is not true if “maximal” is replaced with “minimal.” This is because there may be justifications that the more liberal DM never feels the need to use. He may agree that these are acceptable justifications, but he doesn't appeal to them because something better is always available. Thus, his minimal set of justifiable preferences doesn't include them.

Corollary 2 is connected to the random justification model of Section 1.5. There, we consider a population of DMs that can be totally ordered by the strength of their standards for acceptable behavior. The maximal/minimal distinction does not arise in the stochastic setup because heterogeneity in primary preferences eliminates non-uniqueness in sets of justifiable preferences. Intuitively, non-uniqueness persists when the primary preference is fixed because the sets of justifiable preferences are identified through conflict with the primary preference. If the DM prefers  $p$  to  $q$  and chooses  $p$  over  $q$ , we simply cannot tell whether he felt able to choose  $q$ . This problem goes away when the dataset contains multiple DMs with the same justifications but different preferences.

#### 1.4 PRIMARY PREFERENCE UNOBSERVABLE

This section characterizes the Justification model when the primary preference is Unobservable (JU). For simplicity, we now take  $c$  to be a choice function rather than a choice correspondence. A

JU representation consists of a strict preference  $\succ$  as well as a set of preferences  $\mathcal{M}$ .

This section provides two complementary behavioral characterizations of JU. The first characterization is a straightforward modification of Theorem 1. It shows that a single easy-to-check axiom is necessary and sufficient for a JU representation. It also provides insight into the DM's primary preference: it tells the researcher exactly which primary preferences are consistent with the DM's behavior, and which are not. However, it does not provide direct insight into the preferences the DM considers justifiable. The second characterization fills this gap. It shows how to identify the preferences that the DM may consider justifiable, and how to rule out the rest. Conveniently, the constraints on the justifiable preferences come from simple, easily recognized patterns of choice.

#### 1.4.1 FIRST CHARACTERIZATION

This section builds directly on Theorem 1. Recall that IUA is a necessary and sufficient condition for a JO representation when  $c$  is a choice function. It says that  $a$  is irrelevant when  $A$  is present if  $a \succ c(B \cup \{a\})$  for some subset  $B$  of  $A$ . We can flip IUA backward to derive conditions on the primary preference. To see how this works, suppose that choice from set  $A$  changes when item  $a$  is added. Clearly, it is not the case that  $a$  is irrelevant when  $A$  is present. Thus, it cannot be that  $a \succ c(B \cup \{a\})$  for any subset  $B$  of  $A$ . Intuitively, if the addition of  $a$  affects choice on set  $A$ , then it is justifiable to choose  $a$  over  $A$ . If the DM fails to choose  $a$  over some subset  $B$  of  $A$ , he must be acting out of inclination rather than obligation. We can repeat this argument to obtain a full set of restrictions on  $\succ$ . If the resulting restrictions form a cycle, no strict preference can satisfy them, so there is no representation. But if the restrictions do not form a cycle,  $c$  satisfies IUA conditional on any strict preference that obeys the restrictions, so Theorem 1 delivers a representation. This is exactly what Proposition 2 says.

**Definition 12** (Revealed Preference). *If  $a \neq c(B \cup \{a\})$  and, for some  $A \supseteq B$ ,  $c(A) \neq c(A \cup \{a\})$ ,*

then  $c(B \cup \{a\})$  is revealed preferred to  $a$ .

**Axiom 11** (Acyclicity). *The revealed preference relation for  $c$  is acyclic.*

**Proposition 2.**  *$c$  satisfies Acyclicity if and only if it has a JU representation. Moreover, a preference  $\succ$  extends the revealed preference relation for  $c$  if and only if there is some  $\mathcal{M}$  such that  $(\succ, \mathcal{M})$  represents  $c$ .*

Corollary 2 is helpful not just because it delivers a representation, but because it reveals the set of primary preferences consistent with choice behavior. Sometimes, it pins down a unique preference. Example 3 illustrates.

**Example 3.** *Recall the example from Exley (2016) in Section 1.1:*

$$\begin{aligned} a &= \text{experimenter pays } \$2.50 \text{ to DM} \\ b &= 50\% \text{ chance experimenter donates } \$10 \text{ to charity} \\ d &= \text{experimenter donates } \$4 \text{ to charity} \\ c(\{a, b\}) &= a \quad c(\{b, d\}) = b \quad c(\{a, d\}) = d. \end{aligned}$$

*Exley (2016) considered only binary menus, but choice from menus with more than two alternatives are important for identification in JU. We assume  $b = c(\{a, b, d\})$  since we have already seen a plausible  $(\succ, \mathcal{M})$  that generates this choice.*

*Notice that choice from the grand set changes when  $d$  or  $b$  is removed. Since  $b$  affects choice on a menu containing  $a$ , but  $b$  is not chosen over  $a$ , we must have  $a \succ b$ . Since  $d$  affects choice on a menu containing  $d$ , but  $d$  is not chosen over  $b$ , we must have  $b \succ d$ . Putting these two restrictions together, we get  $a \succ b \succ d$ . The DM prefers the risky donation over the safe donation, but above all prefers to keep more for himself.*

Conditional on any primary preference  $\succ$ , we can work out the set of justifiable preferences. We saw how to do this following Theorem 1. First, find each menu  $A$  and item  $a \in A$  such that  $a \succ A$  but  $a \notin c(A \cup \{a\})$ . Translate each of these pairs into a restriction on the justifiable preferences: every justifiable preference strictly prefers something in  $A$  to  $a$ . Finally, take the set of justifiable preferences to be the preferences that satisfy all these restrictions.

This process leaves something to be desired. It would be better to learn about the set of justifiable preferences just by looking at the data, not by constructing a set of primary preferences and then working out the constraints associated with each one. The next section explains how to do that, via an alternative characterization of the same model. The second characterization is more involved than the first one, but should still be of interest to the empirically inclined reader.

#### 1.4.2 SECOND CHARACTERIZATION

Two patterns of choice behavior are key to understanding the set of justifiable preferences. We define these patterns, explain their implications for the primary preference and the justifiable preferences, and leverage them to obtain another representation theorem for JU.

The first key pattern is a three-element cycle. As suggested in Example 3, the primary preference on any three-element cycle is uniquely pinned down. The item chosen from the full three-element set is middle-ranked, and the item that beats it is top-ranked. Since the bottom-ranked item beats the top-ranked item, it must be unjustifiable to choose the latter over the former. (In Example 3, it is unjustifiable to choose the safe payment  $a$  over the safe donation  $d$ .) Further restrictions on the justifications may be obtained by chaining together multiple cycles. For instance: if one cycle reveals  $a$  to be better than  $b$ , and another reveals  $b$  to be better than  $d$ , but  $d$  is chosen over  $a$ , it must be unjustifiable to choose  $a$  over  $d$ .

**Definition 13** (Cycle/Chain).  $(a_1, a_2, a_3)$  is a cycle if

$$c(\{a_1, a_2\}) = a_1 \quad c(\{a_1, a_2, a_3\}) = a_2 \quad c(\{a_1, a_3\}) = a_3.$$

For  $k \geq 3$ ,  $(a_1, \dots, a_k)$  is a chain if for each  $i \in \{2, \dots, k-1\}$ ,  $(a_{i-1}, a_i, a_{i+1})$  is a cycle and/or both  $(a_{i-2}, a_{i-1}, a_i)$  and  $(a_i, a_{i+1}, a_{i+2})$  are cycles.

The second key pattern of choice is an almost-WARP set. (The reasons for the name will soon become clear.) Fix some set  $A$ , and suppose that choice satisfies WARP on all its proper subsets. If  $|A| = 3$ , suppose further that pairwise choice is not cyclic. (Pairwise choice cannot be cyclic if  $|A| > 3$ .) Then, there is a unique preference on  $A$  that is maximized by choice from each proper subset. This preference is pinned down by pairwise choice. If choice on  $A$  violates WARP, so  $c(A)$  is pairwise-defeated by some other item in  $A$ , the primary preference on  $A$  is uniquely pinned down. In fact, it is the preference given by pairwise choice. (This is not immediately obvious, but it is straightforward to prove.) Like cycles, almost-WARP sets are informative about the justifications as well as the primary preference. Since the item that pairwise-beats everything else in  $A$  is better than everything else in  $A$ , but is not chosen from  $A$ , it must be unjustifiable to choose that item from  $A$ .

**Example 4.** *This example is a continuation of Example 1, and uses the same notation. Suppose that the DM chooses the selfish allocation  $a$  over the performance-based allocation  $b$ , and separately chooses  $a$  over the equitable allocation  $d$ . Suppose further that the DM chooses  $b$  over  $d$ . If the DM selects  $b$  from the grand set  $\{a, b, d\}$ , then the grand set is almost-WARP. The DM's primary preference must be  $a \succ b \succ d$ , which is the preference given by pairwise choice. Moreover,  $a$  must be unjustifiable in the presence of  $\{b, d\}$ . Intuitively, the DM can choose  $a$  over  $b$  by pretending that he is averse to inequity, and he can choose  $a$  over  $d$  by pretending that he believes in rewarding performance, but he cannot do both at once.*

**Definition 14** (Almost-WARP set). *Suppose that  $A$  is not a cycle.  $A$  is an almost-WARP set if choice violates WARP on  $\mathcal{F}(A)$ , but satisfies WARP on  $\mathcal{F}(A) \setminus A$ .*

The implications of cycles and almost-WARP sets for justifiable preferences are summed up in Definition 15.

**Definition 15** (Revealed exclusion). *An item  $a$  is revealed excluded by a menu  $B$  if:*

1. *For  $|B| > 1$ :  $B \cup \{a\}$  is an almost-WARP set, and  $a$  pairwise-defeats  $c(B \cup \{a\})$ .*
2. *For  $B = \{b\}$ :  $b = c(\{a, b\})$ , and  $a$  comes before  $b$  in a chain.*

Since three-element cycles and almost-WARP sets are easy to spot, so is revealed exclusion. One may wonder whether more restrictions on the justifiable preferences could be obtained from more complicated patterns of choice. The answer is no: cycles and almost-WARP sets tell the analyst all she could hope to know about the preferences the DM considers justifiable. Every preference that respects revealed exclusion appears in some JU representation. (In fact, there is a JU representation in which the set of justifiable preferences is *precisely* the set of preferences that respect revealed exclusion. We return to this point after the next representation result.) Proposition 3 summarizes.

**Proposition 3.** *Suppose  $c$  has a JU representation. For any menu  $A$  and item  $a$  such that  $a \notin A$ , the following are equivalent:*

1.  *$a$  is revealed excluded by a subset of  $A$ .*
2. *No justifiable preference in any representation ranks  $a$  above  $A$ .*

The next axiom is the analogue of IUA for the unknown-primary-preference case. It says that an item  $a$  is irrelevant for choice from  $B$  if  $a$  is revealed excluded by a subset of  $B$ .

**Axiom 12** (Irrelevance of Excluded Alternatives (IEA)). *If each  $a \in A \subset B$  is revealed excluded by a subset of  $B$ , then  $c(B) = c(B \setminus A)$ .*



To understand the restrictions imposed by IEA, consider an almost-WARP set  $A$ . By definition, there is a unique preference maximized by choice on the proper subsets of  $A$ . We can index the items in  $A$  from best to worst according to this preference:  $a_1 \succ \cdots \succ a_n$ . (As noted above, the primary preference agrees with the WARP-implied preference, hence the notation.) Since choice on  $A$  violates WARP, we can't have  $c(A) = a_1$ . It turns out we can only have  $c(A) = a_2$ —the item chosen from  $A$  is the DM's second-favorite item.<sup>7</sup> This means the DM's choice cannot deteriorate too quickly as we expand the choice set. He can move from always choosing his favorite item to choosing his second-favorite, but not to his third-favorite or worse. For instance, the DM cannot choose  $d$  from  $\{a, b, d\}$  in Example 4.

Of course, IUA has no bite if nothing is revealed excluded. The reader may wonder how often cycles and almost-WARP sets actually arise. The answer is reassuring: unless choice satisfies WARP (in which case a standard preference-maximization model is perfectly adequate), there will be at least one cycle or almost-WARP set, so at least one item will be revealed excluded. This will provide an opportunity to falsify the model.

**Theorem 4.**  *$c$  satisfies IEA if and only if  $c$  has a JU representation.*

The proof of Theorem 4 constructs a particular JU representation for  $c$ , which we call the “canonical representation.” This representation is notable because the set of justifiable preferences is precisely the set of preferences that respect revealed exclusion. Thus, the set of justifiable preferences is maximal: it includes the set of justifiable preferences from every other representation. Maximality of justifications corresponds to minimality of constraints: the more justifications are available to the DM, the fewer constraints he faces in making his decision. Therefore, the canonical representation for  $c$  is the most parsimonious model of constrained decision-making that explains  $c$ .

---

<sup>7</sup>Suppose  $c(A) = a_3$ . Then,  $a_2$  is revealed excluded by  $A \setminus \{a_2\}$ . By IEA, we can remove  $a_2$  from  $A$  without changing choice. But we know that  $c(A \setminus \{a_2\}) = a_1 \neq c(A)$ , so we have a contradiction. This argument generalizes to  $i > 3$ .

The primary preference in the canonical representation is easily constructed. First, impose  $a \succ b$  whenever  $a$  comes before  $b$  in a chain. Then, if  $a$  and  $b$  have not yet been ranked, impose  $a \succ b$  if  $a = c(\{a, b\})$ , and  $b \succ a$  otherwise. Although this may not be the only preference consistent with behavior, it is the only preference consistent with the *maximal* set of justifiable preferences. Corollary 3 summarizes the properties of the canonical representation.

**Definition 16** (Canonical representation). *JU representation  $(\succ^*, \mathcal{M}^*)$  is canonical if (1)  $\mathcal{M}^*$  is the set of preferences that respects revealed exclusion, and (2)  $\succ^*$  is the preference that has  $a \succ b$  whenever  $a$  comes before  $b$  in a chain, and that agrees with pairwise choice on pairs not connected by any chain.*

**Corollary 3.** *Suppose  $c$  has a JU representation  $(\succ, \mathcal{M})$ . Then, it has a unique canonical representation  $(\succ^*, \mathcal{M}^*)$ , and  $\mathcal{M}^* \supseteq \mathcal{M}$ .*

### 1.4.3 EXTENSION: COMPARING DECISION ENVIRONMENTS

Sections 1.4.1 and 1.4.2 cover identification in a single, static setting. Intuitively, we should be able to learn more about the primary preference and/or set of justifiable preferences by varying the pressure to find a good justification. For instance, we would expect the DM's choices to be closer to his primary preference when he chooses anonymously than when he must announce his choice to some ethically conscious peers. Several experiments find evidence that the choice setting matters in this way. As mentioned in Section 1.2, [Hamman et al. \(2010\)](#) find more selfish behavior when decisions are implemented by an intermediary, and [Dana et al. \(2006\)](#) find more selfish behavior when decisions are unobserved by those affected. [Falk \(2017\)](#) finds less selfish behavior when subjects are forced to watch themselves in a mirror, and [Haley and Fessler \(2005\)](#) find less selfish behavior among subjects who are “watched” by a pair of stylized eyespots.

To formalize the two-setting case, let  $c_L$  be the choice function corresponding to the low-pressure setting, and let  $c_H$  be the choice function corresponding to the high-pressure setting. We are now

looking for a pair of JU representations with the same primary preference and nested sets of justifiable preferences. We accomplish this by building on Theorem 4. First, we ensure that  $c_L$  has a JU representation by requiring it to satisfy IEA. Second, we impose consistency between  $c_H$  and  $c_L$ . This consistency condition, IREA, is essentially a stronger version of IEA. It says that any item revealed excluded in the low-pressure setting is irrelevant in the high-pressure setting. This is clearly necessary for the set of justifications to be smaller in the high-pressure case. IREA also says that anything the DM chose in the low-pressure setting, but not in the high-pressure setting, is irrelevant in the high-pressure setting. Intuitively, this is because an item is replaced in the high-pressure setting only if it no longer counts as justifiable.

**Definition 17** (Replacement).  *$a$  is replaced in  $A$  if  $a = c_L(A) \neq c_H(A)$ .*

**Axiom 13** (Irrelevance of Replaced or Excluded Alternatives (IREA)). *If each  $a \in A \subset B$  is revealed excluded in  $L$  by, or replaced in, a subset of  $B$ , then  $c_H(B) = c_H(B \setminus A)$ .*

**Example 5.** *This example is very loosely based on Norton et al. (2004). Let*

$m_1 = \text{educated male applicant}$

$m_2 = \text{experienced male applicant}$

$f_1 = \text{educated female applicant}$

*Suppose there are two treatments, a low-pressure setting ( $L$ ) in which the DM's choice is observed only by the experimenter, and a high-pressure setting ( $H$ ) in which the DM's choice is observed by a female peer.*

*Suppose that choice in the low-pressure setting is given by*

$$c_L(\{m_1, m_2\}) = m_1 \quad c_L(\{m_2, f_1\}) = c_L(\{m_1, m_2, f_1\}) = m_2 \quad c_L(\{f_1, m_1\}) = f_1.$$

*This is a cycle, so the DM must feel unable to choose the educated male applicant over the educated*

female applicant. Suppose further that  $c_H(\{m_1, m_2, f_1\}) = f_1$ . Since  $m_2$  is replaced in  $\{m_1, m_2, f_1\}$ , IREA implies that

$$c(\{m_1, f_1\}) = c(\{m_2, f_1\}) = c(\{m_1, m_2, f_1\}) = f_1.$$

The DM now feels unable to choose either of the male applicants over the educated female applicant.

Proposition 4 is the representation result for the two-setting case. At the expense of additional notation, it could easily be extended to more than two settings.

**Proposition 4.**  $c_L$  and  $c_H$  have JU representations  $(\succ, \mathcal{M}^L)$  and  $(\succ, \mathcal{M}^H)$  such that  $\mathcal{M}^H \subseteq \mathcal{M}^L$  if and only if  $(c_L, c_H)$  satisfies IREA and  $c_L$  satisfies IEA.

Proposition 4 also helps to interpret JO, the justification model with observable primary preference. It may seem mysterious for the analyst to observe a component of the representation. Corollary 4 explains what is really going on. Rather than directly observing the primary preference, the analyst observes choice behavior in a low-pressure situation. Provided this behavior satisfies WARP, she identifies the WARP-implied preference with the primary preference. Corollary 4 says there is little harm in this: if there is any JU representation, there is one in which the primary preference is the WARP-implied preference.

**Corollary 4.** Suppose that  $c_L$  satisfies WARP, so the restriction of  $c_L$  to binary menus pins down a unique preference  $\succ$ .  $c_H$  has a JU representation if and only if it satisfies IUA conditional on  $\succ$ .

## 1.5 RANDOM JUSTIFICATION MODEL

In Section 1.4, WARP violations are used to reveal inconsistencies between the DM's preferences and his notion of acceptable behavior. Thus, the DM is ultimately unable to maintain the illusion that his preferences are beyond reproach. A sophisticated DM may recognize this danger and adjust

his behavior accordingly, reducing or eliminating the WARP violations needed to identify the representation. This section addresses the problem by moving to a stochastic setup, in which each data point can be collected from a different DM.

The Random Justification model (RJ) builds on the deterministic EU model in Section 1.3.4. An RJ representation has two components: a distribution  $\mu$  over primary preferences, and a distribution  $\nu$  over sets of justifications. (Primary preferences and justifications are assumed to be drawn independently.) Just as in Section 1.3.4, both primary preferences and justifications have an EU form, and sets of justifications are assumed to be closed and convex. Formally, let  $\mathcal{U}$  be the set of EU preferences on  $\Delta(Z)$ . Say that  $U \subset \mathcal{U}$  is convex (closed) if

$$U_R := \bigcup_{\succsim \in U} \{u \in \mathbb{R}^Z : u \text{ represents } \succsim\}$$

is convex (closed). Let  $\mathfrak{U}$  be the set of nonempty closed, convex subsets of  $\mathcal{U}$ .

Following [Gul and Pesendorfer \(2006\)](#), we assume that ties in the primary preferences happen with zero probability. Some of the results in this section also require  $\mu$  to have full support.

**Definition 18** (Preference distribution).  $\mu \in \Delta(\mathcal{U})$  is a preference distribution if, for any distinct  $x, y \in \Delta(Z)$ ,  $\mu(\{\succsim : x \sim y\}) = 0$ .

For tractability, we restrict attention to a particular form of heterogeneity in justifications. The key assumption is that the support of  $\nu$  is ordered by set inclusion, so the sets of justifications can be ranked from “most permissive” to “most strict.” This assumption is appropriate for populations of decision-makers who agree on the values to be promoted, but disagree on the level of commitment needed to effectively promote those values.

**Example 6.** This example is based on [Fong and Oberholzer-Gee \(2011\)](#). Subjects can make a small transfer to a poor person who suffers from a physical disability or a drug addiction. Subjects with any

degree of generosity may choose a transfer to a disabled recipient over keeping the money. Especially generous subjects may choose the transfer regardless of recipient type. The remaining subjects may choose the transfer only if the recipient is sufficiently likely to be disabled.

For technical reasons, we also assume that (1) a positive mass of DMs can justify anything, (2) the constraint distribution has no other mass points, and (3) the support of the constraint distribution has no gaps. These restrictions are formalized in Definition 19.

**Definition 19** (Constraint distribution).  $\nu \in \Delta(\mathfrak{U})$  is a constraint distribution if it satisfies the following conditions:

1. For any distinct  $\mathcal{M}_1, \mathcal{M}_2 \in \text{supp}(\nu)$ , there exists  $\mathcal{M}_3 \in \text{supp}(\nu)$  such that

$$\mathcal{M}_1 \subset \text{int}(\mathcal{M}_3) \subset \mathcal{M}_3 \subset \text{int}(\mathcal{M}_2) \quad \text{or} \quad \mathcal{M}_2 \subset \text{int}(\mathcal{M}_3) \subset \mathcal{M}_3 \subset \text{int}(\mathcal{M}_1).$$

2.  $\nu(\mathfrak{U}) > 0$ , and for each  $t \in (\nu(\mathfrak{U}), 1]$ , there exists a unique  $\mathcal{N} \in \text{supp}(\nu)$  such that

$$\nu(\mathcal{M} : \mathcal{M} \supset \mathcal{N}) = t.$$

Once the primary preference  $\succsim$  and set of justifications  $\mathcal{M}$  are realized, a choice is made according to the deterministic EU model in Section 1.3.4.

**Definition 20** (Random justification (RJ) representation).  $(\mu, \nu)$  is an RJ representation for  $\rho$  if  $\mu$  is a preference distribution,  $\nu$  is a constraint distribution, and

$$\rho(p|A) = \int_{\mathfrak{U}} \int_{\mathfrak{U}} \mathbf{1}\{p \in \arg \max(\succsim_m, \mathcal{M}(A))\} d\nu(\mathcal{M}) d\mu(u)$$

$$\text{where } \mathcal{M}(A) := \bigcup_{\succsim_m \in \mathcal{M}} \arg \max(\succsim_m, A).$$

### 1.5.1 PROPERTIES

Recall from Sections 1.3 and 1.4 that the deterministic justification models violate WARP: unchosen alternatives can affect choice by restricting the set of justifiable alternatives. It will not be surprising that RJ violates Regularity, the stochastic analogue of WARP.

**Definition 21** (Regularity). *For all  $p \in \Delta(Z)$  and all  $A, B \in \mathcal{F}(\Delta(Z))$ ,*

$$\rho(p|A) \geq \rho(p|A \cup B).$$

A simple type of Regularity violation recurs systematically within RJ. Example 7 illustrates.

**Example 7.** *Suppose that the DM can keep \$10 ( $p$ ) or donate it to Charity A ( $q$ ) or Charity B ( $\tilde{q}$ ). Both charities represent good causes, neither of which is obviously more pressing than the other. In RJ, two groups of DMs choose  $q$  from  $\{p, q, \tilde{q}\}$ : DMs with  $q \succ p, \tilde{q}$ , and DMs with  $p \succ q \succ \tilde{q}$  who feel unable to choose  $p$ . Now suppose that the donation to Charity B is implemented with error, so the DM is sometimes able to keep the money. Formally,  $\tilde{q}$  is replaced with  $\alpha\tilde{q} + (1 - \alpha)p$ , for  $\alpha$  close to 1. DMs with  $q \succ p, \tilde{q}$  continue to choose  $q$ , but some DMs with  $p \succ q \succ \tilde{q}$  now choose  $\alpha\tilde{q} + (1 - \alpha)p$  when they feel unable to choose  $p$ . Intuitively, mixing virtuous option  $\tilde{q}$  with  $p$  makes it more attractive to some DMs with less-than-virtuous preferences, reducing the probability of choosing the other virtuous option  $q$ .*

*This is not quite a Regularity violation because  $\alpha\tilde{q} + (1 - \alpha)p$  does not belong to the original menu  $\{p, q, \tilde{q}\}$ . However,  $\alpha\tilde{q} + (1 - \alpha)p$  can be added to that menu without changing choice. (In RJ, adding items in the convex hull of a menu never affects choice on that menu.) This delivers the required Regularity violation.*

Definition 22 and the first part of Proposition 5 show that Example 7 generalizes. Say that  $q$  is more virtuous than  $p$  if it is sometimes unjustifiable to choose  $p$  over  $q$ , i.e.  $\nu(\{\mathcal{M} : p \notin \mathcal{M}(\{p, q\})\}) > 0$ . (Going forward, we abbreviate this as  $\nu(p \notin \mathcal{M}(\{p, q\})) > 0$ .) For any two lotteries  $p$  and

$q$  such that  $q$  is more virtuous than  $p$ , there is a pair of nested menus along the lines of Example 7 such that  $q$  is chosen more often from the larger menu. The intuition is always the same. When one virtuous option is implemented with error (so a DM who selects it sometimes ends up with a less virtuous option), constrained DMs are more likely to select it. They substitute away from the competing virtuous option, causing a Regularity violation.

The second part of Proposition 5 says that Regularity violations along the lines of Example 7 can be used to tell which of two lotteries is the more virtuous. Lottery  $q$  is more virtuous than lottery  $p$  if and only if the smallest  $\mathcal{M}$  in the support of  $\nu$  does not justify choosing  $p$  over  $q$ . Thus, Regularity violations can be used to identify this  $\mathcal{M}$ , which represents the strictest notion of virtue present in the population being studied.

**Definition 22** (Anomalous).  $(p, q) \in \Delta(Z)^2$  is anomalous if, for every  $\varepsilon > 0$ , there exist  $\tilde{q} \in B_\varepsilon(q)$  and  $\alpha, \lambda \in (0, 1)$  such that

$$\rho(q | \{p, q, \alpha\tilde{q} + (1 - \alpha)p, \lambda\tilde{q} + (1 - \lambda)q\}) < \rho(q | \{p, q, \tilde{q}\}). \quad (1.2)$$

**Proposition 5.** Suppose that  $\rho$  has an RJ representation  $(\mu, \nu)$ , where  $\mu$  has full support. For any  $(p, q) \in \text{int}(\Delta(Z))^2$ :

1. If  $\nu(p \notin \mathcal{M}(\{p, q\})) > 0$ , then  $(p, q)$  is anomalous.
2. If there exists  $\varepsilon > 0$  s.t.  $(p, \tilde{q})$  is anomalous for all  $\tilde{q} \in B_\varepsilon(q)$ , then  $\nu(p \notin \mathcal{M}(\{p, q\})) > 0$ .

The second part of Proposition 5 is used as a Lemma in the proof of Theorem 5, which establishes the uniqueness of the RJ representation.

**Theorem 5.** Any RJ representation  $(\mu, \nu)$  with full-support  $\mu$  is unique.

In the absence of Theorem 5, an analyst who suspects that her data is consistent with RJ can do only two things. She can show that the data includes Regularity violations, so it is inconsistent with



random EU. Second, she can try to assess the extent of the violations by computing

$$\rho(p|B) - \rho(p|A)$$

whenever this quantity is positive and  $p \in A \subset B$ . Theorem 5 shows that she can go beyond these simple bounds. For each menu, she can recover precisely the proportion of DMs who are maximizing their primary preferences, and she can tell what choice would be in a hypothetical world in which all DMs maximize their primary preferences.

The proof of Theorem 5 provides an explicit, relatively simple procedure for recovering the components of the RJ representation. It proceeds in three steps. First, for each  $p \in \Delta(Z)$ , we identify

$$D(p) := \arg \max_{x \in \Delta(Z)} \nu(x \notin \mathcal{M}(\{p, x\})).$$

The set  $D(p)$  reveals the most permissive notion of acceptable behavior present within the society being studied, just as  $\{x \in \Delta(Z) : \nu(x \notin \mathcal{M}(\{p, x\})) > 0\}$  reveals the strictest notion of acceptable behavior. It is easier to identify the rest of  $\nu$  once these bounds are known.

We exploit Regularity violations to identify  $D(p)$ . If  $q \in D(p)$ , then  $q$  is the first item to become unjustifiable in any menu containing both  $p$  and  $q$ . Thus, a constrained DM can never choose  $q$  over  $p$ . Since Regularity violations come from DMs who face binding constraints, it is not possible to find menus  $A, B$  such that  $\{p, q\} \subset A \subset B$  and  $\rho(q|A) < \rho(q|B)$ . This is always possible if  $q \notin D(p)$ , as we can construct menus containing  $p$  and  $q$  in which  $q$  is selected by some constrained DMs.

The second step identifies  $\nu$  given the “bounds” from the first step. Like Proposition 5, it uses simple three- or four-element menus. The trick is to identify two groups of DMs who have exactly the same (distribution of) preferences, but different sets of justifications. Since preferences are held

constant, any difference in behavior across the two groups must be attributed to justifications. This allows us to pin down  $\nu(q \notin \mathcal{M}(\{p, q\}))$  for each pair of lotteries  $p, q$ . Since the support of  $\nu$  is ordered by set inclusion, these probabilities pin down  $\nu$  itself.

The final step identifies  $\mu$  given  $\nu$ . This is done by “correcting”  $\rho(p|A)$  for the DMs who choose  $p$  because they are constrained, leaving only the DMs who genuinely prefer  $p$ . To see how this works, fix a three-element menu  $A$ , and index the elements from hardest-to-justify to easiest-to-justify. (Again, this is possible because the support of  $\nu$  is ordered by set inclusion.) Identifying  $\mu(p_1 \succsim A)$  is easy because the only DMs who choose  $p_1$  from  $A$  are those who like  $p_1$  best and are able to justify it. We have already identified the probability that  $p_1$  is justifiable, so we can identify the probability that it is best. Identifying  $\mu(p_2 \succsim A)$  is slightly more complicated because two groups of DMs choose  $p_2$  from  $A$ . The first group consists of DMs who like  $p_2$  best and are able to justify it. The second group consists of DMs who like  $p_1$  better than  $p_2$  better than  $p_3$ , and are able to justify  $p_2$  but not  $p_1$ . We already know the probability that  $p_2$  is justifiable but  $p_1$  is not. To identify  $\mu(p_1 \succ p_2 \succ p_3)$ , notice that DMs with this preference pattern substitute away from  $p_2$  when  $p_1$  is added to the choice set:

$$\rho(p_2|\{p_2, p_3\}) - \rho(p_2|A) = \mu(p_1 \succ p_2 \succ p_3)\nu(\{\mathcal{M} : p_1 \in \mathcal{M}(A)\}).$$

Deducting the DMs in the second group from  $\rho(p_2|A)$  (and adjusting for the probability that  $p_2$  is justifiable), we recover  $\mu(p_2 \succsim A)$ . We now know the full distribution of preferences on  $A$ . The procedure is essentially the same for larger menus, although it must account for additional groups of DMs.

## 1.5.2 EXTENSION: INFORMATION CHOICE

### MODEL OF INFORMATION

Some of the best-known evidence that people exploit “moral wiggle room” comes from experiments that allow subjects to acquire (or avoid) information before making a decision. Several experiments find that subjects fail to acquire, or even pay to avoid, information that would have affected their ultimate decision. The classic example is [Dana et al. \(2007\)](#), in which each DM was told that his interests could be aligned with or opposed to the interests of another subject. DMs who could not avoid learning the state, and found that the conflict state had been realized, nearly always chose unselfishly. DMs who could avoid learning the state did so about 40% of the time, and nearly always chose selfishly.

The results of [Dana et al. \(2007\)](#) raise two questions about the information demand of justifiers. First, can information be used to bring people’s behavior into alignment with their own notions of virtue? [Dana et al. \(2007\)](#) found that it could, and the effects were diminished but not eliminated when people could choose not to look at the information. Proposition 7 says that these results generalize: in RJ, there are almost always opportunities to improve average behavior through information, although some individuals will avoid virtue-promoting information when they have the opportunity. Second, is virtue best promoted by encouraging people to take as much information as possible? Proposition 8 says that the answer is no: an injudicious choice of information may make behavior unambiguously worse by providing new justifications for bad behavior. This possibility is typically ignored in the literature on information and moral wiggle room, although there is empirical evidence of it. The generality of RJ is helpful here; a model tailored to a particular situation would not be able to capture both the “moral suasion” and “excuse” functions of information.

We now formally extend RJ to information choice. We restrict attention to a simple kind of information, designed to match a typical experimental setup. The DM faces a binary choice set. In the

first stage, he is offered a single (binary) signal about one of the items in the set. If he accepts the signal, he observes it before making his selection in the second stage. Definitions 23 and 24 formalize this setup.

**Definition 23** (Signal).  $(q_1, q_2, \alpha) \in \Delta(Z)^2 \times [0, 1]$  is a signal about  $q \in \Delta(Z)$  if  $q = \alpha q_1 + (1 - \alpha)q_2$ .

**Definition 24** (Information choice problem). An information problem is

$$\{\delta_{\{p,q\}}, \alpha\delta_{\{p,q_1\}} + (1 - \alpha)\delta_{\{p,q_2\}}\}, \quad (1.3)$$

where  $(q_1, q_2, \alpha)$  is a signal about  $q$ , and  $p, q \in \Delta(Z)$ .

Perhaps the most natural way to extend RJ to information choice is as follows. Recall that a DM endowed with menu  $\{p, q\}$  limits himself to

$$\mathcal{M}(\{p, q\}) = \bigcup_{m \in \mathcal{M}} \arg \max_{\{p,q\}} m.$$

This allows him to pool with DMs who have preferences in  $\mathcal{M}$ . Now consider a DM who faces information problem (1.3). If the DM chooses  $\delta_{\{p,q\}}$ , then he can only pool with DMs who have preferences in

$$\mathcal{M}^{\text{avoid}} := \left\{ m \in \mathcal{M} : \arg \max_{\{p,q\}} m = \alpha \arg \max_{\{p,q_1\}} m + (1 - \alpha) \arg \max_{\{p,q_2\}} m \right\}. \quad (1.4)$$

He must therefore limit himself to  $\mathcal{M}^{\text{avoid}}(\{p, q\})$  at the second stage. (If  $\mathcal{M}^{\text{avoid}} = \emptyset$ , the DM cannot choose  $\delta_{\{p,q\}}$  in the first place, as doing so would automatically separate him from all DMs with preferences in  $\mathcal{M}$ .)

Unfortunately, Proposition 6 shows that this version of RJ cannot generate a strict preference for information avoidance. In some cases, the DM feels compelled to become informed because he believes every virtuous person would do so. In all other cases, he is free to forego information, but does not gain by it. Intuitively, the DM cannot use information to get closer to his second-stage ideal if information choice and second-stage choice are subject to the same constraints.

**Proposition 6.** *For any justification representation  $(u, \mathcal{M})$  and any information problem  $(p, q_1, q_2, \alpha)$ : either the set  $\mathcal{M}^{avoid}$  given by (1.4) is empty, or*

$$\alpha \max_{\mathcal{M}(\{p, q_1\})} u + (1 - \alpha) \max_{\mathcal{M}(\{p, q_2\})} u \geq \max_{\mathcal{M}^{avoid}(\{p, q\})} u.$$

To account for the empirical prevalence of information avoidance, the constraints on information choice must be relaxed. In the body of the paper, we adopt the simplifying assumption that information choice is fully unconstrained: the DM chooses information to maximize his expected utility, given his beliefs about his own second-stage behavior. Appendix A.2.4 presents a more general model that allows information choice to be *less* constrained than second-stage choice without being fully unconstrained. As shown in the Appendix, the results in this section carry over to the more general model.

To complete the model, we still need to specify the DM's first-stage beliefs about his own second-stage behavior. We assume that he knows his primary utility in both stages, but may have imperfect information about his justifications until he reaches the second stage. This allows for the possibility that the set of justifications is partly influenced by passing sentiments, such as feelings of sympathy. We refer to DMs' first-stage information about their second-stage justifications as "self-knowledge."

**Definition 25** (Self-knowledge). *Self-knowledge about constraint distribution  $\nu$  is  $N \in \Delta^2(\mathfrak{U})$ , where*

$$\int_{\tilde{\nu}} \tilde{\nu}(\cdot) dN_{\nu}(\tilde{\nu}) = \nu(\cdot).$$

Ties are inevitable in the first stage because the DM must be indifferent to information that does not affect his behavior. Since we are interested in a strict preference for information avoidance, we assume that each DM breaks ties in favor of becoming informed. Information “avoidance” that results from indifference is thereby ruled out. This makes no difference to the results.

We are now ready to extend RJ to information choice. Fix a preference distribution  $\mu$  and a self-knowledge  $N$ . Let  $U_{(u,\tilde{v})}$  denote the expected utility of a DM with primary utility  $u \in \text{supp}(\mu)$  and beliefs  $\tilde{v} \in \text{supp}(N)$ . We use these expected utilities to define a stochastic choice function  $\rho_{(\mu,N)}$  on the set of information choice problems:

$$\begin{aligned} \rho_{(\mu,N)}(\delta_{\{p,q\}} | \{\delta_{\{p,q\}}, \alpha\delta_{\{p,q_1\}} + (1-\alpha)\delta_{\{p,q_2\}}\}) \\ := \int_u \int_{\tilde{v}} 1 \{U_{(u,\tilde{v})}(\delta_{\{p,q\}}) > U_{(u,\tilde{v})}(\alpha\delta_{\{p,q_1\}} + (1-\alpha)\delta_{\{p,q_2\}})\} dN(\tilde{v}) d\mu(u). \end{aligned}$$

To assess the effects of information, it is helpful to have a standard for desirable behavior. To this end, we introduce a social planner who must decide how much information to provide to her society. The planner has her own Bernoulli utility  $s \in \mathbb{R}^Z$ , which need not match that of any DM. For instance, the planner could seek to maximize the expected value of charitable donations from a society of selfish DMs. For each dilemma  $\{p, q\}$  confronting members of her society, the planner has three options: provide no information, design a signal about  $q$  and provide it for free, or design a signal about  $q$  and require everyone to observe it. The planner knows the preference distribution  $\mu$  and self-knowledge  $N$  that characterize her society, and maximizes her expected utility  $S_{(\mu,N)}$  given her (rational) expectations of behavior.

## INFORMATION DEMAND AND AVOIDANCE

Despite the simplicity of the setup, RJ allows for complex attitudes to information. We use Example 8, which will recur throughout this section, to illustrate.

**Example 8.** *This is a simplified version of the experiment in Fong and Oberholzer-Gee (2011). Let*

$p = \text{experimenter pays \$10 to DM}$

$q_1 = \text{experimenter pays \$10 to physically disabled poor person}$

$q_2 = \text{experimenter pays \$10 to poor person with drug addiction}$

$$q = \frac{1}{2}q_1 + \frac{1}{2}q_2.$$

*Consider a DM who prefers  $p$  to  $q$ , but may feel that it is unjustifiable to choose  $p$  over  $q$ .*

Information may affect the DM in Example 8 in two ways. First, she may feel able to keep the money more or less often on average. Second, the relative weights on the two recipient groups may change, as the DM may feel compelled to donate to one type of recipient more frequently than the other. The directions of these two effects may be hard to predict. Even if the directions are known, it may be hard to predict which effect will dominate and, by extension, whether the DM will choose to become informed. Finally, it may not be clear whether a virtuous social planner would want the DM to become informed in the first place. For instance, information might reduce total donations but raise donations to a favored recipient group.

In fact, experimental subjects do exhibit complex attitudes to information, and informational interventions often have mixed effects. The goal of this section is not to provide sharp predictions for all or even most information choice problems, but to identify two broad patterns of information choice within RJ. The first pattern is moral suasion: a social planner for a population of RJ agents can further her objectives just by informing people about an option in their choice sets. More specifically, a planner who prefers more virtuous option  $p$  to less virtuous option  $q$  can design a binary signal about  $q$  that she would like every DM who prefers  $q$  to  $p$  to acquire. Surprisingly, some DMs will voluntarily acquire appropriate information even though they disagree with the planner on the

original choice set. The planner benefits from providing this information even if she cannot require anyone to pay attention to it.

**Proposition 7.** *Fix a constraint distribution  $\nu$ , lotteries  $p, q \in \text{int}(\Delta(Z))$  such that  $\nu(\mathcal{U}) < \nu(q \in \mathcal{M}(\{p, q\})) < 1$ , and a social planner  $s$  such that  $s(p) > s(q)$ . There exists a signal  $(q_1, q_2, \alpha)$  about  $q$  such that*

$$S_{(\mu, N)}(\alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}}) > S_{(\mu, N)}(\{\delta_{\{p, q\}}, \alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}}\}) > S_{(\mu, N)}(\delta_{\{p, q\}})$$

for any preference distribution  $\mu$  such that  $\text{supp}(\mu) = \{u \in \mathcal{U} : u(q) \geq u(p)\}$  and any self-knowledge  $N$ .

We use Example 8 to illustrate. Consider a social planner who cares only about expected donations to disabled recipients, so  $s(q_1) > s(q) > s(q_2) = s(p)$ . Plausibly, the obligation to donate to a disabled person is stronger than the obligation to donate to a drug addict, so

$$\nu(p \notin \mathcal{M}(\{p, q_1\})) > \nu(p \notin \mathcal{M}(\{p, q\})) > \nu(p \notin \mathcal{M}(\{p, q_2\})).$$

It is easy to see that the social planner is better off if the DM learns the recipient's type. If  $q_2$  is realized, the planner doesn't care what the DM does; if  $q_1$  is realized, the planner benefits because the DM is more likely to make a transfer than he would if he were ignorant. To see why some DMs will learn the recipient's type voluntarily, consider a DM who has a slight preference for  $p$  over  $q_1$ , but a strong preference for  $q_1$  over  $q_2$ . Since this DM is primarily concerned with avoiding a transfer to a drug addict, he prefers to become informed even if he must sacrifice some chance of keeping the money. By contrast, a DM who does not care about the recipient's type will actively avoid any information that reduces his chance of keeping the money.

The findings of [Fong and Oberholzer-Gee \(2011\)](#) are aligned with the above, although they use a



more complex setup in which DMs can transfer any part of \$10 to the recipient. Expected transfers to disabled recipients do rise from the no-information treatment (\$3) to the full-information treatment (\$4.30). In a third treatment, DMs could pay \$1 to learn the recipient's type before making a transfer decision. This setup does not quite match Proposition 7, in which information is freely available. Fong and Oberholzer-Gee (2011) found that two-thirds of DMs declined to pay for information and transferred \$2 on average, while the remaining one-third paid for information and transferred \$4.50. As a result, average transfers to disabled recipients were \$2.80, slightly lower than in the no-information case. This was likely because DMs passed most of the information cost on to the recipient. Fong and Oberholzer-Gee (2011) estimated that having \$9 rather than \$10 reduced transfers by \$0.80 on average. Correcting for this raises average transfers to disabled recipients to \$3.10, slightly higher than the no-information case. This is probably an underestimate because more DMs would have acquired information (and raised their donations to disabled recipients) if information had been free.

Many other experiments, including Ehrich and Irwin (2005) on ethical consumption, Dana et al. (2007) and Grossman and Van Der Weele (2017) on sharing with others, Serra-Garcia and Szech (2019) on donations to charity, Kajackaite (2015) on donations to a lobbying organization, and Woolley and Risen (2018) on task choice, study the effects of information about a possible negative externality to an appealing action. The informational interventions in these papers encourage virtuous behavior: the externality is mitigated in full-information treatments relative to no-information treatments.<sup>8</sup> As predicted in Proposition 7, a substantial fraction of subjects seem to anticipate the effects of information, and avoid it when they have the chance. This erodes the benefits of the informational intervention.

Generalizing from these findings, one might conclude that information is a force for virtue: peo-

---

<sup>8</sup>Kajackaite (2015) is an exception: she finds that subjects are very conservative when they do not know whether the externality has been realized, so they act like subjects who know there is a negative externality.

ple behave better when they are better informed about the consequences of their actions. However, it is important to remember that the informational interventions discussed above were *designed* to promote virtuous behavior and provoke avoidance. The next result, Proposition 8, shows that information can have exactly the opposite effect. There are signals that every DM with less-than-virtuous preferences would like to acquire, but every social planner with virtuous preferences would like to withhold. This type of information is not limited to contrived examples; it exists whenever virtuous people disagree. Example 9 illustrates the formal notion of disagreement, which is given in Definition 26.

**Example 9.** *This example is based on Norton et al. (2004). The DM must hire someone for a managerial role at a construction company. Candidates are evaluated on both education and experience. Let*

*$p$  = male candidate*

*$q_1$  = female candidate with more education, less experience than  $p$*

*$q_2$  = female candidate with more experience, less education than  $p$*

$$q = \frac{1}{2}q_1 + \frac{1}{2}q_2.$$

*Since women have been historically underrepresented in this field, the DM may believe that an unbiased person would give female candidates a slight edge. This makes it unjustifiable to choose  $p$  over  $q$ , but not to choose  $p$  over  $q_1$  or  $q_2$ . An unbiased person who considers education more important than experience might well choose  $p$  over  $q_2$ , while one who considers experience more important than education might choose  $p$  over  $q_1$ .*

**Definition 26** (Disagreement). *Constraint distribution  $v$  exhibits disagreement about  $(q_1, q_2, p)$  if*

there exists  $\mathcal{M} \in \text{supp}(\nu)$  such that the sets

$$\{m \in \mathcal{M} : m(q_1) \geq m(p)\} \text{ and } \{m \in \mathcal{M} : m(q_2) \geq m(p)\}$$

are nonempty and disjoint.

**Proposition 8.** *If constraint distribution  $\nu$  exhibits disagreement about lotteries  $(q_1, q_2, p)$ , there exist  $\alpha \in (0, 1)$  and  $q \in \Delta(Z)$  such that  $(q_1, q_2, \alpha)$  is a signal about  $q$ , and*

$$S_{(\mu, N)}(\delta_{\{p, q\}}) > S_{(\mu, N)}(\{\delta_{\{p, q\}}, \alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}}\}) = S_{(\mu, N)}(\alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}})$$

for any social planner  $s$  such that  $s(p) > \max\{s(q_1), s(q_2)\}$ , any preference distribution  $\mu$  such that  $\mu(\{u \in \mathcal{U} : \min\{u(q_1), u(q_2)\} > u(p)\}) = 1$ , and any self-knowledge  $N$ .

In Example 9, it is not difficult to see why information might have a pernicious effect. A DM who wishes to hire a male candidate regardless of qualifications may choose  $q$  over  $p$  to avoid revealing his bias. However, he can justify choosing  $p$  over  $q_1$  by arguing that experience is more important, and  $p$  over  $q_2$  by arguing that education is more important. This is exactly what Norton et al. (2004) find in their study of hypothetical hiring decisions. Subjects picked the male candidate 66% of the time (75% when the male candidate was more educated, and 57% when the male candidate was more experienced). Norton et al. (2004) also asked subjects to explain their decisions. When the male candidate was more educated (and in a control condition without gender), half of subjects mentioned that they considered education more important than experience. This proportion dropped below one-quarter when the female candidate was more educated, suggesting that some subjects used experience as an excuse to choose their favored candidate.

## 1.6 RELATED MODELS

The justification model with unknown primary preference (JU) is closely connected to the model of attention in Masatlioglu et al. (2012) (MNO), and to the model of rationalization in Cherepanov et al. (2013b) (CFS). Both models use a two-tier structure in which a preference breaks ties on a consideration set. MNO interpret the consideration set as the subset of items to which the DM pays attention. The key restriction in their model is that removal of an item outside the consideration set does not change choice. This property holds in JU as well, so JU is a special case of MNO. CFS interpret the consideration set as the subset of items the DM can rationalize. The difference between CFS and JU is that CFS take rationales to be unstructured binary relations rather than preferences. Thus, JU is a special case of CFS. Since CFS can impose transitivity of rationales without loss of generality, completeness is really the distinguishing factor. While this may seem like a technical distinction, the requirement that justifications be preferences underpins the primary interpretation of JU: people justify their decisions by pretending to be better versions of themselves.

For a more formal comparison, consider the revealed preference relation from each model.

**Definition 27** (Revealed preference in CFS, MNO).

1. If  $a \neq c(A \cup \{a\})$  and, for some  $B \supseteq A$ ,  $c(B) \neq c(B \cup \{a\})$ , then  $c(A \cup \{a\}) R^{JU} a$ .
2. If  $a \neq c(A \cup \{a\})$  and, for some  $B \supset A$ ,  $a = c(B \cup \{a\})$ , then  $c(A \cup \{a\}) R^{CFS} a$ .
3. If  $a \neq c(A \cup \{a\})$  and  $c(A \cup \{a\}) \neq c(A)$ , then  $c(A) \cup \{a\} R^{MNO} a$ .

The revealed preference relation has the same interpretation in all three models:  $(a, b)$  belongs to the transitive closure of the revealed preference relation if and only if  $(a, b)$  belongs to the primary preference relation in all representations for  $c$ . It is not difficult to see that  $R^{JU} \supseteq R^{MNO}$  and  $R^{JU} \supseteq R^{CFS}$ . We claim that the inclusion is strict on any dataset that violates WARP. Since any such dataset

contains a cycle or an almost-WARP set, it suffices to show that neither MNO nor CFS delivers full identification on either pattern. Consider a cycle:

$$a = c(\{a, b\}) \quad b = c(\{a, b, d\}) = c(\{b, d\}) \quad d = c(\{a, d\}).$$

The three revealed preference relations are given by

$$a R^{CFS} b$$

$$b R^{MNO} d$$

$$a R^{JU} b R^{JU} d$$

Now consider an almost-WARP set, where the items are indexed from pairwise-best ( $a_1$ ) to pairwise-worst ( $a_n$ ). Suppose that  $a_i = c(\{a_1, \dots, a_n\})$ , where  $i > 1$ . Then,  $R^{CFS}$  and  $R^{MNO}$  are given by

$$a_j R^{CFS} a_i \text{ for all } j < i$$

$$\text{if } i \neq 2 \quad a_i R^{MNO} a_j \text{ for all } j \neq i$$

$$\text{if } i = 2 \quad a_2 R^{MNO} a_j \text{ for all } j > 2$$

Since  $R^{JU} \supseteq R^{MNO} \cup R^{CFS}$ , the only case in which  $R^{JU}$  is acyclic is  $i = 2$ , in which case

$$a_1 R^{JU} a_2 \cdots a_{n-1} R^{JU} a_n$$

Thus, JU delivers strictly stronger identification on any dataset that is inconsistent with preference maximization.

In addition to providing stronger identification, JU places stronger restrictions on the data than

CFS and MNO. JU is falsifiable with as few as three alternatives, but CFS and MNO are not. To see why, consider an almost-WARP set with three alternatives in which  $a_3 = c(\{a_1, \dots, a_3\})$ . CFS and MNO give opposing interpretations of this case, so JU rules it out.

From this example, the reader may wonder whether JU is simply the intersection of CFS and MNO. This is not the case: on almost-WARP sets with more than three elements, JU delivers stronger identification than CFS and MNO put together. Moreover, on choice domains with more than four elements, JU imposes stronger restrictions on the data than CFS and MNO put together. Thus, the collection of datasets consistent with JU is the intersection of the collections consistent with CFS and MNO for  $|\mathcal{A}| \leq 4$ , but is a strict subset of the intersection for larger domains.

There is an obvious connection between Corollary 3 of this paper and Proposition 3 of CFS. Both results provide unique “canonical” representations in which the constraints on the DM are minimized. Both state that the primary preference in the canonical representation is pinned down by pairwise choice when choice data is acyclic. The key distinction is that Proposition 3 of CFS *requires* choice data to be acyclic, while Corollary 3 applies to all choice data consistent with JU. As shown in Section 1.4.2, cycles play an important role in JU: they occur whenever one item is revealed unjustifiable in the presence of another.<sup>9</sup> Moreover, the leading empirical example in CFS is a cycle. For these reasons, the assumption of acyclicity is not innocuous.

JU is also a special case of the model of rationalization in Kalai et al. (2002). That model defines the consideration set in the same way as JU, but does not restrict choice from the consideration set. Kalai et al. (2002) show that this model has no empirical content. Thus, the empirical content of JU comes from the tiebreaking “primary” preference. Kalai et al. (2002) consider another way of introducing empirical content by restricting the number of rationales. That route is not taken in this paper, as there is nothing particularly implausible about a justification model with many justi-

---

<sup>9</sup>Almost-WARP sets occur when one item is revealed unjustifiable in the presence of a non-singleton *set* of items.

fications. First, wherever “good and reasonable people” exhibit subtle differences of opinion, many preferences will count as justifiable. Second, more justifications correspond to fewer constraints on the DM, so a representation with large  $\mathcal{M}$  is arguably more parsimonious.

Several existing models are special cases of JU/JO. The dynamic choice model in [Gul and Pesendorfer \(2005\)](#) is JU with a single justification (or a set of justifications that can be collapsed into a single weak preference).<sup>10</sup> The model in [Aizerman and Malishevski \(1981\)](#) is JU with a constant tiebreaker, so  $c(A) = \mathcal{M}(A)$ . [Aizerman and Malishevski \(1981\)](#) offer a behavioral characterization of their model, which may be reinterpreted as a characterization of the justifiable sets in JU. The model of willpower in [Masatlioglu et al. \(2020\)](#) is a special case of JO. Like JO, the willpower model uses  $(\succsim, c)$  as the primitive. [Masatlioglu et al. \(2020\)](#) show that their model is characterized by IUA and other conditions. Given the structural differences between the two models, this commonality is surprising. The consideration set in the willpower model is characterized by a temptation utility and a willpower cutoff, and it is not immediately obvious how to translate these into justifications.

Several other models of two-tier decision making are related to, but not nested with, JU. The model of sequential rationalization in [Manzini and Mariotti \(2007\)](#) uses one asymmetric binary relation to pin down a consideration set, and another to break ties. [Manzini and Mariotti \(2007\)](#) show that their model satisfies a property they call “Always Chosen:”  $a$  is chosen from  $A$  if it pairwise-defeats every other item in  $A$ . JU violates this property when  $|\mathcal{M}| > 1$  because different binary choices may appeal to different justifications.

The model of selfishness in [Dillenberger and Sadowski \(2012\)](#) pits the DM’s primary preference against a subjective norm, but allows the DM to override the norm when his preference is strong enough. The model satisfies Always Chosen because the morally *best* item in the choice set is the only impediment to utility maximization. It does not matter whether the choice set contains many

---

<sup>10</sup>The interpretation is the opposite of JU, though. [Gul and Pesendorfer \(2005\)](#) focus on situations in which the “justification” is shortsighted/tempted and the “primary preference” is forward-looking/rational.

morally good items, or only one. [Masatlioglu et al. \(2020\)](#) satisfies Always Chosen for a similar reason.

[Cunningham and de Quidt \(2015\)](#) (CD) consider a domain in which each item is characterized by a bundle of binary attributes. They distinguish “explicit” preferences over the domain from “implicit” preferences over the attributes. Explicit preferences are independent of the choice set, while implicit preferences may be activated to varying degrees by different choice sets. Both types of preferences are aggregated in a linear utility function. Note that the explicit preference in CD is more like a justification than a primary preference: the DM departs from it more when he can conceal his motivations for doing so. In fact, there is no primary preference in CD. While the signs of the implicit preferences are fixed, the magnitudes are context-dependent. CD provide a series of tests that an analyst can use to recover the signs of the implicit preferences. These tests are similar in spirit to the patterns of choice highlighted in Section 1.4.2, but CD do not offer a characterization result analogous to Theorem 4.

The proliferation of deterministic models of two-tier decision making is not replicated on the stochastic side. [Cattaneo et al. \(2020\)](#) study a population of decision makers who pay attention to different sets of alternatives. Attention satisfies a monotonicity property: alternatives are more likely to be attended to in smaller sets than in larger sets. Justifiable sets in RJ satisfy monotonicity. The key distinction between RJ and [Cattaneo et al. \(2020\)](#) (as well as [Manzini and Mariotti 2014](#) and [Brady and Rehbeck 2016](#)) is that RJ allows for preference heterogeneity as well as heterogeneity in justifications, and fully recovers both. These improvements do not come for free, though: they require the domain to be a set of lotteries, and for preferences to have an EU structure.



## 1.7 CONCLUSION

The family of models in this paper make several points that may be of interest to empirical researchers. First, the puzzling behaviors well documented in moral domains are likely to extend to some non-moral domains. In fact, the same patterns of behavior should be present in any choice setting that invites conflict between primary preferences and social image or self-image. Social image is determined by many factors besides altruism, including rationality, courage, self control, good taste, work ethic, and other virtues. To the author's knowledge, [Woolley and Risen \(2018\)](#) is the only paper to look for "wiggling" behavior in some of these domains.

Second, caution is warranted when interpreting "direct" measures of pro-sociality. A subject who reports a high willingness to donate to charity may be highly constrained rather than highly altruistic. The same is true of a subject who reports that fairness or generosity is central to her sense of self. Putting too much trust in these measures may lead to misleading conclusions. For instance, subjects might seek out information in hopes that it will free them from pressing obligations, not because they place a high value on doing good. A more reliable picture of subjects' pro-sociality can be obtained by offering opportunities to "wiggle" out of obligations. Even so, subjects with a particularly strong sense of obligation may resist these opportunities.

Third, little is known about the way justifiers respond to bundles of choices. [Gneezy et al. \(2019\)](#) and [Haisley and Weber \(2010\)](#) demonstrate that subjects exhibit consistency motives, at least when inconsistencies are sufficiently easy to spot. Subjects find it harder to inflict a bad outcome on someone else when they have already formed or expressed a negative opinion of that outcome. On the other hand, [Exley \(2016\)](#) demonstrates that significant inconsistencies can persist in within-subject data. It would be interesting to see how quickly consistency motives decay over time, and whether they are mitigated or eliminated by anonymity.

We conclude with a note on welfare implications. Although we have referred to the tiebreaker in

the justification model as the “primary preference,” we do not take it to be a reliable measure of welfare. One can always draw a distinction between choice behavior and welfare—even when choices maximize a standard preference relation—but the point seems particularly important when moral principles are involved. It is entirely plausible that people are better off when they choose in accordance with their principles rather than their baser inclinations. Even those who lack lofty moral principles may benefit from having a virtuous social image. The “primary preference” should therefore be interpreted as an interesting feature of a decision maker’s psychology, not the preference of a benevolent agent acting on his behalf.

# 2

## A Dynamic Model of Peer Effects

### 2.1 INTRODUCTION

Outside situations with a clear strategic component, a decision maker (DM) has two good reasons for attending to the choices of others.<sup>1</sup> First, the DM may use his observations to infer which op-

---

<sup>1</sup>For suggestions that substantially improved this paper, I am indebted to Krishna Dasaratha, Ben Golub, Jerry Green, Shengwu Li, Matthew Rabin, Debraj Ray, Collin Raymond, Marciano Siniscalchi, Tomasz Strzalecki, and several anonymous referees.

tion is best, or most closely aligned with his personal taste. Second, he may favor options that will give him a certain status among his peers. The DM might wish to appear particularly successful, responsible, or generous; to avoid appearing poor, lazy or selfish; or to simply fit a norm. He may be motivated by explicit rewards or sanctions, others' admiration or disapproval, or simply self image.

The main model in this paper captures all these mechanisms in reduced form. I model a DM who—for informational reasons, status concerns, or both—would like to condition his choice on the distribution of peers' choices. This is not straightforward because the DM does not know the distribution. He learns about it by observing, one by one, the choices of individual peers and updating his beliefs according to Bayes' rule. Although he cares ultimately about the peer group as a whole, not about the choice of any particular peer, his preferences respond to individual peers' choices through belief updating.

To motivate this setup, consider a diner who is paying his bill at a counter-serve restaurant. When his card is swiped, he is presented with several possible tip amounts. The diner's tip may be influenced by the tips of the diners who preceded him in line, particularly if he is unfamiliar with local tipping culture. This is not because he cares about the good opinion of those particular diners (who are probably not observing him anyway), but because their behavior is informative about what most people do. Notice that the diner can be influenced by others' choices without copying them. An undergraduate may raise his intended tip only slightly if he observes office workers leaving high tips. More starkly, a generous person may raise his intended tip if others' low tips make him feel sorry for the staff.

Although the tipping example provides an especially clean illustration of the setup, the model's applicability is much broader. The main model does not restrict the dependence of preferences on others' choices, so it can be used to study any (non-strategic) situation in which information about others' choices affects behavior. Section 2.2 reviews empirical evidence for these effects in diverse choice settings. Some evidence pertains to decisions that are minor from an individual point of view,

but consequential on aggregate, such as electricity conservation or voting. Other evidence pertains to decisions that are consequential even for the individual, such as retirement plan enrollment and stock market participation. Section 2.2 also contrasts the model in this paper with existing models of conformity, social learning, and belief updating.

Section 2.3 provides an axiomatic characterization of the main model. The broad strategy is to use bets on a randomly selected peer's choices to elicit beliefs about the true distribution of peers' choices. If the elicited beliefs turn out to be exchangeable, de Finetti's theorem and a martingale convergence theorem can be used to obtain the components of the representation. The key axiom is Forward Exchangeability, which effectively requires the DM to place equal probability on any permutation of a fixed sequence of future observations. The catch is that the choice domain is limited to bets on individual choices, not sequences. Imposing invariance to permutations is not straightforward when the DM has no obvious way to express beliefs about order. The solution is a form of dynamic consistency or no-arbitrage, but not one that has appeared elsewhere in the literature (to the author's knowledge). The main representation result is Theorem 6.

The model in Section 2.3 imposes no substantive restrictions beyond expected utility and Bayesian updating. Section 2.4 studies some non-parametric special cases that correspond to particular types of peer effects. All these special cases restrict the way that the perceived marginal return to some activity depends on the activity levels chosen by others. Proposition 9 provides conditions for the DM to raise (or lower) his action when he comes to believe that his peers are choosing higher actions than he originally thought. Similarly, Proposition 9 provides conditions for the DM to raise (or lower) his action when he becomes more certain about the actions chosen by his peers. Although these results are fairly straightforward applications of standard machinery, they help interpret some of the empirical evidence reviewed in Section 2.2. Corollary 6 translates these comparative statics on beliefs to predictions about the evolution of behavior over time.

Section 2.5 is an effort to distinguish the two interpretations mentioned above: information

extraction from others' choices vs. direct regard for others' choices. First, the model is extended to give formal meaning to both interpretations. In the pure learning version, the DM cares directly about maximizing his utility conditional on his preference parameter (type). He cares indirectly about the distribution of peers' choices because it reveals the mapping from types to optima. In the social version, the DM may care both indirectly and directly about the distribution of peers' choices.

Formally distinguishing the two models does not solve the problem of distinguishing them in practice. In fact, Proposition 11 shows that this problem is insoluble when the analyst can only observe preferences on the domain studied in Section 2.3. For any social model, it is possible to construct a pure-learning model that generates exactly the same behavior on this domain. The models can only be distinguished by appealing to additional information. The final result, Proposition 12, shows that the analyst can falsify the pure learning model if she can elicit (or manipulate) beliefs about the distribution of peers' types. The pure learning model, but not the social model, requires the sensitivity to peers' choices to decrease in uncertainty about the type distribution.

## 2.2 RELATED LITERATURE

### 2.2.1 EMPIRICAL WORK

This paper is motivated by empirical work that investigates the effect of information about peers' choices on subjects' own choices. This work covers a wide range of choice settings. It can roughly be divided into two groups: papers where subjects are directly provided with information about a peer group, and papers where subjects gather piecemeal information by observing others. The second group of papers is most closely aligned with the setup introduced in Section 2.3. The first group is useful because it provides direct evidence about the dependence of preferences on the true distribution of others' choices.

Numerous papers in the first group find evidence that subjects adjust their behavior to match

(or approach) the typical behavior of peers. Examples include Frey and Meier (2004) and Shang and Croson (2009) on small charitable donations, Bicchieri and Xiao (2009) on sharing in dictator games, Goldstein et al. (2008) on towel usage in hotels, Cai et al. (2009) on selections from a restaurant menu, Gerber and Rogers (2009) on voting, Chen et al. (2010) on voluntary provision of movie ratings, and Coffman et al. (2017) on acceptance of Teach for America job offers.

Of even greater interest are papers that find asymmetric and/or surprising responses to information about peer behavior, since these indicate a need for a flexible model of peer effects that can accommodate heterogeneous reactions to information. Allcott (2011) found that households with high baseline electricity usage substantially reduced it when they received a report comparing their usage to that of 100 neighboring households. Households with low baseline electricity usage were largely unaffected by the report. Card et al. (2012) found that university employees who learned that their salaries were below median for their pay unit and occupation reported lower job satisfaction and were more likely to look for a new job, while above-median earners reported no change. Relatedly, Beshears et al. (2015) found that a group of low-earning workers were *less* likely to participate in a retirement plan when they were told that a large fraction of similarly aged coworkers were already participating. The effect of peer information on other groups of workers was slightly positive or zero. Finally, in a study of Greek high school students, Goulas and Megalokonomou (2015) found that information about peers' scores on a nationwide exam harmed students who scored poorly, while benefiting students who did well. Specifically, low-scoring students who were able to observe peers' scores after the first round of the exam did worse on the second round than low-scoring students who were unable to observe peers' scores. The pattern was reversed for high-scoring students.

Turning now to the second group, several papers find evidence that people are influenced by individual neighbors' or coworkers' choices. Examples include Duflo and Saez (2003) on retirement plan enrollment, Brown et al. (2008) on stock market participation, Godlonton and Thornton

(2012) on HIV testing in Malawi, and [Ouimet and Tate \(2020\)](#) on participation in employee stock purchase plans. The study of asset purchases by [Bursztyn et al. \(2014\)](#) is of particular interest because it distinguishes pure learning from other types of peer effects, as in Section 2.5 of this paper. [Bursztyn et al.](#) accomplished this by rationing the asset, so that some individuals who wished to purchase the asset were ultimately unable to do so. Since the intention to purchase is what matters for inference, rationing leaves the pure learning channel unaffected. The study found that the two types of peer effects were of roughly equal importance.

The peer effects literature extends well beyond the two groups reviewed above. Many papers show a link between peer groups and outcomes (e.g. test scores or crime) but do not emphasize any particular decision or piece of information. In the interest of space, these papers will not be discussed here; see [Sacerdote \(2014\)](#) for a broader review.

#### 2.2.2 THEORETICAL WORK

The main theorem in this paper is essentially a characterization of Bayesian updating in a particular setting—specifically, when beliefs are elicited through bets and preferences over the prize space are allowed to depend on the distribution being learned. Other characterizations of Bayesian updating exist in different settings. [Majumdar \(2004\)](#) provides one when beliefs are directly observable, and new evidence serves to rule out some states of the world. [Majumdar](#) also states an axiom that requires beliefs to be independent of the order of signals, but the axioms in this paper do not resemble his because the settings are so different. [Cripps \(2018\)](#) provides a characterization of a broader class of updating rules that includes Bayesian updating.

[Billot et al. \(2005\)](#) studies the formulation of a Bayesian prior rather than the updating process. Beliefs are built up from individual observations, which may be more or less relevant to the situation at hand. [Billot et al.](#) also use the term “dataset,” but datasets have a different function in their paper and in this one. Here, the DM already has a Bayesian prior, which he uses to interpret all the obser-



vations he makes. Although he may find some observations more informative than others, there is no need for an explicit similarity relation on datasets.

[Safonov \(2017\)](#) studies Bayesian updating about preferences. The key difference there is the nature of signals. In [Safonov's](#) model, the framing of a decision problem triggers a signal that informs the DM about some feature of his own preferences. Both the frames and the signals are unobservable to the analyst, who sees only the distribution of choices. [Frick et al. \(2019\)](#) provide a different characterization of learning about one's own preferences, where choice does not depend on framing but may depend on (one's own) previous choices. The key axiom says that a DM who currently prefers  $x$  to  $y$  does not expect to flip to  $y$  in future. This type of condition does not appear here because the DM does not make choices over consumption streams. Another paper in this tradition, [Natenzon \(2019\)](#) studies a tractable special case of a learning model and applies it to menu effects. Peers' choices do not play an explicit role in any of these papers; they are models of introspection rather than learning about/from others.

The classic references on learning from others are [Banerjee \(1992\)](#) and [Bikhchandani et al. \(1992\)](#). These papers do not provide axiomatic characterizations of an individual DM's learning process. Instead, they focus on the way that diffusion of information shapes the behavior of a large group. The two approaches are suitable in different settings. This paper studies the way that an individual (or relatively small group) responds to information about a large group that is already in equilibrium. Others' choices are drawn from a stable distribution; the DM does not need to worry about the information that others have or might gather in future. This strong assumption is needed to impose some structure on an otherwise highly general problem. That being said, the "pure learning" version of the model in this paper could be used as a building block for future social learning models.

[Kamenica \(2008\)](#) presents a special case of a social learning model in which consumers learn about a "global preference parameter" by observing the range of available items. Learning dynamics are simplified because the product range is the only signal. [Kamenica's](#) goal is not an axiomatic

characterization of the learning process, but a Bayesian explanation for menu effects.

Unlike a DM in a social learning model, the DM in this paper may care about his relative status as well as the “intrinsic” value of the option he selects. The classic reference in this area is [Bernheim \(1994\)](#), which studies the tension between personal taste and conformist motives. [Bernheim’s](#) model does not include a pure learning channel or allow other social motivations (e.g. wanting to look better than others). Moreover, his focus is on equilibrium rather than learning dynamics.

[Chambers et al. \(2019\)](#) also study social influence in an equilibrium setting. There is no need for Bayesian updating in their model because agents directly observe the choice frequencies of their peers. Since the primary objective of their paper is to disentangle the ways in which multiple DMs influence one another, they impose strong parametric assumptions on social influence. In particular, they model social influence as a distortion from a default “isolation” preference. There is no isolation preference in this paper, as others’ choices are taken to be an integral part of utility and/or a source of valuable information. [Chambers et al.](#) do not distinguish between pure learning and status channels of social influence. In fact, the author is not aware of any preexisting models that elucidate this distinction.

## 2.3 MODEL AND CHARACTERIZATION

### 2.3.1 PRIMITIVES

Let  $X = \{x_1, \dots, x_n\}$  be a finite set of items or actions. We study the evolution of the decision-maker’s preferences over  $X$  as he learns about others’ choices from  $X$ . To that end, we must keep track of the DM’s information about others’ choices, and the beliefs informed by those choices. Information about others’ choices is captured by a dataset  $D \in \mathcal{D} := \{0, 1, \dots\}^X$ .  $D(i) = k$  means the DM has observed  $k$  peers choosing  $x_i$  from  $X$ . Notice that  $D$  does not record the order in which the observations were made. Since we are modeling a Bayesian DM, whose beliefs could not possibly

depend on the order of signals, information about order is not relevant.

The DM does not care directly about  $D$ , the choices from  $X$  that he happened to observe. Instead, he cares about the true distribution of others' choices. Since the DM's beliefs about that distribution are not observed by the analyst, we will need a richer domain than  $X$  to elicit them. Specifically, we allow the DM to place bets on the choice of a randomly selected peer. The payoffs of these bets will be objective lotteries over  $X \cup \{\diamond\}$ . The objective lotteries are a (completely standard) technical convenience. It is well known that adding objective risk to subjective uncertainty improves tractability, as in the Anscombe-Aumann model.  $\diamond$  should be interpreted as "getting nothing," or some other unambiguously bad outcome. This is also a technical convenience.

The only non-standard feature of the setup is the double role of  $X$ . Since the DM is betting on a randomly selected peer's choice from  $X$ ,  $X$  is the state space. With the addition of  $\diamond$ , it is also the prize space. Although we could have introduced a separate prize space (money), this would have been undesirable for two reasons. First, preferences over some less interesting space would clutter the model. Second, we would need an additional strong assumption: that preferences over the prize space are independent of  $D$ , so that the prize space is a numeraire good. We avoid all this by restricting the model to  $X$ , even though state-dependent preferences over prizes bring their own difficulties.

Formally, the analyst observes preferences indexed by datasets  $D$ : the primitive is  $\{\succsim_D\}_{D \in \mathcal{D}}$ . The domain of each  $\succsim_D$  is  $\mathcal{F} := \Delta(X \cup \{\diamond\})^X$ , with generic element  $f$ .<sup>2</sup> For any  $x \in X$ ,  $f(x) \in \Delta(X)$  is the objective lottery that  $f$  delivers if the randomly drawn peer chooses  $x$ . Notice that preferences over  $X$  can be elicited by choosing  $f$  that delivers a given item with certainty no matter what the randomly drawn peer does. This is the virtue of  $\mathcal{F}$ : it provides a neat, unified way to elicit beliefs about peers' choices and the preferences that depend on those beliefs.

**Example 10.** *As suggested in the introduction, let  $X = \{H, L\}$  be the set of possible tip amounts*

---

<sup>2</sup> $f$  is often called a "horse lottery."

displayed on a tablet at a counter-serve restaurant. The following are elements of  $\mathcal{F}$ :

$$\begin{aligned}
 f(x) &= x \\
 g(x) &= \begin{cases} \frac{1}{2}H + \frac{1}{2}L & \text{if } x = H \\ L & \text{if } x = L \end{cases} \\
 h(x) &= L
 \end{aligned}$$

Consider a DM who wants to follow the local norm, but has no information about it. Since his next observation will be highly informative, he will have  $f \succ g \succ h$ . If he later observes many people choosing  $L$ , he will become confident that the norm is  $L$ , and will have  $h \succ g \succ f$ .

### 2.3.2 REPRESENTATION

We model a DM who maximizes expected utility given his current dataset. The DM's beliefs conditional on any dataset are the Bayesian updates of his prior. Thus, the representation requires two components: Bernoulli utility function  $u$  and prior  $P$ . Since the DM ultimately cares about the true distribution of others' choices, not the data points he happens to observe,  $u$  takes the true distribution as a parameter. It does not depend on the dataset  $D$ .  $D$  influences behavior because it informs beliefs about the true distribution.

Formally,  $u$  is a function from  $\Delta(X \cup \{\diamond\}) \times \Delta(X)$  to  $\mathbb{R}$ . The first argument is an objective lottery over the prize space, and the second is the true distribution of others' choices.  $u$  is linear in its first argument; it takes an expected utility form. The dependence of preferences on the true distribution of others' choices is unrestricted, with one exception. It is convenient to have all the  $u(\cdot, \mu)$  agree that  $\diamond$  is at least as bad as anything from  $X$ .

**Definition 28** (Peer-effect utility). *A peer-effect utility is bounded  $u : \Delta(X \cup \{\diamond\}) \times \Delta(X) \rightarrow \mathbb{R}$*

such that, for any  $\mu \in \Delta(X)$ ,  $u(\cdot, \mu)$  is linear and non-constant, and minimized at  $\diamond$ .

The prior  $P$  belongs to  $\Delta(\Delta(X))$ ; it is a distribution over distributions—specifically, distributions over others' choices. To ensure well-defined conditional probabilities given any dataset  $D$ , we require the DM to place positive probability on full-support distributions.

**Definition 29** (Prior). *A prior is a probability measure  $P$  on  $\Delta(X)$  that places positive probability on full-support  $\mu \in \Delta(X)$ .*

The posterior on  $\Delta(X)$  for any given dataset  $D$  is given by

$$p(\mu \in E|D) = \frac{\int_E \prod_{d \in D} \mu(d) P(d\mu)}{\int \prod_{d \in D} \mu(d) P(d\mu)}.$$

The posterior can be used to compute the conditional probability that a randomly selected peer will choose item  $x$ :

$$P(x|D) = \int \nu(x) P(d\mu|D).$$

These conditional probabilities can be used to compute the expected utility of any bet  $f \in \mathcal{F}$ :

$$\begin{aligned} U_D(f) &= \sum_{x \in X} P(x|D) U_{Dx}(f(x)) \\ &= \sum_{x \in X} P(x|D) \mathbb{E}[u(f(x), \mu) | Dx] \\ &= \mathbb{E}_P[u(f(x), \mu) | D]. \end{aligned} \tag{2.1}$$

This is the main representation.

**Definition 30** (Peer-effect representation). *A peer-effect representation is a prior-utility pair  $(P, u)$  such that, for each  $D \in \mathcal{D}$ ,  $\succsim_D$  is represented by (2.1).*

**Example 11.** *This is a continuation of Example 10. Suppose that the value of tipping  $L$ , given the local tip distribution  $\mu$ , is  $-2\sqrt{\mu(H)}$ ; leaving a low tip subjects the DM to the restaurant staff's disapproval if many patrons tip well. The value of  $H$  is  $-1$ ; tipping well is costly, but tipping too much is not punished. Suppose further that the DM believes that the local tipping culture is either strong ( $\mu_S$ ) or weak ( $\mu_W$ ).*

*To compute the value of a bet  $f$  to a DM with dataset  $D$ , we proceed in two steps. First, for each  $x \in \{H, L\}$ , we compute the value of  $f(x)$  to a DM with dataset  $Dx$ . Then, we take an expectation over  $x$ , conditional on dataset  $D$ .*

*For  $f(x) = x$ , we have  $U_{DH}(f(H)) = -1$ , and*

$$U_{DL}(f(L)) = \Pr(\mu = \mu_S | DL) \left(-2\sqrt{\mu_S(H)}\right) + \Pr(\mu = \mu_W | DL) \left(-2\sqrt{\mu_W(H)}\right).$$

*That is, the value of  $f(L) = L$  to a DM with dataset  $DL$  is a weighted average of the value of  $L$  when tipping culture is strong, and the value of  $L$  when tipping culture is weak. The dataset  $DL$  (along with the DM's prior) determines the weights. Putting everything together, the total value of  $f$  is*

$$U_D(f) = \Pr(H|D)U_{DH}(f(H)) + \Pr(L|D)U_{DL}(f(L)).$$

### 2.3.3 BEHAVIORAL CHARACTERIZATION

This section provides a full behavioral characterization of the peer-effect model. This characterization poses two challenges. First, unlike standard expected-utility models, there is no neat separation of beliefs and preferences. A DM placing a bet on a peer's choice from  $X$  expects to learn something from the peer's choice that may affect his own preferences over  $X$  and, by extension, his preferences over the bet he is currently placing. This entanglement makes it harder to obtain valuations  $U_D(\cdot)$  and beliefs  $P(\cdot|D)$  about the choice of a randomly selected peer. It turns out that beliefs and prefer-

ences can be partially, but not totally, disentangled. The “intensity” of the DM’s preferences across different distributions of others’ choices cannot be separated from the probabilities placed on these distributions. Fortunately, this is the only source of non-uniqueness: the product of the intensity and the probability is unique, so a tight characterization of possible belief-preference pairs can be offered.

The second challenge arises in moving from the one-step-ahead beliefs  $P(\cdot|D)$  to a prior over distributions of others’ choices. De Finetti’s theorem says that a set of one-step-ahead beliefs can be generated by an appropriate prior if and only if they can be chained together into an exchangeable probability measure on  $X^\infty$ . That is, for any sequence of observations  $x_1, \dots, x_m$ :

$$Pr(x_m|x_1, \dots, x_{m-1}) \cdots Pr(x_2|x_1)Pr(x_1) \tag{2.2}$$

must not depend on the choice of indices. One might expect this condition to be automatically satisfied when the DM cares only about his dataset  $D$ , not about the order of its elements. This intuition is incorrect. Exchangeability of beliefs has an additional behavioral implication, which we call Forward Exchangeability (FE).

Fully characterizing the model requires five axioms. Quantifiers are omitted where not needed for clarity. The first four axioms are entirely familiar. Free Disposal says that  $\diamond$  is weakly worse than any other choice, and strictly worse than some other choice (not necessarily the same one for all  $D$ ). This minimal level of consistency across datasets is technically convenient. It also prevents uncertainty about (non-trivial) future preferences from causing exact indifference today.

**Axiom 14** (Free Disposal). *For any  $D$ :  $x \succsim_D \diamond$  for all  $x \in X$ , and  $x \succ_D \diamond$  for some  $x$ .*

Independence and Mixture Continuity deliver a linear (expected-utility) representation for each  $\succsim_D$ .

**Axiom 15** (Independence).

$$f \sim g \implies \frac{1}{2}f + \frac{1}{2}b \sim \frac{1}{2}g + \frac{1}{2}b.$$

**Axiom 16** (Mixture Continuity). *The sets*

$$\{\alpha \in [0, 1] : \alpha f + (1 - \alpha)g \succ b\}$$

$$\{\alpha \in [0, 1] : \alpha f + (1 - \alpha)g \prec b\}$$

*are open in*  $[0, 1]$ .

Dynamic Consistency provides a link across datasets. If  $f$  is better than  $g$  for every dataset that extends  $D$  by one data point, then  $f$  is better than  $g$  for  $D$ .

**Axiom 17** (Dynamic Consistency). *If*

$$f(x) \succ_{Dx} g(x) \tag{2.3}$$

*for all*  $x \in X$ , *then*  $f \succ_D g$ . *Also, if* (2.3) *holds with strict preference for some*  $x$ , *then*  $f \succ_D g$ .

Combining the first four axioms, we get a representation for initial preferences ( $D = \emptyset$ ) in terms of preferences following a single signal.

$$U_{\emptyset}(f) = \sum_{x \in X} P(x) U_x(f(x)). \tag{2.4}$$

This part is easy because exchangeability places no restrictions on  $P(x)$ . To extend this construction to larger datasets, we need FE. FE resembles Dynamic Consistency, but with mixing across pairs of states rather than state-by-state comparisons. Mixing across states  $x$  and  $y$  forces the DM to treat



signal streams  $(y, x)$  and  $(x, y)$  as equiprobable given his current dataset  $D$ . This restriction has to be imposed indirectly because the domain of choice does not include bets on pairs of signals; the DM is only offered bets on the choice of *one* randomly selected peer.

**Axiom 18** (Forward Exchangeability (FE)). *If*

$$\frac{1}{2}f(x) + \frac{1}{2}f(y) \succsim_{Dxy} \frac{1}{2}g(x) + \frac{1}{2}g(y) \quad (2.5)$$

for all  $x, y \in X$ , then  $f \succsim_D g$ . Also, if (2.5) holds with strict preference for some  $x, y$ , then  $f \succ_D g$ .

**Example 12.** Take the same domain and utility function of Example 11, and let

$$\begin{aligned} f(x) &= x \\ g(x) &= \frac{1}{2}H + \frac{1}{2}L. \end{aligned}$$

$f$  prescribes matching the choice of the patron preceding the DM, while  $g$  prescribes randomizing no matter what the other patron does. Suppose that the DM places equal probability on  $\mu(H) = 0.8$  and  $\mu(L) = 0.2$ .  $HH$  is a strong signal that high tippers are the majority, so  $f(H) \succ_{HH} g(H)$ . Similarly,  $LL$  is a strong signal that low tippers are the majority, so  $f(L) \succ_{LL} g(L)$ . The mixed-signal case is more subtle. Although the DM's preferences after receiving two signals do not depend on the order of those signals, the payoffs from bets  $f$  and  $g$  depend only on the first signal. Since  $H$  is better than  $L$  given a mixed signal,  $g$  is actually better than  $f$  in state  $LH$ . But  $f$  has an equally large advantage in state  $HL$ , so

$$\frac{1}{2}f(H) + \frac{1}{2}f(L) \sim_{HL} \frac{1}{2}g(H) + \frac{1}{2}g(L).$$

FE says we don't have to compute  $U_\emptyset(f) - U_\emptyset(g)$  to know that  $f$  is better than  $g$  when  $D = \emptyset$ . Intuitively,  $HL$  and  $LH$  must be equally likely, so  $g$ 's advantage in state  $LH$  is fully compensated by  $f$ 's

advantage in state HL. Since  $f$  is strictly better in the other possible states,  $f$  is strictly better overall.

FE allows us to write (2.4) for arbitrary  $D$ :

$$U_D(f) = \sum_{x \in X} P(x|D) U_{Dx}(f(x)), \quad (2.6)$$

where the  $P(x|D)$  satisfy (2.2). Chaining together these one-step-ahead probabilities, we get an exchangeable measure  $P$  over sequences  $(x_1, x_2, \dots) \in X^\infty$ . (2.6) says that the value of any item  $x$  is a martingale under  $P$ , so we use a martingale convergence theorem to write the value of  $x$  given any dataset  $D$  as a conditional expectation over limiting values. We then use the de Finetti theorem to write the limiting values in  $(x_1, x_2, \dots)$  as a function of the limiting frequencies with which each item appears in  $(x_1, x_2, \dots)$ . This is the true distribution of others' choices, so the limiting valuations are the realizations of a peer-effect utility function. The de Finetti theorem also allows us to move from an exchangeable measure over  $X^\infty$  to a prior over  $\Delta(X)$ . We conclude that the value of  $x$  given  $D$  can be written as the conditional expectation under prior  $P$  of peer-effect utility  $u$ . This is the peer-effect representation.

**Theorem 6.**  $\{\succsim_D\}_{D \in \mathcal{D}}$  satisfies Axioms 14 - 18 if and only if it has a peer-effect representation.

A peer-effect representation can be non-unique for two reasons. First, two different signals might happen to induce the same preferences, making it impossible to determine the weight placed on each signal. This case is not very interesting, so we assume it away.

**Definition 31 (Full Rank).**  $\{\succsim_D\}_{D \in \mathcal{D}}$  has full rank if, for all distinct  $x, y \in X$  and all  $D \in \mathcal{D}$ ,

$$\succsim_{Dx} \neq \succsim_{Dy}.$$

Corollary 5 shows that Full Rank allows the analyst to pin down the weights the DM places on each signal (and, by extension, the weights he places on each distribution of others' choices). These

weights incorporate both beliefs and the “intensity” of preferences, or the scaling of the Bernoulli utility function. Choice data on  $\mathcal{F}$  is never informative about beliefs or intensities alone, so one representation can be transformed into another by lowering the probability of some states and raising the intensity of preferences in those states to compensate.

**Corollary 5.** *Suppose that  $\{\succsim_D\}_{D \in \mathcal{D}}$  has full rank. If  $(P, u)$  and  $(Q, v)$  are peer-effect representations with  $u$  and  $v$  continuous in the second argument, then*

$$v(\cdot, \mu) - v(\diamond, \mu) = A(\mu) [u(\cdot, \mu) - u(\diamond, \mu)]$$

$$\left. \frac{dQ}{dP} \right|_{\mu} = \frac{\mathbb{E}_Q A}{A(\mu)}$$

for some  $A : \Delta(X) \rightarrow \mathbb{R}_{++}$ .

#### 2.4 TYPES OF PEER EFFECTS

Thus far, the shape of  $u$ —the dependence of preferences on others’ choices—has been entirely unrestricted. This section proposes some natural special cases and sets out the behavioral implications.

Many of the examples discussed in Section 2.2 have a directional nature. The DM must decide how much to donate, invest, study, etc., and cares about how much others do. These examples are well captured by a simplified setting in which the domain is an interval in  $\mathbb{R}$  and the support of the prior has the Monotone Likelihood Ratio Property (MLRP). MLRP is convenient because it makes belief updating very intuitive: observing  $x + \Delta$  rather than  $x$  leads to unambiguously higher beliefs about others’ choices. It is also convenient to assume that the return to increasing one’s action is decreasing sufficiently quickly in  $x$ . This assumption rules out multiple optima and leads to neat results about behavior over time. These restrictions are collected in Assumption 1. The marginal

benefit of action is denoted

$$\Delta u(x, \mu) := u(x, \mu) - u(x - \Delta, \mu).$$

f

**Assumption 1.**

1.  $X$  is an open interval in  $\mathbb{R}$ .
2.  $u$  is strictly concave in  $x$ , and  $\Delta u$  is concave in  $x$ .
3.  $\text{supp}(P)$  has the Monotone Likelihood Ratio Property: for all distinct  $\mu, \nu \in \text{supp}(P)$ , the ratio  $\mu(x)/\nu(x)$  is monotone in  $x$ .

All the results in this section are predictions about the optimal selection from  $X$ . We denote the optimum<sup>3</sup> as

$$x^*(P) = \arg \max_{x \in X} \mathbb{E}_P[u(x, \mu)].$$

The first result provides a condition under which higher beliefs induce higher actions, where the ordering on beliefs is first-order stochastic dominance (FOSD). This condition is a familiar increasing differences property. It says that the optimum is increasing in beliefs provided the benefits to increasing one's action from  $x$  to  $x + \Delta$  are greater if others are doing more. As with the other results in this section, the sign can be flipped to obtain the opposite result. If the benefits of increasing one's action are decreasing in beliefs, so is the optimal action.

**Proposition 9** (Effect of shift in beliefs). *If  $\Delta u(x, \mu)$  is increasing (decreasing) in  $\mu$ , then*

$$x^*(P) \geq (\leq) x^*(Q)$$

---

<sup>3</sup> $x^*(P)$  is guaranteed to be unique when it exists. All results in this section should be read with the caveat “when the optimum exists.”

for  $P \succ_{FOSD} Q$ .

**Example 13.** If  $X = (1, 2)$  and

$$u(x, \mu) = -(\mathbb{E}_\mu x)(x - \mathbb{E}_\mu x)^2,$$

then  $\Delta u(x, \mu)$  is increasing in  $\mu$ . This utility function is maximized at

$$x^*(\mu) = \mathbb{E}_\mu x.$$

Clearly, the optimal action is increasing in peers' actions when there is no uncertainty. Proposition 9 says that this is still true in the presence of uncertainty.

This result is theoretically straightforward, but empirically useful. Several papers have documented the effects of informing low performers (in varied domains) that their peers are doing better than they are. In [Allcott \(2011\)](#), households using large amounts of electricity reduced their consumption when told about their neighbors' electricity usage. In [Card et al. \(2012\)](#), below-median earners reported lower job satisfaction and were more likely to search for a new job when informed about coworkers' salaries. In [Goulas and Megalokonomou \(2015\)](#), students who learned that they had scored poorly on an exam did worse on a subsequent exam. Notice that this effect goes in the opposite direction as the previous two. Apparently, low-achieving students were demotivated rather than energized by peers' superior scores.

Proposition 9 provides one way to think about these changes in beliefs. It is plausible that the return to decreasing one's energy usage is higher if one's peers are using less. This could be because low peer usage signals low costs of energy-saving investments, or because it makes one feel irresponsible or wasteful. Either way, a homeowner should cut his energy usage if he learns that others' usage is lower than he thought. Similarly, the return to looking for a new job is probably higher for an

employee whose coworkers are more valued and better compensated than she is. The exam score example is more complicated; the effect could reasonably go either way. A student who learns that he has performed poorly might study more to avoid feeling lazy. Alternatively, he might conclude that his peers are more academically gifted than he is, so he should shift time from academics to activities in which he has a comparative advantage. The results in [Goulas and Megalokonomou \(2015\)](#) suggest that the second effect dominates, although other mechanisms could be at work.

The next result suggests one alternative mechanism: changes in uncertainty rather than changes in level. People can be motivated by the possibility that their peers are doing more than they are, even if it is also possible that they are doing more than their peers. Formally, if higher beliefs induce higher actions, then uncertainty also induces higher actions provided marginal benefits exhibit decreasing differences:

$$\Delta u(x_H, \mu_L) - \Delta u(x_L, \mu_L) \geq \Delta u(x_H, \mu_H) - \Delta u(x_L, \mu_H)$$

where  $x_H > x_L$  and  $\mu_H >_{FOSD} \mu_L$ . To see why this works, suppose that  $x_H$  is optimal for  $\mu_H$  and  $x_L$  is optimal for  $\mu_L$ , so  $\Delta u(x_H, \mu_H)$  and  $\Delta u(x_L, \mu_L)$  are approximately zero for small  $\Delta$ . The decreasing differences condition becomes

$$\Delta u(x_L, \mu_H) \geq -\Delta u(x_H, \mu_L).$$

This says the benefit of increasing one's action from the low-belief optimum, given high beliefs, exceeds the cost of increasing one's action from the high-belief optimum, given low beliefs. Intuitively, the possibility of falling behind one's peers provides pressure that is not fully offset by the possibility of being ahead, so the optimum is increasing in uncertainty.

**Proposition 10** (Effect of uncertainty). *If  $\Delta u(x, \mu)$  is increasing (decreasing) in  $\mu$  and exhibits*

decreasing (increasing) differences, then

$$x^* \left( \sum_i \lambda_i P_i \right) \geq \sum_i \lambda_i x^*(P_i).$$

Example 13 illustrates Proposition 10 as well as Proposition 9:  $\Delta u$  exhibits decreasing differences.

Proposition 10 provides an alternative mechanism for the demotivating effect of peers' test scores. Some low scorers may not have been surprised by their results, given their experience of the exam and their grades in school. In fact, they may have lost the motivation provided by their worst fears (provided they were not at the very bottom of the grade distribution).

This reasoning generalizes beyond students and exams. When attempting to motivate the lower tail (in any activity) it is important to determine what people in the lower tail already know about their relative position. If they think they are closer to the middle than they actually are, e.g. because they have attributed negative signals partly to chance, they are likely to raise their effort when they learn about peer performance. This is Proposition 9. If they are pretty well calibrated, but still uncertain, information may be ineffective or outright harmful. Perverse effects are likely when the risk of being further behind than one thinks is a strong motivator, e.g. when low performers face penalties like firing or expulsion.

Propositions 9 and 10 are most useful when the analyst has good information about, and can directly influence, subjects' beliefs. Corollary 6 shows how these results translate into the setting from Section 2.3. There, beliefs are not observed or directly controlled; the DM learns by collecting data about individual peers' choices. Fortunately, the above restrictions on  $u$  still have simple predictions.

**Corollary 6.**

1. If  $\Delta u$  is increasing (decreasing) in  $\mu$ , then

$$x^*(P|Dx) \geq (\leq) x^*(P|Dy)$$

for any  $x > y$ .

2. If  $\Delta u$  is increasing (decreasing) in  $\mu$  and exhibits decreasing (increasing) differences, then

$$x^*(P|D) \geq (\leq) \mathbb{E}[x^*(P|Dx)].$$

The first part of Corollary 6, which is the analogue of Proposition 9, says that the DM's optimal choice is increasing (decreasing) in peer actions if  $u$  has increasing (decreasing) differences. The study by [Rogers and Feller \(2016\)](#) provides an illustration. Students in an online course were asked to rate essays written by a few randomly assigned fellow students. Although grading in the course was not relative, students who were assigned high-quality essays were more likely to drop out. This corresponds to the decreasing differences case. Some students may have believed that “drop out” and “produce low quality work” were the only choices available to them. Those students preferred to drop out if they came to believe that most of their peers were producing high quality work.

The second part of Corollary 6 is the analogue of Proposition 10. It says that a DM who is motivated by uncertainty will reduce his action on average as he learns more about his peers' actions. Two empirical papers on pay-what-you-want pricing schemes find evidence of this effect. Both [Schons et al. \(2014\)](#) (within-subject) and [Riener and Traxler \(2012\)](#) (across-subject) find that consumers' chosen prices decline over time, although they converge to a positive level. [Riener and Traxler \(2012\)](#) interpret this exactly as in Corollary 6: consumers pay more when they are less familiar with the pricing scheme because paying less than the socially acceptable amount is very unattractive. This result may also explain why [Johnson and Cui \(2013\)](#) find that external reference prices, including minimum and maximum prices, tend to reduce consumers' chosen prices. Reference prices provide information, which makes it less risky for a consumer to select a lower price.



## 2.5 INTERPRETING PEER EFFECTS: PURE LEARNING VS. SOCIAL IMAGE

All the analysis thus far has been consistent with two different interpretations. The DM might care directly about his relative standing within his peer group (social interpretation), or he might take peers' choices as a source of information about  $X$  (pure learning interpretation). This section formalizes each interpretation through an extension of the main model. It turns out that choice data on  $\mathcal{F}$  is not enough to distinguish the two interpretations. Even though the pure learning model is a special case of the social model, data generated by a social model can always be replicated by an appropriate pure learning model. However, the pure learning interpretation is falsifiable in an enriched environment, where the analyst can observe or control beliefs about the characteristics (not just choices) of the peer group. The pure learning model predicts less sensitivity to peer choices when there is more uncertainty about the distribution of characteristics. This need not be true in the broader social model.

As in Section 2.4, we take the domain of choice to be an interval in  $\mathbb{R}$ . The informational value of others' choices comes from uncertainty about the utility function  $u \in \mathcal{U}$ . To account for heterogeneity across people, we allow utility to depend on a type parameter  $\theta \in \Theta$ , where  $\Theta$  is an interval in  $\mathbb{R}$ . Higher types are assumed to prefer higher actions;  $u_x$  is increasing in  $\theta$ . The DM may be uncertain about the distribution of others' types as well as the utility function. We write  $T \in \mathcal{T}$  for the cdf of the type distribution. Uncertainty about both  $u$  and  $T$  is captured by a prior  $P \in \Delta(\mathcal{U} \times \mathcal{T})$ . To simplify belief updating, we assume that each  $T$  has strictly positive density on  $\Theta$ , and that each  $u$  will generate a bijective mapping between types and optima. The formal conditions are stated in Assumption 2 below.

**Example 14.**  $X = \Theta = \mathbb{R}$ . Types are normally distributed with mean  $\bar{\theta}$ , which is itself normally distributed. Utility is given by

$$u(x, \theta) = -(x - (\theta + k))^2$$

where  $k$  is normally distributed with zero mean.  $k$  captures the common value of activity  $x$ , while  $\theta$  captures personal taste. Beliefs about  $k$  capture uncertainty about utility. Beliefs about  $\bar{\theta}$  capture uncertainty about the type distribution.

**Assumption 2.**

1.  $X$  and  $\Theta$  are open intervals in  $\mathbb{R}$ .
2. Each  $T \in \mathcal{T}$  has strictly positive density on  $\Theta$ .
3. For each  $u \in \mathcal{U}$ ,  $u_x$  is strictly decreasing in  $x$ , strictly increasing in  $\theta$ , and satisfies

$$\lim_{x \rightarrow \sup(X)} \sup_{\theta} u_x = \lim_{x \rightarrow \inf(X)} \inf_{\theta} u_x = 0$$

As in the main model, we abstract away from inference about the learning of others. The DM gathers choice data from agents who already know the model parameters. These agents always maximize their utilities conditional on their types. Types are private information; the DM knows his own type, but does not observe anyone else's.

In the pure learning model, utilities are defined on  $X \times \Theta$ , so the choice of informed agents is pinned down by

$$x^*(\theta; u) := \arg \max_{x \in X} u(x, \theta).$$

For the utility function in Example 14, we have  $x^*(\theta; u) = \theta + k$ .

**Definition 32** (Pure learning model). *A pure learning model is  $P \in \Delta(\mathcal{T}, \mathcal{U})$  such that  $\mathcal{U} \subset \mathbb{R}^{X \times \Theta}$ .*

The pure learning model does not allow  $u$  to depend directly on the distribution of others' choices. Others' choices matter because they are informative about the realization of  $u$ . To see how

this works, consider the DM's utility function

$$\mathbb{E}_P[u(x, \theta) | D] \propto \int_{T, u} u(x, \theta) \prod_{d \in D} \frac{d}{dd} T \circ (x^*)^{-1}(d; u) dP(T, u).$$

$T \circ (x^*)^{-1}(d; u)$  is the probability, given the type distribution  $T$  and utility  $u$  that a randomly selected agent has an optimum below  $d$ . Thus,

$$\prod_{d \in D} \frac{d}{dd} T \circ (x^*)^{-1}(d; u)$$

measures the likelihood of getting dataset  $D$  from model parameters  $(u, T)$ . The DM puts more weight on  $(u, T)$  pairs that are more likely to generate his dataset.

The social model is a generalization of the pure learning model in which  $u$  takes  $\Delta(\bar{X})$  as an additional argument. Here, the DM is allowed to care directly about others' choices as well as their informational content. Since the other people in the model also care about their peers' choices, equilibrium concerns arise. We handle this by specifying the equilibrium  $x^*$  as one of the parameters of the model. We assume that the DM knows  $x^*$ . That is, for each possible  $(u, T)$ , the DM knows the equilibrium the informed agents would play if  $(u, T)$  were realized. This strong assumption can be avoided by choosing  $u$  to guarantee a unique equilibrium, as in Example 15 below.

**Example 15.**  $X, \Theta$  and  $\mathcal{T}$  are as in Example 14. Utility is non-random and given by

$$u(x, \theta, T) = -(\lambda\theta + (1 - \lambda)\mathbb{E}_T[x] - x)^2,$$

where  $\lambda \geq 0$ . Agents are conformist; they like being close to the mean.  $x^*$  is uniquely pinned down:

$$x^*(\theta; T) = \lambda\theta + (1 - \lambda)\bar{\theta}.$$

**Definition 33** (Social model). *A social model is  $P \in \Delta(\mathcal{T}, \mathcal{U})$  such that  $\mathcal{U} \subset \mathbb{R}^{X \times \Theta \times \Delta(x)}$ , and  $x^* : \Theta \times \text{supp}(P) \rightarrow \mathbb{R}$  such that*

$$x^*(\theta; T, u) = \arg \max_{x \in X} u(x, \theta, T \circ (x^*)^{-1}).$$

In the social model, the DM maximizes

$$\mathbb{E}_P [u(x, \theta, T \circ (x^*)^{-1})] \propto \int_{T, u} u(x, \theta, T \circ (x^*)^{-1}) \prod_{d \in D} \frac{d}{dd} T \circ (x^*)^{-1}(d; u) dP(T, u).$$

Examples 14 and 15 hint at the difficulty of distinguishing pure learning and social models. In both cases, the DM will appear to imitate peer behavior: a dataset that contains mostly high values will induce a higher action than a dataset that contains mostly low values. Empirical researchers are already aware of this. An informal version of this point is found in [Manski \(2000\)](#), and several empirical papers on peer effects, including [Duflo and Saez \(2003\)](#) and [Brown et al. \(2008\)](#), offer both information spillovers and social pressure as possible explanations for their findings.

The natural question is whether the difficulty of distinguishing the two types of models is limited to imitative behavior. Could any pattern of choice behavior provide a conclusive rejection of the pure learning model? Proposition 11 says the answer is no. When the analyst can only observe choice on  $\mathcal{F}$ , the domain studied in Section 2.3, she cannot reject the pure learning model without also rejecting the social model. The additional generality of the social model is only apparent. All choices on  $\mathcal{F}$  that can be interpreted as conformist, envious, defiant, and so on can be recast as the product of rational inference about a utility function that does not depend on peer choices.<sup>4</sup>

**Proposition 11** (Observational equivalence). *For any social model  $(P, x^*)$ , there is a pure learning*

---

<sup>4</sup>This is not because  $\mathcal{F}$  is overly restricted; it subsumes preferences on  $\Delta(X)$ , the set of all lotteries on the domain of interest.

model  $\hat{P}$  that generates the same preferences over  $\mathcal{F}$  for every type  $\theta$  and dataset  $D$ :

$$\mathbb{E}_P [u(x, \theta, T \circ (x^*)^{-1}) | D] = \mathbb{E}_{\hat{P}} [\hat{u}(x, \theta) | D].$$

Proposition 11 should not be read as an argument for restricting attention to pure learning models of peer effects. Although it is always possible to construct an inference-based alternative to compete with any social model, the inference-based model may not offer better understanding or predictions outside the domain at hand. The remainder of this section shows that the two types of models *can* issue divergent predictions in an enriched environment. The key parameter is uncertainty about the type distribution  $T$ . The analyst can either elicit beliefs about  $T$  directly, or vary the information that the DM receives about the composition of the group he observes. To make this more concrete, consider the charitable donation case, in which the DM observes the amount donated by individuals randomly selected from a pool of potential donors. The experimenter can give or withhold information about the socioeconomic, political, or religious composition of the donor pool, which may affect the way the DM processes the data he receives.

Intuitively, a DM who is concerned only with the effectiveness of the charity should place less weight on the donations he observes if he does not know anything about the donor pool. High donations may signal simply that the donor pool is rich, or shares the ideology of the charity. To formalize this point, we need some way of measuring the DM's sensitivity to his dataset. We consider two measures. The first,  $\mathcal{M}$ , captures movement in the valuation of a fixed item  $x$ . The second,  $\mathcal{R}$ , captures reversals in preference. A reversal occurs if the DM initially prefers  $x$  over  $y$ , but comes to prefer  $y$  to  $x$  after observing  $D$ .  $\mathcal{R}$  increases in the size of the reversal as well as the frequency; flipping from one strong preference to another has a bigger effect than flipping from one weak preference to

another. Formally, for any  $x, y \in X$ , any  $\theta \in \Theta$ , and any finite  $n$ , let

$$M_P(x, \theta) := \mathbb{E}_P[|U_D(x, \theta) - U_\emptyset(x, \theta)|]$$

$$R_P(x, y, \theta) := \mathbb{E}_P[1_{\Delta U_D \Delta U_\emptyset \leq 0} |\Delta U_D \Delta U_\emptyset|]$$

where  $\Delta U_D := U_D(x, \theta) - U_D(y, \theta)$

and all expectations are taken over datasets of size  $n$ .

Proposition 12 confirms the intuition above. In the pure learning model, both measures of sensitivity to data decrease when the DM has less information about the type distribution.

**Proposition 12.** *Let  $P = \sum_i \lambda_i P_i$ , where  $P_i(u \in U|T) = P_j(u \in U|T)$  for all  $i, j$ . For any  $x, y \in X$ , any  $\theta \in \Theta$ , and any finite dataset size,*

$$M_P \leq \sum_i \lambda_i M_{P_i}$$

$$R_P \leq \sum_i \lambda_i R_{P_i}.$$

This is not the case for all social models. This is easiest to see with a family of social models in which  $T$  is the only source of uncertainty ( $u$  is known). As uncertainty about  $T$  vanishes, neither a pure learning DM nor a social DM learns anything from others' behavior, so preferences do not depend on  $D$ . In fact, the pure learning DM's preferences are independent of  $D$  no matter how much uncertainty he faces about  $T$ .  $T$  is useful only insofar as it helps him interpret evidence about  $u$ , which he already knows. By contrast, the social DM is responsive to uncertainty about  $T$ . When the DM does not know  $T$ ,  $D$  is informative about the distribution of others' choices, which feeds back into preferences.

To fix ideas, consider the charitable donation case with utility

$$u(x, \theta, T) = -(\theta + \lambda T(x) - x)^2. \quad (2.7)$$

Here, the relevant uncertainty is not about the value of the charity, but others' generosity. The DM likes appearing (or feeling) generous relative to others. If he faces higher uncertainty about the generosity of the potential donor pool, his beliefs will move more in response to the donations he observes. High donations send a strong signal that others are generous, so the DM needs to donate more to avoid looking (or feeling) cheap. This is not the case if the DM is well informed about the composition of the donor pool. Then, high donations are attributed mostly to chance, so the DM remains confident that most people will donate less. The total effect is exactly opposite that of the pure learning model: more uncertainty about the type distribution generates more reversals.

## 2.6 CONCLUSION

A large and growing empirical literature finds that people's choices are often sensitive to information about peers' choices. These effects are present in surprisingly important decisions, including savings (Duflo and Saez, 2003), investment (Bursztyn et al., 2014; Brown et al., 2008), and job choice (Coffman et al., 2017), and they sometimes go in surprising directions (Beshears et al., 2015; Goulas and Megalokonomou, 2015). This paper provides a unified theoretical framework for interpreting these findings and designing new studies. The first contribution is a general Bayesian model that is flexible enough to accommodate the interesting patterns found in the evidence. The second is a set of special cases that deliver neat predictions about different types of informational interventions. The third is a careful analysis of "pure learning" and "social" mechanisms behind peer effects, and a novel strategy for distinguishing the two.

# 3

## Fantasies and Worries as Distorted Probabilities

### 3.1 INTRODUCTION

In settings conducive to anticipatory emotions like fantasy or worry, the length of time before a given risk resolves affects its value. For instance, people may find it costly to live with the thought of



something bad happening in future. They may respond by convincing themselves that the bad event is unlikely, or they may simply spend more time thinking about better outcomes. In either case, they may place too little weight on the bad outcome when making decisions. This problem goes away if the risk is resolved immediately, so there is never any need to live with (and try to avoid) unpleasant thoughts. Thus, fixed downside risks will be evaluated less favorably as they are moved from the future into the present.

To provide a concrete example, it is easier to take a rosier view of a major life event like marriage or retirement when it is far away. Some of this apparent optimism may be attributed to lack of knowledge, but arguably not all of it. People may be well aware of rates of divorce or dementia but still fail to plan for either, suggesting that they do not genuinely expect those things to happen to them.

This paper makes room for fantasy and worry within standard models of preferences over consumption trees. The resulting representation agrees with expected utility (EU) for risks that are resolved immediately. For risks resolved in future, the representation is expected utility with distorted probabilities. The distorted probabilities obey a certain structure: distortions for lotteries with more than two outcomes are built up from distortions for two-outcome lotteries. Nevertheless, this structure turns out to be quite flexible. It is significantly more flexible than the standard model of probability weighting.

This flexibility is desirable. It seems clear that people do fantasize and worry about future risks, but the precise patterns of fantasy and worry are less clear. Do people generally fantasize in some situations and worry in others, or do certain individuals predominantly fantasize while others predominantly worry? Beyond the absence of fantasy and worry for present risks, how does the tendency to fantasize or worry depend on the horizon? When there is a big spread in possible outcomes, does this make the bad outcome look worse or the good outcome look better (or does it depend on the probabilities)? The model presented below does not impose answers to these questions, although

particular specifications of the model will. It should remain useful as researchers learn more about the tendency to fantasize or worry.

The generality of the domain is also desirable. It consists of trees that capture an agent's evolving consumption prospects. Such trees will be familiar from [Epstein and Zin \(1989\)](#) and [Chew and Epstein \(1991\)](#) or from applications in asset pricing and household finance. The trees can be built up from discrete or continuous distributions with bounded or unbounded payoffs. They can capture risks that are resolved gradually over several periods, or in one shot. Trees extend into the infinite future, but it is easy to impose a finite horizon by specifying a payoff of 0 in every period after some fixed date. Thus, the model applies in a variety of settings.

The remainder of this section connects the model to existing literature. Section 3.2 describes the domain of choice and presents the main representation. Section 3.3 provides some convenient parametric specifications. Section 3.4 explains what it means for one person to fantasize (or worry) more than another and connects this to the representation. Section 3.5 presents axioms for the representation. Section 3.6 gives an additional axiom that guarantees monotonicity in first-order stochastic dominance, thereby ruling out some plainly irrational behavior. Section 3.7 concludes with a discussion of future work.

### 3.1.1 RELATED LITERATURE

The starting point for this paper is the recursive model introduced by [Kreps and Porteus \(1978\)](#) and extended by [Epstein and Zin \(1989\)](#) and [Chew and Epstein \(1991\)](#). The domain of choice is drawn from the infinite “consumption programs” of Chew and Epstein, who extend the finite “temporal lotteries” of Kreps and Porteus. Kreps-Porteus preferences have a recursive structure with expected utility (not necessarily the same Bernoulli utility) in each period. Chew-Epstein preferences depart from EU but use the same certainty-equivalent function in each period. In this paper, the decision maker's attitude to a given risk depends on the period in which it will be resolved. His certainty-

equivalent function for present risks, which do not inspire fantasies or worries, is EU; for future risks, which do, it is EU with distorted probabilities. Probability distortions do not feature explicitly in Chew and Epstein or Kreps and Porteus.

Probability distortions do feature in a large literature on rank-dependent expected utility (RDEU). RDEU was introduced in [Kahneman and Tversky \(1979\)](#) and further developed in [Tversky and Kahneman \(1992\)](#). [Quiggin \(1982\)](#), [Yaari \(1987\)](#), [Wakker \(1994\)](#), and [Abdellaoui \(2002\)](#) all provided axiomatizations. In RDEU, the weight the decision maker places on getting more than  $c$  depends only on the probability of getting more than  $c$ . (This differs from EU in that the weight need not *equal* the probability.) The probability distortions in this paper are more general: the weight on getting more than  $c$  from a given distribution may (but does not have to) depend on the value of its conditional distributions below  $c$  and above  $c$ . For instance, a decision maker could place a lower probability on winning a fifty-fifty bet when the payoff from losing is particularly bad. This additional generality comes at the cost of imposing EU preferences over present risks. EU at  $t = 0$  provides separate identification for the Bernoulli utility function and the probability distortion function, thereby avoiding the main difficulty with the RDEU approach.

Probability distortions are not the only way to model fantasies and worries. [Caplin and Leahy \(2001\)](#) construct a model of expected utility over psychological states. Although they discuss the mapping between physical lotteries and psychological states, their key assumptions are on preferences over psychological lotteries. All the axioms in this paper pertain to preferences over physical lotteries. [Epstein \(2008\)](#) studies the demand for commitment of a decision maker with anticipatory emotions. For Epstein's purposes, it is enough that the decision maker evaluates risks differently when they resolve at different times. This paper goes further by specifying *how* future risks should be evaluated: by EU with distorted probabilities.

[Dillenberger et al. \(2017\)](#) provide a model in which the probability distributions a decision maker uses to evaluate subjective acts depend on the payoffs associated with those acts. The deci-

sion maker’s preferences over objective lotteries are simply EU. In this paper, the probability distributions the decision maker uses to evaluate *objective lotteries* depends on the payoffs associated with those lotteries. This approach complements Dillenberger et al. by modeling optimism and pessimism without the use of subjective acts.

### 3.2 MODEL

#### 3.2.1 DOMAIN OF CHOICE

The domain of choice is very similar to the space of consumption trees in [Chew and Epstein \(1991\)](#). A more detailed description of the domain is in Appendix C.1; an intuitive characterization is below.

An agent who faces a consumption tree  $d$  draws a pair  $(c_t, d_{t+1})$  in each period  $t \in \{0, 1, \dots\}$ .  $c_t$  belongs to an interval  $C$  of  $\mathbb{R}$  (open, closed, or neither) with  $0 \in C$ .<sup>1</sup>  $d_{t+1}$  is a continuation tree; it captures the agent’s consumption prospects from period  $t + 1$  onward.

A consumption tree  $d$  can be built up from continuous or discrete distributions, or it can combine the two. Today’s payoff  $c_t$  can completely determine tomorrow’s prospects  $d_{t+1}$ , give no information about  $d_{t+1}$ , or act as a noisy signal about  $d_{t+1}$ . All uncertainty about payoffs can be resolved by some future date  $T$ , or uncertainty can be resolved only asymptotically.

Chew & Epstein use the full space  $\bar{D}$  of consumption trees as the domain of choice. That is an option here as well, but it has the undesirable consequence of a bounded Bernoulli utility function.<sup>2</sup> This would rule out some convenient special cases of the model, such as an agent with CRRA utility and a probability distortion function that preserves lognormality.

The solution is to stipulate that the domain  $D$  is a large enough subset of  $\bar{D}$ , where “large enough”

---

<sup>1</sup>Chew & Epstein allow for any separable metric space, but the additional generality is not needed here.

<sup>2</sup>The reason for this is familiar. If  $u$  were unbounded, then some risks would have infinite expected utility and no finite certainty equivalent.

means “satisfying several richness conditions.” ( $\bar{D}$  is of course large enough, so  $D = \bar{D}$  is permitted.) Since these conditions are not particularly interesting in themselves, they are relegated to Appendix C.1.

### 3.2.2 PROBABILITY DISTORTIONS

The representation presented in the next section uses probability distortions to capture anticipatory emotions. An agent who fantasizes about a future risk maps the true distribution over payoffs into a more favorable distribution, then evaluates the result by expected utility.

The structure of probability distortions used in this paper comes from a very simple idea: an agent who fantasizes subtracts weight from bad outcomes and adds it to good ones, while an agent who worries does the opposite. The nature of the outcomes in question should be allowed to matter. For instance, if the outcomes in a lottery are all very bad, or if the lottery is very spread out, it may be especially conducive to worry (large pessimistic probability distortions). The implications of this for binary lotteries are obvious: the agent should act as if he misperceives the probability of the bad outcome, where his misperception may depend on the values of both outcomes as well as the probabilities. Extending this idea to a more complicated lottery, we may imagine the agent dividing the lottery into better and worse parts and adjusting their relative probabilities just as he would if facing a binary lottery. There are many ways to divide a lottery into better and worse parts, and no clear reason to privilege one over another. Thus, we assume that the agent treats each possible outcome of the lottery as a breakpoint and readjusts the probabilities at all breakpoints.

The benefit of this approach is that all the agent’s probability distortions are pinned down by his distortions of binary lotteries. Such distortions are simple and easy to work with. The downside is that the distortions of binary lotteries must be chosen carefully so that the distortions of non-binary lotteries are always well defined. The key idea is that the probability of getting the bad payoff cannot fall too quickly as the values of the good and bad payoffs rise. This is because the values of a given

lottery below quantile  $q$  and above  $q$  both rise with  $q$ . If this improvement in prospects made the agent very optimistic, the distorted “cdf” would have some decreasing regions (so would not be a cdf at all).

Definition 34 states the condition that distortions of binary lotteries must satisfy to avoid the problem above.

**Definition 34** (Probability distortion function). *Fix a continuous and strictly increasing Bernoulli utility  $u : C \rightarrow \mathbb{R}$ .*

$$\pi_t : \{(q, c_L, c_H) \in [0, 1] \times [C \cup \emptyset]^2 : \\ c_L = \emptyset \Rightarrow q = 0, c_H = \emptyset \Rightarrow q = 1, c_L, c_H \neq \emptyset \Rightarrow c_L < c_H\} \rightarrow [0, 1]$$

*is a probability distortion function for  $u$  if, for all  $F$ ,*

$$\pi_t (F(x), u^{-1}\mathbb{E}_F[u(\tilde{x})|\tilde{x} \leq x], u^{-1}\mathbb{E}_F[u(\tilde{x})|\tilde{x} > x])$$

*is weakly increasing in  $x$  over  $\{x : F(x) \in (0, 1)\}$ , and*

$$\pi_t(0, \emptyset, c_H) = 0$$

$$\pi_t(1, c_L, \emptyset) = 1$$

$$\forall q \in (0, 1) \pi_t(q, c_L, c_H) \in (0, 1).$$

The continuity condition for a probability distortion function is tricky because  $\pi_t$  sometimes takes  $\emptyset$  as an argument. We would like  $\pi_t$  to be a continuous function over the subset of the domain that excludes  $\emptyset$ . We would also like  $\pi_t$  to approach 0 when the weight  $q$  on the worse outcome  $c_L$  vanishes, and to approach 1 when the weight  $1 - q$  on the better outcome  $c_H$  vanishes. The following

definition captures these requirements.

**Definition 35** (Continuous  $\pi_t$ ). *A probability distortion function  $\pi_t$  is continuous if for any sequence  $(q^n, c_L^n, c_H^n)$  of vectors in its domain that converges to  $(q, c_L, c_H)$  with  $c_L < c_H$ ,*

$$\lim \pi_t(q^n, c_L^n, c_H^n) = \pi_t(q, c_L^*, c_H^*)$$

$$\text{where } c_L^* = \begin{cases} c_L & \text{if } q \neq 0 \\ \emptyset & \text{if } q = 0 \end{cases}, \quad c_H^* = \begin{cases} c_H & \text{if } q \neq 1 \\ \emptyset & \text{if } q = 1 \end{cases}$$

**Lemma 1.**

$$\hat{F}_t(x) \equiv \begin{cases} \pi_t(F(x), u^{-1}\mathbb{E}_F[u(\tilde{x})|\tilde{x} \leq x], u^{-1}\mathbb{E}_F[u(\tilde{x})|\tilde{x} > x]) & \text{if } F(x) \in (0, 1) \\ F(x) & \text{if } F(x) \in \{0, 1\} \end{cases}$$

is a cdf if  $\pi_t$  is a continuous probability distortion function.

*Proof.* The only requirement that is not obvious is right-continuity. Since  $\pi_t$  is a continuous function over the subset of its domain that excludes  $\emptyset$ ,  $\hat{F}_t$  can only violate right-continuity at some  $x$  such that  $F(x) \in (0, 1)$  if one or more of its arguments does. But all three of its arguments are right-continuous. Since  $\hat{F}_t$  is flat at 0 for all  $x$  such that  $F(x) = 0$  and flat at 1 for all  $x$  such that  $F(x) = 1$ , it cannot violate right-continuity there either.  $\square$

The main condition for  $\pi_t$  to be a probability distortion function (and  $\hat{F}_t$  to be a cdf) is not easy to interpret. It is not immediately obvious what a function satisfying this condition might look like, or how to find a function that does. Of course,  $\pi_t(q, c_L, c_H) = f(q)$  with  $f$  increasing will work. Thus, the model nests rank-dependent expected utility (RDEU).

One easy option is to get dependence on  $c_L, c_H$  by making  $\pi_t$  increase in  $q(u(c_H) - u(c_L))$  or

decrease in  $(1 - q)(u(c_H) - u(c_L))$ . This works because

$$F(x) (\mathbb{E}[u(\tilde{x})|\tilde{x} \leq x] - \mathbb{E}[u(\tilde{x})|\tilde{x} > x])$$

and  $(1 - F(x)) (\mathbb{E}[u(\tilde{x})|\tilde{x} \leq x] - \mathbb{E}[u(\tilde{x})|\tilde{x} > x])$

are increasing in  $x$  for any lottery  $F$ . These specifications have a nice intuitive interpretation. If  $\pi_t$  is increasing in  $q(u(c_H) - u(c_L))$ , the agent puts a greater weight on a bad outcome in a binary lottery when it is probable and when it is very bad compared to the good outcome. If  $\pi_t$  is decreasing in  $(1 - q)(u(c_H) - u(c_L))$ , the agent puts a greater weight on the bad outcome when it is probable and when it is *not* too bad compared to the good outcome.

Another convenient specification is

$$\frac{\pi_t(q, c_L, c_H)}{1 - \pi_t(q, c_L, c_H)} = g(c_L, c_H) \frac{q}{1 - q}$$

with  $g$  defined as follows. If  $u$  is bounded below by  $\circ$ , set

$$g(c_L, c_H) = (1 + au(c_H) + bu(c_L))^{-\alpha}$$

with  $\alpha \in (0, 1]$  and  $a \geq b \geq 0$ . Since  $g \leq 1$ , an agent with this  $g$  is always optimistic: he downweights the probability of worse outcomes relative to better ones. Moreover, he becomes more optimistic as the value of the worse outcome and/or the better outcome improves. This ensures that his preferences satisfy FOSD-monotonicity.

If  $u$  is bounded above by  $\circ$ , set

$$g(c_L, c_H) = (1 - au(c_H) - bu(c_L))^\alpha$$



with  $\alpha \in (0, 1]$  and  $b \geq a \geq 0$ . Since  $g \geq 1$ , an agent with this  $g$  is always pessimistic. He becomes less pessimistic as the value of the worse outcome and/or the better outcome improves, so his preferences satisfy FOSD-monotonicity.<sup>3</sup>

It may be useful to see this  $\pi_t$  in action, turning a physical lottery  $F$  into a less favorable distorted lottery  $\hat{F}_t$ . Let  $u = -\exp(-x)$  and  $\alpha = a = b = 1$ . Consider a lottery that puts weight  $1/3$  on each of 0, 1, and 2. We have

$$\begin{aligned}\hat{F}_t(0) &= \pi_t \left( \frac{1}{3}, 0, -\ln \left( \frac{1}{2} \exp(-1) + \frac{1}{2} \exp(-2) \right) \right) \\ &= 0.53 \\ \hat{F}_t(1) &= \pi_t \left( \frac{2}{3}, -\ln \left( \frac{1}{2} \exp(0) + \frac{1}{2} \exp(-1) \right), 1 \right) \\ &= 0.78.\end{aligned}$$

Thus, the distorted lottery puts weight over  $1/2$  on 0, about  $1/4$  on 1, and about  $1/5$  on 2.

### 3.2.3 MAIN REPRESENTATION

This section places the probability distortion function within the broader representation. The notation for the representation is easier to grasp with an intuitive understanding of its structure. Each tree  $d$  will be identified with a payoff such that the agent is indifferent between getting  $d$  and getting an infinite stream of that payoff. The payoff is the expected-utility certainty equivalent of the lottery with cdf  $F_d$ . It is denoted  $u^{-1}\mathbb{E}_{F_d}u(\tilde{x})$ , and it is called the “value” of  $d$ .

The same idea can be applied to any  $t$ -subtree  $d_t$  of a tree  $d$ . Each  $d_t$  is identified with a payoff that can be substituted for  $d_t$  without changing the value of the overall tree. This payoff is the expected-utility certainty equivalent of the lottery with cdf  $\hat{F}_{d_t}$ . It is denoted  $u^{-1}\mathbb{E}_{\hat{F}_{d_t}}u(\tilde{x})$ , and it is called the

---

<sup>3</sup>To understand where these functional forms come from, or for guidance on generating additional functional forms), see Appendix C.2.

“value” of  $d_t$  at  $t$ .

$F_d$  and the  $\hat{F}_{d_t}$  are obtained by working backwards through the tree. The idea is to transform  $d$  into the lottery  $F_d$  by making a series of modifications that do not affect its value. This process is easiest to understand when  $d$  is a lottery-terminating tree: there is some  $T$  such that each  $T$ -subtree  $d_T$  of  $d$  is a lottery. In this case,  $F_{d_T}$  is the cdf associated with  $d_T$ , and  $\hat{F}_{d_T}$  is the distorted cdf obtained by applying the probability distortion function  $\pi_T$  to  $F_{d_T}$ .

We use each  $\hat{F}_{d_T}$  to compute the value  $u^{-1}\mathbb{E}_{\hat{F}_{d_T}} u(\tilde{x})$  of the corresponding  $d_T$ , and then we replace each  $d_T$  with an infinite sequence of its value. In the resulting tree, each  $(T-1)$ -subtree  $d_{T-1}$  is a measure over streams of the form  $(a, b^\infty)$ . We use the time-invariant aggregator  $W$  to replace each  $(a, b^\infty)$  with the constant stream  $W(a, b)^\infty$ . Now each  $(T-1)$ -subtree is a lottery and can be associated with a cdf  $F_{d_{T-1}}$ . We iterate the process until we are left with a single lottery resolving in period 0. This is  $F_d$ .

**Definition 36** (Aggregator).  $W : C^2 \rightarrow C$  is an aggregator if it is continuous, strictly increasing in both arguments, and satisfies  $W(c, c) = c$  for all  $c \in C$ .

**Definition 37** (Distorted-probability representation). Let  $\{\pi_t\}_{t=1}^\infty$  be a set of continuous probability distortion functions for a strictly increasing and continuous Bernoulli utility  $u$ . Let  $W$  be an aggregator.  $(\{\pi_t\}_{t=1}^\infty, u, W)$  is a distorted-probability representation for  $\succeq$  if

$$\forall d \in D \quad d \sim \left( u^{-1}\mathbb{E}_{F_d} [u(\tilde{x})] \right)^\infty,$$

where  $F_d$  is recursively defined by

$$\begin{aligned} F_{d_t}(x) &= d_t \left\{ (c_t, d_{t+1}) : W \left( c_t, u^{-1}\mathbb{E}_{\hat{F}_{d_{t+1}}} u(\tilde{x}) \right) \leq x \right\} \\ \hat{F}_{d_t}(x) &= \pi_t \left( F_{d_t}(x), u^{-1}\mathbb{E}_{F_{d_t}} [u(\tilde{x}) | \tilde{x} \leq x], u^{-1}\mathbb{E}_{F_{d_t}} [u(\tilde{x}) | \tilde{x} > x] \right) \end{aligned}$$

The structure of the representation, aside from the probability distortions, should already be familiar from [Epstein and Zin \(1989\)](#). Epstein and Zin specialize to a CES aggregator  $W$  and a CRRA Bernoulli utility  $u$ . Plugging these into the representation above gives

$$\begin{aligned}
 d &\sim \left( (\mathbb{E}_{F_d} \tilde{x}^\alpha)^{1/\alpha} \right)^\infty \\
 F_{d_t}(x) &= d_t \left\{ (c_t, d_{t+1}) : \left( (1 - \delta)c_t^\rho + \delta \left( \mathbb{E}_{\hat{F}_{d_{t+1}}} \tilde{x}^\alpha \right)^{\rho/\alpha} \right)^{1/\rho} \leq x \right\} \\
 \hat{F}_{d_t}(x) &= \pi_t \left( F_{d_t}(x), \mathbb{E}_{\hat{F}_{d_{t+1}}} [\tilde{x}^\alpha | \tilde{x} \leq x]^{1/\alpha}, \mathbb{E}_{\hat{F}_{d_{t+1}}} [\tilde{x}^\alpha | \tilde{x} > x]^{1/\alpha} \right)
 \end{aligned}$$

where  $\delta \in [0, 1)$ ,  $\alpha \in (-\infty, 1]$ , and  $\rho \in (0, 1]$ . This is Epstein-Zin utility with probability distortions. Epstein-Zin utility itself is the special case  $\pi_t(q, c_L, c_H) = q$ .

### 3.3 PARAMETRIC SPECIAL CASES

This section discusses a convenient feature of the probability distortions introduced in this paper: it is often possible to choose the distortion function  $\pi_t$  so that it preserves a particular class of distribution (e.g. lognormal, normal, or exponential). An agent with such a  $\pi_t$  will behave as if he misperceives the parameters of distributions in the relevant class. However, he will still be able to evaluate distributions outside that class. This is an improvement over models in which a new distortion must be chosen for each class of distribution the agent might face.

The main difficulty in choosing  $\pi_t$  is not achieving the correct mapping between cdfs in a given class, but ensuring that  $\pi_t$  always produces a valid cdf when confronted with distributions *outside* that class. This restricts the range of possible parameter distortions. Illustrative restrictions appear in [Propositions 13 and 15](#).

### 3.3.1 LOGNORMAL/NORMAL DISTRIBUTIONS

Consider an agent or group of agents who will face lognormally distributed payoffs at some time  $t$ . Suppose we don't wish to restrict the possible parameters of the distribution; perhaps they are initially uncertain, or can differ from agent to agent. In this case, we may want to find a  $\pi_t$  that always preserves lognormality. Whenever an agent with this  $\pi_t$  faces a given lognormal distribution, he should behave like an EU agent who faces a lognormal distribution with different parameters. Proposition 13 provides a way to achieve this when Bernoulli utility is CRRA.

**Proposition 13.** *Suppose*

$$u(x) = \frac{x^{1-\gamma}}{1-\gamma}$$

for  $\gamma \in [0, 1) \cup (1, \infty)$ . Let

$$\begin{aligned} \pi_t(q, c_H, c_L) &= \Phi \left( \Phi^{-1}(q) - \frac{\hat{\mu}(\mu, \sigma) - \mu}{\sigma} \right) \\ \text{where } \mu &= \frac{1}{1-\gamma} \ln \left( (1-\gamma) (qu(c_L) + (1-q)u(c_H)) \right) - \frac{1}{2}(1-\gamma)\sigma^2 \\ \sigma &= \frac{1}{1-\gamma} \left( \Phi^{-1}(q) - \Phi^{-1} \left( \frac{qu(c_L)}{qu(c_L) + (1-q)u(c_H)} \right) \right) \end{aligned}$$

for some differentiable  $\hat{\mu}(\mu, \sigma)$  that satisfies

$$\hat{\mu}_1 - \frac{\hat{\mu}_2}{(1-\gamma)\sigma} + \frac{\hat{\mu} - \mu}{(1-\gamma)\sigma^2} \in [0, 1].$$

Then  $\pi_t$  is a probability distortion function, and an agent with this  $\pi_t$  will map  $\log \mathcal{N}(\mu, \sigma^2)$  into  $\log \mathcal{N}(\hat{\mu}, \sigma^2)$ .

It is worth pointing out two convenient  $\hat{\mu}$  that satisfy the restriction above. The first is

$$\hat{\mu}(\mu, \sigma) = \mu + k\sigma$$

for any  $k \in \mathbb{R}$ . This works because it makes  $\pi_t$  independent of  $c_L$  and  $c_H$  (so this is a case of RDEU).

The second is

$$\hat{\mu}(\mu, \sigma) = \mu + k(1 - \gamma)\sigma^2$$

with  $k \in [0, 1]$ . This specification turns out to have a nice property:  $k$  is a measure of the tendency to fantasize (for agents with  $\gamma < 1$ ) or worry (for agents with  $\gamma > 1$ ). It may prove useful in applications to have a parameter that governs anticipatory emotions and can move independently of risk aversion  $\gamma$ .<sup>4</sup> This point will be formalized in Section 3.4.

A very similar result can be obtained for normal distributions when utility is CARA.

**Proposition 14.** *Suppose*

$$u(x) = -\exp(-\alpha x)$$

for  $\alpha \in (0, \infty)$ . Let

$$\begin{aligned} \pi_t(q, c_H, c_L) &= \Phi \left( \Phi^{-1}(q) - \frac{\hat{\mu}(\mu, \sigma) - \mu}{\sigma} \right) \\ \text{where } \mu &= -\frac{1}{\alpha} \ln(-qu(c_L) - (1-q)u(c_H)) + \frac{1}{2}\alpha\sigma^2 \\ \sigma &= \frac{1}{\alpha} \left( \Phi^{-1} \left( \frac{qu(c_L)}{qu(c_L) + (1-q)u(c_H)} \right) - \Phi^{-1}(q) \right) \end{aligned}$$

for some differentiable  $\hat{\mu}(\mu, \sigma)$  that

$$\hat{\mu}_1 + \frac{\hat{\mu}_2}{\alpha\sigma} - \frac{\hat{\mu} - \mu}{\alpha\sigma^2} \in [0, 1].$$

Then  $\pi_t$  is a probability distortion function, and an agent with this  $\pi_t$  will map  $\mathcal{N}(\mu, \sigma^2)$  into  $\mathcal{N}(\hat{\mu}, \sigma^2)$ .

The convenient specifications are  $\hat{\mu}(\mu, \sigma) = \mu + k\sigma$  for  $k \in \mathbb{R}$ , which is a case of RDEU, and

---

<sup>4</sup>Since risk resolved at  $t = 0$  is evaluated by expected utility with parameter  $\gamma$ ,  $\gamma$  can still be interpreted as aversion to risk in the absence of anticipatory emotions.

$\hat{\mu}(\mu, \sigma) = \mu - k\alpha\sigma^2$  with  $k \in [0, 1]$ .  $k$  turns out to be a measure of the tendency to worry.

### 3.3.2 EXPONENTIAL DISTRIBUTIONS

Consider an agent who will face exponentially distributed payoffs at some time  $t$ . Again, we look for a  $\pi_t$  that preserves exponential distributions while changing  $\lambda$ . We assume CRRA utility.

**Proposition 15.** *Suppose*

$$u(x) = \frac{x^{1-\gamma}}{1-\gamma}$$

for  $\gamma \in [0, 1) \cup (1, \infty)$ . Let

$$\pi_t(q, c_H) = 1 - \exp\left(\frac{\hat{\lambda}(\lambda)}{\lambda} \ln(1-q)\right)$$

$$\text{where } \lambda = \frac{\left(\frac{1}{1-q} \int_q^1 (-\ln(1-x))^{1-\gamma} dx\right)^{\frac{1}{1-\gamma}}}{c_H}$$

for some continuous, increasing  $\hat{\lambda}(\lambda)$  that satisfies

$$\hat{\lambda} \geq \lambda \hat{\lambda}'.$$

Then  $\hat{\pi}_t$  is a probability distortion function, and an agent with this probability distortion function will map  $\text{Exp}(\lambda)$  into  $\text{Exp}(\hat{\lambda})$ .

Again, it is worth noting one convenient case:  $\hat{\lambda} = k\lambda$  for any  $k \in \mathbb{R}^+$ . This works because it makes  $\pi_t$  independent of  $c_H$  (so this is a case of RDEU).

### 3.4 COMPARATIVE STATICS

This section provides meaning to the statements “ $\succsim$  fantasizes” and “ $\succsim^A$  fantasizes more than  $\succsim^B$ .” In defining these concepts, I assume that there is such a thing as “the value  $I_t(d_t)$  that  $\succsim$  places on subtree  $d_t$  at time  $t$ .” That is, I assume there is a unique function  $I_t$  from  $t$ -subtrees to payoffs that satisfies the following: for all  $d \in D$ , if  $d'$  is constructed from  $d$  by replacing each  $d_t$  with an infinite sequence of  $I_t(d_t)$ , then  $d' \sim d$ . I also assume that  $\succsim$  prefers larger payoffs to smaller ones, so  $I_t(d_t) > I_t(f_t)$  indicates a preference for  $d_t$  over  $f_t$ .

The main representation satisfies these conditions. It has

$$I_0(d) = u^{-1}\mathbb{E}_{F_d}u(\tilde{x})$$

$$\forall t > 0 \ I_t(d_t) = u^{-1}\mathbb{E}_{\tilde{F}_{d_t}}u(\tilde{x}).$$

The assumptions required to

These assumptions are much weaker than those needed for a full distorted-probability representation. Still, some generality is lost in assuming that the value of a  $t$ -subtree  $d_t$  depends on  $d_t$  and  $t$ , but not on other features of the overall tree  $d$ . Generality could be regained by allowing the  $I_t$  to take more than one argument and conditioning the statements below on that argument. For instance, the value of a subtree  $d_t$  could depend on the probability of reaching it. Then, one could make statements like “A fantasizes more than B about  $t$ -subtrees with probability  $\alpha$ .” I do not take this discussion further because it does not contribute much to an understanding of a distorted-probability representation.

Note that definitions and results for “ $\succsim$  worries” and “ $\succsim^A$  worries more than  $\succsim^B$ ” can easily be obtained by reversing the inequalities below. In this model, there is nothing special about fantasy as distinct from worry.

**Definition 38** (Fantasizes).  $\succeq$  fantasizes about  $t$  if

$$I_0(F) \leq I_t(F).$$

Definition 38 should be intuitive: for an agent who fantasizes, the value of a given lottery rises as it is shifted from the present (in which there is no opportunity to fantasize) into the future. Proposition 16 says that this translates into a natural condition on the representation: distortions must reduce the probability of doing worse than  $c$  relative to the probability of doing better.

**Proposition 16.** *Suppose  $\succeq$  has a distorted-probability representation.  $\succeq$  fantasizes about  $t$  if and only if*

$$\forall q, c_L, c_H \quad \pi_t(q, c_L, c_H) \leq q.$$

The difficulty in formulating a definition of “A fantasizes more than B” is that A and B may have different preferences over lotteries at time 0. Intuitively, if A’s valuation of  $F$  goes up a lot between 0 and  $t$ , and B’s only goes up a little, there are two possible explanations. One is that A has a greater predisposition to fantasize. Another is that A and B are in different situations, and A’s situation is more conducive to fantasy (holding the predisposition to fantasize constant). Specifically, A may be more comfortable with the risk in  $F$  than B is, and his greater optimism may come from his (subjectively) superior prospects.

There are two possible solutions. The simplest is to require that A and B have the same preferences over lotteries at time 0. Then if A likes  $F$  at time  $t > 0$  more than B does, it must be because he has a greater predisposition to fantasize. This is formalized in Definition 39. Proposition 17 shows that A fantasizes more than B if and only if A always pushes up the probability of doing better than  $c$  more than B does.



**Definition 39** (Fantasizes more than).  $\succeq^A$  fantasizes more than  $\succeq^B$  about  $t$  if

$$\forall F \quad I_0^A(F) = I_0^B(F)$$

and

$$\forall F \quad I_t^A(F) \geq I_t^B(F).$$

**Proposition 17.** Suppose  $\succeq^A$  and  $\succeq^B$  have distorted-probability representations.  $\succeq^A$  fantasizes more than  $\succeq^B$  about  $t$  if and only if  $u_A = \alpha u_B + \beta$  (where  $\alpha > 0$ ) and

$$\forall q, c_L, c_H \quad \pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, c_L, c_H).$$

The second possible solution is to allow A and B to have different risk preferences, but to compare their valuations over lotteries on different domains. Suppose A faces lottery  $F^A$ . The idea is to construct  $F^B$  so that B's preferences over lotteries with payoffs in  $\text{supp}(F^B)$  perfectly mirror A's preferences over lotteries with payoffs in  $\text{supp}(F^A)$ . If B is less risk averse than A, for instance, the bad payoffs in  $F^B$  can be pushed down and the good ones pushed up so that B's situation becomes as "subjectively risky" as A's. This idea is formalized below.

**Definition 40** (Comparable/comparison function).  $\succeq^A$  and  $\succeq^B$  are comparable if there exists some  $\varphi : C \rightarrow C$  such that

$$\forall F \quad I_0^A(F) = \varphi^{-1}(I_0^B(\varphi(F))) \text{ and } I_0^B(F) = \varphi(I_0^A(\varphi^{-1}(F))).$$

In this case,  $\varphi$  is an A-to-B comparison function.

**Definition 41** (Fantasizes more than).  $\succeq^A$  fantasizes more than  $\succeq^B$  about  $t$  if A and B are compara-

ble and, for every A-to-B comparison function  $\varphi$ ,

$$\forall F \quad I_t^A(F) \geq \varphi^{-1}(I_t^B(\varphi(F))).$$

Unfortunately, there may be many A-to-B comparison functions, and they may not have the same implications for probability distortions. This problem goes away if both Bernoulli utility functions are bounded on both sides. There is no further loss of generality in assuming that they have the same range (since this can always be accomplished through rescaling).

**Proposition 18.** *If  $\succeq^A$  and  $\succeq^B$  have distorted-probability representations with  $u^A$  and  $u_B$  bounded and  $\text{range}(u_A) = \text{range}(u_B)$ , then  $u_B^{-1}(u_A)$  is the unique A-to-B comparison function, and  $\succeq^A$  fantasizes more than  $\succeq^B$  about  $t$  iff*

$$\forall q, c_L, c_H \quad \pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, u_B^{-1}(u_A(c_L)), u_B^{-1}(u_A(c_H))).$$

If A and B are unbounded on the same side, there will still be many A-to-B comparison functions, but we can force all of them to have the same implications for probability distortions. Thus, we can still do comparative statics for CRRA and CARA utility functions (among others) if we choose  $\pi_t$  appropriately.

**Proposition 19.** *If  $\succeq^A$  and  $\succeq^B$  have distorted-probability representations with  $u_A$  and  $u_B$  unbounded on one side by 0 and unbounded on the other, and*

$$\forall \alpha > 0 \quad \pi_t^I(q, c_L, c_H) = \pi_t^I(q, u_I^{-1}(\alpha u_I(c_L)), u_I^{-1}(\alpha u_I(c_H))),$$

*then  $\succeq^A$  fantasizes more than  $\succeq^B$  about  $t$  iff*

$$\forall q, c_L, c_H \quad \pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, u_B^{-1}(u_A(c_L)), u_B^{-1}(u_A(c_H))).$$

Proposition 19 is not just a mathematical curiosity. An appropriate  $\pi_t$  arose naturally in Section 3.3.1. Corollary 7 applies Proposition 19 to this case.

**Corollary 7.** *Suppose that  $\succeq^A$  and  $\succeq^B$  have distorted-probability representations where  $u_I$  and  $\pi_t^I$  take the forms in Proposition 13. Suppose that either  $\gamma_A, \gamma_B > 1$  or  $\gamma_A, \gamma_B < 1$ . If*

$$\frac{\hat{\mu}^I(\mu, \sigma) - \mu}{\sigma} = f^I(\sigma),$$

*then  $\succeq^A$  fantasizes more than  $\succeq^B$  if and only if*

$$f^A(\sigma) \geq f^B\left(\frac{1 - \gamma^A}{1 - \gamma^B} \sigma\right).$$

*In the special case*

$$\hat{\mu}^I(\mu, \sigma) = \mu + k^I(1 - \gamma^I)\sigma^2,$$

*$\succeq^A$  fantasizes more than  $\succeq^B$  if  $\gamma_A, \gamma_B < 1$  and  $k_A \geq k_B$ , or  $\gamma_A, \gamma_B > 1$  and  $k_A \leq k_B$ .*

### 3.5 AXIOMS

As usual, we assume that  $\succeq$  is a weak order on  $D$ .

A standard certainty equivalent axiom (adapted to this setting) says that, for any tree, there exists a payoff such that the decision maker is indifferent between getting the tree and getting the infinitely repeated payoff. The below is somewhat stronger: it says that, for any tree starting in future, there exists a payoff such that the decision maker is indifferent between getting the tree in future and getting the infinitely repeated payoff in future. This payoff is the recursive certainty equivalent (loosely, the value) of the future tree. Note that the value of a tree need not stay the same when the tree is shifted forward. This is a key feature of the model: the opportunity for fantasy/worry depends on

the period in which a risk is resolved, so the value of the risk does too.

**Axiom 19** (Recursive Certainty Equivalence). *For any  $(0^t, d_t) \in D$ ,*

$$\exists c \in C \quad (0^t, c^\infty) \sim (0^t, d_t).$$

The next axiom says that the value of a subtree depends only on the period in which the subtree starts. It does not depend on the payoffs that precede it (history) or other subtrees that do not overlap with it (counterfactuals).

**Axiom 20** (History/Counterfactual Irrelevance). *If tree  $d'$  is constructed from  $d$  by replacing each  $t$ -subtree  $d_t$  with  $c(d_t)^\infty$ , where*

$$(0^t, c(d_t)^\infty) \sim (0^t, d_t),$$

*then  $d' \sim d$ .*

No Savoring looks much like History/Counterfactual Irrelevance, but with two important differences. First, it deals with deterministic consumption streams rather than general subtrees. Second, the “value”  $c(c_t, c_{t+1}, \dots)$  of a stream is assessed by moving the stream back to period 0. (Compare to History/Counterfactual Irrelevance, in which the “value”  $c(d_t)$  of a subtree is assessed without moving  $d_t$  out of period  $t$ .) Thus, No Savoring says that the value of a stream is independent of the period in which it takes place as well as its history and counterfactuals.<sup>5</sup>

**Axiom 21** (No Savoring). *Fix  $d = (0^t, d_t)$  where  $d_t$  is a measure over deterministic sequences. If  $d'$  is constructed from  $d$  by replacing each  $(c_t, c_{t+1}, \dots)$  in the support of  $d_t$  with  $c(c_t, c_{t+1}, \dots)^\infty$ , where*

$$c(c_t, c_{t+1}, \dots)^\infty \sim (c_t, c_{t+1}, \dots),$$

---

<sup>5</sup>Note that No Savoring does not rule out discounting. It says that  $y^\infty \sim (c_0, c_1, \dots)$  implies  $(0^t, y^\infty) \sim (0^t, c_0, c_1, \dots)$ . It does not say that  $(c_0, c_1, \dots) \sim (0^t, c_0, c_1, \dots)$ . The latter rules out discounting; the former is perfectly consistent with it.

then  $d' \sim d$ .

Monotonicity says that the value of a measure over consumption streams is between the value of (1) a stream consisting of the worst possible payoffs period by period and (2) a stream consisting of the best possible payoffs period by period. This ensures that a large certain payoff is preferred to a small certain payoff and that the value of a lottery is between its worst outcome and its best outcome. It does not ensure that the decision-maker's preferences are monotonic in first-order stochastic dominance. FOSD-monotonicity is treated in Section 3.6.1.

**Axiom 22** (Monotonicity). *Fix any  $t \in \{0, 1, \dots\}$  and any  $d = (0^t, d_t)$  where  $d_t$  is a nondegenerate measure over ultimately constant deterministic sequences. If*

$$\forall (c_t, c_{t+1}, \dots) \in \text{supp}(d_t) \quad (c_t^L, c_{t+1}^L, \dots) \leq (c_t, c_{t+1}, \dots) \leq (c_t^H, c_{t+1}^H, \dots),$$

then

$$(0^t, c_t^L, c_{t+1}^L, \dots) \prec d \prec (0^t, c_t^H, c_{t+1}^H, \dots).$$

Recall that a tree  $d$  is a lottery if it consists of a probability measure on constant streams, and that a lottery can be identified with a cdf  $F$ . Section 3.2.1 mentioned the truncations  $F^a$  and  $F_a$  of  $F$ . Since they also appear in the third part of the Continuity axiom, we define them more formally below. For any  $a \in C$ , let

$$F^a(c) = \begin{cases} F(c) & \text{if } c < a \\ 1 & \text{if } c \geq a. \end{cases} \quad F_a(c) = \begin{cases} F(c) & \text{if } c \geq a \\ 0 & \text{if } c < a. \end{cases}$$

Truncation Continuity is from [Wakker \(1993\)](#). It helps secure an EU representation without imposing bounded  $u$ . Deterministic Continuity ensures a continuous time aggregator, and Binary Lottery Continuity ensures a continuous probability distortion function.

**Axiom 23** (Continuity).

1. (Deterministic Continuity) For any open  $B$ ,  $\{(a, b) : \exists c \in B (a, b^\infty) \sim c\}$  is open in  $C \times C$ .
2. (Binary Lottery Continuity) For any open  $B$ ,  $\{pa \oplus (1-p)b : \exists c \in B (0^t, pa^\infty \oplus (1-p)b^\infty) \sim (0^t, c^\infty)\}$  is open.
3. (Truncation Continuity) For any lottery  $F$ ,

$$F \succ c^\infty \Rightarrow \exists a F^a \succ c^\infty$$

$$F \prec c^\infty \Rightarrow \exists a F_a \prec c^\infty$$

The below is a standard independence axiom. Note that it only applies to risk resolved immediately. Preferences over risks resolved in future need not satisfy independence because they need not have an EU representation.

**Axiom 24** (Initial Independence). For any lotteries  $F, F', F''$ ,

$$\forall \alpha \in (0, 1) F \succ (\sim)F' \Rightarrow \alpha F + (1 - \alpha)F'' \succ (\sim)F' + (1 - \alpha)F''.$$

A new definition simplifies the exposition of the next axiom. The first part of the definition applies only to binary lotteries. We have already seen that the value of a binary lottery  $pc_1 \oplus (1-p)c_2$  may depend on the period in which it takes place. Say the value of this lottery at  $t$  is  $y$ . A distortion of  $pc_1 \oplus (1-p)c_2$  is a binary lottery  $\hat{p}c_1 \oplus (1-\hat{p})c_2$  that has value  $y$  at 0. Intuitively, the distorted lottery at period 0 is as good as the original lottery at period  $t$  because the probabilities have been altered to make up for any fantasies or worries associated with the latter.

**Definition 42** (Distortion (binary lottery)).  $\hat{p}_t c_1 \oplus (1 - \hat{p}_t) c_2$  is a distortion of  $p c_1 \oplus (1 - p) c_2$  at  $t$  if

$$(0^t, p c_1^\infty \oplus (1 - p) c_2^\infty) \sim (0^t, y^\infty) \Rightarrow \hat{p} c_1^\infty \oplus (1 - \hat{p}) c_2^\infty \sim y^\infty.$$

The second part of the definition extends the notion of a distortion to non-binary lotteries. This is done by dividing a given lottery  $F$  into a “worse part” with time-0 value  $z^c$  and a “better part” with time-0 value  $z_c$ . A distortion is obtained for this binary lottery. Then the process is repeated for all  $c$ , i.e. for all ways of dividing  $F$  into better and worse parts. The resulting function  $\hat{F}$  may seem like an odd thing to define, but the next axiom establishes its utility.

Let  $F^c(F_c)$  denote the lottery obtained by conditioning lottery  $F$  on a payoff less than or equal to (greater than)  $c$ .

**Definition 43** (Distortion (general lottery)).  $\hat{F} : C \rightarrow [0, 1]$  is a distortion of  $F$  at  $t$  if, for all  $c$ , for some  $z^c$  and  $z_c$  such that

$$z^c \sim F^c \text{ and } z_c \sim F_c,$$

$\hat{F}(c) z^c \oplus (1 - \hat{F}(c)) z_c$  is a distortion of  $F(c) z^c \oplus (1 - F(c)) z_c$  at  $t$ .

**Axiom 25** (Binary Distortion). If  $\hat{F}$  is a distortion of  $F$  at  $t$ , then  $\hat{F} \in D$ , and

$$(0^t, F) \sim (0^t, y^\infty) \Rightarrow \hat{F} \sim y^\infty.$$

This is the key axiom of the model. It says that the value  $y$  of lottery  $F$  at  $t$  can be obtained by finding the distortion of  $F$  at  $t$ . This distortion will be a lottery in its own right, and it will have value  $y$ .

Intuitively, Binary Distortion says that the agent acts as if he misperceives distributions of future payoffs. These misperceptions follow a particular pattern. The agent divides a given distribution into a better part and a worse part, and adjusts the probabilities of the two parts just as he would if

he were facing a binary lottery between the certainty equivalents of the better part and the worse part. Said another way, the agent’s weight on “doing better than  $c$ ” doesn’t depend on the shape of the conditional distribution below  $c$  or above  $c$ , as long as the value of these distributions (absent fantasies and worries) is held fixed.

Notice that the agent considers each possible way to cut the distribution, and he distorts the probabilities at each cut. The distortions are “binary” because they can be identified from binary lotteries, not because the agent cuts each distribution in exactly one place.

**Proposition 20.** *Axioms 1-7 are necessary and sufficient for  $\succeq$  to have a distorted-probability representation.*

**Corollary 8.** *In a distorted-probability representation,  $u$  is unique up to a positive affine transformation.  $W$  and all the  $\pi_t$  are unique.*

## 3.6 EXTENSIONS

### 3.6.1 FIRST-ORDER STOCHASTIC DOMINANCE

Thus far, we have not required the agent’s probability distortions to be FOSD-monotonic. This generality is desirable because there are intuitively plausible cases that do not satisfy this condition. For instance, an agent may place more weight on getting the bad payoff in a binary lottery when it is *very* bad relative to the good payoff. This means the agent will place less weight on the bad payoff if the good payoff is reduced a bit. The new lottery will be strictly worse than the old lottery (in the sense of FOSD), but this will not be true of the distorted lotteries “perceived” by the agent.

Unfortunately, there is no clear way to impose FOSD-monotonicity of the agent’s preferences without imposing FOSD-monotonicity of his probability distortion function. We must choose between (1) allowing more flexibility in probability distortions while admitting the possibility a



dominated lottery may be chosen, and (2) ruling out this possibility along with some interesting specifications. The correct choice may depend on the application. The following axiom and proposition explain how to do (2).

The axiom is a condition on pairs of binary lotteries with the same probabilities but different payoffs, chosen so one lottery is FOSD-better than the other. Suppose that the agent doesn't like the good lottery at time  $t$  as much as he likes a given reweighted version of that lottery at time 0. Then, he must not like the *bad* lottery at time  $t$  as much as he likes the reweighted version of that lottery at time 0. Intuitively, the axiom says that the opportunity to fantasize is greater for a good lottery than a bad lottery. If a given reweighting is beneficial enough to offset the lost opportunity to fantasize about the good lottery, it will certainly offset the lost opportunity to fantasize about the bad one.

**Axiom 26** (Distortion Monotonicity). *If  $c_L < c'_L < c_H < c'_H$ , then*

$$Z_0(qc'_L \oplus (1 - q)c'_H) \geq Z_t(pc'_L \oplus (1 - p)c'_H) \Rightarrow Z_0(qc_L \oplus (1 - q)c_H) \geq Z_t(pc_L \oplus (1 - p)c_H).$$

**Proposition 21.** *Suppose  $\succeq$  has a distorted-probability representation. Distortion Monotonicity is necessary and sufficient for  $\pi_t(q, c_L, c_H)$  to be (weakly) increasing in  $(c_L, c_H)$ . It is sufficient for*

$$F \geq_{FOSD} G \Rightarrow \hat{F}_t(F) \geq_{FOSD} \hat{F}_t(G) \text{ and } I_t(F) \geq I_t(G)$$

$$F >_{FOSD} G \Rightarrow \hat{F}_t(F) >_{FOSD} \hat{F}_t(G) \text{ and } I_t(F) > I_t(G).$$

To see why the proposition works, consider a pair of lotteries  $F, G$  such that  $F \geq_{FOSD} G$ . Cut  $F$  and  $G$  at quantile  $q$  and rescale the parts above  $q$  so they become cdfs. The result for  $F$  will still be FOSD-better, so will have a higher EU certainty equivalent, than the result for  $G$ . The same is true of the parts below  $q$ . These certainty equivalents are inputs to the distortion process. The axiom says that better certainty equivalents lead to more optimism (equivalently, a lower cdf). This ensures

that the distorted version of  $F$  remains below the distorted version of  $G$ .

### 3.6.2 CONTEXT DEPENDENCE

Thus far, we have assumed that the value of a given subtree is pinned down by the structure of that subtree and the period in which it starts. To make this more concrete, consider an agent who will become a doctor if admitted to medical school, and a nurse otherwise. Once in medical school, he will become a surgeon if he does very well, and a GP otherwise. In the main representation, the value of becoming a surgeon does not depend on any of the agent's other possible career paths; it depends only on (1) the career path of a surgeon and (2) the amount of time the agent has to fantasize about becoming one. This section presents an extended representation in which (3) the value of becoming a GP may matter too. Intuitively, an agent who strongly dislikes the idea of becoming a GP may find that the juxtaposition of GP with surgery makes surgery look even more appealing (and GP even less). This may prompt him to fantasize even more about becoming a surgeon (and worry more about becoming a GP) than he otherwise would.

The extended representation will not allow the value of becoming a surgeon to depend on the career path of a nurse. This is because the surgeon part of the tree diverges from the nurse part early on, but from the GP part at the last minute. Since surgeon and GP appear in the same lottery but surgeon and nurse do not, it is plausible that the surgeon/GP contrast will have a greater effect on the agent's perception of a surgeon's life than the surgeon/nurse contrast. We make the simplifying assumption that the surgeon/nurse contrast does not matter at all.

In this example, GP is the "close counterfactual" of surgeon, while nurse is a "distant counterfactual." More formally, fix a  $t$ -subtree  $d_t$  of a tree  $d$  and select a  $(t + 1)$ -subtree  $d_{t+1}$  that emanates from  $d$ . The close counterfactual  $d_{t+1}^c$  of  $d_{t+1}$  is the consumption prospect faced by an agent who knows that he was in  $d_t$  last period and is not in  $d_{t+1}$  now. It is possible to draw  $(\tilde{c}_{t+1}, \tilde{d}_{t+2})$  from  $d_{t+1}^c$  as follows. (1) Draw  $(\tilde{c}_t, \tilde{d}_{t+1})$  from  $d_t$ , discarding and redrawing if  $d_{t+1}$  is drawn. (2) Draw

$(\tilde{c}_{t+1}, \tilde{d}_{t+2})$  from the result of (1). Notice that  $d_{t+1}$  does not have a close counterfactual if  $d_t$  is degenerate. In this case,  $d_{t+1}$  takes on its “standalone” value.

Close counterfactuals enter the extended representation in a simple way. For any  $t$ -subtree  $d_t$  with close counterfactual  $d_t^c$ , the probability distortions of  $d_t$  are allowed to depend on the standalone value of  $d_t^c$  and the probability that  $d_t$  will happen rather than something in  $d_t^c$ . Due to complications of dealing with zero-probability subtrees, we restrict attention to trees built up from simple lotteries (i.e. lotteries with finitely many outcomes).

**Definition 44** (Context-dependent distortion function). *Fix a continuous and strictly increasing Bernoulli utility  $u : C \rightarrow \mathbb{R}$ . A continuous function*

$$\pi_t : \{(q, p, c_L, c_H, y) \in [0, 1] \times (0, 1) \times C^3 : c_L < c_H\} \rightarrow [0, 1]$$

*is a context-dependent distortion function for  $u$  if*

$$\pi_t(0, p, c_L, c_H, y) = 0, \pi_t(1, p, c_L, c_H, y) = 1,$$

*and for any cdf  $F$  associated with any simple lottery, any  $y \in C$ , and any  $p \in (0, 1)$ , the following is weakly increasing in  $x$ :*

$$\pi_t(F(x), p, u^{-1}\mathbb{E}_F[u(\tilde{x}) | \tilde{x} \leq x], u^{-1}\mathbb{E}_F[u(\tilde{x}) | \tilde{x} > x], z).$$

**Definition 45** (Context-dependent representation). *Let  $\{\pi_t\}_{t=1}^\infty$  be a set of context-dependent distortion functions for Bernoulli utility  $u$ . Let  $W : C \times C \rightarrow C$  be continuous and strictly increasing in both arguments, and suppose that  $W(c, c) = c$ .  $(\{\pi_t\}_{t=1}^\infty, u, W)$  is a context-dependent representation for  $\succeq$  if*

$$\forall d \in D \ d \sim I_0(d)^\infty,$$

where

$$\begin{aligned}
I_0(d) &= u^{-1} \mathbb{E}_{F_d} [u(\tilde{x})] \\
I_t(d_t; P_{t-1}(d_t), d_t^c) &= u^{-1} \mathbb{E}_{\hat{F}_{d_t}} [u(\tilde{x})] \\
F_{d_t}(x) &= d_t \{ (c_t, d_{t+1}) : W(c_t, I_{t+1}(d_{t+1}; P_t(d_{t+1}), d_{t+1}^c)) \leq x \} \\
\hat{F}_{d_t}(x) &= \pi_t(F_{d_t}(x), P_{t-1}(d_t), u^{-1} \mathbb{E}_F [u(\tilde{x}) | \tilde{x} \leq x], u^{-1} \mathbb{E}_F [u(\tilde{x}) | \tilde{x} > x], I_t(d_t^c)).
\end{aligned}$$

### 3.7 FUTURE WORK

There are two main avenues for future research. The first is an exploration of the effects of fantasy and worry on the demand for non-instrumental information. This is natural because people often explain a desire to obtain or avoid information in terms of anticipatory emotions. There is no need to build into the model an opportunity to buy information; rather, trees that provide early signals about future payoffs can be compared to trees that do not. The results of this exercise are not obvious. Although an agent who fantasizes will prefer to avoid total resolution of uncertainty, which eliminates the opportunity to fantasize, he will not necessarily be averse to partial resolution, which could provide even greater opportunities to fantasize. Moreover, an agent may like information structures that provide more precise information following a good signal but not a bad one (or the other way around).

The second avenue is an extension of the model from a static to a dynamic choice problem. An agent who fantasizes or worries will be dynamically inconsistent because (as has been emphasized many times) his valuation of a given tree changes as it is brought forward into the present. Thus, an extension to the dynamic case must take a stand on the sophistication or naivety of the agent. The naive case is technically simple because the agent will behave in each period as if he were making a choice once and for all (the environment modeled here). This case may still make some interesting

predictions, such as an optimistic agent's failure to plan appropriately for the pragmatic choices he will ultimately make.

# A

## Appendix to Chapter I

### A.1 PROOFS OF RESULTS IN TEXT

#### A.1.1 PROOF OF THEOREM I

First, we show necessity. Necessity of Optimization follows because the items that maximize a preference over a set must all be indifferent. For IUA, fix  $A$ ,  $a$  such that  $a \in A$ . Suppose  $a \succsim c(A)$  and  $a \notin c(A)$ , and fix  $B \supseteq A$ . For all  $\succsim_m \in \mathcal{M}$ , we have  $a \not\prec_m A$ , so  $a \not\prec_m B$ . To confirm that

$c(B) = c(B \setminus \{a\})$ , it suffices to show that  $\mathcal{M}(B) = \mathcal{M}(B \setminus \{a\})$ . Take any  $b \in \mathcal{M}(B)$ . Since  $b \succsim_m B$  for some  $\succsim_m \in \mathcal{M}$ ,  $b \succsim_m A$  for some  $\succsim_m \in \mathcal{M}$ , so  $b \neq a$ . Since  $b \succsim_m B$  implies  $b \succsim_m B \setminus \{a\}$ , we have  $b \in \mathcal{M}(B \setminus \{a\})$ . Now take any  $b \in \mathcal{M}(B \setminus \{a\})$ . There exists  $\succsim_m \in \mathcal{M}$  such that  $b \succsim_m B \setminus \{a\}$ . Since it cannot be that  $a \succsim_m b \succsim_m B \setminus \{a\}$ , we have  $b \succsim_m B$ , so  $b \in \mathcal{M}(B)$ .

Now we show sufficiency. To define  $\mathcal{M}$ , we need the notion of exclusion from below.

**Definition 46** (Exclusion from below ( $\triangleright$ )). *Say  $X \in \mathcal{F}(\mathcal{A})$  excludes  $y \notin X$  from below (written  $X \triangleright y$ ) if  $y \succsim X$  and  $y \notin c(X \cup \{y\})$ .*

Say that a strict preference  $\succsim_m$  respects exclusion from below if

$$X \triangleright y \implies \exists x \in X \ x \succ_m y.$$

We take  $\mathcal{M}$  to be the set of strict preferences on  $\mathcal{A}$  that respect exclusion from below.

Fix  $A \in \mathcal{F}(\mathcal{A})$  and  $b \notin A$  such that  $b \succ c(A \cup \{b\})$  and  $b \notin c(A \cup \{b\})$ . We show that  $b \notin \mathcal{M}(A \cup \{b\})$ .

**Lemma 2.** *Fix  $X \in \mathcal{F}(\mathcal{A})$  and  $y \notin X$ . If  $y \succ c(X \cup \{y\})$  and  $y \notin c(X \cup \{y\})$ , then*

$$c(X \cup \{y\}) \cup \{x \in X : y \succ x\} \triangleright y.$$

*Proof.* Take any  $x \in X$  such that  $x \succ y$  and  $x \notin c(X \cup \{y\})$ . Since  $y \succ c(X \cup \{y\})$ , we have  $x \succ c(X \cup \{y\})$ . By IUA,  $c(X \cup \{y\}) = c((X \cup \{y\}) \setminus \{x\})$ , so  $y \notin c((X \cup \{y\}) \setminus \{x\})$  and  $y \succ c((X \cup \{y\}) \setminus \{x\})$ . Iterating this argument, we can remove every  $x \in X$  such that  $x \succ y$  and  $x \notin c(X \cup \{y\})$  without changing choice. We end up with

$$y \notin c(\{y\} \cup c(X \cup \{y\}) \cup \{x \in X : y \succ x\}).$$

This implies  $c(X \cup \{y\}) \cup \{x \in X : y \succ x\} \triangleright y$ .  $\square$

Lemma 2 implies that  $c(A \cup \{b\}) \cup \{a \in A : b \succ a\} \triangleright b$ . Since each  $\succ_m \in \mathcal{M}$  respects exclusion from below, we have  $b \not\succeq_m c(A \cup \{b\}) \cup \{a \in A : b \succ a\}$  for all  $\succ_m \in \mathcal{M}$ . This implies  $b \not\succeq_m A$  for all  $\succ_m \in \mathcal{M}$ , so  $b \notin \mathcal{M}(A \cup \{b\})$ .

Now fix  $A \in \mathcal{F}(A)$  and  $b \notin A$  such that  $b \in c(A \cup \{b\})$ . We show that  $b \in \mathcal{M}(A \cup \{b\})$ : there is a strict preference  $\succ_m$  that respects exclusion from below and has  $b \succ_m A$ .

We will construct an appropriate  $\succ_m$  by extending  $\triangleright$ . We define several useful properties of  $\triangleright$ .

**Definition 47** (Menu-item relation). *A menu-item relation is a subset of  $(\mathcal{F}(A) \cup \{\emptyset\}) \times A$ .*

**Definition 48** (Transitivity). *A menu-item relation  $R$  is transitive if*

$$(X R x, Y R y \text{ and } x \in Y) \implies (X \cup Y) \setminus \{x, y\} R y.$$

We denote the transitive closure of a menu-item relation  $R$  by  $\text{tr}(R)$ .

To see why  $\triangleright$  is transitive, suppose  $X \triangleright x$ ,  $Y \triangleright y$ , and  $x \in Y$ . We have  $x \succ X$  and  $y \succ Y$ . Since  $x \in Y$ , we have  $y \succ x \succ X$ . IUA implies  $y$  is irrelevant for choice on any superset of  $Y$ , so  $y \notin c(X \cup Y \cup \{y\})$ . IUA also implies  $x$  is irrelevant for choice on any superset of  $X$ , so  $c(X \cup Y \cup \{y\}) = c((X \cup Y \cup \{y\}) \setminus \{x\})$ . We have  $y \notin c((X \cup Y \cup \{y\}) \setminus \{x\})$  as well as  $y \succ X \cup Y \setminus \{x, y\}$ , so  $X \cup Y \setminus \{x, y\} \triangleright y$ .

**Definition 49** (Properness). *A menu-item relation  $R$  is proper if  $X R x \implies X \neq \emptyset$ .*

**Definition 50** (Irreflexivity). *A menu-item relation  $R$  is irreflexive if  $X R x \implies x \notin X$ .*

**Definition 51** (Consistency with  $b \succ A$ ). *A menu-item relation  $R$  is consistent with  $b \succ A$  if it is not the case that  $A' R b$  for any  $A' \in \mathcal{F}(A)$ .*

To see why  $\triangleright$  is consistent with  $b \succ A$ , suppose  $A' \triangleright b$  for some  $A' \subset A$ . By IUA,  $b$  is irrelevant for choice on any superset of  $A'$ , including  $A$ . This contradicts  $b \notin c(A \cup \{b\})$ .



The following two lemmas will be useful for extending  $\triangleright$ .

**Lemma 3.** *Fix an irreflexive, transitive and proper menu-item relation  $R$ . Fix distinct  $x, y \in \mathcal{A}$  such that  $\neg(\{y\} R x)$ . Then,  $\text{tr}(R \cup (\{x\}, y))$  is irreflexive and proper.*

*Proof.* Let  $R^0 := R$ . For  $i > 0$ , let  $R^i$  be the extension of  $R^{i-1}$  obtained by imposing

$$\left( \bigcup_{j=1}^k X_j \cup \{y_{k+1}, \dots, y_n\} \right) \setminus \{y\} R^i y$$

whenever

$$\{y_1, \dots, y_n\} R^0 y \text{ and, for all } j \leq k, X_j R^{i-1} y_j.$$

Then, the transitive closure of  $R$  is  $\bigcup_{i=0}^{\infty} R^i$ . This is a standard result about the transitive closure.

The usual proof goes through with the version of transitivity used here.

Repeated applications of transitivity will not lead to a violation of irreflexivity, so we only need to check whether  $\text{tr}(R \cup (\{x\}, y))$  is proper. To keep track of repeated applications of transitivity, we introduce the notion of a tree. To simplify notation, we write  $z^k$  instead of  $(z_0, \dots, z_k)$  and  $\{z^k\}$  instead of  $\{z_0, \dots, z_k\}$ .

**Definition 52 (Q-tree).** *For a menu-item relation  $Q$ , a Q-tree from  $W \in (\mathcal{F}(\mathcal{A}) \cup \{\emptyset\})$  to  $w \in \mathcal{A}$  is inductively defined as follows:*

- *The level-0 node  $z_0 := w$  is mapped to a parent set  $Z_1(z_0)$  such that  $Z_1(z_0) Q z_0$ . A generic member of  $Z_1(z_0)$  is denoted  $z_1(z_0)$ .*
- *For  $k > 0$ : each level- $k$  node  $z_k(z^{k-1}) \notin W \cup \{z^{k-1}\}$  is mapped to a parent set  $Z_{k+1}(z^k)$  such that  $Z_{k+1}(z^k) Q z_k$ . A generic member of  $Z_{k+1}(z^k)$  is denoted  $z_{k+1}(z^k)$ .*
- *For some finite  $K > 0$ : each level- $K$  node  $z_K(z^{K-1})$  belongs to  $W \cup \{z^{K-1}\}$ .*

We refer to nodes that do not have parents as top nodes. A branch of a tree is a sequence  $(z_0, z_1(z_0), z_2(z^1), \dots, z_k(z^{k-1}))$  where  $z_k(z^{k-1})$  is a top node. We refer to  $(z_0, \dots, z_{i-1}(z^{i-2}))$  as descendants of  $z_i(z^{i-1})$ , and  $(z_{i+1}(z^i), \dots, z_k(z^{k-1}))$  as ancestors of  $z_i$ .

It is not difficult to see that  $(W, w)$  belongs to  $\text{tr}(Q)$  if and only if there is a  $Q$ -tree from  $W$  to  $w$ . Suppose that  $\text{tr}(R \cup (\{x\}, y))$  is improper, so there is a  $R \cup (\{x\}, y)$ -tree from  $\emptyset$  to  $w$  for some  $w \in \mathcal{A}$ . Notice that there must be at least one point in the tree in which  $x$  is the sole parent of  $y$ . Otherwise, there would be an  $R$ -tree from  $\emptyset$  to  $w$ , contradicting properness of  $R$ .

Construct a new tree by removing all the ancestors of  $y$  wherever  $x$  is the sole parent of  $y$ . The result is an  $R$ -tree. Let  $V$  be the set of items that descend from any instance of  $y$  that had  $x$  as its sole parent in the original tree. Fix any  $v \in V$ , and drop all the items from the  $R$ -tree that are not ancestors of  $v$ . The result is an  $R$ -tree from  $\{y\}$  to  $v$ , so it must be that  $\{y\} R v$ . Now return to the original tree. Take any point in the tree where  $x$  is the sole parent of  $y$ . Construct a new tree by removing everything except this instance of  $x$  and its ancestors. The result is an  $R$ -tree from a subset of  $V \cup \{y\}$  to  $x$ . To see why, recall that every top node in the original tree is a duplicate of one of its descendants. Fix any top node  $z_k$  of the new tree, and consider the branch of the original tree running through it:  $(w, \dots, y, x, \dots, z_k)$ . If  $z_k$  is not duplicated in  $(x, \dots, z_{k-1})$ , it must be duplicated in  $(w, \dots, y)$ —so it must belong to  $V \cup \{y\}$ . Since  $R$  is transitive, we have  $V' R x$  for some  $V' \subseteq V \cup \{y\}$ . Since  $\{y\} R v$  for all  $v \in V$ , applying transitivity once more gives  $\{y\} R x$ —contradiction. □

**Lemma 4.** *Fix an irreflexive, transitive, proper and  $(b, A)$ -consistent menu-item relation  $R$ . For any  $a \in A$ ,  $\text{tr}(R \cup (\{b\}, a))$  is consistent with  $b \succ A$ .*

*Proof.* Suppose that  $\text{tr}(R \cup (\{b\}, a))$  is inconsistent with  $b \succ A$ , so  $(A', b) \in \text{tr}(R \cup (\{b\}, a))$  for some  $A' \in \mathcal{F}(A)$ . Then, there must be an  $R \cup (\{b\}, a)$  tree from  $A'$  to  $b$ . Construct a new tree by removing all the ancestors of  $a$  wherever  $b$  is the sole parent of  $a$ . The result is an  $R$ -tree from a

subset of  $A' \cup \{a\}$  to  $b$ . Since  $R$  is transitive, we have  $A'' R b$  for some  $A'' \subseteq A' \cup \{a\} \subseteq A$ . This contradicts consistency of  $R$  with  $b \succ A$ .  $\square$

Let  $A = \{a_1, \dots, a_n\}$ . Let  $\triangleright^0 := \triangleright$ . For  $i \in \{1, \dots, n\}$ , let  $\triangleright^i = \text{tr}(\triangleright_{i-1} \cup (\{b\}, a_i))$ . Since  $\triangleright$  is irreflexive, proper, transitive, and consistent with  $b \succ A$ , we can use Lemmas 3 and 4 to show that the same is true of each  $\triangleright^i$ . Notice that  $\{b\} \triangleright^n a$  for all  $a \in A$ .

Now we use Lemma 3 to show that  $\triangleright^n$  can be extended to an irreflexive, proper and transitive relation  $\triangleright^+$  such that, for all distinct  $x, y \in \mathcal{A}$ ,  $\{x\} \triangleright^+ y$  or  $\{y\} \triangleright^+ x$ . The proof is similar to that of the Szpilrajn Extension Theorem. Consider the set of irreflexive, proper and transitive relations that extend  $\triangleright^n$ , ordered by set inclusion. Take any chain in the partially ordered set. The union of its elements is clearly irreflexive, proper and transitive, so it is an upper bound for the chain. By Zorn's Lemma, the partially ordered set must have a maximal element  $\triangleright^+$ . Suppose that, for some distinct  $x, y$ , neither  $\{x\} \triangleright^+ y$  nor  $\{y\} \triangleright^+ x$ . By Lemma 3,  $\triangleright^+$  can be extended to another irreflexive, proper and transitive relation containing  $(\{x\}, y)$ . Then  $\triangleright^+$  cannot be maximal, a contradiction. Moreover, for each  $X \in \mathcal{F}(\mathcal{A})$  and  $y \in \mathcal{A}$ ,  $\triangleright^+$  must satisfy

$$X \triangleright y \implies (\exists x \in X \text{ s.t. } \{x\} \triangleright^+ y). \quad (\text{A.1})$$

Suppose not. Then  $\{y\} \triangleright^+ x$  for all  $x \in X$ , as well as  $X \triangleright^+ y$ . Since  $\triangleright^+$  is transitive,  $\emptyset \triangleright^+ y$ . Since  $\triangleright^+$  is proper, this is a contradiction. Similarly, suppose that  $\{x\} \triangleright^+ y$  and  $\{y\} \triangleright^+ x$ . By transitivity,  $\emptyset \triangleright^+ x$ , a contradiction.

We can use  $\triangleright^+$  to define a strict preference  $\succ_m$ :

$$\{x\} \triangleright^+ y \iff x \succ_m y.$$

It is easy to see that  $\succ_m$  is antisymmetric, complete and transitive. It respects exclusion from below

because of (A.1), so it is indeed in  $\mathcal{M}$ . It also satisfies  $b \succ_m A$  because  $\triangleright^+$  extends  $\triangleright^n$ , and  $\{b\} \triangleright^n a$  for all  $a \in A$ .

### A.1.2 PROOF OF PROPOSITION 1

First, we show necessity of ISA. Fix  $B \in \mathcal{F}(\mathcal{A})$  and  $A \subseteq S(B)$ . For all  $\succ_m \in \mathcal{M}$ , we have  $a \not\prec_m A$ , so  $a \not\prec_m B$ . To confirm that  $c(B) = c(B \setminus A)$ , it suffices to show that  $\mathcal{M}(B) = \mathcal{M}(B \setminus A)$ . Take any  $b \in \mathcal{M}(B)$ . Since there exists  $\succ_m \in \mathcal{M}$  such that  $b \prec_m B$ , it cannot be that  $b \in A$ . Since  $b \prec_m B$  implies  $b \prec_m B \setminus A$ , we have  $b \in \mathcal{M}(B \setminus A)$ . Now suppose  $b \in \mathcal{M}(B \setminus A)$ , so  $b \prec_m B \setminus A$  for some  $\succ_m \in \mathcal{M}$ . Suppose that  $a \prec_m b \prec_m B \setminus A$  for some  $a \in A$ . For the  $\succ_m$ -best such  $a$ , we must have  $a \prec_m B$ —contradiction. Conclude that  $b \prec_m B$ , so  $b \in \mathcal{M}(B)$ .

Now we show sufficiency. The proof is similar to that of Theorem 1. Let  $D$  be the menu-item relation such that

$$\{x\} D y \iff x \succ_D y.$$

For any  $y \in \mathcal{A}$  and  $X \in \mathcal{F}(\mathcal{A})$ , let  $D(y, X)$  be the subset of items in  $X$  that are not dominated by  $y$ :

$$D(y, X) := X \setminus \{x \in X : y \succ_D x\}.$$

Define  $D$ -exclusion from below as follows:

$$X \triangleright_D y \iff (X \triangleright y \text{ and } X = D(y, X)).$$

We will take  $\mathcal{M}$  to be the set of  $D$ -monotone strict preferences that respect  $\triangleright_D$ . Formally, a  $D$ -monotone strict preference  $\succ_m$  belongs to  $\mathcal{M}$  if and only if

$$X \triangleright_D y \implies \exists x \in X \ x \succ_m y.$$

Fix  $A \in \mathcal{F}(\mathcal{A})$  and  $b \notin A$  such that  $b \succsim c(A \cup \{b\})$  and  $b \notin c(A \cup \{b\})$ . We show that  $b \notin \mathcal{M}(A \cup \{b\})$ .

**Lemma 5.** Fix  $X \in \mathcal{F}(\mathcal{A})$  and  $y \notin X$ . If  $y \succsim c(X \cup \{y\})$  and  $y \notin c(X \cup \{y\})$ , then

$$D(y, c(X \cup \{y\}) \cup \{x \in X : y \succ x\}) \triangleright_D y.$$

*Proof.* By Lemma 2,  $c(X \cup \{y\}) \cup \{x \in X : y \succ x\} \triangleright y$ . This implies  $y \notin c(\{y\} \cup c(X \cup \{y\}) \cup \{x \in X : y \succ x\})$ . By ISA, every item in  $c(X \cup \{y\}) \cup \{x \in X : y \succ x\}$  that is dominated by  $y$  can be removed without changing choice. We get

$$y \notin c(\{y\} \cup D(y, c(X \cup \{y\}) \cup \{x \in X : y \succ x\})).$$

This implies

$$D(y, c(X \cup \{y\}) \cup \{x \in X : y \succ x\}) \triangleright y,$$

so

$$D(y, c(X \cup \{y\}) \cup \{x \in X : y \succ x\}) \triangleright_D y,$$

□

Lemma 5 implies that  $D(b, c(A \cup \{b\}) \cup \{a \in A : b \succ a\}) \triangleright_D b$ . Since all the preferences in  $\mathcal{M}$  respect  $D$ -exclusion from below, we have  $b \not\succeq_m D(b, c(A \cup \{b\}) \cup \{a \in A : b \succ a\})$  for all  $\succ_m \in \mathcal{M}$ . This implies  $b \not\succeq_m A$  for all  $\succ_m \in \mathcal{M}$ , so  $b \notin \mathcal{M}(A \cup \{b\})$ .

Now fix  $A \in \mathcal{F}(\mathcal{A})$  and  $b \notin A$  such that  $b \in c(A \cup \{b\})$ . We show that  $b \in \mathcal{M}(A \cup \{b\})$ : there is a strict preference  $\succ_m$  that respects  $D$ -exclusion from below and has  $b \succ_m A$ .

We will construct an appropriate  $\succ_m$  by extending  $\triangleright_D$ . First, we define two useful properties of menu-item relations.

**Definition 53** (*D-transitivity*). *A menu-item relation  $R$  is  $D$ -transitive if*

$$(X R x, Y R y \text{ and } x \in Y) \implies D(y, X \cup Y \setminus \{x, y\}) R y.$$

For any menu-item relation  $R$ , let  $D\text{-tr}(R)$  denote the  $D$ -transitive closure of  $R$ .

**Definition 54** (*D-monotonicity*). *A menu-item relation  $R$  is  $D$ -monotone if  $R$  extends  $D$  and*

$$X R y \implies X = D(y, X).$$

**Lemma 6.**  *$D\text{-tr}(\triangleright_D \cup D)$  is irreflexive,  $D$ -monotone and proper, and consistent with  $b \succ A$ .*

*Proof.* Since  $\triangleright_D \cup D$  is irreflexive and  $D$ -monotone, and application of  $D$ -transitivity preserves irreflexivity and  $D$ -monotonicity, it is clear that  $D\text{-tr}(\triangleright_D \cup D)$  is irreflexive and  $D$ -monotone. As in the proof of Lemma 3, we check properness via a tree. The relevant construction is very similar to Definition 52. The only difference is as follows. In Definition 52, each top node is in  $W$  or identical to one of its descendants. Now, each top node is in  $W$  or identical to *or dominated by* one of its descendants.

Suppose there is a  $(\triangleright_D \cup D)$ -tree from  $\emptyset$  to  $w$ . Consider a menu  $Z$  that consists of all the items in the tree. Take any  $z \in c(Z)$ . By assumption, there exists  $Z' \subset Z$  such that  $Z' \triangleright_D z$  or  $Z' D z$ . By ISA,  $c(Z) = c(Z \setminus \{z\})$ , which contradicts the assumption that  $z \in c(Z)$ .

Suppose that there is a  $(\triangleright_D \cup D)$ -tree from  $A' \in \mathcal{F}(A)$  to  $b$ . Consider a menu  $Z$  that consists of all the items in the tree as well as  $A \setminus A'$ . By assumption, for each  $z \in Z \setminus A$ , there exists  $Z' \subset Z$  such that  $Z' \triangleright_D z$  or  $Z' D z$ . By ISA,  $b \notin c(Z)$ . Also by ISA,  $c(Z) = c(A \cup \{b\})$ . This contradicts the assumption that  $b \in c(A \cup \{b\})$ .  $\square$

**Lemma 7.** *Fix an irreflexive, proper,  $D$ -transitive and  $D$ -monotone menu-item relation  $R$ . Fix distinct  $x, y \in A$  such that  $\neg(\{y\} R x)$ . The  $D$ -transitive closure of  $R \cup (\{x\}, y)$  is irreflexive,  $D$ -*

*monotone and proper.*

*Proof.* The proof is very similar to that of Lemma 3, but using the modified notion of a tree from the proof of Lemma 6. A pair  $(W, w)$  belongs to  $D\text{-tr}(R \cup (\{x\}, y))$  if and only if there is an  $R \cup (\{x\}, y)$ -tree from  $W$  to  $w$ . Repeated application of  $D$ -transitivity will not cause a violation of irreflexivity or  $D$ -monotonicity, so we only need to check properness.

Suppose there is an  $R \cup (\{x\}, y)$ -tree from  $\emptyset$  to  $w$ . Construct a new tree by removing all the ancestors of  $y$  wherever  $x$  is the sole parent of  $y$ . Let  $V$  be the set of items that descend from any instance of  $y$  that had  $x$  as its sole parent in the original tree. Just as in the proof of Lemma 3, we have  $\{y\} R v$  for each  $v \in V$ . Let  $W$  be the set of items in the original tree that are dominated by a member of  $\{y\} \cup V$ . Since  $R$  is  $D$ -monotone and  $D$ -transitive,  $\{y\} R w$  for each  $w \in W$ .

Now return to the original tree. Take any point in the tree where  $x$  is the sole parent of  $y$ . Construct a new tree by removing everything except this instance of  $x$  and its ancestors. The result is an  $R$ -tree from a subset of  $W \cup V \cup \{y\}$  to  $x$ . To see why, notice that every top node in the original tree is a duplicate of one of its descendants or dominated by one of its descendants. Fix any top node  $z_k$  of the new tree, and consider the branch of the original tree running through it:  $(w, \dots, y, x, \dots, z_k)$ . If  $z_k$  is not duplicated in, or dominated by something in,  $(x, \dots, z_{k-1})$ , then it must be duplicated in, or dominated by something in,  $(w, \dots, y)$ . Since each item in  $(w, \dots, y)$  belongs to  $V \cup \{y\}$ ,  $z_k$  belongs to  $W \cup V \cup \{y\}$ . Since  $R$  is  $D$ -transitive, we have  $X R x$  for some  $X \subseteq W \cup V \cup \{y\}$ . Since  $\{y\} R v$  for all  $v \in V$ , and  $\{y\} R w$  for all  $w \in W$ , applying  $D$ -transitivity once more gives  $\{y\} R x$ —contradiction.  $\square$

**Lemma 8.** *Fix an irreflexive, proper,  $D$ -transitive,  $D$ -monotone menu-item relation  $R$  consistent with  $b \succ A$ . For any  $a \in A$ ,  $D\text{-tr}(R \cup (\{b\}, a))$  is consistent with  $b \succ A$ .*

*Proof.* The proof is exactly the same as that of Lemma 4, but using  $D$ -transitivity in place of transitivity and using the modified notion of a tree from the proof of Lemma 6.  $\square$

We can define  $\triangleright_D^i$  for  $i \in \{0, \dots, n\}$  exactly as in the proof of Theorem 1, but using  $D\text{-tr}(\triangleright_D \cup D)$  in place of  $\triangleright$ , and  $D\text{-tr}$  instead of  $tr$ . For each  $i$ ,  $\triangleright_D^i$  will be irreflexive, proper,  $D$ -transitive,  $D$ -monotone, and consistent with  $b \succ A$ . We will have  $\{b\} \triangleright_D^n a$  for all  $a \in A$ .

As in the proof of Theorem 1, we can use Lemma 7 to extend  $\triangleright_D^n$  to an irreflexive, proper,  $D$ -transitive and  $D$ -monotone relation  $\triangleright_D^+$  such that, for all distinct  $x, y \in \mathcal{A}$ ,  $\{x\} \triangleright_D^+ y$  or  $\{y\} \triangleright_D^+ x$ . We will have

$$X \triangleright_D y \implies (\exists x \in X \text{ s.t. } \{x\} \triangleright_D^+ y). \quad (\text{A.2})$$

We can use  $\triangleright_D^+$  to define a strict preference  $\succ_m$ :

$$\{x\} \triangleright_D^+ y \iff x \succ_m y.$$

It is easy to see that  $\succ_m$  is antisymmetric, complete and transitive. It extends  $\succ_D$  because  $\triangleright_D^+$  extends  $D$ . It respects exclusion from below because of (A.2). Finally, it satisfies  $b \succ_m A$  because  $\triangleright_D^+$  extends  $\triangleright^n$ , and  $\{b\} \triangleright^n a$  for all  $a \in A$ .

### A.1.3 PROOF OF THEOREM 2

**Lemma 9.** *C-IUA implies IUA.*

*Proof.* We first show that  $B \subset \bar{W}(B)$  for all  $B \in \mathcal{F}(\mathcal{A})$ . Take any  $b \in B$ . By Improvability, we can find a sequence  $B_i \rightarrow b$  such that  $b \in W(B_i)$  for all  $i$ . Let  $\hat{B}_i := B_i \cup B \setminus \{b\}$ . We have  $\hat{B}_i \rightarrow B$ , and  $b \in W(\hat{B}_i)$  for all  $i$ , so  $b \in \bar{W}(B)$ .

Take any  $A \in \mathcal{F}(\mathcal{A})$ , and index the items in  $A$  from best (1) to worst ( $|A|$ ). Break ties arbitrarily, with one exception: everything in  $\{a \in A : a \sim c(A), a \notin c(A)\}$  must have a lower index than everything in  $c(A)$ .

If  $a_1 \in c(A)$ , then there is no  $a \in A$  such that  $a \succsim c(A)$  and  $a \notin c(A)$ , so IUA has no bite.



Suppose that  $a_1 \notin c(A)$ . We have  $a_1 \succsim A \setminus \{a_1\}$ , so  $a_1 \in W(A \setminus \{a_1\})$ , so  $a_1 \in W(A)$ . We showed above that  $W(A) \in \mathcal{F}(\bar{W}(A))$ , so C-IUA implies  $c(A) = c(A \setminus \{a_1\})$ . Now suppose  $a_2 \succsim c(A)$  but  $a_2 \notin c(A)$ , so  $a_2 \notin c(A \setminus \{a_1\})$ . We have  $a_2 \succsim A \setminus \{a_1, a_2\}$ , so  $a_2 \in W(A \setminus \{a_1, a_2\})$ , so  $a_2 \in W(A)$ . C-IUA implies  $c(A) = c(A \setminus \{a_2\}) = c(A \setminus \{a_1, a_2\})$ . Iterating the argument, we get  $c(A) = c(A \setminus \{a_i\})$  for all  $a_i$  such that  $a_i \succsim c(A)$  and  $a_i \notin c(A)$ . This is IUA.  $\square$

**Definition 55** (Respects exclusion).  $m : \mathcal{A} \rightarrow \mathbb{R}$  respects exclusion if, for all  $(A, b) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $b \in W(A)$ , there exists  $a \in A$  such that  $m(a) > m(b)$ .

Let

$$\mathcal{M} := \{m \in C(\mathcal{A}, \mathbb{R}) : m \text{ respects exclusion}\}.$$

Later, we will confirm that  $\mathcal{M}$  is nonempty.

For any  $A \in \mathcal{F}(\mathcal{A})$ , let

$$\mathcal{M}(A) := \bigcup_{m \in \mathcal{M}} \arg \max_{a \in A} m(a).$$

We show that, for any  $A \in \mathcal{F}(\mathcal{A})$ ,  $a \succsim c(A)$  and  $a \notin c(A)$  implies  $a \notin \mathcal{M}(A)$ . By Lemma 2,  $a \succsim c(A)$  and  $a \notin c(A)$  implies  $c(A) \cup \{a \in A : c(A) \succ a\} \triangleright a$ . By definition of  $W$ ,  $a \in W(c(A) \cup \{a \in A : c(A) \succ a\})$ , so  $a \in W(A)$ . For any  $m : \mathcal{A} \rightarrow \mathbb{R}$  that respects exclusion,  $a \notin \arg \max_{\tilde{a} \in A} m(\tilde{a})$ . This implies  $a \notin \mathcal{M}(A)$ .

The remainder of the proof establishes that, for any  $A \in \mathcal{A}$ ,  $c(A) \subseteq \mathcal{M}(A)$ . To this end, take any  $A, b \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $b \in c(A \cup \{b\})$ . We will show that there exists  $m^* \in \mathcal{M}$  such that  $b \in \arg \max_{a \in A} m^*(a)$ , so  $b \in \mathcal{M}(A)$ .

If there is no  $(X, y) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $y \in W(X)$ , then  $\mathcal{M} = C(\mathcal{A}, \mathbb{R})$ . This implies  $\mathcal{M}(X) = X$  for all  $X \in \mathcal{F}(\mathcal{A})$ . We can therefore assume non-triviality when constructing  $m^*$ .

**Definition 56** (Non-triviality). There exists  $(X, y) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $y \in W(X)$ .

Under non-triviality, we can apply Theorem 3.4 in Herden and Pallack (2002) (HP). The theorem uses the following definition, adapted to our notation.

**Definition 57** (HP-system). *Let  $R$  be a binary relation on  $\mathcal{F}(\mathcal{A})$ . A family  $\{E_i\}_{i=0}^\infty$  of open subsets of  $\mathcal{F}(\mathcal{A})$  is an HP-system for  $R$  if it satisfies the following conditions:*

1. *There exist  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $d(E) \subset E'$ .*
2. *For all  $E, E' \in \{E_i\}_{i=0}^\infty$ ,  $E \subseteq E'$  or  $E' \subseteq E$ .*
3. *For all  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $d(E) \subset E'$ , there exists some  $E'' \in \{E_i\}_{i=0}^\infty$  such that  $d(E) \subset E'' \subset d(E'') \subset E'$ .*
4. *For any  $X, Y \in \mathcal{F}(\mathcal{A})$  and any  $E \in \{E_i\}_{i=0}^\infty$ : if  $X \in E$  and  $X R Y$ , then  $Y \in E$ .*
5. *For each  $X, Y \in \mathcal{F}(\mathcal{A})$  such that  $X R Y$  and  $\neg(Y R X)$ , there exist  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $d(E) \subset E'$ ,  $Y \in E$  and  $X \notin E'$ .*

We will construct an HP system  $\{E_i\}_{i=0}^\infty$  for the relation  $R$  given by

$$X R Y \iff Y \subset W(X). \tag{A.3}$$

The following Lemma will be used repeatedly.

**Lemma 10.** *For any  $z \in \mathcal{A}$  and any  $X, Y \in \mathcal{F}(\mathcal{A})$ : if  $z \in W(Y)$  and  $Y \subset \bar{W}(X)$ , or if  $z \in \bar{W}(Y)$  and  $Y \subset W(X)$ , then  $z \in W(X)$ .*

*Proof.* Suppose  $z \in W(Y)$  and  $Y \subset \bar{W}(X)$ . By definition of  $W$ , we have  $Y' \subseteq Y$  such that  $z \succsim Y'$  and  $z \notin c(\{z\} \cup Y')$ . By definition of  $\bar{W}$ , we have  $X_i \rightarrow X$  and  $y_i \rightarrow y$  such that  $y_i \in W(X_i)$  for all  $i$ .  $y_i \in W(X_i)$  means there is some  $X_i(y) \subseteq X_i$  such that  $y_i \succsim X_i(y)$  and  $y_i \notin c(\{y_i\} \cup X_i(y))$ . Passing to a subsequence if necessary, let  $X(y) := \lim_{i \rightarrow \infty} X_i(y)$ . Let  $X' := \bigcup_{y \in Y'} X(y)$ . We have

$X_i(y) \cup X' \setminus X(y) \rightarrow X'$ , Since  $y_i \in W(X_i(y))$  for all  $i$ , we also have  $y_i \in W(X_i(y) \cup X' \setminus X(y))$  for all  $i$ . We conclude that  $Y' \subset \bar{W}(X')$ . By C-IUA,  $z \notin c(\{z\} \cup X')$ . Since  $y_i \succsim X_i(y)$  for all  $i$  and all  $y \in Y'$ , and since  $\succsim$  is continuous,  $y \succsim X(y)$  for all  $y \in Y'$ . Since  $Y'$  is finite, there must be some  $y \in Y'$  such that  $y \succsim X'$ . Since  $z \succsim Y'$ , we have  $z \succsim X'$ . We conclude that  $z \in W(X')$ , so  $z \in W(X)$ .

A parallel argument covers the second case:  $z \in \bar{W}(Y)$  and  $Y \subset W(X)$ . Suppose  $z \in \bar{W}(Y)$ . Just as above, we can find  $Y'$  such that  $z \in \bar{W}(Y')$  and  $z \succsim Y'$ . For each  $y \in Y'$ , we can find  $X(y)$  such that  $y \in W(X(y))$  and  $y \succsim X(y)$ . Letting  $X' := \bigcup_{y \in Y'} X(y)$ , we get  $Y' \subset W(X')$  and  $y \succsim X'$  for some  $y \in Y'$ . Since  $z \succsim Y'$ ,  $z \succsim X'$ . C-IUA gives  $z \notin c(\{z\} \cup X')$ . We conclude that  $z \in W(X')$ , so  $z \in W(X)$ .  $\square$

Now we construct  $\{E_i\}_{i=0}^\infty$ .

1. Let

$$E_0 := \mathcal{F}(W(\{b\} \cup A)).$$

Continuity says that  $W(\{b\} \cup A)$  is open, so  $\mathcal{F}(W(\{b\} \cup A))$  is open. By IUA,  $b \notin E_0$ .

Let

$$\bar{E}_0 := \mathcal{F}(\bar{W}(\{b\} \cup A)).$$

We show that  $\bar{E}_0$  is closed. Take  $\{X_i \in \mathcal{F}(\bar{W}(\{b\} \cup A))\}_{i=1}^\infty$  such that  $X_i \rightarrow X \in \mathcal{F}(A)$ .

For each  $i$ , there exist  $\{X_{ij} \in \mathcal{F}(A)\}_{j=1}^\infty \rightarrow X_i$  and  $\{Y_{ij} \in \mathcal{F}(A)\}_{j=1}^\infty \rightarrow \{b\} \cup A$  such that  $X_{ij} \in W(Y_{ij})$  for all  $j$ . Notice that  $X_{ii} \rightarrow X$  and  $Y_{ii} \rightarrow \{b\} \cup A$ . Since  $X_{ii} \in W(Y_{ii})$  for all  $i$ ,  $X \in \mathcal{F}(\bar{W}(\{b\} \cup A))$ .

By the definitions of  $W$  and  $\bar{W}$ ,  $E_0 \subset \bar{E}_0$ . Since  $\bar{E}_0$  is closed,  $\text{cl}(E_0) \subset \bar{E}_0$ . We showed in the proof of Lemma 9 that, for any  $X \in \mathcal{F}(A)$ ,  $X \subset \bar{W}(X)$ . Thus,  $\{b\} \cup A \subset \bar{E}_0$ . Since  $b \notin E_0$ ,  $b \in \bar{E}_0 \setminus E_0$ .

2. The inductive hypothesis is as follows. For each  $j \in \{0, \dots, i-1\}$ , suppose that we have already constructed  $E_j, \bar{E}_j \in \mathcal{F}(\mathcal{A})$  such that

$$E_j = \mathcal{F}(W(Z(E_j))) \quad (\text{A.4})$$

$$\bar{E}_j = \mathcal{F}(\bar{W}(Z(E_j))) \quad (\text{A.5})$$

for some  $Z(E_j) \in \mathcal{F}(\mathcal{A})$ . Suppose that, for each  $j$ ,  $E_j$  is open and  $\bar{E}_j$  is closed, and that  $\{z_j\} \in \bar{E}_j \setminus E_j$ . Finally, suppose that there is a permutation  $\pi$  of  $\{0, \dots, i-1\}$  such that

$$E_{\pi(0)} \subset \bar{E}_{\pi(0)} \subset E_{\pi(1)} \subset \bar{E}_{\pi(1)} \subset \dots \subset E_{\pi(n-1)} \subset \bar{E}_{\pi(n-1)}. \quad (\text{A.6})$$

- (a) Suppose that  $\{z_i\} \in \bar{E} \setminus E$  for some  $E \in \{E_1, \dots, E_{i-1}\}$ . Set  $Z(E_i) := Z(E)$ , and set

$$E_i = \mathcal{F}(W(Z(E_i))) = E \quad (\text{A.7})$$

$$\bar{E}_i = \mathcal{F}(\bar{W}(Z(E_i))) = \bar{E}. \quad (\text{A.8})$$

- (b) Suppose that  $\{z_i\} \in \cap_{j=0}^{i-1} E_j$ . Set  $Z(E_i) := \{z_i\}$ , and set

$$E_i := \mathcal{F}(W(Z(E_i)))$$

$$\bar{E}_i := \mathcal{F}(\bar{W}(Z(E_i))).$$

Since  $z_i \in \bar{W}(z_i) \setminus W(z_i)$ ,  $\{z_i\} \in \bar{E}_i \setminus E_i$ .

Suppose that  $E$  is the smallest member of  $\{E_j\}_{j=0}^{i-1}$ . We show that  $\bar{E}_i \subset E$ . Since  $\{z_i\} \in E$ ,  $z_i \in W(Z(E))$ . Fix any  $X \in \mathcal{F}(\bar{W}(z_i))$ . By Lemma 10,  $X \in \mathcal{F}(W(Z(E))) = E$ .

- (c) Suppose that  $\{z_i\} \notin \cup_{j=0}^{i-1} \bar{E}_j$ . Suppose that  $E$  is the largest member of  $\{E_k\}_{k=0}^{i-1}$ . By Improvability, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $Z(E)$  such that  $Z(E) \subset$

$W(Z)$ . If we choose  $Z$  sufficiently close to  $Z(E)$ , we will have  $z_i \notin W(Z)$ . Suppose not. Then we have a sequence  $\{Z_j \in \mathcal{F}(\mathcal{A})\}_{j=1}^{\infty} \rightarrow Z(E)$  such that  $z_i \in W(Z_j)$  for all  $j$ . This implies  $z_i \in \bar{W}(Z(E))$ —contradiction.

Let  $Z(E_i) := Z \cup \{z_i\}$ , and set  $E_i, \bar{E}_i$  according to (A.7). Since  $z_i \in \bar{W}(Z \cup \{z_i\})$ , and since  $z_i \notin W(Z)$ ,  $\{z_i\} \in \bar{E}_i \setminus E_i$ . We show that  $\bar{E} \subset E_i$ . Suppose  $X \in \bar{E}$ , so  $X \in \mathcal{F}(\bar{W}(Z(E)))$ . Since  $Z(E) \in \mathcal{F}(W(Z))$ , Lemma 10 implies  $X \in \mathcal{F}(W(Z))$ , so  $X \in \mathcal{F}(W(Z \cup \{z_i\})) = E_i$ .

- (d) Suppose that  $z_i \in E_j$  for some  $j \in \{0, \dots, i-1\}$  and that  $z_i \notin \bar{E}_{j'}$  for some  $j' \in \{0, \dots, i-1\}$ . Suppose that  $E$  is the largest element of  $\{E_j : \{z_i\} \notin \bar{E}_j\}_{j=0}^{i-1}$ , and that  $E'$  is the smallest element of  $\{E_j : \{z_i\} \in E_j\}_{j=0}^{i-1}$ . We must have  $\bar{E} \subset E'$ .

As in Step 2c, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $Z(E)$  such that  $Z(E) \subset W(Z)$  and  $z_i \notin W(Z)$ . Since  $Z(E) \in \bar{E} \subset E'$ , and since  $E'$  is open, we can ensure  $Z \in E'$  by choosing  $Z$  sufficiently close to  $Z(E)$ . Let  $Z(E_i) := Z \cup \{z_i\}$ , and set  $E_i, \bar{E}_i$  according to (A.7).

As in Step 2c,  $\{z_i\} \in \bar{E}_i \setminus E_i$ , and  $\bar{E} \subset E_i$ . We show that  $\bar{E}_i \subset E'$ . Suppose that  $X \in \bar{E}_i$ , so  $X \in \mathcal{F}(\bar{W}(Z \cup \{z_i\}))$ . Since  $\{z_i\} \cup Z' \in E'$ ,  $\{z_i\} \cup Z' \in \mathcal{F}(W(Z(E')))$ . By Lemma 10,  $X \in \mathcal{F}(W(Z(E'))) = E'$ .

Rather than proving that  $\{E_i\}_{i=0}^{\infty}$  is an HP-system for  $R$  given in (A.3), we will prove a stronger result. Define a new relation  $S$  by

$$X S Y \iff \nexists E, E' \in \{E_i\}_{i=0}^{\infty} \text{ s.t. } X \in E, \bar{E} \subset E', Y \notin E'.$$

Notice that  $b S A$  because  $\{b\} \in \bar{E}_0 \setminus E_0$  and  $A \in \bar{E}_0$ . We show that

$$X R Y \implies X S Y \text{ and } \neg(Y S X).$$

Suppose that  $X R Y$ , so  $Y \in \mathcal{F}(W(X))$ . Suppose that  $\neg(X S Y)$ , so there exist  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $X \in E, \bar{E} \subset E', Y \notin E'$ . We have  $X \in \mathcal{F}(W(Z(E))) = E$ . By Lemma 10,  $Y \in \mathcal{F}(W(Z(E))) = E$ , which contradicts  $Y \notin E' \supset E$ . Thus,  $X R Y$  implies  $X S Y$ .

Now we show that  $\neg(Y S X)$ . We need to find  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $Y \in E, \bar{E} \subset E'$ , and  $X \notin E'$ . By Improvability, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $Y$  such that  $Y \subset W(Z)$ . Since  $Y \subset W(X)$ , we ensure  $Z \subset W(X)$  by choosing  $Z$  sufficiently close to  $Y$ . Since  $Z \subset \mathcal{Z}$ , for each  $z \in Z$ , there is some  $E(z) \in \{E_i\}_{i=0}^\infty$  such that  $z \in \bar{E}(z) \setminus E(z)$ . Choose the  $z \in Z$  for which  $E(z)$  is largest; call it  $z^*$  and write  $E^*$  instead of  $E(z^*)$ . We have  $Z \subset \bar{E}^*$ . Suppose that  $X \in \bar{E}^* = \mathcal{F}(\bar{W}(Z(E^*)))$ . Since  $z^* \in W(X)$ , Lemma 10 implies  $z^* \in W(Z(E^*))$ , so  $\{z^*\} \in E^*$ . This contradicts the definition of  $E^*$ . Conclude that  $X \notin \bar{E}^*$ . Since  $Y \subset W(Z)$  and  $Z \subset \bar{E}^* = \mathcal{F}(\bar{W}(Z(E^*)))$ , Lemma 10 implies  $Y \in \mathcal{F}(W(Z(E^*))) = E^*$ . We can repeat the same arguments with  $Z(E^*)$  in place of  $Y$  to get  $Z \subset \mathcal{F}(\mathcal{Z})$  such that  $Z(E^*) \subset E(Z)$  and  $X \notin \bar{E}(Z)$ .  $Z(E^*) \subset E(Z)$  implies  $\bar{E}^* \subset E(Z)$ . Putting everything together, we have  $Y \in E^*, \bar{E}^* \subset E(Z)$ , and  $X \notin E(Z)$  as desired. Conclude that  $S$  extends  $R$ .

Now we show that  $\{E_i\}_{i=0}^\infty$  is indeed an HP-system for  $S$ . We check the requirements of Definition 57 in order.

1. This is implied by requirement 4. By non-triviality, there is some  $(X, Y) \in \mathcal{F}(\mathcal{A})^2$  such that  $Y \subset W(X)$ , i.e.  $X R Y$ . Notice that  $Y \subset W(X)$  implies  $\neg(X \subset W(Y))$ . Otherwise, we would have  $Y \subset W(Y)$  by Lemma 10. Applying IUA gives  $c(Y) = \emptyset$ —contradiction. Thus, we have  $(X, Y)$  such that  $(X R Y)$  and  $\neg(Y R X)$ . Under this condition, requirement 4 says that there exist  $E, E' \in \{E_i\}_{i=1}^\infty$  such that  $\text{cl}(E) \subset E'$ .
2. From (A.6), it is clear that  $E_i \subseteq E_j$  or  $E_i \supseteq E_j$  for all  $E_i, E_j \in \{E_i\}_{i=0}^\infty$ .
3. For any  $E_i, E_j$  such that  $\text{cl}(E_i) \subset E_j$ , we need to find  $E_k$  such that  $\text{cl}(E_i) \subset E_k \subset \text{cl}(E_k) \subset E_j$ . We first show there is  $E_k$  such that  $\bar{E}_i \subset E_k \subset \bar{E}_k \subset E_j$ . Since  $E_j$  is open and  $\bar{E}_i$  is closed,

there must be some  $z_k \in \mathcal{Z}$  such that  $k > \max_{i,j}$  and  $\{z_k\} \in E_j \setminus \bar{E}_i$ . We constructed the  $\{E_j\}_{j=0}^\infty$  so that  $\{z_k\} \notin \bar{E}_i$  implies  $\bar{E}_i \subset E_k$  and  $\{z_k\} \in E_j$  implies  $\bar{E}_k \subset E_j$ . To see why, notice that  $k > i$  implies  $\bar{E}_i$  was already present when  $E_k$  was defined. Since  $\{z_k\} \notin \bar{E}_i$ , step 2d requires  $\bar{E}_i \subset E_k$ . Similarly,  $k > j$  implies  $E_j$  was already present when  $\bar{E}_k$  was defined. Since  $\{z_k\} \in E_j$ , step 2d requires  $\bar{E}_k \subset E_j$ .

Since  $\text{cl}(E_i) \subseteq \bar{E}_i$ ,  $\text{cl}(E_i) \subset E_k$  follows from  $\bar{E}_i \subset E_k$ . Similarly,  $\text{cl}(E_k) \subset E_j$  follows from  $\bar{E}_k \subset E_j$ .

4. Take any  $(X, \gamma) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $X \mathcal{S} \{\gamma\}$ . Suppose that  $X \in E_i$  for some  $E_i \in \{E_i\}_{i=0}^\infty$ . We show that  $\{\gamma\} \in E_i$ . By definition of  $\mathcal{S}$ , there cannot be  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $X \in E, \bar{E} \subset E'$ , and  $\{\gamma\} \notin E'$ . It suffices to show that, for all  $E \in \{E_i\}_{i=0}^\infty, X \in E$  implies that there exists  $E' \in \{E_i\}_{i=0}^\infty$  such that  $X \in E' \subset \bar{E}' \subset E$ .

The argument is very similar to the one we used to show that  $\mathcal{S}$  extends  $\mathcal{R}$ .  $X \in E$  means  $X \in \mathcal{F}(W(Z(E)))$ . By Improvability, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $X$  such that  $X \in W(Z)$ . By choosing  $Z$  sufficiently close to  $X$ , we will have  $Z \in \mathcal{F}(W(Z(E))) = E$ . Since  $Z \subset \mathcal{Z}$ , for each  $z \in Z$ , there is some  $E(z) \in \{E_i\}_{i=0}^\infty$  such that  $z \in \bar{E}(z) \setminus E(z)$ . Choose the  $z \in Z$  for which  $E(z)$  is largest; call it  $z^*$  and write  $E^*$  instead of  $E(z^*)$ . We have  $Z \subset \bar{E}^*$ . Suppose that  $\bar{E} \subset \bar{E}^*$ , so  $Z(E) \subset \bar{W}(Z(E^*))$ . Since  $z^* \in Z \subset W(Z(E))$ , Lemma 10 implies  $z^* \in W(Z(E^*))$ , so  $\{z^*\} \in E^*$ —contradiction. Conclude that  $\bar{E}$  is a strict superset of  $\bar{E}^*$ , so  $\bar{E}^* \subset E$ . Since  $X \subset W(Z)$  and  $Z \subset \bar{E}^* = \mathcal{F}(\bar{W}(Z(E^*)))$ , Lemma 10 implies  $X \in \mathcal{F}(W(Z(E^*))) = E^*$ . Putting everything together,  $X \in E^* \subset \bar{E}^* \subset E$  as desired.

5. Take any  $(X, Y) \in \mathcal{F}(\mathcal{A})^2$  such that  $X \mathcal{S} Y$  and  $\neg(Y \mathcal{S} X)$ .  $\neg(Y \mathcal{S} X)$  implies that there exist  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $Y \in E, \bar{E} \subset E'$ , and  $X \notin E'$ .

By HP Theorem 3.4 (and the fact that  $\mathcal{A}$  is a separable metric space), there exists a continuous

$f: \mathcal{F}(\mathcal{A}) \rightarrow \mathbb{R}$  such that, for any  $(X, Y) \in \mathcal{F}(\mathcal{A})^2$ ,

$$X S Y \implies f(X) \geq f(Y)$$

$$X S Y \text{ and } \neg(Y S X) \implies f(X) > f(Y).$$

Since  $Y \in \mathcal{F}(W(X))$  implies  $X S Y$  and  $\neg(Y S X)$ ,  $Y \in \mathcal{F}(W(X))$  implies  $f(X) > f(Y)$ . Since  $b S A$ ,  $f(\{b\}) \geq f(A)$ .

Take any  $X \in \mathcal{F}(\mathcal{A})$ . It is easy to see that  $X S \{x\}$  for all  $x \in X$ . For some  $x \in X$ , we will also have  $\{x\} S X$ . To see why, notice that there must be some  $x \in X$ , denoted  $x^*$ , such that  $\{x^*\} \in E$  implies  $X \in E$  for all  $E \in \{E_i\}_{i=0}^\infty$ . Suppose that  $\neg(\{x^*\} S X)$ , so there exist  $E, E' \in \{E_i\}_{i=0}^\infty$  such that  $X \notin E'$ ,  $\bar{E} \subset E'$ , and  $\{x^*\} \in E$ . By definition of  $x^*$ ,  $\{x^*\} \in E$  implies  $X \in E$ , so  $X \in E'$ . We conclude that  $\{x^*\} S X$ . We must have  $f(\{x^*\}) \geq f(X) \geq \max_{x \in X} f(\{x\})$ . This implies  $f(X) = \max_{x \in X} f(\{x\})$ . Define  $m^* : \mathcal{A} \rightarrow \mathbb{R}$  by

$$m^*(\cdot) := f(\{\cdot\}).$$

Since  $f(\{b\}) \geq f(A) = \max_{a \in A} f(\{a\})$ , we have  $m^*(b) \geq \max_{a \in A} m^*(a)$ . Since, for any  $(X, y) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$ ,  $y \in W(X)$  implies  $\max_{x \in X} f(\{x\}) = f(X) > f(y) = \max_{y \in Y} f(\{y\})$ ,  $m^*$  respects exclusion. Since  $m^*$  inherits continuity from  $f$ ,  $m^* \in \mathcal{M}$  as desired. This completes the construction of  $m^*$ .

Since  $(A, b)$  was chosen arbitrarily from  $\{\mathcal{F}(\mathcal{A}) \times \mathcal{A} : b \in c(\{b\} \cup A)\}$ , we conclude that  $c(A) \subseteq \mathcal{M}(A)$ .

Since  $\succsim$  is continuous, there is some  $u \in C(A, \mathbb{R})$  such that  $u$  represents  $\succsim$ . Since  $c(A) \subseteq \mathcal{M}(A)$  and since

$$a \succsim c(A) \text{ and } a \notin c(A) \implies a \notin \mathcal{M}(A),$$



$(u, \mathcal{M})$  represents  $(\succsim, c)$ . We show that  $(u, \mathcal{M})$  satisfies local non-satiation, recoverability, and closedness. For local non-satiation: take any  $a \in \mathcal{A}$ . By Improvability, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $a$  such that  $a \in W(Z)$ . Since the preferences in  $\mathcal{M}$  respect exclusion,  $Z$  is strictly preferred to  $a$  by  $\mathcal{M}$ .

For recoverability: take any  $(B, a) \in \mathcal{F}(\mathcal{A}) \times \mathcal{A}$  such that  $\max_{b \in B} m(b) > m(a)$  for all  $m \in \mathcal{M}$ . We show that  $\max_{b \in B: u(b) \leq u(a)} m(b) > m(a)$  for all  $m \in \mathcal{M}$ . It suffices to show that

$$\exists m \in \mathcal{M} \text{ s.t. } m(a) \geq \max_{b \in B: u(b) \leq u(a)} m(b) \implies \exists m \in \mathcal{M} \text{ s.t. } m(a) \geq \max_{b \in B} m(b).$$

Let  $B' := \{b \in B : u(b) \leq u(a)\}$ , and suppose  $m(a) \geq \max_{b \in B'} m(b)$  for some  $m \in \mathcal{M}$ . Since  $(u, \mathcal{M})$  represents  $(\succsim, c)$ ,  $a \succsim B'$  and  $a \in c(\{a\} \cup B')$ , so  $a \notin W(B')$ . Since  $B \setminus B' \succ a$ ,  $a \notin W(B)$ . We can construct  $\hat{m} \in \mathcal{M}$  such that  $\hat{m}(a) \geq \max_{b \in B} \hat{m}(b)$  in exactly the same way that we constructed  $m^*$  above. (Notice that the construction does not require  $a \in c(\{a\} \cup B)$ , which might not hold; it only requires  $a \notin W(B)$ , which we have just shown.)

For closedness: since  $W(B)$  is open for all  $B \in \mathcal{F}(\mathcal{A})$ , it suffices to show that

$$\bigcap_{m \in \mathcal{M}} \{a \in \mathcal{A} : m(a) < \max_{b \in B} m(b)\} = W(B).$$

Since each  $m \in \mathcal{M}$  respects exclusion,  $a \in W(B)$  implies  $m(a) < \max_{b \in B} m(b)$ . Now suppose  $m(a) < \max_{b \in B} m(b)$  for all  $m \in \mathcal{M}$ . Since recoverability holds, we have  $m(a) < \max_{b \in B'} m(b)$  where  $B' := \{b \in B : u(a) \geq u(b)\}$  for all  $m \in \mathcal{M}$ . Since  $(u, \mathcal{M})$  represents  $(\succsim, c)$ ,  $a \succ B'$  and  $a \notin c(\{a\} \cup B')$ . By definition of  $W$ ,  $a \in W(B')$ , so  $a \in W(B)$ .

#### A.1.4 PROOF OF THEOREM 3

The following piece of notation is useful for both necessity and sufficiency. Recall the definition of  $B(p)$  in (1.1). Let

$$NB(p) := \{q \in \Delta(Z) : p \succsim q\} \setminus B(p) = \{q \in \Delta(Z) : p \succsim q \text{ and } p \in \mathcal{M}(\{p, q\})\}.$$

First, we show necessity. It is well known that the first parts of Continuity and Independence are necessary (and sufficient) for  $\succsim$  to have an EU representation. The second part of Continuity says that  $NB(p)$  is closed. Take any convergent sequence  $(q_n)$  such that each  $q_n \in NB(p)$ . By definition,  $q_n \succsim p$  and  $p \in c(\{q_n, p\})$  for all  $n$ . Since  $\succsim$  is continuous,  $q \succsim p$ . For each  $n$ , we have  $(m_n)'p \geq (m_n)'q$  for some  $m_n \in \mathcal{M}$ . Since  $\mathcal{M}$  is compact, some subsequence of  $m_n$  has a limit  $m \in \mathcal{M}$ . We will have  $m'p \geq m'q$ , so  $p \in c(\{q, p\})$ . We conclude that  $q \in NB(p)$ , so  $NB(p)$  is closed.

Now consider Convexity. Suppose that  $A$  excludes  $p$ . Let  $\underline{A} := c(A \cup \{p\}) \cup \{a \in A : p \succ a\}$ . By Lemma 2,  $p \notin c(\underline{A} \cup \{p\})$ . For each  $m \in \mathcal{M}$ , we have some  $a \in \underline{A}$  such that  $m'p < m'a$ . It is without loss to assume that  $m'p = 0$  for all  $m \in \mathcal{M}$ . We want to find a set of weights  $\alpha$  such that  $p \succ \sum_{a \in \underline{A}} \alpha(a) \delta_a$  and

$$\sum_{a \in \underline{A}} \alpha(a) m'a > 0$$

for all  $m \in \mathcal{M}$ . The first part is easy—it will hold for any  $\alpha$  since  $\succsim$  is EU—so we focus on the second. For each  $m \in \mathcal{M}$ , let  $m_{\underline{A}} := (m'a)_{a \in \underline{A}}$ . Let  $\mathcal{M}_{\underline{A}}$  be the set of the  $m_{\underline{A}}$ . Like  $\mathcal{M}$ ,  $\mathcal{M}_{\underline{A}}$  is nonempty, compact and convex. Let  $N := \mathbb{R}_{-}^{|\underline{A}|}$ , which is nonempty, closed and convex. Notice that no element of  $\mathcal{M}_{\underline{A}}$  can be weakly negative (otherwise, some  $m \in \mathcal{M}$  would rank  $p$  weakly higher than each member of  $\underline{A}$ ). Thus,  $\mathcal{M}_{\underline{A}}$  and  $N$  are disjoint, and we can apply the Separating Hyperplane Theorem. This delivers a nonzero  $\alpha \in \mathbb{R}^{|\underline{A}|}$  and  $c \in \mathbb{R}$  such that  $\alpha'n < c < \alpha'm_A$  for all  $n \in N, m_A \in \mathcal{M}_{\underline{A}}$ . Since the zero vector belongs to  $N$ , we must have  $c > 0$ . Suppose the  $i$ th

element of  $\alpha$  is strictly negative. By choosing  $n$  with a sufficiently negative number in  $i$ th position and zeros elsewhere, we get  $\alpha'n > c$ , a contradiction. Thus, each element of  $\alpha$  is weakly positive. If we rescale  $\alpha$  to a unit sum, we still have  $\alpha'm_A > 0$  for all  $m_A \in \mathcal{M}_A$ . We can rewrite this as  $\sum_{a \in A} \alpha(a)m'a > 0$  for all  $m \in \mathcal{M}$ , which is exactly what we needed.

For the other direction of Convexity, suppose  $p \in c(A \cup \{p\})$ . For some  $m \in \mathcal{M}$ , we have  $m'p \geq m'a$  for all  $a \in A$ . For this  $m$ , we clearly have  $m'p \geq m'a$  for all  $a \in \text{co}(A)$ . No  $a \in \text{co}(A)$  can belong to  $B(p)$ .

Consider the second part of Monotonicity. Suppose that  $p >_{\text{FOSD}} q$  and that  $A \supset \{p, q\}$ . First, we show that  $q \notin c(A)$ . Suppose otherwise, so  $m'q = \max_{a \in A} m'a$  for some  $m \in \mathcal{M}$ . Since each  $m \in \mathcal{M}$  is weakly FOSD-monotone, we have  $m'p = \max_{a \in A} m'a$ , so  $p \in \mathcal{M}(A)$ . Since  $\succsim$  is strictly FOSD-monotone and  $p >_{\text{FOSD}} q$ , it cannot be that  $q \in c(A)$ . To confirm that  $c(A) = c(A \setminus \{q\})$ , it now suffices to show that  $\mathcal{M}(A) \setminus \{q\} = \mathcal{M}(A \setminus \{q\})$ . Take any  $r \in \mathcal{M}(A \setminus \{q\})$ , so  $m'r = \max_{a \in A \setminus \{q\}} m'a$ . Since  $p \in A \setminus \{q\}$  and each  $m \in \mathcal{M}$  is weakly FOSD-monotone, we have  $m'r = \max_{a \in A} m'a$ , so  $r \in \mathcal{M}(A)$ . Now take any  $r \in \mathcal{M}(A) \setminus \{q\}$ , so  $m'r = \max_{a \in A} m'a \geq \max_{a \in A \setminus \{q\}} m'a$ . Since  $r \neq q$ , we have  $r \in \mathcal{M}(A \setminus \{q\})$ .

Now we show sufficiency.

**Lemma 11.**  $B(p)$  is a convex cone.

*Proof.* First, suppose that  $q \in B(p)$ , so  $p \succsim q$  and  $\{q\} = c(\{p, q\})$ . By  $\succsim$ -Independence,  $p \succsim \alpha p + (1 - \alpha)q$  for all  $\alpha \in (0, 1)$ . By  $c$ -Independence,  $\{\alpha p + (1 - \alpha)q\} = c(\{p, \alpha p + (1 - \alpha)q\})$ , so  $\alpha p + (1 - \alpha)q \in B(p)$ . Similarly, suppose that  $\alpha p + (1 - \alpha)q \in B(p)$ , so  $p \succsim \alpha p + (1 - \alpha)q$  and  $\{\alpha p + (1 - \alpha)q\} = c(\{\alpha p + (1 - \alpha)q, p\})$ . By  $\succsim$ -Independence,  $p \succsim q$ . By  $c$ -Independence,  $\{q\} = c(\{q, p\})$ , so  $q \in B(p)$ . Now suppose that  $q, r \in B(p)$ , so  $p \succsim q, r$  and  $\{q\} = c(\{p, q\})$ ,  $\{r\} = c(\{p, r\})$ . By

$\succsim$ -Independence,  $p \succsim \alpha q + (1 - \alpha)r$  for all  $\alpha \in (0, 1)$ . By  $\succsim$ - and  $c$ -Independence,

$$\begin{aligned} \alpha p + (1 - \alpha)r &\succsim \alpha q + (1 - \alpha)r \\ \{\alpha q + (1 - \alpha)r\} &= c(\{\alpha q + (1 - \alpha)r, \alpha p + (1 - \alpha)r\}), \end{aligned}$$

so  $\alpha q + (1 - \alpha)r$  excludes  $\alpha p + (1 - \alpha)r$  from below. Also by  $\succsim$ - and  $c$ -Independence,

$$\begin{aligned} p &\succsim \alpha p + (1 - \alpha)r \\ \{\alpha p + (1 - \alpha)r\} &= c(\{p, \alpha p + (1 - \alpha)r\}), \end{aligned}$$

so  $\alpha p + (1 - \alpha)r$  excludes  $p$  from below. By IUA, we can add  $p$  to a set containing  $\alpha p + (1 - \alpha)r$  without affecting choice, so

$$\{\alpha q + (1 - \alpha)r\} = c(\{\alpha p + (1 - \alpha)r, \alpha q + (1 - \alpha)r, p\}).$$

Also by IUA, we can remove  $\alpha p + (1 - \alpha)r$  from a set containing  $\alpha q + (1 - \alpha)r$  without affecting choice, so

$$\{\alpha q + (1 - \alpha)r\} = c(\{\alpha q + (1 - \alpha)r, p\})$$

so  $\alpha q + (1 - \alpha)r \in B(p)$ . □

**Lemma 12.** Fix  $p, p', q, q' \in \Delta(Z)$  such that  $q - p = q' - p'$ . If  $q \in B(p)$ , then  $q' \in B(p')$ .

*Proof.* We have

$$\frac{1}{2}p + \frac{1}{2}q' = \frac{1}{2}p' + \frac{1}{2}q.$$

Since  $q \in B(p)$ ,  $p \succsim q$  and  $\{q\} = c(\{p, q\})$ . By  $c$ -Independence,

$$\begin{aligned} \left\{ \frac{1}{2}q + \frac{1}{2}p' \right\} &= c \left( \left\{ \frac{1}{2}q + \frac{1}{2}p', \frac{1}{2}p + \frac{1}{2}p' \right\} \right) \\ \left\{ \frac{1}{2}p + \frac{1}{2}q' \right\} &= c \left( \left\{ \frac{1}{2}p + \frac{1}{2}q', \frac{1}{2}p + \frac{1}{2}p' \right\} \right) \\ \{q'\} &= c(\{p', q'\}). \end{aligned}$$

Similarly, by  $\succsim$ -Independence,

$$\begin{aligned} \frac{1}{2}p + \frac{1}{2}p' &\succsim \frac{1}{2}q + \frac{1}{2}p' \\ \frac{1}{2}p + \frac{1}{2}p' &\succsim \frac{1}{2}p + \frac{1}{2}q' \\ p' &\succsim q'. \end{aligned}$$

We conclude that  $q' \in B(p')$ . □

We are now ready to define  $\mathcal{M}$ . Take  $p$  in the interior of  $\Delta(Z)$ . Take any supporting hyperplane  $H$  of  $B(p)$  that passes through some boundary point  $b$  of  $B(p)$  with  $p \succ b$ . Since  $B(p)$  is a cone with vertex  $p$ ,  $H$  will also pass through  $p$ . Since  $B(p)$  is open in  $\{q \in \Delta(Z) : p \succsim q\}$ ,  $H$  cannot include any point in  $B(p)$ . Let  $\mathcal{H}$  be the set with generic member  $H$ . For each  $H$ , take a unit-norm  $m \in \mathbb{R}^{|Z|}$  such that  $m'b = 0$  for all  $b \in H$  (including  $p$ ) and  $m'q > 0$  for all  $q \in B(p)$ . Collect these  $m$ , and take the closed convex hull. This is  $\mathcal{M}$ .

We show that everything in  $\mathcal{M}$  is weakly  $D$ -monotone. It suffices to show that any preference with an indifference curve in  $\mathcal{H}$  that has  $B(p) \succ p$  is weakly  $D$ -monotone. (All the other preferences are combinations and/or limits of these, so will inherit weak  $D$ -monotonicity.) Suppose that some preference  $\succsim_b$  with an indifference curve in  $\mathcal{H}$  has  $r \succ_b q$  even though  $q >_{FOSD} r$ . By definition of  $\mathcal{H}$ , we have  $b \sim_b p$  for some boundary point  $b$  of  $B(p)$  such that  $p \succ b$ . It is without loss to assume

that  $b$  is interior, so there exists  $\lambda > 0$  such that  $b + \lambda(q - r) \in \Delta(Z)$ . Since  $p \succ b$ , we will have  $p \succ b + \lambda(q - r)$  for  $\lambda$  small enough. We also have  $b \succ_b b + \lambda(q - r)$  and  $b + \lambda(q - r) \succ_{FOSD} b$ . Since  $b$  is a boundary point of  $B(p)$ , we can find  $\tilde{b}$  arbitrarily close to  $b$  such that  $\tilde{b} \in B(p)$ . Since  $p \sim_b b$ ,  $b \succ_b b + \lambda(q - r)$ , we can ensure  $p \succ_b \tilde{b} + \lambda(q - r)$  by choosing  $\tilde{b}$  sufficiently close to  $b$ . This implies  $\tilde{b} + \lambda(q - r) \notin B(p)$ . Since  $p \succ b + \lambda(q - r)$ , we can also ensure  $p \succ \tilde{b} + \lambda(q - r)$  by choosing  $\tilde{b}$  sufficiently close to  $b$ . This implies  $\tilde{b} + \lambda(q - r) \in NB(p)$ . Now we can derive a contradiction. To simplify the notation, let  $b^* = \tilde{b} + \lambda(q - r)$ . Notice that  $b^* \succ_{FOSD} \tilde{b}$ , so  $c(\{p, b^*, \tilde{b}\}) = c(\{p, b^*\})$  by the second part of Monotonicity. Since  $b^* \notin NB(p)$ , we have  $p \in c(\{p, b^*\})$ , so  $p \in c(\{p, b^*, \tilde{b}\})$ . But since  $\tilde{b} \in B(p)$ , IUA implies  $p \notin c(\{p, b^*, \tilde{b}\})$ . Conclude that  $\succ_b$  is weakly FOSSD-monotone, so every preference in  $\mathcal{M}$  is weakly  $D$ -monotone.

We now show

$$B(p) = \bigcap_{m \in \mathcal{M}} \{q \in \Delta(Z) : m'q > m'p\} \cap \{q \in \Delta(Z) : p \succsim q\}. \quad (\text{A.9})$$

Suppose  $b \notin B(p)$ , but  $b \in \bigcap_{m \in \mathcal{M}} \{q \in \Delta(Z) : m'q > m'p\}$  and  $p \succsim b$ . By construction of  $\mathcal{M}$ ,  $b$  must be a boundary point of  $B(p)$  that is not a limit of points in  $NB(p)$  that are strictly worse than  $p$ . Clearly,  $p \sim b$ . Moreover, for any  $q \prec p$ , there must be some  $\alpha$  sufficiently close to 1 such that

$$\alpha b + (1 - \alpha)q \in B(p).$$

(If this were not the case, then  $b$  could be written as a limit of points in  $NB(p)$  that are strictly worse than  $p$ .) Take  $r$  such that  $p \succ_{FOSSD} r$ . ( $p$  is interior, so some such  $r$  must exist.) Since the primary preference is strictly FOSSD-monotone, we have  $p \succ r$ . We also have  $\alpha \in (0, 1)$  such that

$$\alpha b + (1 - \alpha)r \in B(p).$$

That is,  $p \notin c(\{p, \alpha b + (1 - \alpha)r\})$  even though  $p \succsim \alpha b + (1 - \alpha)r$ . By the second part of Convexity,  $p \notin c(\{p, b, r\})$ . By the second part of Monotonicity and  $p \succ_{FOSD} r$ ,  $c(\{p, b, r\}) = c(\{p, b\})$ , so  $p \notin c(\{p, b\})$ . Since  $p \sim b$ , this implies  $b \in B(p)$ —contradiction.

Now suppose  $B(p) \not\subseteq \bigcap_{m \in \mathcal{M}} \{q \in \Delta(Z) : m'q > m'p\} \cap \{q \in \Delta(Z) : p \succsim q\}$ . By construction of  $\mathcal{M}$ , this can only happen if  $b \sim p$ , and there is a sequence  $\{H\}_{n=1}^{\infty}$  of hyperplanes in  $\mathcal{H}$  converging to  $\{q \in \Delta(Z) : p \sim q\}$ . Recall that each hyperplane in  $\mathcal{H}$  passes through some boundary point of  $B(p)$  that is strictly worse than  $p$ . Take the sequence of such points corresponding to  $\{H\}_{n=1}^{\infty}$ . Passing to a subsequence if necessary, let  $\tilde{b}$  be the limit of this sequence of points. For any sufficiently small perturbation  $\tilde{b}$  of  $b$  with  $\tilde{b} \prec p$ , we must have  $\tilde{b} \in B(p)$ . (Otherwise,  $\{q \in \Delta(Z) : p \sim q\}$  could not be the limit of  $\{H\}_{n=1}^{\infty}$ .) We can now apply the argument in the previous paragraph. Take any  $r$  such that  $p \succ_{FOSD} r$ . We must have  $\alpha b + (1 - \alpha)r \in B(p)$  for some  $\alpha \in (0, 1)$ . Applying the second part of Convexity,  $p \notin c(\{p, b, r\})$ . Applying the second part of Monotonicity,  $c(\{p, b, r\}) = c(\{p, b\})$ , so  $p \notin c(\{p, b\})$ , so  $b \notin NB(p)$ . Since  $b$  is a limit point of  $NB(p)$ , this contradicts the second part of Continuity.

Notice that it does not matter which  $p$  we use to define  $\mathcal{M}$ , since Lemma 12 ensures that we will get the same set of utilities (up to an irrelevant additive constant) for any interior  $p$ . To finish the proof, we have to show that  $\mathcal{M}$  satisfies two conditions. First, no utility in  $\mathcal{M}$  would justify choosing an item that the DM doesn't choose, but likes as much as anything he does choose. Second, for any item the DM chooses, some utility in  $\mathcal{M}$  justifies it.

Consider the first part. Suppose  $q \succsim c(A \cup \{q\})$  and  $q \notin c(A \cup \{q\})$ . Let  $\underline{A} = c(A \cup \{q\}) \cup \{a \in A : q \succ a\}$ . By Lemma 2,  $q \notin c(\underline{A} \cup \{q\})$ . By Convexity, we can find  $a^* \in \text{co}(\underline{A})$  such that  $q \notin c(\{q, a^*\})$ . Since  $q \succsim a^*$ ,  $a^* \in B(q)$ . By (A.9),  $m'a^* > m'q$  for all  $m \in \mathcal{M}$ . For each  $m \in \mathcal{M}$ , we must have  $a \in \underline{A}$  such that  $m'a > m'q$ . This is exactly what we needed.

For the second part, suppose  $q \in c(A \cup \{q\})$ . To start, suppose  $q \succsim A$ . Suppose that we cannot find  $m \in \mathcal{M}$  so that  $m'q \geq m'a$  for all  $a \in A$ . Recall the argument we used to show necessity of

Convexity. Since  $\mathcal{M}$  is compact and convex, we can use the same argument to find an  $a^* \in \text{co}(A)$  such that

$$a^* \in \bigcap_{m \in \mathcal{M}} \{r \in \Delta(Z) : m'r > m'q\} \cap \{r \in \Delta(Z) : q \succsim r.\}$$

By (A.9),  $a^* \in \text{co}(A) \cap B(q)$ . By Convexity,  $q \notin c(A \cup \{q\})$ , a contradiction.

Now we relax the assumption that  $q \succsim A$ . Let  $\underline{A} = \{a \in A : q \succsim a\}$ . We know  $q \in c(\underline{A} \cup \{q\})$ . (Suppose not. Then  $q \notin c(A \cup \{q\})$  by IUA, a contradiction.) By the previous argument, we can find  $m^* \in \mathcal{M}$  such that  $(m^*)'q \geq \max_{a \in \underline{A}} (m^*)'a$ . Now take any item  $\bar{a} \in A \setminus \underline{A}$ . Since  $q \in c(\{q\} \cup A)$  and  $\bar{a} \succ q$ ,  $\bar{a} \notin c(\{q\} \cup A)$ . By Lemma 2,  $\bar{a} \notin c(\{\bar{a}, q\} \cup \underline{A})$  even though  $\bar{a} \succ \{q\} \cup \underline{A}$ . We already showed that there is no  $m \in \mathcal{M}$  such that  $m'\bar{a} \geq \max_{x \in \underline{A} \cup \{q\}} m'x$ . In particular, it cannot be that  $(m^*)'\bar{a} \geq (m^*)'q \geq \max_{a \in \underline{A}} (m^*)'a$ , so  $(m^*)'q \geq \max_{a \in \underline{A} \cup \bar{a}} (m^*)'a$ . Since  $\bar{a}$  was an arbitrary selection from  $A \setminus \underline{A}$ , we have  $(m^*)'q \geq \max_{a \in A} m'a$ .

#### A.1.5 PROOF OF COROLLARY I

For the minimal set, recall the construction in the proof of Theorem 3. We start with  $p$  in the interior of  $\Delta(Z)$ . Then we take the supporting hyperplanes of  $B(p)$  that pass through some boundary point of  $B(p)$  that is strictly worse than  $p$ .  $\mathcal{H}$  is the set of such hyperplanes. Let  $\bar{c}(\mathcal{H})$  be the closed convex hull of  $\mathcal{H}$ . Let  $\mathcal{M}_\succ$  denote the set of preferences with representations in  $\mathcal{M}$ .  $\mathcal{M}_\succ$  is precisely the set of EU preferences that have indifference curves in  $\bar{c}(\mathcal{H})$  and that strictly prefer  $B(p)$  to  $p$ .

We show that  $\mathcal{M}_\succ$  is minimal. Take any boundary point  $b$  of  $B(p)$  such that  $p \succ b$ . Since  $B(p)$  is open in  $\{q \in \Delta(Z) : p \succsim q\}$ ,  $b \notin B(p)$ , so  $p \in c(\{p, b\})$ . Clearly, every set of justifiable preferences in every representation must contain some preference  $\succsim_m$  such that  $p \succsim_m b$ . Take any closed, convex proper subset of  $\mathcal{M}_\succ$ , and call it  $\mathcal{N}_\succ$ . Recall that  $\mathcal{M}_\succ$  is the closed, convex hull of the set of EU preferences that are indifferent between  $p$  and a boundary point of  $B(p)$  that is strictly worse



than  $p$ . Suppose that, for each boundary point  $b$  of  $B(p)$  such that  $b \succ p$ , there exists  $\succsim_n \in \mathcal{N}_\succ$  such that  $p \sim_n b$ . Since  $\mathcal{N}_\succ$  is convex and closed, this implies  $\mathcal{N}_\succ \supseteq \mathcal{M}_\succ$ —contradiction.

Now we turn to the maximal set. We will modify the construction in the proof of Theorem 3. As before, take  $p$  in the interior of  $\Delta(Z)$ . Let  $\mathcal{M}_\succ^{max}$  be the set of weakly  $D$ -monotone EU preferences that strictly prefer  $B(p)$  to  $p$ . If  $\mathcal{M}_\succ^{max}$  is indeed a representation, it is obviously maximal. It is easy to see that  $\mathcal{M}_\succ$  is convex. Suppose that it is not closed, so there exists a point  $b \in B(p)$  such that  $p \succsim_m b$  for some limit of  $D$ -monotone preferences that strictly prefer  $B(p)$  to  $p$ . The argument will be familiar from the proof of Theorem 3. Notice that  $b$  must be on the boundary of  $B(p)$ . Also, since  $b \notin NB(p)$  and  $NB(p)$  is closed,  $b$  cannot be a limit of points in  $NB(p)$ . In this situation, only one preference with  $p \succsim_m b$  could possibly be in  $\mathcal{M}_\succ^{max}$ : the preference exactly opposite  $\succ$ . But this preference is not weakly  $D$ -monotone, so it is not in  $\mathcal{M}_\succ^{max}$ .

We can define a compact, convex set of utility representations for the maximal set of EU preferences just as we did for the minimal set. Call the result  $\mathcal{M}^{max}$ . To confirm that (A.9) holds, suppose that  $b \notin B(p)$  but  $p \succsim b$  and  $m'b \succ m'p$  for all  $m \in \mathcal{M}^{max}$ . That is, every  $D$ -monotone EU preference that strictly prefers  $B(p)$  to  $p$  also prefers  $b$  to  $p$ . We already know that the preferences in the minimal set are  $D$ -monotone and prefer  $B(p)$  to  $p$ , so they must prefer  $b$  to  $p$  as well. But since the minimal set of EU preferences represents  $(\succsim, c)$ , we must have  $\{b\} = c(\{p, b\})$ , which contradicts  $b \notin B(p)$ . The rest of the proof of Theorem 3 relies only on (A.9), so it goes through as before.

#### A.1.6 PROOF OF COROLLARY 2

Let  $\mathcal{M}_\succ^1$  be the maximal set of EU preferences corresponding to  $(\succsim, c_1)$ . Recall that  $\mathcal{M}_\succ^2$  is the set of  $D$ -monotone EU preferences that strictly prefer  $B^2(p)$  to  $p$ . If  $B_1(p) \supset B_2(p)$ , then  $B^1(p) \succ_m p$  implies  $B^2(p) \succ_m p$ . Thus,  $\mathcal{M}_\succ^2$  includes all the  $D$ -monotone preferences that strictly prefer  $B^2(p)$  to  $p$ , so  $\mathcal{M}_\succ^2 \supseteq \mathcal{M}_\succ^1$ . To confirm that the inclusion is strict, take any point  $b \in B_1(p)$  such that  $b \notin B_2(p)$ . Since  $b \in B_1(p)$ , we have  $p \notin c_1(\{p, b\})$ , so there does not exist  $\succsim_m \in \mathcal{M}_\succ^1$  such that

$p \succsim_m b$ . Since  $b \notin B_2(p)$ , we have  $p \in c_2(\{p, b\})$ , so there exists  $\succsim_m \in \mathcal{M}_\succ^2$  such that  $p \succsim_m b$ .

For the converse, suppose  $\mathcal{M}_\succ^1 \subset \mathcal{M}_\succ^2$ . Suppose that there exists  $b \in B^2(p)$  such that  $b \notin B^1(p)$ . We must have  $p \notin c_2(\{p, b\})$ , so we must not have any  $\succsim_m \in \mathcal{M}_\succ^2$  such that  $p \succsim_m b$ . On the other hand, we must have  $p \in c_1(\{p, b\})$ , so we must have  $\succsim_m \in \mathcal{M}_\succ^1$  such that  $p \succsim_m b$ . Since this preference cannot belong to  $\mathcal{M}_\succ^2$ , we have a contradiction. Conclude that  $B^2(p) \subseteq B^1(p)$ . To confirm that the inclusion is strict, suppose that  $B^2(p) = B^1(p)$ . Since the maximal set of EU preferences is precisely the set of  $D$ -monotone EU preferences that prefer  $B(p)$  to  $p$ , we have  $\mathcal{M}_\succ^1 = \mathcal{M}_\succ^2$ —contradiction.

#### A.1.7 PROOF OF COROLLARY 2

If the revealed preference relation for  $c$  is acyclic, then there exists a strict preference  $\succ$  that extends it. For any such  $\succ$ ,  $(\succ, c)$  satisfies IUA. To see why, suppose that  $b \succ c(\{b\} \cup A)$ . We need to show that  $c(\{b\} \cup B) = c(B)$  for any  $B \supseteq A$ . Suppose that  $c(\{b\} \cup B) \neq c(B)$  for some  $B \supseteq A$ . Then,  $c(A \cup \{b\})$  is revealed preferred to  $b$ , so  $c(A \cup \{b\}) \succ b$ —contradiction. By Theorem 1,  $(\succ, c)$  has a justification representation  $\mathcal{M}$ —so  $c$  has a justification representation  $(\succ, \mathcal{M})$ .

Suppose that  $c$  has a justification representation  $(\succ, \mathcal{M})$ , so  $(\succ, c)$  has a justification representation  $\mathcal{M}$ . Fix any  $A, b$  such that  $b \neq c(A \cup \{b\})$  and, for some  $B \supseteq A$ ,  $c(B) \neq c(B \cup \{b\})$ . Since IUA is necessary, it must be that  $c(A \cup \{b\}) \succ b$ . That is,  $\succ$  must extend revealed preference, so revealed preference must be acyclic.

#### A.1.8 PROOF OF PROPOSITION 3

This proof builds on that of Theorem 4.

Suppose that no justifiable preference in any representation has  $a \succ_m B$ . Consider the  $\mathcal{M}$  constructed in the proof of Theorem 4. In that proof, we showed that there exists  $\succ_m$  such that

$a \succ_m B$  unless there is a revealed-exclusion-tree from a subset of  $A$  to  $b$ . Suppose there is a revealed-exclusion-tree from  $A' \subset A$  to  $b$ .

Consider the  $\succ$  constructed in the proof of Theorem 4. By Lemma 14,  $y \succ X$  whenever  $y$  is revealed excluded by  $X$ . Thus, each node in a revealed-exclusion-tree must be strictly  $\succ$ -better than all its parents. Since we have a revealed-exclusion-tree from  $A'$  to  $b$ , and since  $\succ$  is transitive, we must have  $b \succ A'$ . As shown in the proof of Theorem 4, we must also have  $b \neq c(A' \cup \{b\})$ .

Take any  $A'' \subset A'$  such that  $b \neq c(A'' \cup \{b\})$  but  $b = c(A''' \cup \{b\})$  for any strict subset  $A'''$  of  $A''$ . Since  $b \neq c(A' \cup \{b\})$ , there must be some such  $A''$ . We show that  $b$  is revealed excluded by  $A''$ .

Suppose that  $A''$  is not a singleton, and choice on the proper subsets of  $A'' \cup \{b\}$  violates WARP. Then, a subset of  $A'' \cup \{b\}$  is a cycle or almost-WARP set. We conclude that some item in  $A'' \cup \{b\}$  is revealed excluded by a subset of  $A'' \cup \{b\}$ . Since  $b = c(\{b\} \cup A''')$  for every strict subset  $A'''$  of  $A''$ ,  $b$  cannot be revealed excluded. Suppose  $a \neq b$  is revealed excluded. By IEA,  $b \neq c(\{b\} \cup A'') = c(\{b\} \cup A'' \setminus \{a\}) = b$ —contradiction. Conclude that choice on the proper subsets of  $A'' \cup \{b\}$  satisfies WARP. Since  $b$  is chosen over every proper subset of  $A''$ , but not over  $A''$  itself, choice on  $A''$  violates WARP. Conclude that  $A''$  is almost-WARP, and that  $b$  is revealed excluded by  $A''$ .

Now suppose that  $a$  is revealed excluded by  $B$ . Suppose  $b = \{b\}$ , so there is a chain from  $a$  to  $b$ , and  $b = c(\{a, b\})$ . Recall that  $a \succ b \succ d$  for any cycle  $(a, b, d)$  and any  $\succ$  in any representation. This implies  $x_1 \succ x_2 \succ \dots \succ x_n$  for any chain  $(x_1, \dots, x_n)$  and any  $\succ$  in any representation. Since there is a chain from  $a$  to  $b$ , we have  $a \succ b$  for any  $\succ$  in any representation. Since  $b = c(\{a, b\})$ , we have  $b \notin \mathcal{M}(\{a, b\})$  for any  $\mathcal{M}$  in any representation.

Now suppose  $|B| > 1$ , so  $\{a\} \cup B$  is an almost-WARP set, and  $a = c(\{a, b\})$  for all  $b \in B$ . We show that every  $\succ$  in every representation agrees with pairwise choice on  $\{a\} \cup B$ . Index the items in  $\{a\} \cup B$  from pairwise-best to pairwise-worst:  $x_1, \dots, x_n$ , where  $x_1 = a$ . Now suppose there is a representation in which  $x_j \succ x_i$ , for  $j > i$ . Since  $x_i = c(\{x_i, x_j\})$ ,  $x_i \succ_m x_j$  for all  $\succ_m \in \mathcal{M}$ . This implies  $c(\{a\} \cup B) = c(\{a\} \cup B \setminus \{x_j\})$ . Since  $\{a\} \cup B$  is almost-WARP,  $a \neq c(\{a\} \cup B)$

but  $a = c(\{a\} \cup B')$  for all proper subsets  $B'$  of  $B$ . Thus,  $c(\{a\} \cup B) = c(\{a\} \cup B \setminus \{x_j\})$  holds only if  $x_j = a$ . But  $a = x_1$  and  $j > i \geq 1$ —contradiction. Conclude that every  $\succ$  in every representation agrees with pairwise choice on  $\{a\} \cup B$ . In particular, every  $\succ$  has  $a \succ B$ . This implies  $a \notin \mathcal{M}(\{a\} \cup B)$  for every  $\mathcal{M}$  in every representation.

#### A.1.9 PROOF OF THEOREM 4

Write  $x C y$  if there is some  $z$  such that  $(x, y, z)$  is a cycle, or  $(z, x, y)$  is a cycle.

**Lemma 13.** *Define a binary relation  $\succ$  by  $a \succ b$  if (1)  $(a, b) \in \text{tr}(C)$  or (2)  $(a, b) \notin \text{tr}(C)$ ,  $(a, b) \notin \text{tr}(C)$ , and  $a = c(\{a, b\})$ . Then,  $\succ$  is a strict preference.*

*Proof.* Clearly,  $\succ$  is complete. Suppose it contains a cycle:  $x_1 \succ x_2 \succ \cdots \succ x_n$  where  $x_1 = x_n$ . For each adjacent pair  $(x_i, x_{i+1})$ , either (1)  $(x_i, x_{i+1}) \in \text{tr}(C)$ , or (2)  $(x_i, x_{i+1}) \notin \text{tr}(C)$ ,  $(x_{i+1}, x_i) \notin \text{tr}(C)$ , and  $x_i = c(\{x_i, x_{i+1}\})$ .

We show that  $(x_i, x_{i+1}) \in \text{tr}(C)$  for some  $i$ . Suppose not. Then we have  $x_i = c(\{x_i, x_{i+1}\})$  for each  $i$ . If  $x_1 = c(\{x_{n-2}, x_1\})$ , then one of the following is a cycle:  $(x_{n-2}, x_{n-1}, x_1)$ ,  $(x_{n-1}, x_1, x_{n-2})$ ,  $(x_1, x_{n-2}, x_{n-1})$ . Then,  $x_{n-2} C x_{n-1}$  or  $x_{n-1} C x_1$ , which contradicts the assumption that no adjacent pair is in  $\text{tr}(C)$ . Conclude that  $x_{n-2} = c(\{x_{n-2}, x_1\})$ . But then we can remove  $x_{n-1}$  without breaking the cycle in  $\succ$ . We can iterate this argument, removing one item at each step, until we end up with  $x_1 = c(\{x_1, x_2\})$ ,  $x_2 = c(\{x_2, x_3\})$ , and  $x_3 = c(\{x_1, x_3\})$ . One of the following must be a cycle:  $(x_1, x_2, x_3)$ ,  $(x_2, x_3, x_1)$ ,  $(x_3, x_1, x_2)$ . In any case,  $x_1 C x_2$  or  $x_2 C x_3$ , which contradicts the assumption that no adjacent pair is in  $\text{tr}(C)$ .

Suppose there is some  $(x_{i-1}, x_i) \notin \text{tr}(C)$ . By the previous step, there is some  $(x_{i-1}, x_i) \notin \text{tr}(C)$  such that  $(x_i, x_{i+1}) \in \text{tr}(C)$ . We can temporarily expand the cycle by adding  $(y_1, \dots, y_k)$  between  $x_i, x_{i+1}$ , where  $x_i C y_1 C \dots C y_k C x_{i+1}$ . Since  $(x_{i-1}, x_i) \notin \text{tr}(C)$ , we must have  $x_{i-1} = c(\{x_{i-1}, x_i\})$ . Since  $x_i C y_1$ , we must have  $x_i = c(\{x_i, y_1\})$ . Now suppose  $y_1 = c(\{y_1, x_{i-1}\})$ .

One of the following must be a cycle:  $(x_i, y_1, x_{i-1}), (y_1, x_{i-1}, x_i), (x_{i-1}, x_i, y_1)$ . The second and third cases are ruled out because they contradict  $(x_{i-1}, x_i) \notin \text{tr}(C)$ , so we must have  $x_i \succ C y_1 \succ C x_{i-1}$ . But then  $(x_i, x_{i-1}) \in \text{tr}(C)$ , which contradicts  $x_{i-1} \succ x_i$ . Conclude that  $x_{i-1} = c(\{y_1, x_{i-1}\})$ . If  $y_1 \succ x_{i-1}$ , it must be that  $(y_1, x_{i-1}) \in \text{tr}(C)$ . Since  $x_i \succ C y_1$ , we must have  $(x_i, x_{i-1}) \in \text{tr}(C)$ , which is a contradiction. Conclude that  $x_{i-1} \succ y_1$ . This means we can remove  $x_i$  while preserving the cycle. Suppose  $(x_{i-1}, y_1) \in \text{tr}(C)$ . Since  $(y_1, x_{i+1}) \in \text{tr}(C)$ , we must have  $(x_{i-1}, x_{i+1}) \in \text{tr}(C)$ . Conclude that  $x_1 \succ \dots \succ x_{i-1} \succ x_{i+1} \succ \dots \succ x_n$ . We have shortened the original cycle by one item. Now suppose  $(x_{i-1}, y_1) \notin \text{tr}(C)$ . We can repeat the argument above to remove  $y_1$  while preserving the cycle. If  $(x_{i-1}, y_2) \in \text{tr}(C)$ , we have again shortened the cycle by one item. If  $(x_{i-1}, y_2) \notin \text{tr}(C)$ , we can repeat the argument once more. We can keep repeating it until we have either shortened the original cycle by one item, or all the  $y$ s have been removed. In that case, we will have  $x_{i-1} \succ x_{i+1}$ —so we will still have shortened the cycle by one item.

We can repeat the procedure above until we have removed all the  $(x_{i-1}, x_i) \notin \text{tr}(C)$ . Re-indexing the elements, we now have a cycle in which each  $(x_{i-1}, x_i) \in \text{tr}(C)$ . For each  $(x_{i-1}, x_i)$ , we can find a finite sequence  $(y_1, \dots, y_k)$  such that  $x_i \succ C y_1 \succ C \dots \succ y_k \succ C x_{i+1}$ . Re-indexing the elements, we now have a cycle in which  $x_{i-1} \succ C x_i$  for each  $i$ . This implies both (1)  $(x_i, x_{i-1}) \in \text{tr}(C)$  for each  $i$ , and (2)  $x_{i-1} = c(\{x_{i-1}, x_i\})$ . Putting (1) and (2) together,  $x_i$  is revealed excluded by  $x_{i-1}$  for each  $i$ . By IEA,  $x_i \notin c(\{x_1, \dots, x_{n-1}\})$  for all  $i$ . This is a contradiction. Conclude that  $\succ$  does not contain a cycle.  $\square$

We let  $\mathcal{M}$  be the set of strict preferences that respect revealed exclusion. That is,  $\succ_m$  belongs to  $\mathcal{M}$  if and only if

$$a \text{ is revealed excluded by } B \implies b \succ_m a \text{ for some } b \in B.$$

It remains to show that  $(\succ, \mathcal{M})$  deliver the correct predictions. First, suppose that  $b = c(\{b\}) \cup$

A). We show that  $b \succ_m A$  for some  $\succ_m \in \mathcal{M}$ . We construct an appropriate  $\succ_m$  as follows. Let

$$\begin{aligned}
B_0 &:= A \\
B_i &= B_{i-1} \cup \bigcup_{B' \in \mathcal{F}(B_i)} \{a \in A : a \text{ is revealed excluded by } B'\} \quad \text{for } i > 0 \\
B &:= \bigcup_{i \geq 0} B_i \\
T &= A \setminus B
\end{aligned}$$

For any distinct  $x, y \in A$ , impose  $x \succ_m y$  if (1)  $\{x, y\} \subset B$  or  $\{x, y\} \subset T$  and  $y \succ x$ , or (2)  $x \in T$  and  $y \in B$ . Notice that  $\succ_m$  is a strict preference.

We show that  $b \in T$ , so  $b \succ_m A$ . Suppose not. Then (using the definition of a tree from the proof of Theorem 1), there is a revealed-exclusion-tree from a subset of  $A$  to  $b$ . Consider the menu  $Z$  that consists of everything in the tree along with any other items in  $A$ . Since  $b$  is revealed excluded in  $Z$ ,  $b \neq c(Z)$ . Since every item in  $Z \setminus A$  is revealed excluded in  $Z$ ,  $c(Z) = c(A \cup \{b\})$ , so  $b \neq c(A \cup \{b\})$ —contradiction.

Now we show that  $\succ_m$  respects revealed exclusion.

**Lemma 14.** *If  $y$  is revealed excluded by  $X$ , then  $y \succ X$ .*

*Proof.* Suppose  $X = \{x\}$ . Then there is a chain from  $y$  to  $x$ , so  $y \succ x$ . Now suppose  $|X| > 1$ . Then  $X \cup \{y\}$  is almost-WARP, so  $y \neq c(X)$  but  $y = c(X' \cup \{y\})$  for any strict subset  $X'$  of  $X$ . Suppose that there is a chain from  $x$  to  $y$  for some  $x \in X$ . Since  $y = c(\{x, y\})$ ,  $x$  is revealed excluded by  $y$ . Since IEA is necessary,  $c(\{y\} \cup X' \setminus \{x\}) = c(\{y\} \cup X)$ —contradiction. Conclude that there is no chain from  $x$  to  $y$  for any  $x \in X$ . Since  $y = c(\{x, y\})$  for each  $x \in X$ , we have  $y \succ X$ .  $\square$

Suppose  $y$  is revealed excluded by  $X$ , but  $y \succ_m X$ . We can write  $X = T' \cup B'$  where  $B' \subseteq B$  and  $T' \subseteq T$ . Suppose  $y \in B$ . If  $T'$  is nonempty, we have  $t \succ_m y$  for all  $t \in T'$ . Thus,  $T'$  must be empty.

It must be that  $y \succ_m B'$  but  $y$  is revealed excluded by  $B'$ . By Lemma 14, we have  $y \succ B'$ . Since  $\{y\} \cup B' \subseteq B$ , we have  $B' \succ_m y$  by definition of  $\mathcal{M}$ —contradiction. Now suppose  $y \in T$ . If  $T'$  is empty, then  $y$  is revealed excluded by  $B' \subseteq B$ , which contradicts the assumption that  $y \in T$ . Thus,  $T' \neq \emptyset$ . By Lemma 14, we have  $y \succ B' \cup T'$ . Since  $\{y\} \cup T' \subseteq T$ , we have  $T' \succ_m y$ —contradiction. Conclude that  $\succ_m$  respects revealed exclusion. There is a preference  $\succ_m \in \mathcal{M}$  such that  $b \succ_m A$ .

Now suppose that  $b \succ c(\{b\} \cup A)$ . We show that there is no  $\succ_m \in \mathcal{M}$  such that  $b \succ_m A$ .

**Lemma 15.** *If  $x \in X$  is not revealed excluded by any subset of  $X$ , and if  $x \neq c(X)$ , then  $c(X) \succ x$ .*

*Proof.* Let

$$X^* := \{x \in X : x \text{ is not revealed excluded by any } X' \subset X\}.$$

By assumption,  $x \in X^*$ . By IEA,  $c(X) = c(X^*)$ , so  $x \neq c(X^*)$ . Take any  $X' \subseteq X^*$  such that  $|X'| = 3$ . If choice on  $X'$  violates WARP, then  $X'$  is an almost-WARP set or a cycle, so something in  $X'$  is revealed excluded. This contradicts the definition of  $X^*$ . By induction on the size of  $X'$ , we can show that choice on  $X^*$  satisfies WARP. Since  $x \in X^*$ , we must have  $c(X^*) = c(\{x, c(X^*)\})$ . This implies  $c(X^*) \succ x$  unless there is a chain from  $x \rightarrow c(X^*)$ . In that case,  $x$  is revealed excluded by  $c(X^*)$ , which contradicts  $x \in X^*$ . Conclude that  $c(X^*) \succ x$ . Since  $c(X^*) = c(X)$ , we have  $c(X) \succ x$ . □

Suppose that  $b$  is not revealed excluded by a subset of  $A$ . Lemma 15 implies  $c(\{b\} \cup A) \succ b$ , which contradicts the assumption that  $b \succ c(\{b\} \cup A)$ . Conclude that  $b$  is revealed excluded by a subset of  $A$ , so there is no  $\succ_m \in \mathcal{M}$  such that  $b \succ_m A$ .

#### A.1.10 PROOF OF COROLLARY 3

The proof of Theorem 4 constructs precisely this representation. Suppose the  $\mathcal{M}$  constructed in that proof is not maximal. Then, there is some  $\succ_m$  in some representation that does not respect

revealed exclusion. That is,  $b \succ_m A$  for some menu  $A$  and item  $b$  such that  $b$  is revealed excluded by  $A$ . Proposition 3 says that this cannot be the case.

For uniqueness of  $\succ$ , consider some representation  $(\succ', \mathcal{M})$  such that  $a \succ b$  but  $b \succ' a$ . Suppose that there is a chain from  $a$  to  $b$ . Then,  $a$  is revealed preferred to  $b$ , which contradicts  $b \succ' a$ . Now suppose that there is a chain from  $b$  to  $a$ . This implies  $b \succ a$ , which contradicts  $a \succ b$ . Finally, suppose that  $a$  and  $b$  are not linked by a chain. By definition of  $\succ$ , we have  $a = c(\{a, b\})$ . Since  $b \succ' a$  and  $(\succ', \mathcal{M})$  is a representation, it must be that  $a \succ_m b$  for all  $\succ_m \in \mathcal{M}$ . Consider a preference  $\succ_{bad}$  that is exactly opposite  $\succ$ . We show that  $\succ_{bad} \in \mathcal{M}$ . Recall from the proof of Proposition 3 that  $y \succ X$ , so  $X \succ_{bad} y$ , whenever  $y$  is revealed excluded by  $X$ . Conclude that  $\succ_{bad}$  respects revealed exclusion, so  $\succ_{bad} \in \mathcal{M}$ . Since  $a \succ b$ , we have  $b \succ_{bad} a$ . This contradicts the assumption that  $a \succ_m b$  for all  $\succ_m \in \mathcal{M}$ . Conclude that  $(\succ', \mathcal{M})$  is not a representation.

#### A.1.1.1 PROOF OF PROPOSITION 4

We show sufficiency. Since  $c_L$  satisfies IEA, we construct  $\succ$  in accordance with Lemma 13. We then let  $\mathcal{M}^L$  be the set of strict preferences consistent with revealed exclusion in  $L$ . That is,  $\succ_m \in \mathcal{M}^L$  if and only if

$$a \text{ is revealed excluded by } B \text{ in } L \implies b \succ_m a \text{ for some } b \in B.$$

By Theorem 4,  $(\succ, \mathcal{M}^L)$  represents  $c_L$ .

For  $\mathcal{M}^H$ , we need to define a new relation  $R$  that captures replacement as well as revealed exclusion.

**Definition 58** (Relation  $R$ ). *Say that  $Z R z$  if either of the following holds:*

1.  $z$  is revealed excluded in  $L$  by  $Z$ .
2.  $z$  is replaced in  $Z \cup \{z\}$ , and no item in  $Z$  is revealed excluded in  $L$  by any subset of  $Z \cup \{z\}$ .



Let  $\mathcal{M}^H$  be the set of strict preferences that respect  $R$ . That is,  $\succ_m \in \mathcal{M}^H$  if and only if

$$B R a \implies b \succ_m a \text{ for some } b \in B.$$

Since  $R$  extends revealed exclusion in  $L$ ,  $\mathcal{M}^H \subseteq \mathcal{M}^L$ .

**Lemma 16.** *If  $a$  is replaced in  $B \cup \{a\}$ , then  $a$  is replaced in*

$$B^* = \{b \in B : b \text{ is not revealed excluded in } L \text{ by any } B' \subset B\} \cup \{a\},$$

so  $B^* R a$ .

*Proof.* Suppose  $a = c_L(B \cup \{a\}) \neq c_H(B \cup \{a\})$ . Since  $c_L$  satisfies IEA and  $c_H$  satisfies IREA, we have  $c_L(B \cup \{a\}) = c_L(B^* \cup \{a\})$  and  $c_H(B \cup \{a\}) = c_H(B^* \cup \{a\})$ . Thus,  $a = c_L(B^* \cup \{a\}) \neq c_H(B^* \cup \{a\})$ , so  $a$  is replaced in  $B^* \cup \{a\}$ .  $\square$

Lemma 16 implies that each  $\succ_m \in \mathcal{M}^H$  respects replacement:

$$a \text{ is replaced in } B \cup \{a\} \implies b \succ_m a \text{ for some } b \in B \setminus \{a\}.$$

It remains to show that  $(\succ, \mathcal{M}^H)$  delivers the correct predictions. First, suppose that  $b = c_H(\{b\} \cup A)$ . We show that  $b \succ_m A$  for some  $\succ_m \in \mathcal{M}_H$ . We can construct an appropriate  $\succ_m$  following the approach of Theorem 4. The only difference is that we use  $R$  instead of revealed exclusion, so

$$B_i = B_{i-1} \cup \bigcup_{B' \in \mathcal{F}(B_i)} \{a \in \mathcal{A} : B' R a\} \text{ for } i > 0.$$

We can then define  $\succ_m$  as before. To use the argument that  $\succ_m \in \mathcal{M}^H$ , we need to show that  $Z R z$  implies  $z \succ Z$ . From Lemma 14, we know that  $z \succ Z$  if  $z$  is revealed excluded in  $L$  by  $Z$ . Suppose

instead that  $z$  is replaced in  $Z \cup \{z\}$ , so  $z = c_L(Z \cup \{z\}) \neq c_H(Z \cup \{z\})$ . Suppose further that no item in  $Z$  is revealed excluded in  $L$  by any subset of  $Z \cup \{z\}$ . Since  $z = c_L(Z \cup \{z\})$ ,  $z$  is not revealed excluded in  $L$  by any subset of  $Z$ .

Toward a contradiction, suppose  $z' \succ z$  for some  $z' \in Z$ . There are two possibilities: (1)  $z = c_L(\{z', z\})$  and  $z'$  comes before  $z$  in a chain in  $L$ , or (2)  $z' = c_L(\{z', z\})$  and  $z, z'$  are not linked by a chain in  $L$ . In case (1),  $z'$  is revealed excluded in  $L$  by  $z$ , which contradicts our assumption about  $Z$ . To rule out case (2), we show that  $c_L(\{z, z'\}) = c_L(Z \cup \{z\}) = z$ . In the proof of Theorem 4, we showed that the restriction of  $c$  to a set in which nothing is revealed excluded satisfies WARP. Since nothing in  $Z \cup \{z\}$  is revealed excluded in  $L$ ,  $c_L$  satisfies WARP on  $Z \cup \{z\}$ . WARP and  $z = c_L(Z \cup \{z\})$  imply  $z = c_L(\{z, z'\})$ . This completes the proof that  $Z R z$  implies  $z \succ Z$ . We can now use the arguments from Theorem 4, with  $R$  in place of revealed exclusion, to show that  $\succ_m$  respects  $R$ .

Now suppose that  $b \succ c_H(\{b\} \cup A)$ . We show that there is no  $\succ_m \in \mathcal{M}_H$  such that  $b \succ_m A$ . Toward a contradiction, suppose that  $\neg(A' R b)$  for all  $A' \subseteq A$ . Let

$$A^* := \{a \in A : \neg(A'' R a) \text{ for all } A'' \subseteq A\}.$$

Suppose  $c_L(A^* \cup \{b\}) \neq c_H(A^* \cup \{b\})$ , so  $c_L(A^* \cup \{b\})$  is replaced in  $A^* \cup \{b\}$ . By Lemma 16, there exists  $X \subset A^* \cup \{b\}$  such that  $X R c_L(A^* \cup \{b\})$ . By definition of  $A^*$ ,  $c_L(A^* \cup \{b\}) \notin A^* \cup \{b\}$ —contradiction. Conclude that  $c_L(A^* \cup \{b\}) = c_H(A^* \cup \{b\})$ . Since  $c_H$  satisfies IREA, we have  $c_H(A^* \cup \{b\}) = c_H(A \cup \{b\})$ . Putting both equalities together, we have  $b \succ c_L(A^* \cup \{b\})$ . But since  $b$  is not revealed excluded in  $L$  by any subset of  $A^* \cup \{b\}$ , Lemma 15 says that  $c_L(A^* \cup \{b\})$ —contradiction. Conclude that  $A' R b$  for some  $A' \subseteq A$ , so  $\neg(b \succ_m A)$  for all  $\succ_m$  that respect  $R$ .

A.I.12 LEMMAS USED IN PROOFS FOR RJ

For any  $A \in \mathcal{F}(\Delta(Z))$  and any  $\mathcal{M} \in \text{supp}(\nu)$ , let

$$W_{\mathcal{M}}(A) := \bigcap_{\mathcal{Z} \in \mathcal{M}} \{x \in \Delta(Z) : \exists a \in A \text{ s.t. } a \succ x\}.$$

For brevity, we typically write  $\nu(x \in W(A))$  instead of  $\nu(\{\mathcal{M} : x \in W_{\mathcal{M}}(A)\})$ .

**Lemma 17.**

1. For any  $\mathcal{M}, \mathcal{N} \in \text{supp}(\nu)$  such that  $\mathcal{M} \subset \mathcal{N}$  and any  $p \in \Delta(Z)$ ,  $\text{cl}(W_{\mathcal{N}}(p)) \setminus \{p\} \subset W_{\mathcal{M}}(p)$ .
2. For any  $p, q \in \Delta(Z)$  such that  $\nu(q \notin W(p)) \in (\nu(\mathcal{U}), 1)$ , we can find  $\mathcal{M} \in \text{supp}(\nu)$  such that  $q$  is on the boundary of  $W_{\mathcal{M}}(p)$ .

*Proof.*

1. Fix any  $q \in \text{cl}(W_{\mathcal{N}}(p)) \setminus \{p\}$ . We have  $n'p \geq n'q$  for all  $n \in \mathcal{N}_R$ . Thus, any  $n \in \mathcal{N}_R$  such that  $n'p = n'q$  is on the boundary of  $\mathcal{N}_R$ . By the first property of  $\nu$ , we have  $\mathcal{M} \subset \text{int}(\mathcal{N})$ , so  $\mathcal{M}$  cannot contain any boundary point of  $\mathcal{N}$ . In particular,  $m'p > m'q$  for all  $m \in \mathcal{M}$ . Thus,  $q \in W_{\mathcal{M}}(p)$ .
2. Let  $\mathcal{M}^*$  be the element of  $\text{supp}(\nu)$  such that

$$\nu(\{\mathcal{M} : \mathcal{M} \subseteq \mathcal{M}^*\}) = \nu(q \in W(p)).$$

By the second property of  $\nu$ ,  $\mathcal{M}^*$  exists and is unique.

We show that  $q$  is on the boundary of  $W_{\mathcal{M}^*}(p)$ . Suppose  $q \notin \text{cl}(W_{\mathcal{M}^*}(p))$ . Let  $\underline{\mathcal{M}}$  be the largest element of  $\text{supp}(\nu)$  such that  $q \in \text{cl}(W_{\underline{\mathcal{M}}}(p))$ . ( $\underline{\mathcal{M}}$  must exist because  $\text{supp}(\nu)$  is

closed.) By the first property of  $\nu$ , there exists  $\tilde{\mathcal{M}} \in \text{supp}(\nu)$  such that  $\underline{\mathcal{M}} \subset \text{int}(\tilde{\mathcal{M}}) \subset \tilde{\mathcal{M}} \subset \text{int}(\mathcal{M}^*)$ . By definition of  $\underline{\mathcal{M}}$ , we must have  $q \notin \text{cl}(W_{\tilde{\mathcal{M}}})$ , so

$$\nu(q \in W(p)) < \nu(\{\mathcal{M} : \mathcal{M} \subseteq \tilde{\mathcal{M}}\}) < \nu(\{\mathcal{M} : \mathcal{M} \subseteq \mathcal{M}^*\}),$$

which contradicts the definition of  $\mathcal{M}^*$ .

Now suppose  $q \in W_{\mathcal{M}^*}(p)$ . Let  $\bar{\mathcal{M}}$  be the smallest element of  $\text{supp}(\nu)$  such that  $q \notin W_{\bar{\mathcal{M}}}(p)$ . (Again,  $\bar{\mathcal{M}}$  must exist because  $\text{supp}(\nu)$  is closed.) By the first property of  $\nu$ , there exists  $\tilde{\mathcal{M}} \in \text{supp}(\nu)$  such that  $\mathcal{M}^* \subset \text{int}(\tilde{\mathcal{M}}) \subset \tilde{\mathcal{M}} \subset \text{int}(\bar{\mathcal{M}})$ . Since  $q \in W_{\tilde{\mathcal{M}}}(p)$ , we have

$$\nu(q \in W(p)) \geq \nu(\{\mathcal{M} : \mathcal{M} \subseteq \tilde{\mathcal{M}}\}) > \nu(\{\mathcal{M} : \mathcal{M} \subseteq \mathcal{M}^*\}),$$

which contradicts the definition of  $\mathcal{M}^*$ .

□

**Lemma 18.** *For any  $p, q, r \in \Delta(Z)$ :*

1.  $\nu(p \in W(r)) \geq \min\{\nu(p \in W(q)), \nu(q \in W(r))\}$ .
2. *If  $\nu(p \in W(q)) \neq \nu(q \in W(r))$  and  $\nu(p \in W(r)) > 0$ , then*

$$\nu(p \in W(r)) > \min\{\nu(p \in W(q)), \nu(q \in W(r))\}.$$

*Proof.*

1. Since  $\text{supp}(\nu)$  can be ordered by set inclusion,

$$\nu(p \in W(q) \text{ and } q \in W(r)) = \min\{\nu(p \in W(q)), \nu(q \in W(r))\}.$$

$p \in W(q)$  and  $q \in W(r)$  implies  $p \in W(r)$ , so

$$\nu(p \in W(r)) \geq \min\{\nu(p \in W(q)), \nu(q \in W(r))\}.$$

2. The non-trivial case is  $\min\{\nu(p \in W(q)), \nu(q \in W(r))\} > 0$ . Suppose  $\nu(p \in W(q)) > \nu(q \in W(r))$ . If  $\nu(p \in W(q)) < \max_{x \in \Delta(Z)} \nu(x \in W(q))$ , we can use the second part of Lemma 17 to obtain  $\mathcal{M}_1$  such that  $p$  is on the boundary of  $W_{\mathcal{M}_1}(q)$ . If  $\nu(p \in W(q)) = \max_{x \in \Delta(Z)} \nu(x \in W(q))$ , set  $\mathcal{M}_1 = \mathcal{U}$ . Since the first part of Lemma 17 says  $\text{cl}(W_{\mathcal{M}_1}(q)) \setminus \{q\} \subset W_{\mathcal{M}}(q)$  for all  $\mathcal{M} \subset \mathcal{M}_1$ , we have  $p \in W_{\mathcal{M}}(q)$  for all  $\mathcal{M} \subset \mathcal{M}_1$ . By assumption,  $\nu(q \in W(r)) < \max_{x \in \Delta(Z)} \nu(x \in W(r))$ , so we can find  $\mathcal{M}_2$  such that  $q$  is on the boundary of  $W_{\mathcal{M}_2}(r)$ . We have  $r \succsim q$  for all  $\succsim \in \mathcal{M}_2$ . We also have  $\mathcal{M}_2 \subset \mathcal{M}_1$ , so  $p \in W_{\mathcal{M}_2}(q)$ . That is,  $q \succ p$  for all  $\succ \in \mathcal{M}_2$ . By transitivity,  $r \succ p$  for all  $\succ \in \mathcal{M}_2$ , so  $p \in W_{\mathcal{M}_2}(r)$ . We can find  $\mathcal{M}_3 \supset \mathcal{M}_2$  such that  $p \in W_{\mathcal{M}_3}(r)$ , so  $\nu(p \in W(r)) > \nu(q \in W(r))$ . A parallel argument covers the case  $\nu(p \in W(q)) < \nu(q \in W(r))$ .

□

**Lemma 19.**

1. Suppose that  $\nu(\mathcal{U}) < 1$ . For any  $p \in \text{int}(\Delta(Z))$ , the sets

$$\{x \in \Delta(Z) : \nu(p \in W(x)) > 0\} \text{ and } \{x \in \Delta(Z) : \nu(x \in W(p)) > 0\}$$

are disjoint, convex, open and nonempty.

2. For any  $p, q, r \in \Delta(Z)$ : if  $\nu(p \in W(q)) > 0$  and  $\nu(r \in W(p)) > 0$ , then  $\nu(p \in W(x)) = \nu(x \in W(p)) = 0$  for some  $x \in \text{co}(\{q, r\})$ .

*Proof.*

- i. Since  $\rho$  does not have an REU representation,  $\nu(\mathcal{U}) < 1$ , so there exist  $x, y \in \Delta(Z)$  such that  $\nu(y \in W(x)) > 0$ . By Independence,  $\nu(p + \lambda(y - x) \in W(p)) > 0$  for any  $\lambda > 0$  such that  $p + \lambda(y - x) \in \Delta(Z)$ , and  $\nu(p \in W(p + \lambda(x - y))) > 0$  for any  $\lambda > 0$  such that  $p + \lambda(x - y) \in \Delta(Z)$ . Since  $p$  is interior,  $p + \lambda(y - x)$  and  $p + \lambda(x - y)$  will belong to  $\Delta(Z)$  for  $\lambda$  small enough.

Suppose that  $\nu(p \in W(q)) > 0$  and  $\nu(q \in W(p)) > 0$ . If  $p \in W(q)$ , then  $W(p) \subset W(q)$ . If in addition  $q \in W(p)$ , then  $q \in W(q)$ . Since  $\text{supp}(\nu)$  can be ordered by set inclusion,

$$\nu(p \in W(q) \text{ and } q \in W(p)) = \min\{\nu(p \in W(q)), \nu(q \in W(p))\} > 0$$

so  $\nu(q \in W(q)) > 0$ . This implies  $\rho(q|\{q\}) < 1$ —contradiction. Conclude that  $\nu(q \in W(p)) = 0$  whenever  $\nu(p \in W(q)) > 0$ .

For convexity, suppose that  $\nu(q \in W(p)) \geq \nu(r \in W(p)) > 0$ . We have  $\nu(\{q, r\} \subset W(p)) = \nu(r \in W(p))$ . Since each realization of  $W(p)$  is convex,  $\nu(\text{co}(\{q, r\}) \subset W(p)) = \nu(r \in W(p)) > 0$ . A parallel argument establishes that  $\nu(q \in B(p)) \geq \nu(r \in B(p)) > 0$  implies  $\nu(\text{co}(\{q, r\}) \subset B(p)) > 0$ .

For any  $\mathcal{M} \in \text{supp}(\mu)$ , the sets  $W_{\mathcal{M}}(p)$  and  $\{x \in \Delta(Z) : p \in W_{\mathcal{M}}(x)\}$  are open because  $\mathcal{M}$  is closed. If  $\nu(x \in W(p)) > 0$ , then  $x \in W_{\mathcal{M}}(p)$  for some  $\mathcal{M} \neq \underline{\mathcal{M}} := \min(\text{supp}(\nu), \supset)$ . That implies  $x \in W_{\underline{\mathcal{M}}}(p)$ . Conversely, if  $x \in W_{\underline{\mathcal{M}}}(p)$ ,  $x \in W_{\mathcal{M}}(p)$  for some  $\mathcal{M} \neq \underline{\mathcal{M}}$ , so  $\nu(x \in W(p)) > 0$ . Thus,  $\nu(x \in W(p)) > 0$  if and only if  $x \in W_{\underline{\mathcal{M}}}(p)$ .

We have

$$\{x \in \Delta(Z) : \nu(x \in W(p)) > 0\} = W_{\underline{\mathcal{M}}}(p).$$

A similar argument establishes that  $\nu(p \in W(x)) > 0$  if and only if  $p \in W_{\underline{\mathcal{M}}}(x)$ , so

$$\{x \in \Delta(Z) : \nu(p \in W(x)) > 0\} = \{x \in \Delta(Z) : p \in W_{\underline{\mathcal{M}}}(x)\}.$$

2. Since  $\{x \in \Delta(Z) : \nu(p \in W(x)) > 0\} \cap \text{co}(\{q, r\})$  and  $\{x \in \Delta(Z) : \nu(x \in W(p)) > 0\} \cap \text{co}(\{q, r\})$  are disjoint, nonempty, and open in  $\text{co}(\{q, r\})$ , the set  $\{x \in \Delta(Z) : \nu(p \in W(x)) = \nu(x \in W(p)) = 0\} \cap \text{co}(\{q, r\})$  must be nonempty.

□

**Lemma 20.** For any  $p \in \Delta(Z)$  and any  $a \in \mathcal{F}(\Delta(Z))$ ,

$$\nu(p \in W(A)) = \max_{a \in \text{co}(A)} \nu(p \in W(a)).$$

*Proof.* Fix any  $\mathcal{M} \in \text{supp}(\nu) \setminus \{\mathcal{U}\}$ . We show that  $p \in \text{cl}(W_{\mathcal{M}}(A))$  if and only if  $p \in \text{cl}(W_{\mathcal{M}}(a))$  for some  $a \in \text{co}(A)$ .  $p \in \text{cl}(W_{\mathcal{M}}(A))$  means that for each  $\succsim \in \mathcal{A}$ , there exists  $a \in A$  such that  $a \succsim p$ . Let

$$\begin{aligned} \mathcal{M}_R &= \bigcup_{\succsim \in \mathcal{M}} \{(m(a_1), \dots, m(a_{|A|})) : m \text{ represents } \succsim \text{ and } m(p) = 0\} \\ \mathcal{N}_R &= \mathbb{R}_{--}^{|A|} \end{aligned}$$

Since  $\mathcal{M}$  is convex and nonempty, so is  $\mathcal{M}_R$ . Clearly,  $\mathcal{N}_R$  is convex and nonempty, and  $\mathcal{N}_R \cap \mathcal{M}_R = \emptyset$ . By the Separating Hyperplane Theorem, there exist nonzero  $v \in \mathbb{R}^{|A|}$  and  $c$  such that  $v'm \geq c$  for all  $m \in \mathcal{M}_R$  and  $v'n \leq c$  for all  $n \in \mathcal{N}$ . Suppose that  $v'n = c$  for some  $n$ . Since  $\mathcal{N}_R$  is open and  $v$  is nonzero, we can find  $\tilde{n} \in \mathcal{N}_R$  such that  $v'\tilde{n} > c$ —contradiction. Thus,  $v'n < c$  for all  $n \in \mathcal{N}_R$ . Suppose that  $c < 0$ . By choosing  $n \in \mathcal{N}_R$  sufficiently close to  $o$ , we get  $v'n > c$ —contradiction. Thus,  $c \geq 0$ . Suppose  $c > 0$ . By choosing  $m \in \mathcal{M}_R$  sufficiently close to  $o$ , we get  $v'm < c$ —

contradiction. Thus,  $c = 0$ . Suppose  $v(i) < 0$  for some  $i$ . By choosing  $n(i)$  sufficiently negative and  $n(j)$  sufficiently close to 0 for  $j \neq i$ , we get  $v'n > 0$ —contradiction. Thus,  $v(i) \geq 0$  for all  $i$ , and (since  $v$  is nonzero)  $\sum_i v(i) > 0$ . Let

$$a^* = \sum_i \frac{v(i)}{\sum_j v(j)} a_i.$$

Since  $\sum_j v(j) > 0$ , we have

$$\sum_i \frac{v(i)}{\sum_j v(j)} m(i) \geq 0$$

for all  $m \in \mathcal{M}_R$ . This implies  $m(a^*) \geq m(p)$  for every  $m$  that represents some  $\zeta \in \mathcal{M}$ , so  $a^* \succeq p$  for all  $\zeta \in \mathcal{M}$ .

For each  $\mathcal{M} \in \text{supp}(v)$  such that  $p \in \text{cl}(W_{\mathcal{M}}(A))$ , let

$$A_{\mathcal{M}} := \{a \in \text{co}(A) : p \in \text{cl}(W_{\mathcal{M}}(a))\}.$$

We have just shown that  $A_{\mathcal{M}}$  is nonempty; we now show that it is closed. Take any convergent sequence belonging to  $A_{\mathcal{M}}$ . Since

$$\text{cl}(W_{\mathcal{M}}(x)) = (x - y) + \text{cl}(W_{\mathcal{M}}(y))$$

for any  $x, y \in \Delta(Z)$ , we have  $\text{cl}(W_{\mathcal{M}}(a_i)) \rightarrow \text{cl}(W_{\mathcal{M}}(\lim_i a_i))$ . Since  $p \in \text{cl}(W_{\mathcal{M}}(a_i))$  for all  $i$ ,  $p \in \text{cl}(W_{\mathcal{M}}(\lim_i a_i))$ . Thus,  $\lim_i a_i \in A_{\mathcal{M}}$  as desired.

Note that  $A_{\mathcal{M}} \supset A_{\mathcal{M}'}$  if  $\mathcal{M} \subset \mathcal{M}'$ . Since  $\text{supp}(v)$  can be ordered by set inclusion, so can the  $A_{\mathcal{M}}$ .



Since each  $A_{\mathcal{M}}$  is closed and the  $A_{\mathcal{M}}$  can be ordered by set inclusion,

$$\bigcap_{\{\mathcal{M} \in \text{supp}(\nu) : p \in \text{cl}(W_{\mathcal{M}}(A))\}} A_{\mathcal{M}}$$

is nonempty. For any  $a^*$  belonging to this set,  $p \in \text{cl}(W_{\mathcal{M}}(a^*))$  for all  $\mathcal{M} \in \text{supp}(\nu)$  such that  $p \in \text{cl}(W_{\mathcal{M}}(A))$ . That is,

$$p \in \text{cl}(W_{\mathcal{M}}(a^*)) \iff p \in \text{cl}(W_{\mathcal{M}}(A)).$$

Fix any  $\mathcal{M}$  such that  $p \in W_{\mathcal{M}}(A)$ . Suppose  $p \notin W_{\mathcal{M}}(a^*)$ . We can find  $\mathcal{N} \supset \mathcal{M}$  such that  $p \in W_{\mathcal{N}}(A)$ . Since  $\text{cl}(W_{\mathcal{N}}(a^*)) \setminus \{a^*\} \subset W_{\mathcal{M}}(a^*)$ ,  $p \notin \text{cl}(W_{\mathcal{N}}(a^*))$ . This contradicts the definition of  $a^*$ . We conclude that

$$p \in W_{\mathcal{M}}(a^*) \iff p \in W_{\mathcal{M}}(A).$$

This implies

$$\nu(p \in W(a^*)) = \nu(p \in W(A)).$$

Since  $W(a) \subset W(A)$  for all  $a \in \text{co}(A)$ , we must have  $\nu(p \in W(a)) \leq \nu(p \in W(A))$  for all  $a \in \text{co}(A)$ . We conclude that

$$\nu(p \in W(a^*)) = \max_{a \in \text{co}(A)} \nu(p \in W(a)).$$

□

**Lemma 21.** For any  $p, q \in \Delta(Z)$  and any  $A \in \mathcal{F}(\Delta(Z))$ :

$$I. \nu(p \in W(A)) = \nu(p \in W(A \setminus \{p\})).$$

2. If  $\nu(p \in W(A)) \leq \nu(q \in W(A))$ , then  $\nu(p \in W(A)) = \nu(p \in W(A \setminus \{q\}))$ .

*Proof.*

1. By Lemma 20, there exists some  $a \in \text{co}(A)$  such that  $p \in W(A)$  if and only if  $p \in W(a)$ . Suppose that  $a = \lambda p + (1 - \lambda)a'$  for some  $a' \in \text{co}(A \setminus \{p\})$  and some  $\lambda \in (0, 1)$ . By Independence,  $p \in W(\lambda p + (1 - \lambda)a')$  if and only if  $p \in W(a')$ . Thus,  $\nu(p \in W(a')) = \nu(p \in W(A))$ , so

$$\nu(p \in W(A \setminus \{p\})) \geq \nu(p \in W(a')) = \nu(p \in W(A)).$$

Since  $\text{co}(A \setminus \{p\}) \subseteq \text{co}(A)$ , it cannot be that  $\nu(p \in W(A \setminus \{p\})) > \nu(p \in W(A))$ , so we have  $\nu(p \in W(A \setminus \{p\})) = \nu(p \in W(A))$  as desired.

2. Suppose that  $\nu(p \in W(A)) \leq \nu(q \in W(A))$ . Since  $\text{supp}(\nu)$  can be ordered by set inclusion,  $p \in W(A)$  implies  $q \in W(A)$ . By the first part of this Lemma,  $q \in W(A)$  implies  $q \in W(A \setminus \{q\})$ . By Independence,  $q \in W(A \setminus \{q\})$  implies  $\lambda q + (1 - \lambda)a \in W(\lambda A \setminus \{q\} + (1 - \lambda)a)$  for all  $\lambda \in (0, 1]$  and all  $a \in \text{co}(A \setminus \{q\})$ . Since  $\lambda A \setminus \{q\} + (1 - \lambda)a \subset A \setminus \{q\}$ ,  $\lambda q + (1 - \lambda)a \in W(\lambda A \setminus \{q\})$  implies  $\lambda q + (1 - \lambda)a \in W(A \setminus \{q\})$ .

Suppose that  $p \in W(\lambda q + (1 - \lambda)a)$  for some  $\lambda \in (0, 1]$  and  $a \in \text{co}(A \setminus \{q\})$ . Since  $\lambda q + (1 - \lambda)a \in \text{co}(A)$ ,  $p \in W(A)$ . We have just seen that this implies  $\lambda q + (1 - \lambda)a \in W(A \setminus \{q\})$ , which in turn implies  $W(\lambda q + (1 - \lambda)a) \subset W(A \setminus \{q\})$ . Since  $p \in W(\lambda q + (1 - \lambda)a)$  implies  $q \in W(A \setminus \{q\})$ , it cannot be that

$$\nu(p \in W(\lambda q + (1 - \lambda)a)) > \nu(p \in W(A \setminus \{q\})).$$

This implies

$$\nu(p \in W(A)) = \max_{x \in \text{co}(A)} \nu(p \in W(x)) = \nu(p \in W(A \setminus \{q\})).$$

□

**Lemma 22.** For any  $p \in \text{int}(\Delta(Z))$ :

1. The set  $D(p) := \arg \max_{x \in \Delta(Z)} \nu(x \in W(p))$  is nonempty.
2. For any  $d \in D(p)$ ,  $\nu(d \in W(p)) \geq \nu(q \in W(A))$  for any  $q \in \Delta(Z)$  and  $A \in \mathcal{F}(\Delta(Z))$ .
3. For any  $q \in \Delta(Z) \setminus D(p)$  such that  $\nu(q \in W(p)) > 0$  and any  $\lambda \in [0, 1)$ :

$$d \in D(p) \implies \nu(q \in W(\lambda p + (1 - \lambda)d)) < \nu(q \in W(p))$$

$$q \in D(b) \implies \nu(\lambda b + (1 - \lambda)q \in W(p)) < \nu(q \in W(p))$$

*Proof.*

1. Notice that  $D(p) = \bigcap_{\mathcal{M} \in \text{supp}(\nu) \setminus \{\mathcal{U}\}} W_{\mathcal{M}}(p)$ . To see why  $D(p)$  is nonempty, fix any increasing sequence  $\{\mathcal{M}_i \in \text{supp}(\nu) : \mathcal{M}_i \neq \mathcal{U}\}_{i=1}^{\infty}$  such that, for each  $t \in (\nu(\mathcal{U}), 1)$ ,

$$\exists i \in \{1, 2, \dots\} \text{ s.t. } \nu(\{\mathcal{M} : \mathcal{M} \supset \mathcal{M}_i\}) < t.$$

Fix  $\varepsilon$  such that  $p + \varepsilon(q - r) \in \Delta(Z)$  for any  $q, r \in \Delta(Z)$ . For each  $i$ , choose a point  $w_i$  on the boundary of  $W_{\mathcal{M}_i}(p) \cap (B_\varepsilon(p) \setminus B_{\varepsilon/2}(p))$ . Pass to a convergent subsequence if necessary, and let  $w = \lim_i w_i$ . Since  $w \notin B_{\varepsilon/2}(p)$ ,  $w \neq p$ . Since the  $W_{\mathcal{M}_j}(p)$  are decreasing, we must have  $\{w_j, w_{j+1}, \dots\} \subset \text{cl}(W_{\mathcal{M}_j})$  for all  $j$ , so  $w \in \text{cl}(W_{\mathcal{M}_j})$ .

We show that  $w \in \bigcap_i W_{\mathcal{M}_i}(p)$ . Suppose not. Then, there must be some  $i$  such that  $w \notin W_{\mathcal{M}_i}(p)$ . Since  $\text{cl}(W_{\mathcal{M}_{i+1}}) \setminus \{p\} \subset W_{\mathcal{M}_i}(p)$ , and since  $w \neq p, w \notin \text{cl}(W_{\mathcal{M}_{i+1}}(p))$ —contradiction.

Finally, we show that  $w \in D(p)$ . Suppose  $w \notin W_{\mathcal{M}^*}$  where  $\mathcal{M}^* \in \text{supp}(\nu) \setminus \mathcal{U}$ . We can find  $\mathcal{M}_i$  such that

$$\nu(\{\mathcal{M} : \mathcal{M} \supset \mathcal{M}^*\}) > \nu(\{\mathcal{M} : \mathcal{M} \supset \mathcal{M}_i\}).$$

That is,  $\mathcal{M}_i \supset \mathcal{M}^*$ . Since  $w \notin W_{\mathcal{M}}(p), w \notin W_{\mathcal{M}_i}(p)$ —contradiction.

2. Since  $\nu(q \in W(A)) = \max_{x \in \text{co}(A)} \nu(q \in W(x))$  by Lemma 20, it suffices to show that  $\nu(d \in W(p)) \geq \nu(q \in W(r))$  for any  $q, r \in \Delta(Z)$ . Suppose  $\nu(d \in W(p)) < \nu(q \in W(r))$ . By Independence,  $\nu(p + \varepsilon(q - r) \in W(p)) > \nu(d \in W(p))$  for any  $\varepsilon$  small enough that  $p + \varepsilon(q - r) \in W(p)$ . But then  $\nu(d \in W(p)) < \max_{x \in \Delta(Z)} \nu(x \in W(p))$ —contradiction.
3. For the first part, suppose  $\nu(q \in W(\lambda p + (1 - \lambda)d)) = \nu(d \in W(p))$ . Since  $\nu(d \in W(p)) = \nu(\lambda p + (1 - \lambda)d \in W(p))$ , we have  $\nu(q \in W(p)) \geq \nu(d \in W(p))$  by the first part of Lemma 18. Since  $d \in D(p)$ , we have  $q \in D(p)$ , which contradicts the assumption about  $q$ . Conclude that  $\nu(q \in W(\lambda p + (1 - \lambda)d)) < \nu(\lambda p + (1 - \lambda)d \in W(p))$ . By the second part of Lemma 18,  $\nu(q \in W(p)) > \nu(q \in W(\lambda p + (1 - \lambda)d))$ .

For the second part, suppose  $\nu(\lambda b + (1 - \lambda)q \in W(p)) = \nu(q \in W(b))$ . Since  $\nu(q \in W(b)) = \nu(q \in W(\lambda b + (1 - \lambda)q))$ , we have  $\nu(q \in W(p)) \geq \nu(q \in W(b))$  by the first part of Lemma 18. Since  $q \in W(b)$ , we have  $q \in D(p)$ , which contradicts the assumption about  $q$ . Conclude that  $\nu(\lambda b + (1 - \lambda)q \in W(p)) < \nu(q \in W(\lambda b + (1 - \lambda)q))$ . By the second part of Lemma 18,  $\nu(q \in W(p)) > \nu(q \in W(\lambda b + (1 - \lambda)q))$ .

□

**Lemma 2.3.** For any  $p \in \Delta(Z)$  and any  $A, B \in \mathcal{F}(\Delta(Z))$ : if  $B_\varepsilon(p) \cap \text{co}(A \cup \{p\}) = B_\varepsilon(p) \cap \text{co}(B \cup \{p\})$  for some  $\varepsilon > 0$ , then  $v(p \in W(A)) = v(p \in W(B))$ .

*Proof.* By Lemma 2.0, there exists  $a \in \text{co}(A)$  such that  $p \in W(A)$  if and only if  $p \in W(a)$ . For any  $\lambda \in [0, 1)$ ,  $p \in W(a)$  if and only if  $p \in W(\lambda p + (1 - \lambda)a)$ . Since

$$v(p \in W(\lambda p + (1 - \lambda)a)) = v(p \in W(A)) = \max_{\tilde{a} \in \text{co}(A \cup \{p\})} v(p \in W(\tilde{a})),$$

$\lambda p + (1 - \lambda)a \in \arg \max_{\tilde{a} \in \text{co}(A \cup \{p\})} v(p \in W(\tilde{a}))$ . Notice that  $\lambda p + (1 - \lambda)a \in \bar{B}_\varepsilon(p)$  for  $\lambda$  small enough.

Now take any menu  $B$  such that  $\bar{B}_\varepsilon(p) \cap \text{co}(A \cup \{p\}) = \bar{B}_\varepsilon(p) \cap \text{co}(B \cup \{p\})$ . Clearly,  $\lambda p + (1 - \lambda)a \in B_\varepsilon(p) \cap \text{co}(B \cup \{p\})$ . We have

$$\begin{aligned} v(p \in W(B)) &= \max_{\tilde{b} \in \text{co}(B \cup \{p\})} v(p \in W(\tilde{b})) \\ &\geq \max_{\tilde{b} \in \bar{B}_\varepsilon(p) \cap \text{co}(B \cup \{p\})} v(p \in W(\tilde{b})) \\ &= \max_{\tilde{a} \in \bar{B}_\varepsilon(p) \cap \text{co}(A \cup \{p\})} v(p \in W(\tilde{a})) \\ &= v(p \in W(a)) \\ &= v(p \in W(A)). \end{aligned}$$

Now suppose  $v(p \in W(B)) > v(p \in W(A))$ , so  $v(p \in W(b)) > v(p \in W(a))$  for some  $b \in B$ . For any  $\lambda \in [0, 1)$ ,  $p \in W(b)$  if and only if  $p \in W(\lambda p + (1 - \lambda)b)$ . For  $\lambda$  small enough,  $\lambda p + (1 - \lambda)b \in B_\varepsilon(p) \cap \text{co}(B \cup \{p\})$ , so  $\lambda p + (1 - \lambda)b \in A$ . But then

$$v(p \in W(a)) < v(p \in W(\lambda p + (1 - \lambda)b)) \leq \max_{\tilde{a} \in \text{co}(A)} v(p \in W(\tilde{a})) = v(p \in W(A)),$$

which contradicts the definition of  $A$ . Conclude that  $\nu(p \in W(B)) = \nu(p \in W(A))$ .  $\square$

**Lemma 24.** *For any  $q, r \in \Delta(Z)$  and any  $A \in \mathcal{F}(\Delta(Z))$ : if  $\nu(A \subset W(r)) > 0$  and  $\nu(r \in W(q)) = 0$ , then  $\nu(r \in W(A \cup \{q\})) = 0$ .*

*Proof.* Suppose  $\nu(r \in W(A \cup \{q\})) > 0$ . By Lemma 20, there exists  $x \in \text{co}(A \cup \{q\})$  such that  $r \in W(x)$  whenever  $r \in W(A \cup \{q\})$ . Suppose

$$x = \lambda q + (1 - \lambda)a$$

for some  $a \in \text{co}(A)$  and  $\lambda \in (0, 1]$ . Since  $\nu(r \in W(q)) = 0$ , we have

$$\nu(\{\mathcal{M} : \exists \succ \in \mathcal{M} \text{ s.t. } r \succ q\}) = 1.$$

Since  $\nu(A \subset W(r)) > 0$ , we also have

$$\nu(\{\mathcal{M} : \forall \succ \in \mathcal{M} \ r \succ A\}) > 0.$$

Any preference in  $\mathcal{U}$  that has  $r \succ q$  and  $r \succ A$  must have  $r \succ x$ , so

$$\nu(\{\mathcal{M} : \exists \succ \in \mathcal{M} \text{ s.t. } r \succ x\}) > 0.$$

This says  $\nu(x \in W(r)) > 0$ , so  $\nu(r \in W(x)) = 0$ —contradiction.  $\square$

### A.1.13 PROOF OF PROPOSITION 5

#### FIRST PART

Fix  $p, q \in \text{int}(\Delta(Z))$  such that  $\nu(p \in W(q)) > 0$  but  $p \notin D(q)$ . We restrict attention to  $\varepsilon$  small enough that  $p + \varepsilon(x - y), q + \varepsilon(x - y) \in \Delta(Z)$  for all  $x, y \in \Delta(Z)$ .

By Lemma 22,  $D(q) := \arg \max_{x \in \Delta(Z)} \nu(x \in W(q))$  is nonempty. Since

$$d \in D(q) \implies \lambda d + (1 - \lambda)q \in D(q),$$

the set contains points arbitrarily close to  $q$ .

By Independence,  $\nu(q \in W(q + \lambda(q - p))) = \nu(p \in W(q)) > 0$  for all  $\lambda > 0$ . Fix  $\lambda$  such that  $q + \lambda(q - p) \in B_{\varepsilon/2}(q)$ . For any  $x$  close enough to  $q$ , we have  $x + \lambda(x - p) \in B_{\varepsilon}(q)$ ,  $\nu(q \in W(x + \lambda(x - p))) > 0$ , and  $\nu(p \in W(x)) > 0$ . Choose some  $x \in D(q)$  that satisfies these requirements, and call it  $x^*$ .

Suppose

$$\nu(p \in W(\alpha q + (1 - \alpha)(x^* + \lambda(x^* - p)))) \geq \nu(p \in W(q))$$

for some  $\alpha \in [0, 1)$ . By Independence,

$$\nu\left(p \in W\left(\frac{\alpha}{1 + \lambda(1 - \alpha)}q + \frac{(1 - \alpha)(1 + \lambda)}{1 + \lambda(1 - \alpha)}x^*\right)\right) \geq \nu(p \in W(q)).$$

But Lemma 20 and the third part of Lemma 22 imply

$$\nu(p \in W(q)) = \nu(p \in W(\{q, x^*\})) = \max_{x \in \text{co}(\{q, x^*\})} \nu(p \in W(x)),$$

so we have a contradiction. Conclude that

$$\nu(p \in W(q)) = \max_{x \in \text{co}(\{q, x^* + \lambda(x^* - p)\})} \nu(p \in W(x)) = \nu(p \in W(\{q, x^* + \lambda(x^* - p)\})).$$

Suppose  $\nu(q \in W(\{p, x^* + \lambda(x^* - p)\})) \geq \nu(p \in W(\{q, x^* + \lambda(x^* - p)\}))$ . By Lemma 21,

$$\nu(p \in W(\{q, x^* + \lambda(x^* - p)\})) = \nu(p \in W(x^* + \lambda(x^* - p))).$$

But we have just shown that

$$\nu(p \in W(\{q, x^* + \lambda(x^* - p)\})) = \nu(p \in W(q)) > \nu(p \in W(x^* + \lambda(x^* - p))).$$

Conclude that

$$\nu(q \in W(\{p, x^* + \lambda(x^* - p)\})) < \nu(p \in W(\{q, x^* + \lambda(x^* - p)\})).$$

Notice that  $\nu(p \in W(x^* + \lambda(x^* - p))) = \nu(p \in W(x^*)) > 0$ , and recall that  $\nu(q \in W(x^* + \lambda(x^* - p))) > 0$ . By Lemma 24,  $\nu(x^* + \lambda(x^* - p) \in W(\{p, q\})) = 0$ . We are now ready to compute  $\rho(q|\{p, q, x^* + \lambda(x^* - p)\})$ . To simplify notation, let

$$\tilde{q} := x^* + \lambda(x^* - p).$$

Two groups of DMs may choose  $q$ : those who have  $q \succ p, \tilde{q}$ , and those who have  $p \succ q \succ \tilde{q}$ , but



$p \in W(q)$ . We have

$$\begin{aligned} \rho(q|\{p, q, x^* + \lambda(x^* - p)\}) &= \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})) \\ &\quad + \mu(p \succ q \succ \tilde{q}) [\nu(q \notin W(\tilde{q})) - \nu(p \notin W(q))]. \end{aligned}$$

Since  $\nu(\{p, q\} \subset W(\tilde{q})) > 0$ , we can find  $\alpha \in (0, 1)$  such that  $\nu(\{p, q\} \subset W(\alpha\tilde{q} + (1-\alpha)x^*)) > 0$ . Since  $\nu(p \in W(x^*)) > 0$ , and since  $\alpha\tilde{q} + (1-\alpha)x^*$  is a combination of  $p$  and  $x^*$ , we have  $\nu(\alpha\tilde{q} + (1-\alpha)x^* \in W(\tilde{q})) > 0$ . By the second part of Lemma 19, we can find  $\beta \in (0, 1)$  such that

$$\nu(\beta\tilde{q} + (1-\beta)q \in W(\alpha\tilde{q} + (1-\alpha)x^*)) = \nu(\alpha\tilde{q} + (1-\alpha)x^* \in W(\beta\tilde{q} + (1-\beta)q)) = 0.$$

By Lemma 24,

$$\nu(\alpha\tilde{q} + (1-\alpha)x^* \in W(\{p, q, \beta\tilde{q} + (1-\beta)q\})) = 0.$$

Since  $\{y \in \Delta(Z) : \nu(p \in W(y)) > 0\}$  is convex by the first part of Lemma 19 and since  $\min\{\nu(p \in W(q)), \nu(p \in W(\tilde{q}))\} > 0, \nu(p \in W(\beta\tilde{q} + (1-\beta)q)) > 0$ . Since  $\nu(q \in W(\tilde{q})) > 0, \nu(q \in W(\beta\tilde{q} + (1-\beta)q)) > 0$  as well. By Lemma 24,

$$\nu(\beta\tilde{q} + (1-\beta)q \in W(\{p, q, \alpha\tilde{q} + (1-\alpha)x^*\})) = 0.$$

By Lemma 23,

$$\begin{aligned} \nu(p \in W(\{q, \alpha\tilde{q} + (1-\alpha)x^*, \beta\tilde{q} + (1-\beta)q\})) &= \nu(p \in W(\{q, \tilde{q}\})) = \nu(p \in W(q)) \\ \nu(q \in W(\{p, \alpha\tilde{q} + (1-\alpha)x^*, \beta\tilde{q} + (1-\beta)q\})) &= \nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(\tilde{q})). \end{aligned}$$

We are now ready to compute  $\rho(q|\{p, q, \alpha\tilde{q} + (1-\alpha)x^*, \beta\tilde{q} + (1-\beta)q\})$ . Two groups of DMs may

choose  $q$  from this menu: those who have  $q \succ p, \tilde{q}$ , and those who have  $p \succ q \succ \alpha\tilde{q} + (1 - \alpha)x^*$ , but  $p \in W(q)$ . We have

$$\begin{aligned} \rho(q|\alpha\tilde{q} + (1 - \alpha)x^*, \beta\tilde{q} + (1 - \beta)q) &= \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})) + \\ &\mu(p \succ q \succ \alpha\tilde{q} + (1 - \alpha)x^*) [\nu(q \notin W(\tilde{q})) - \nu(p \notin W(q))]. \end{aligned}$$

To show that

$$\rho(q|\alpha\tilde{q} + (1 - \alpha)x^*, \beta\tilde{q} + (1 - \beta)q) < \rho(q|\{p, q, \tilde{q}\}),$$

it suffices to confirm that

$$\mu(p \succ q \succ \alpha\tilde{q} + (1 - \alpha)x^*) < \mu(p \succ q \succ \tilde{q}).$$

Since  $x^*$  is a combination of  $p$  and  $\tilde{q}$ ,  $p \succ q \succ \alpha\tilde{q} + (1 - \alpha)x^*$  implies  $p \succ q \succ \tilde{q}$ . The converse does not hold: a DM who likes  $q$  slightly more than  $\tilde{q}$  and  $p$  a lot more than  $q$  will have  $\tilde{q} + (1 - \alpha)x^* \succ q$ . Since  $\mu$  has full support, the inequality holds.

Now we cover the case  $p \in D(q)$ . The arguments are very similar, but the construction is slightly different. Take any  $\tilde{q} \in B_\varepsilon(q)$  such that  $\nu(q \in W(\tilde{q})) > 0$  but  $q \notin D(\tilde{q})$ . As above, we have

$$\nu(p \in W(q)) = \nu(p \in W(\{q, \tilde{q}\})) > \nu(q \in W(\tilde{q})) = \nu(q \in W(\{p, \tilde{q}\})) > \nu(\tilde{q} \in W(\{p, q\})) = 0,$$

so

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})) + \mu(p \succ q \succ \tilde{q})[\nu(q \notin W(\tilde{q})) - \nu(p \notin W(q))].$$

As in the first case, we can find  $\alpha \in (0, 1)$  such that  $\nu(\{p, q\} \subset W(\alpha\tilde{q} + (1 - \alpha)p)) > 0$  and

$\beta \in (0, 1)$  such that

$$\nu(\beta\tilde{q} + (1 - \beta)q \in W(\alpha\tilde{q} + (1 - \alpha)p)) = \nu(\alpha\tilde{q} + (1 - \alpha)p \in W(\beta\tilde{q} + (1 - \beta)q)) = 0.$$

As in the first case, we have

$$\nu(\alpha\tilde{q} + (1 - \alpha)p \in W(\{p, q, \beta\tilde{q} + (1 - \beta)p\})) = 0$$

$$\nu(\beta\tilde{q} + (1 - \beta)q \in W(\{p, q, \alpha\tilde{q} + (1 - \beta)q\})) = 0$$

$$\nu(q \in W(\{p, \alpha\tilde{q} + (1 - \alpha)p, \beta\tilde{q} + (1 - \beta)p\})) = \nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(\tilde{q}))$$

$$\nu(p \in W(\{q, \alpha\tilde{q} + (1 - \alpha)p, \beta\tilde{q} + (1 - \beta)p\})) = \nu(p \in W(\{q, \tilde{q}\})) = \nu(p \in W(q)).$$

This delivers

$$\begin{aligned} \rho(q|\alpha\tilde{q} + (1 - \alpha)p, \beta\tilde{q} + (1 - \beta)q) &= \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})) + \\ &\quad \mu(p \succ q \succ \alpha\tilde{q} + (1 - \alpha)p) [\nu(q \notin W(\tilde{q})) - \nu(p \notin W(q))]. \end{aligned}$$

Since  $\mu(p \succ q \succ \alpha\tilde{q} + (1 - \alpha)p) < \mu(p \succ q \succ \tilde{q})$ , there is a regularity violation as above.

## SECOND PART

Fix any  $p, q \in \text{int}(\Delta(Z))$  such that  $\nu(p \in W(q)) = 0$ .

Suppose there exists  $\bar{\varepsilon} > 0$  such that  $\nu(p \in W(\tilde{q})) = 0$  for all  $\tilde{q} \in B_{\bar{\varepsilon}}(q)$ . We show that  $(p, q)$  cannot be anomalous for  $\varepsilon$  sufficiently small.

For any plane  $P \in \Delta(Z)$  and any  $x \in P$ , let

$$D_P(x) := \max_{y \in P} \nu(y \in W(x)).$$

If we restrict attention to items in  $P$ , Lemma 22 will still apply with  $D_P$  in place of  $D$ .

Fix any plane  $P$  containing  $p$  and  $q$ . To begin, suppose  $q \in D_P(p)$ . By the second part of Lemma 22,

$$\nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(p)) \geq \max\{\nu(\tilde{q} \in W(\{p, q\})), \nu(p \in W(\{q, \tilde{q}\}))\}.$$

Thus,

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(p)).$$

Similarly, for any  $\beta_1, \beta_2, \delta_1, \delta_2 > 0$  such that  $\beta_1 + \beta_2 = \delta_1 + \delta_2 = 0$ , we have

$$\begin{aligned} \nu(q \in W(p)) &\geq \max\{\nu(p \in W(\{\beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p, q\})), \\ &\quad \nu(\beta_1\tilde{q} + \beta_2q \in W(\{\delta_1\tilde{q} + \delta_2p, p, q\})), \\ &\quad \nu(\delta_1\tilde{q} + \delta_2p \in W(\{\{\beta_1\tilde{q} + \beta_2q, p, q\}\}))\}, \end{aligned}$$

so

$$\begin{aligned} \rho(q|\{p, q, \beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p\}) &= \mu(q \succ p, \beta_1\tilde{q} + \beta_2q)\nu(q \notin W(p)) \\ &= \mu(q \succ p, \tilde{q})\nu(q \notin W(p)) \\ &= \rho(q|\{p, q, \tilde{q}\}). \end{aligned}$$

From now on, we assume  $q \notin D_P(p)$ . For  $\varepsilon$  small enough, any  $\tilde{q} \in B_\varepsilon(q) \cap P$  can be written in one of two ways. The first is

$$\tilde{q} = \alpha_1q + \alpha_2b + (1 - \alpha_1 - \alpha_2)p$$

where  $q \in D_p(b)$  and  $\alpha_1, \alpha_2 \geq 0$  (but  $\alpha_1 + \alpha_2$  may exceed 1).

For  $\varepsilon$  small enough,  $b$  can be taken to be sufficiently close to  $q$  that

$$\nu(p \in W(x)) = 0$$

for all  $x \in \text{co}(\{b, q\})$ . By Independence, the same is true for all  $x \in \text{co}(\{b, q, p\})$ . Since  $\text{co}(q, \tilde{q}) \subset \text{co}(\{b, q, p\})$ , we have

$$\nu(p \in W(\{q, \tilde{q}\})) = 0$$

by Lemma 20.

We show that  $\nu(\tilde{q} \in W(\{p, q\})) \geq \nu(q \in W(\{p, \tilde{q}\}))$ .

$$\begin{aligned} \tilde{q} \in W(\{p, q\}) &= \max_{\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1} \nu(\alpha_1 q + \alpha_2 b + \alpha_3 p \in W(\lambda_1 q + \lambda_2 p)) \\ &= \max \left\{ \max_{\lambda_1 \geq \alpha_1, \lambda_2 \geq \alpha_3} \nu \left( b \in W \left( \frac{\lambda_1 - \alpha_1}{\alpha_2} q + \frac{\lambda_2 - \alpha_3}{\alpha_2} p \right) \right), \right. \\ &\quad \max_{\lambda_1 > \alpha_1, \lambda_2 \leq \alpha_3} \nu \left( \frac{\alpha_2}{\lambda_1 - \alpha_1} b + \frac{\alpha_3 - \lambda_2}{\lambda_1 - \alpha_1} p \in W(q) \right) \\ &\quad \left. \max_{\lambda_1 \leq \alpha_1, \lambda_2 > \alpha_3} \nu \left( \frac{\alpha_1 - \lambda_1}{\lambda_2 - \alpha_3} q + \frac{\alpha_2}{\lambda_2 - \alpha_3} b \in W(p) \right) \right\} \end{aligned}$$

We work through the terms in the max. First, suppose that  $\nu(b \in W(x)) > \nu(q \in W(p))$  for some  $x \in \text{co}(\{p, q\})$ . Since  $b$  was chosen so that  $\nu(q \in W(b)) \geq \nu(b \in W(x))$ , the first part of Lemma 18 implies  $\nu(q \in W(x)) \geq \nu(b \in W(x)) > \nu(q \in W(p))$ . But since  $x \in \text{co}(\{p, q\})$ , Independence implies  $\nu(q \in W(x)) = \nu(q \in W(p))$ —contradiction. Conclude that  $\nu(b \in W(x)) \leq \nu(q \in W(p))$  for all  $x \in \text{co}(\{p, q\})$ , so the first term cannot exceed  $\nu(q \in W(p))$ . Now consider the second term. Since  $\nu(p \in W(q)) = 0$  and  $\nu(q \in W(b)) > 0$ , Lemma 24 gives  $\nu(x \in W(q)) = 0$  for all  $x \in \text{co}(\{b, p\})$ . Thus, the second term is 0. Finally, consider the third term. Suppose that  $\nu(x \in W(p)) > \nu(q \in W(p))$  for some  $x \in \text{co}(\{b, q\})$ . Independence implies

$\nu(q \in W(x)) = \nu(q \in W(b))$ , and  $b$  was chosen so that  $\nu(q \in W(b)) \geq \nu(x \in W(p))$ . We have  $\nu(q \in W(x)) \geq \nu(x \in W(p))$ . By the first part of Lemma 18,  $\nu(q \in W(p)) \geq \nu(x \in W(p)) > \nu(q \in W(p))$ —contradiction. Conclude that  $\nu(x \in W(p)) \leq \nu(q \in W(p))$  for all  $x \in \text{co}(\{b, q\})$ . Thus, the third term cannot exceed  $\nu(q \in W(p))$ . Putting all three terms together, we have

$$\nu(\tilde{q} \in W(\{p, q\})) \leq \nu(q \in W(p)) \leq \nu(q \in W(\{p, \tilde{q}\})).$$

Since

$$\nu(q \in W(\{p, q\})) \geq \max\{\nu(\tilde{q} \in W(\{p, q\})), \nu(p \in W(\{q, \tilde{q}\}))\},$$

we have

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(\{p, \tilde{q}\})).$$

Now we show that

$$\nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\})) \leq \nu(q \in W(\{p, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})) \quad (\text{A.10})$$

for any  $\beta_1, \beta_2, \lambda_1, \lambda_2 > 0$  and  $\beta_1 + \beta_2 = \lambda_1 + \lambda_2 = 1$ . By Lemmas 20 and 23,

$$\begin{aligned} \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\})) &= \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{q, \delta_1 \tilde{q} + \delta_2 p\})) \\ &= \max_{\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1} \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\lambda_1 q + \lambda_2(\delta_1 \tilde{q} + \delta_2 p))) \\ &= \max \left\{ \max_{\lambda_1 \leq \beta_2, \lambda_2 \leq \beta_1/\delta_1} \nu \left( \frac{\beta_1 - \lambda_2 \delta_1}{\delta_2 \lambda_2} \tilde{q} + \frac{\beta_2 - \lambda_1}{\delta_2 \lambda_2} q \in W(p) \right), \right. \\ &\quad \max_{\lambda_1 < \beta_2, \lambda_2 > \beta_1/\delta_1} \nu \left( q \in W \left( \frac{\lambda_2 \delta_1 - \beta_1}{\beta_2 - \lambda_1} \tilde{q} + \frac{\lambda_2 \delta_2}{\beta_2 - \lambda_1} p \right) \right), \\ &\quad \left. \max_{\lambda_1 > \beta_2, \lambda_2 < \beta_1/\delta_1} \nu \left( \tilde{q} \in W \left( \frac{\lambda_1 - \beta_2}{\beta_1 - \lambda_2 \delta_1} q + \frac{\lambda_2 \delta_2}{\beta_1 - \lambda_2 \delta_1} p \right) \right) \right\} \end{aligned}$$

We work through the terms in the max. Expanding  $\tilde{q}$  and rearranging, and using the fact (from the third part of Lemma 22) that

$$\nu(q \in W(p)) = \max_{\lambda \in [0,1]} \nu(\lambda q + (1 - \lambda)b \in W(p)).$$

we can rewrite the first term as  $\nu(q \in W(p))$ . By Lemma 20, the second term is no greater than  $\nu(q \in W(\{p, \tilde{q}\}))$ , which is equal to

$$\nu(q \in W(\{p, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\}))$$

by Lemma 23. Similarly, the third term is no greater than  $\nu(\tilde{q} \in W(\{p, q\}))$ , which we already showed was less than  $\nu(q \in W(\{p, \tilde{q}\}))$ . Combining all three terms, we have (A.10) as desired.

By Lemma 23,

$$\nu(p \in W(\{q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})) = \nu(p \in W(\{q, \tilde{q}\})) = 0.$$

Since

$$\begin{aligned} \nu(q \in W(\{p, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})) &\geq \max\{\nu(p \in W(\{q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})), \\ &\quad \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\}))\}, \end{aligned}$$

the only DMs who would choose  $q$  from  $\{p, q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\}$  are the ones who like  $q$  best.

We have

$$\begin{aligned} \rho(q|\{p, q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\}) &= \mu(q \succ \tilde{q}, p) \nu(q \notin W(\{p, \tilde{q}\})) \\ &= \rho(q|\{p, q, \tilde{q}\}). \end{aligned}$$

This completes the first possibility for  $\tilde{q}$ . The second possibility is

$$\tilde{q} = \alpha_1 q + \alpha_2 w + (1 - \alpha_1 - \alpha_2)p$$

where  $w \in D_P(q)$  and  $\alpha_1, \alpha_2 \geq 0$  (but  $\alpha_1 + \alpha_2$  may exceed 1). We show that  $\nu(p \in W(x)) = 0$  for all  $x \in \text{co}(\{q, \tilde{q}\})$ . Suppose not. By Independence,  $\nu(p \in W(y)) > 0$  for some  $y \in \text{co}(\{q, w\}) \setminus \{q\}$ . Since  $\nu(y \in W(q)) > 0$ ,  $\nu(p \in W(q)) > 0$  by the first part of Lemma 18—contradiction.

For  $\tilde{q}$  sufficiently close to  $q$ , there are two subcases. The first is  $\tilde{q} \in D_P(y)$  for some  $y \in \text{co}(\{p, q\})$ . By Independence, this implies

$$\begin{aligned} \exists y \in \text{co}(\{p, q\}) \quad \beta_1 \tilde{q} + \beta_2 q &\in D_P(y) \\ \exists y \in \text{co}(\{p, q\}) \quad \delta_1 \tilde{q} + \delta_2 p &\in D_P(y) \end{aligned}$$

for all  $\beta_1, \beta_2, \delta_1, \delta_2 > 0$  such that  $\beta_1 + \beta_2 = \delta_1 + \delta_2 = 1$ . We have

$$\begin{aligned} \nu(\tilde{q} \in W(\{p, q\})) &= \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\})) \\ &= \nu(\delta_1 \tilde{q} + \delta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\})). \end{aligned}$$

We also have

$$\nu(\tilde{q} \in W(\{p, q\})) \geq \nu(q \in W(\{p, \tilde{q}\})),$$

so

$$\nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(p))$$

by Lemma 21. There are two groups of DMs that may choose  $q$  from  $\{p, q, \tilde{q}\}$ : those who like  $q$



best, and those who have  $\tilde{q} \succ q \succ p$ , but  $\tilde{q} \in W(\{p, q\})$ . We have

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(p)) + \mu(\tilde{q} \succ q \succ p)[\nu(q \notin W(p)) - \nu(\tilde{q} \notin W(\{p, q\}))].$$

Similarly,

$$\nu(\{\beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p\} \subset W(\{p, q\})) \geq \nu(q \in W(\{p, \beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p\})),$$

so Lemma 2.1 implies

$$\nu(q \in W(\{p, \beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p\})) = \nu(q \in W(p)).$$

We have

$$\begin{aligned} \rho(q|\{p, q, \beta_1\tilde{q} + \beta_2q, \delta_1\tilde{q} + \delta_2p\}) &= \mu(q \succ p, \tilde{q})\nu(q \notin W(p)) \\ &\quad + \mu(\beta_1\tilde{q} + \beta_2q \succ q \succ p) \\ &\quad \times [\nu(q \notin W(p)) - \nu(\beta_1\tilde{q} + \beta_2q \notin W(\{p, q\}))] \\ &= \mu(q \succ p, \tilde{q})\nu(q \notin W(p)) \\ &\quad + \mu(\tilde{q} \succ q \succ p)[\nu(q \notin W(p)) - \nu(\tilde{q} \notin W(\{p, q\}))] \\ &= \rho(q|\{p, q, \tilde{q}\}). \end{aligned}$$

The second subcase is

$$q \in \arg \max_{x \in \text{co}(\{p, q\})} \nu(\tilde{q} \in W(x)).$$

We show that

$$\nu(q \in W(\{\tilde{q}, p\})) = \nu(q \in W(p)).$$

By Lemma 20,

$$\begin{aligned} \nu(q \in W(\{p, \tilde{q}\})) &= \max_{\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1} \nu(q \in W(\lambda_1 \tilde{q} + \lambda_2 p)) \\ &= \max \left\{ \max_{\lambda_1 \alpha_1 < 1, \lambda_1 \alpha_3 + \lambda_2 \geq 0} \nu \left( \frac{\lambda_1 \alpha_2}{1 - \lambda_1 \alpha_1} w + \frac{\lambda_1 \alpha_3 + \lambda_2}{1 - \lambda_1 \alpha_1} p \right), \right. \\ &\quad \left. \max_{\lambda_1 \alpha_1 \geq 1, \lambda_1 \alpha_3 + \lambda_2 < 0} \nu \left( p \in W \left( \frac{\lambda_1 \alpha_1 - 1}{-(\lambda_1 \alpha_3 + \lambda_2)} q + \frac{\lambda_1 \alpha_2}{-(\lambda_1 \alpha_3 + \lambda_2)} w \right) \right) \right\} \end{aligned}$$

For the first term, recall that  $\nu(w \in W(\{p, q\})) = \nu(w \in W(q)) \geq \nu(q \in W(\{w, p\}))$ . By Lemma 21,  $\nu(q \in W(\{w, p\})) = \nu(q \in W(p))$ . By the third part of Lemma 22,

$$\nu(q \in W(p)) = \max_{x \in \text{co}(\{w, p\})} \nu(q \in W(x))$$

so the first term is  $\nu(q \in W(p))$ . The second term is 0 because  $\nu(p \in W(x)) = 0$  for all  $x \in \text{co}(\{q, \tilde{q}\})$ .

Now we show that  $\nu(q \in W(p)) \leq \nu(\tilde{q} \in W(\{p, q\}))$ . Suppose not. By Lemma 21,  $\nu(\tilde{q} \in W(\{p, q\})) = \nu(\tilde{q} \in W(p))$ . Expanding  $\tilde{q}$ , we have

$$\nu(\tilde{q} \in W(p)) = \nu \left( \frac{\alpha_1}{\alpha_1 + \alpha_2} q + \frac{\alpha_2}{\alpha_1 + \alpha_2} w \in W(p) \right).$$

Since  $\text{co}(\{q, w\}) \subset W(p)$  whenever  $q \in W(p)$ , we have

$$\begin{aligned} \nu(q \in W(p)) &\leq \nu\left(\frac{\alpha_1}{\alpha_1 + \alpha_2}q + \frac{\alpha_2}{\alpha_1 + \alpha_2}w \in W(p)\right) \\ &= \nu(\tilde{q} \in W(p)) \\ &= \nu(\tilde{q} \in W(\{p, q\})). \end{aligned}$$

This contradicts the assumption that  $\nu(q \in W(p)) > \nu(\tilde{q} \in W(\{p, q\}))$ .

We can now compute  $\rho(q|\{p, q, \tilde{q}\})$ . We have

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(p)) + \mu(\tilde{q} \succ q \succ p)[\nu(q \notin W(p)) - \nu(\tilde{q} \notin W(\{p, q\}))]$$

exactly as in the first subcase.

We show that

$$\nu(\beta_1\tilde{q} + \beta_2q \in W(\{p, q, \delta_1\tilde{q} + \delta_2p\})) = \nu(\tilde{q} \in W(q)). \quad (\text{A.11})$$

By Lemma 23,

$$\nu(\beta_1\tilde{q} + \beta_2q \in W(\{p, q, \delta_1\tilde{q} + \delta_2p\})) = \nu(\beta_1\tilde{q} + \beta_2q \in W(\{q, \delta_1\tilde{q} + \delta_2p\})).$$

By Lemma 20,

$$\begin{aligned}
\nu(\beta_1 \tilde{q} + \beta_2 q \in W(\{q, \delta_1 \tilde{q} + \delta_2 p\})) &= \max_{\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1} \nu(\beta_1 \tilde{q} + \beta_2 q \in W(\lambda_1(\delta_1 \tilde{q} + \delta_2 p) + \lambda_2 q)) \\
&= \max \left\{ \max_{\beta_1 \geq \lambda_1 \delta_1, \beta_2 \geq \lambda_2} \nu \left( \frac{\beta_1 - \lambda_1 \delta_1}{\lambda_1 \delta_2} \tilde{q} + \frac{\beta_2 - \lambda_2}{\lambda_1 \delta_2} q \in W(p) \right), \right. \\
&\quad \max_{\beta_1 < \lambda_1 \delta_1, \beta_2 > \lambda_2} \nu \left( q \in W \left( \frac{\lambda_1 \delta_1 - \beta_1}{\beta_2 - \lambda_2} \tilde{q} + \frac{\lambda_1 \delta_2}{\beta_2 - \lambda_2} p \right) \right), \\
&\quad \left. \max_{\beta_1 > \lambda_1 \delta_1, \beta_2 < \lambda_2} \nu \left( \tilde{q} \in W \left( \frac{\lambda_2 - \beta_2}{\beta_1 - \lambda_1 \delta_1} q + \frac{\lambda_1 \delta_2}{\beta_1 - \lambda_1 \delta_1} p \right) \right) \right\}
\end{aligned}$$

We work through the terms in the max. For the first term, suppose that  $\nu(x \in W(p)) > \nu(\tilde{q} \in W(q))$  for some  $x \in \text{co}(\{\tilde{q}, q\}) \setminus \{\tilde{q}\}$ . Since  $\{y \in \Delta(Z) : y \in W(p)\}$  is a cone, it cannot be that  $\nu(x \in W(p)) > 0$  and  $\nu(\tilde{q} \in W(p)) = 0$  for  $\varepsilon$  small. Thus,  $\nu(\tilde{q} \in W(p)) > 0$ . By Independence,  $\nu(\tilde{q} \in W(q)) = \nu(\tilde{q} \in W(x))$ , so  $\nu(x \in W(p)) > \nu(\tilde{q} \in W(x))$ . By the second part of Lemma 18,  $\nu(\tilde{q} \in W(p)) > \nu(\tilde{q} \in W(x)) = \nu(\tilde{q} \in W(q))$ . This contradicts the assumption that  $\nu(\tilde{q} \in W(q)) = \nu(\tilde{q} \text{ in } W(\{p, q\}))$ . Conclude that the first term cannot exceed  $\nu(\tilde{q} \in W(q))$ . Now consider the second term. By Lemma 20,  $\nu(q \in W(x)) \leq \nu(q \in W(\{p, \tilde{q}\}))$  for all  $x \in \text{co}(\{p, \tilde{q}\})$ . We already showed that  $\nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(p)) \leq \nu(\tilde{q} \in W(q))$ . Finally, consider the third term. By Lemma 20,  $\nu(\tilde{q} \in W(x)) \leq \nu(\tilde{q} \in W(\{p, q\}))$  for all  $x \in \text{co}(\{p, q\})$ . Since  $\nu(\tilde{q} \in W(\{p, q\})) = \nu(\tilde{q} \in W(q))$ , the third term cannot exceed  $\nu(\tilde{q} \in W(q))$ . This maximum can be achieved by setting  $\lambda_1 = 0$ , so the third term equals  $\nu(\tilde{q} \in W(q))$ . Conclude that (A.11) holds.

By Lemma 20,

$$\nu(\delta_1 \tilde{q} + \delta_2 p \in W(\{p, q, \beta_1 \tilde{q} + \beta_2 q\})) \geq \max_{x \in \text{co}(\{p, q\})} \nu(\delta_1 \tilde{q} + \delta_2 p \in W(x)).$$

By Independence,

$$\nu(\delta_1 \tilde{q} + \delta_2 p \in W(\delta_1 q + \delta_2 p)) = \nu(\tilde{q} \in W(q)).$$

Conclude that

$$\nu(\delta_1 \tilde{q} + \delta_2 p \in W(\{p, q, \beta_1 \tilde{q} + \beta_2 q\})) \geq \nu(\tilde{q} \in W(q)).$$

As in the first subcase, we have

$$\nu(q \in W(\{p, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})) = \nu(q \in W(p, \tilde{q})) = \nu(q \in W(p))$$

$$\nu(p \in W(\{p, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\})) = \nu(p \in W(q, \tilde{q})) = 0.$$

There are two groups of DMs who may choose  $q$  from  $\{p, q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\}$ : those who like  $q$  best, and those who have  $\beta_1 \tilde{q} + \beta_2 q \succ q \succ p$  but  $\beta_1 \tilde{q} + \beta_2 q \in W(\{p, q, \delta_1 \tilde{q} + \delta_2 p\})$ . We have

$$\begin{aligned} \rho(q|\{p, q, \beta_1 \tilde{q} + \beta_2 q, \delta_1 \tilde{q} + \delta_2 p\}) &= \mu(q \succ \beta_1 \tilde{q} + \beta_2 q, p) \nu(q \notin W(p)) \\ &\quad + \mu(\beta_1 \tilde{q} + \beta_2 q \succ q \succ p) \\ &\quad \times [\nu(q \notin W(p)) - \nu(\beta_1 \tilde{q} + \beta_2 q \notin W(\{p, q\}))] \\ &= \mu(q \succ \tilde{q}, p) \nu(q \notin W(p)) \\ &\quad + \mu(\tilde{q} \succ q \succ p) [\nu(q \notin W(p)) - \nu(\tilde{q} \notin W(p))] \\ &= \rho(q|\{p, q, \tilde{q}\}). \end{aligned}$$

This completes the second subcase and, by extension, the second possibility for  $\tilde{q}$ . We have shown that  $(p, q)$  cannot be anomalous for  $\varepsilon$  sufficiently small, provided there exists  $\bar{\varepsilon} > 0$  such that  $\nu(p \in W(\tilde{q})) = 0$  for all  $\tilde{q} \in B_{\bar{\varepsilon}}(q)$ .

Now suppose there exist  $\tilde{q}$  arbitrarily close to  $q$  such that  $\nu(p \in W(\tilde{q})) > 0$ . Since  $\nu(p \in W(q)) = 0$ ,  $q$  must be on the boundary of  $\{x \in \Delta(Z) : \nu(p \in W(x)) > 0\}$ . There are  $q'$

arbitrarily close to  $q$  that are not on this boundary. For any such  $q'$ , there exists  $\varepsilon > 0$  such that  $\nu(p \in W(\tilde{q})) = 0$  for all  $\tilde{q} \in B_\varepsilon(q')$ . We have already seen that  $(p, q')$  cannot be anomalous. Thus, there is no  $\varepsilon > 0$  such that  $(p, q')$  is anomalous for all  $q' \in B_\varepsilon(q)$ .

#### A.1.14 PROOF OF THEOREM 5

If  $\rho$  has an REU representation, the unique RJ representation has  $\nu(\mathcal{U}) = 1$ .  $\mu$  is simply the REU representation. For the remainder of this proof, we assume that  $\rho$  has an REU representation with  $\nu(\mathcal{U}) < 1$ .

#### IDENTIFYING $D(p)$

Fix any interior  $q, \tilde{q}$ . Suppose that  $q \in D(\tilde{q})$ . For any  $p \in \Delta(Z)$ , we have

$$\nu(q \in W(\tilde{q})) \geq \max\{\nu(p \in W(\{q, \tilde{q}\})), \nu(\tilde{q} \in W(\{p, q\}))\}$$

by the second part of Lemma 22. This implies

$$\rho(q|\{p, q, \tilde{q}\}) = \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})).$$

By Independence,  $q \in D(\beta\tilde{q} + (1 - \beta)q)$  for any  $\beta \in (0, 1)$ , so

$$\begin{aligned} \nu(q \in W(\beta\tilde{q} + (1 - \beta)q)) &\geq \max\{\nu(p \in W(\{q, \beta\tilde{q} + (1 - \beta)q, \delta\tilde{q} + (1 - \delta)p\})), \\ &\nu(\beta\tilde{q} + (1 - \beta)q \in W(\{p, q, \delta\tilde{q} + (1 - \delta)p\})), \\ &\nu(\delta\tilde{q} + (1 - \delta)p \in W(\{p, q, \beta\tilde{q} + (1 - \beta)q\}))\} \end{aligned}$$

for any  $\delta \in (0, 1)$ . This implies

$$\begin{aligned} \rho(q|\{p, q, \beta\tilde{q} + (1 - \beta)q, \delta\tilde{q} + (1 - \delta)q\}) &= \mu(q \succ p, \beta\tilde{q} + (1 - \beta)q)\nu(q \notin W(\beta\tilde{q} + (1 - \beta)q)) \\ &= \mu(q \succ p, \tilde{q})\nu(q \notin W(\tilde{q})) \\ &= \rho(q|\{p, q, \tilde{q}\}) \end{aligned}$$

Now suppose that  $\nu(q \in W(\tilde{q})) > 0$ , but  $q \notin D(\tilde{q})$ . (We have already seen how to tell whether  $\nu(q \in W(\tilde{q})) > 0$ .) We showed at the end of Section A.1.13 that

$$\rho(q|\{p, q, \tilde{q}\}) > \rho(q|\{p, q, \beta\tilde{q} + (1 - \beta)q, \delta\tilde{q} + (1 - \delta)p\})$$

for some  $p \in \Delta(Z)$  (in particular, any  $p \in D(q)$  will work) and some  $\beta, \delta \in (0, 1)$ . This allows us to identify  $D(\tilde{q})$  and, by extension,  $D(x)$  for any  $x \in \Delta(Z)$ .

#### IDENTIFYING $\nu$

Take any  $p, q \in \text{int}(\Delta(Z))$  such that  $\nu(p \in W(q)) > 0$  and  $p \notin D(q)$ . Construct  $x^*$  and  $\tilde{q}$  exactly as in Section A.1.13.

Since

$$\nu(p \in W(\{q, \tilde{q}\})) = \nu(p \in W(q)) > \nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(\tilde{q})) > \nu(\tilde{q} \in W(\{p, q\})) = 0,$$

we have

$$\begin{aligned} \rho(p|\{p, q, \tilde{q}\}) &= \mu(p \succ q, \tilde{q})\nu(p \notin W(q)) \\ &= \left[ \mu\left(\frac{1}{2}q + \frac{1}{2}p \succ q, \tilde{q}\right) + \mu\left(p \succ \tilde{q} \succ \frac{1}{2}q + \frac{1}{2}p\right) \right] \nu(p \notin W(q)) \end{aligned}$$

Now consider menu

$$\left\{ \frac{1}{2}q + \frac{1}{2}p, q, \tilde{q} \right\}.$$

By Lemma 23,

$$\nu \left( q \in W \left( \left\{ \frac{1}{2}q + \frac{1}{2}p, q, \tilde{q} \right\} \right) \right) = \nu(q \in W(\{p, \tilde{q}\})) = \nu(q \in W(\tilde{q})).$$

Since  $\nu(\{p, q\} \in W(\tilde{q})) > 0$ ,  $\nu(\tilde{q} \in W(\{p, q\})) = 0$ . By Lemma 20,  $\nu(\tilde{q} \in W(A)) = 0$  for all  $A \subset \text{co}(\{p, q\})$ . In particular,

$$\nu \left( \tilde{q} \in W \left( \left\{ \frac{1}{2}q + \frac{1}{2}p, q \right\} \right) \right) = 0.$$

By Lemma 20,

$$\begin{aligned} \nu \left( \frac{1}{2}q + \frac{1}{2}p \in W(q, \tilde{q}) \right) &= \max_{\lambda_1, \lambda_2 \geq 0, \lambda_1 + \lambda_2 = 1} \nu \left( \frac{1}{2}q + \frac{1}{2}p \in W(\lambda_1 q + \lambda_2 \tilde{q}) \right) \\ &= \max \left\{ \max_{\lambda_1 \geq \frac{1}{2}} \nu(p \in W((2\lambda_1 - 1)q + 2\lambda_2 \tilde{q})), \right. \\ &\quad \left. \max_{\lambda_1 < \frac{1}{2}} \nu \left( \frac{1 - 2\lambda_1}{2\lambda_2} q + \frac{1}{2\lambda_2} p \in W(\tilde{q}) \right) \right\} \end{aligned}$$

We work through the terms in the max. For the first term, recall that

$$\nu(p \in W(q)) = \nu(p \in W(\{q, \tilde{q}\})) = \max_{x \in \text{co}(\{q, \tilde{q}\})} \nu(p \in W(x)).$$

Thus, the first term is maximized at  $\nu(p \in W(q))$  by setting  $\lambda_1 = 1$ . For the second term, suppose  $\nu(x \in W(\tilde{q})) > \nu(p \in W(q))$  for some  $x \in \text{co}(\{p, q\})$ . We showed in Section A.1.13 that  $\nu(p \in W(\tilde{q})) < \nu(p \in W(q))$ , so  $x \neq p$ . We have  $\nu(p \in W(x)) = \nu(p \in W(q))$  by Independence. By the first part of Lemma 18,  $\nu(p \in W(\tilde{q})) \geq \nu(p \in W(q))$ —contradiction. Conclude that the



second term is no greater than  $v(p \in W(q))$ , so

$$v\left(\frac{1}{2}q + \frac{1}{2}p \in W(q, \tilde{q})\right) = v(p \in W(q)).$$

We have

$$\begin{aligned} \rho\left(\frac{1}{2}q + \frac{1}{2}p \mid \left\{\frac{1}{2}q + \frac{1}{2}p, q, \tilde{q}\right\}\right) &= \mu\left(\frac{1}{2}q + \frac{1}{2}p \succ q, \tilde{q}\right) v(p \in W(q)) \\ \rho(p \mid \{p, q, \tilde{q}\}) - \rho\left(\frac{1}{2}q + \frac{1}{2}p \mid \left\{\frac{1}{2}q + \frac{1}{2}p, q, \tilde{q}\right\}\right) &= \mu\left(p \succ \tilde{q} \succ \frac{1}{2}q + \frac{1}{2}p\right) v(p \in W(q)). \end{aligned} \tag{A.12}$$

Finally, consider menu

$$\left\{\frac{1}{2}q + \frac{1}{2}p, x^*, q, \tilde{q}\right\}.$$

By Lemma 23,

$$\begin{aligned} v\left(q \in W\left(\left\{\frac{1}{2}q + \frac{1}{2}p, x^*, \tilde{q}\right\}\right)\right) &= v(q \in W(\{p, \tilde{q}\})) = v(q \in W(\tilde{q})) \\ v\left(\tilde{q} \in W\left(\left\{\frac{1}{2}q + \frac{1}{2}p, x^*, q\right\}\right)\right) &= v(\tilde{q} \in W(\{p, q\})) = 0. \end{aligned}$$

Since  $x^*$  was chosen so that

$$v(x^* \in W(q)) \geq v\left(\frac{1}{2}q + \frac{1}{2}p \in W(\{x^*, q, \tilde{q}\})\right),$$

Lemma 21 gives

$$\begin{aligned}
\nu\left(\frac{1}{2}q + \frac{1}{2}p \in W(\{x^*, q, \tilde{q}\})\right) &= \nu\left(\frac{1}{2}q + \frac{1}{2}p \in W(\{q, \tilde{q}\})\right) \\
&= \nu\left(\frac{1}{2}q + \frac{1}{2}p \in W(q)\right) \\
&= \nu(p \in W(q)).
\end{aligned}$$

We have

$$\begin{aligned}
\rho\left(x^* \left| \left\{ \frac{1}{2}q + \frac{1}{2}p, x^*, q, \tilde{q} \right\} \right.\right) &= \mu\left(x^* \succ \frac{1}{2}p + \frac{1}{2}q, \tilde{q}\right) \nu(x^* \notin W(q)) \\
\rho\left(\frac{1}{2}p + \frac{1}{2}q \left| \left\{ \frac{1}{2}q + \frac{1}{2}p, x^*, q, \tilde{q} \right\} \right.\right) &= \mu\left(\frac{1}{2}p + \frac{1}{2}q \succ x^*, q\right) \nu(p \notin W(q)) \\
&\quad + \mu\left(x^* \succ \frac{1}{2}p + \frac{1}{2}q \succ \tilde{q}\right) \\
&\quad \times [\nu(p \notin W(q)) - \nu(x^* \notin W(q))] \\
\rho\left(\left\{x^*, \frac{1}{2}p + \frac{1}{2}q\right\} \left| \left\{ \frac{1}{2}p + \frac{1}{2}q, x^*, q, r \right\} \right.\right) &= \left[ \mu\left(\frac{1}{2}p + \frac{1}{2}q \succ d^*, q\right) \right. \\
&\quad \left. + \mu\left(d^* \succ \frac{1}{2}p + \frac{1}{2}q \succ \tilde{q}\right) \right] \nu(p \notin W(q)) \\
&\quad + \mu\left(d^* \succ \tilde{q} \succ \frac{1}{2}p + \frac{1}{2}q\right) \nu(x^* \notin W(q)) \\
&= \mu\left(\frac{1}{2}p + \frac{1}{2}q \succ q, \tilde{q}\right) \nu(p \notin W(q)) \\
&\quad + \mu\left(x^* \succ \tilde{q} \succ \frac{1}{2}p + \frac{1}{2}q\right) \nu(x^* \notin W(q)).
\end{aligned}$$

This implies

$$\begin{aligned} & \rho \left( \left\{ x^*, \frac{1}{2}p + \frac{1}{2}q \right\} \middle| \left\{ \frac{1}{2}p + \frac{1}{2}q, x^*, q, r \right\} \right) - \rho \left( \frac{1}{2}q + \frac{1}{2}p \middle| \left\{ \frac{1}{2}q + \frac{1}{2}p, q, \tilde{q} \right\} \right) \\ & = \mu \left( p \succ \tilde{q} \succ \frac{1}{2}p + \frac{1}{2}q \right) \nu(x^* \notin W(q)). \quad (\text{A.13}) \end{aligned}$$

Combining (A.12) and (A.13),

$$\frac{\rho(p|p, q, \tilde{q}) - \rho(\frac{1}{2}p + \frac{1}{2}q | \{\frac{1}{2}p + \frac{1}{2}q, q, \tilde{q}\})}{\rho(\{x^*, \frac{1}{2}p + \frac{1}{2}q\} | \{\frac{1}{2}p + \frac{1}{2}q, x^*, q, \tilde{q}\}) - \rho(\frac{1}{2}p + \frac{1}{2}q | \{\frac{1}{2}p + \frac{1}{2}q, q, \tilde{q}\})} = \frac{\nu(p \notin W(q))}{\nu(x^* \notin W(q))}. \quad (\text{A.14})$$

Now we show how to identify  $\nu(x^* \notin W(q))$ . Choose any point  $p^* \neq q$  on the boundary of  $\{x \in \Delta(Z) : \nu(x \in W(q)) > 0\}$ . We have  $\nu(p^* \in W(q)) = 0$ . but some sequence  $\{p_i\}$  converging to  $p^*$  such that  $\nu(p_i \in W(q)) > 0$  for all  $i$ . For each  $p_i$ , we can construct  $x_i^*$  and  $\tilde{q}_i$  as above, and use them to recover

$$\frac{\nu(p_i \notin W(q))}{\nu(x_i^* \notin W(q))}.$$

Since each  $x_i^* \in D(q)$ ,  $\nu(x_i^* \notin W(q)) = \nu(\mathcal{U})$  for all  $i$ .

Fix any  $t \in [\nu(\mathcal{U}), 1]$ . Let  $\mathcal{M}_t$  be the unique element of  $\text{supp}(\nu)$  such that  $\nu(\{\mathcal{M} : \mathcal{M} \supset t\}) = t$ .  $\mathcal{M}_t$  always exists by the second property of  $\nu$ . Let  $W_t(q) := W_{\mathcal{M}_t}(q)$ . Let  $\varepsilon_t$  be the (Hausdorff) distance between  $W_t(q)$  and  $W_1(q)$ . Notice that  $\lim_i \varepsilon_{t_i} = 0$  by the first property of  $\nu$ . Since no element of  $B_{\varepsilon_t}(p^*)$  can belong to  $W_t(q)$ , we must have  $\nu(\tilde{p} \notin W(q)) \geq t$  for all  $\tilde{p} \in B_{\varepsilon_t}(p^*)$ . Since  $\lim_i p_i = p^*$ , we have  $\lim_i \nu(p_i \notin W(q)) = 1$ .

This implies

$$\lim_i \frac{\nu(p_i \notin W(q))}{\nu(x_i^* \notin W(q))} = \frac{1}{\nu(\mathcal{U})}.$$

Plugging into (A.14) and using  $\nu(x^* \notin W(q)) = \nu(\mathcal{U})$ , we recover  $\nu(p \notin W(q))$ .

Since  $p$  was chosen arbitrarily subject to the requirement  $\nu(p \in W(q)) > 0$ , we can recover

$\nu(x \in W(q))$  for all  $x$  such that  $\nu(x \in W(q)) > 0$ . For each  $t \in (\nu(\mathcal{U}), 1]$ , let

$$W_t(q) := \{x \in \Delta(Z) : \nu(x \notin W(q)) < t\}.$$

Let  $W_{\nu(\mathcal{U})}(q) = \text{int} \left( \bigcap_{t \in (\nu(\mathcal{U}), 1]} W_t(q) \right)$ . (This set may be empty.) For  $t \in [\nu(\mathcal{U}), 1]$ , let

$$\mathcal{M}_t := \{\tilde{\mathcal{L}} \in \mathcal{U} : q \succ W_t(q)\}.$$

$\nu$  is pinned down by  $\nu(\mathcal{U})$  and, for  $t \in (\nu(\mathcal{U}), 1]$ ,

$$\nu(\{\mathcal{M} : \mathcal{M} \supset \mathcal{M}_t\}) = t.$$

#### IDENTIFYING $\mu$

Now we recover  $\mu$ . Fix any  $A \in \mathcal{F}(\Delta(Z))$ . Since we have already recovered  $\nu$ , we can index the elements of  $A$  as follows:

$$\nu(a_1 \in W(A)) \geq \nu(a_2 \in W(A)) \geq \cdots \geq \nu(a_{|A|} \in W(A)).$$

To simplify notation, let

$$A^i := \{a_i, \dots, a_{|A|}\}.$$

For each  $i \in \{1, \dots, |A|\}$ , we have

$$\begin{aligned} \rho(a_i|A) &= \mu(a_i \succsim A) \nu(a_i \notin W(A^{i+1})) \\ &\quad + \sum_{j < i} \mu(a_j \succ a_i \succsim A^{i+1}) [\nu(a_i \notin W(A^{i+1})) - \nu(a_j \notin W(A^{i+1}))] \end{aligned} \quad (\text{A.15})$$

$$= \mu(a_i \succsim A^i) \nu(a_i \notin W(A^{i+1})) - \sum_{j < i} \mu(a_j \succ a_i \succsim A^{i+1}) \nu(a_j \notin W(A^{i+1})). \quad (\text{A.16})$$

By Lemma 21,  $\nu(a_i \in W(A)) = \nu(a_i \in W(A^j))$  for  $j \leq i$ . The ordering we used for  $A$  still works for  $A^j$ :

$$\nu(a_j \in W(A^j)) \geq \dots \geq \nu(a_{|A|} \in W(A^j)).$$

For  $j \leq i$ , we have

$$\begin{aligned} \rho(a_i|A^j) &= \mu(a_i \succsim A^j) \nu(a_i \notin W(A^{j+1})) \\ &\quad - \sum_{k \in \{j, \dots, i-1\}} \mu(a_k \succ a_i \succsim A^{k+1}) \nu(a_k \notin W(A^{k+1})). \end{aligned}$$

Combining this with (A.16),

$$\rho(a_i|A^{j+1}) - \rho(a_i|A^j) = \mu(a_j \succ a_i \succsim A^{j+1}) \nu(a_j \notin W(A^{j+1})).$$

Since we have already recovered  $\nu(a_j \notin W(A^{j+1}))$ , we can use this equation to recover  $\mu(a_j \succ a_i \succsim A^{j+1})$ . Plugging this back into (A.15) and rearranging, we can recover  $\mu(a_i \succsim A)$ . Since  $a_i$  was an arbitrary member of  $A$ , and  $A$  was an arbitrary menu, this is enough to fully recover  $\mu$ .

### A.1.15 PROOF OF PROPOSITION 6

Suppose that  $\mathcal{M}$  contains at least one preference that weakly prefers  $p$  to both  $q_1$  and  $q_2$ , and at least one preference that weakly prefers both  $q_1$  and  $q_2$  to  $p$ . This implies  $\mathcal{M}^{\text{avoid}}(\{p, q\}) = \{p, q\}$ . Since  $\mathcal{M}^{\text{avoid}} \subset \mathcal{M}$ , it also implies  $\mathcal{M}(\{p, q_1\}) = \{p, q_1\}$  and  $\mathcal{M}(\{p, q_2\}) = \{p, q_2\}$ . Since the DM does not face any constraints on any of the feasible menus, he is weakly better off acquiring information.

Now suppose that  $\mathcal{M}$  contains at least one preference that weakly prefers  $p$  to both  $q_1$  and  $q_2$ , but no preference that weakly prefers both  $q_1$  and  $q_2$  to  $p$ . This implies  $\mathcal{M}^{\text{avoid}}(\{p, q\}) = \{p\}$ . It also implies  $p \in \mathcal{M}(\{p, q_1\})$  and  $p \in \mathcal{M}(\{p, q_2\})$ . Since the DM must choose  $p$  if he avoids information, and has the option of choosing  $p$  if he acquires information, he is weakly better off becoming informed.

Finally, suppose that  $\mathcal{M}$  contains at least one preference that weakly prefers both  $q_1$  and  $q_2$  to  $p$ , but no preference that weakly prefers  $p$  to both  $q_1$  and  $q_2$ . This implies  $\mathcal{M}^{\text{avoid}}(\{p, q\}) = \{q\}$ . It also implies  $q_1 \in \mathcal{M}(\{p, q_1\})$  and  $q_2 \in \mathcal{M}(\{p, q_2\})$ . Since  $q = \alpha q_1 + (1 - \alpha)q_2$ , the DM is weakly better off ex ante if he acquires information. This covers all the cases.

### A.1.16 PROOF OF PROPOSITION 7

Let  $\mathcal{M}$  be the element of  $\text{supp}(\nu)$  such that

$$\nu(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \subset \mathcal{M}\}) = \nu(q \in W(p)).$$

We need to find  $q_1 \in \Delta(Z)$  such that  $s(q_1) = s(p)$  and  $m(q_1) > m(p)$  for all  $m \in \mathcal{M} := \{m \in \mathcal{M} : m(q) = m(p)\}$ . Suppose there is no such  $q_1$ . First, suppose there is no  $x$  such that  $m(x) > m(p)$  for all  $m \in \mathcal{M}$ . Then, there is no  $x$  such that  $m(p) > m(p + \lambda(p - x))$  for all  $m \in \mathcal{M}$

and all  $\lambda$  sufficiently small. Since  $M \subset \mathcal{M}$ ,  $W_{\mathcal{M}}(p)$  must be empty—contradiction.

Now suppose  $s(x) > s(p)$  for all  $x$  such that  $m(x) > m(p)$  for all  $m \in M$ . Since  $q$  is on the boundary of  $W_{\mathcal{M}}(p)$ , we can find  $\tilde{q}$  arbitrarily close to  $q$  such that  $m(\tilde{q}) > m(p)$  for all  $m \in M$ . In particular, we can choose  $\tilde{q}$  close enough to  $q$  that  $s(p) > s(\tilde{q})$ —contradiction. Now suppose  $s(x) < s(p)$  for all  $x$  such that  $m(x) > m(p)$  for all  $m \in M$ . Take  $\tilde{q}$  as above. Notice that  $m(p + \lambda(p - \tilde{q})) < m(p)$  for all  $m \in M$  for any  $\lambda$  such that  $p + \lambda(p - \tilde{q}) \in \Delta(Z)$ . Similarly,  $s(p + \lambda(p - \tilde{q})) > s(p)$ —contradiction. There must exist  $x, y$  such that  $s(x) > s(p) > s(y)$  and  $m(x), m(y) > m(p)$  for all  $m \in \mathcal{M}$ . Some combination of  $x$  and  $y$  is the desired  $q_1$ .

Notice that  $m(q + \lambda(q - q_1)) < m(p)$  for all  $m \in M$  and any  $\lambda$  such that  $q + \lambda(q - q_1) \in \Delta(Z)$ . By choosing  $\lambda$  sufficiently small, we can ensure that  $m(q + \lambda(q - q_1)) < m(p)$  for all  $m \in \mathcal{M}$ , i.e.  $q + \lambda(q - q_1) \in W_{\mathcal{M}}(p)$ . Suppose not. Then we have sequences  $\{m_i : m_i \in \mathcal{M} \setminus M\}$  and  $\lambda_i \rightarrow 0$  such that  $m_i(q + \lambda_i(q - q_1)) > m_i(p)$ . We can shift and rescale the  $m_i$  so they belong to a compact subset of  $\mathbb{R}^Z$ , then pass to a convergent subsequence. Call the limit  $m^*$ . We have  $m^*(q) \geq m^*(p)$ . Since  $\mathcal{M}$  is closed,  $m^* \in \mathcal{M}$ , so it cannot be that  $m^*(q) > m^*(p)$ . Conclude that  $m^*(q) = m^*(p)$ , so  $m^* \in M$ . Since each  $m_i \in \mathcal{M} \setminus M$ ,  $m_i(p) > m_i(q)$  for all  $i$ . Since  $m_i(q + \lambda_i(q - q_1)) > m_i(p)$ , we must have  $m_i(q) > m_i(q_1)$  for all  $i$ , so  $m^*(q) \geq m^*(q_1)$  for all  $i$ . This contradicts the definition of  $q_1$ , which requires  $m^*(q_1) > m^*(p) = m^*(q)$ .

Thus, we can find  $\lambda$  such that  $q_2 := q + \lambda(q - q_1) \in W_{\mathcal{M}}(p)$ . This implies  $v(q_2 \in W(p)) > v(q \in W(p))$ . Set  $\alpha = \lambda/(1 + \lambda)$ , so  $(q_1, q_2, \alpha)$  is a signal for  $q$ . We show that

$$S_{(\mu, N)}(\alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}}) > S_{(\mu, N)}(\delta_{\{p, q\}}).$$

for any  $\mu$  such that  $\text{supp}(\mu) = \{u \in \mathcal{U} : u(q) \geq u(p)\}$  and any  $N$ . (Self-knowledge  $N$  does not matter here because the DM does not get to choose whether to acquire information; he is required to be informed, or required to remain ignorant.) To simplify the notation, normalize  $s(p)$  to 0. We

have  $0 = s(p) = s(q_1) > s(q) > s(q_2)$ . Since the social planner does not care whether  $q_1$  or  $p$  is chosen, we can break the support of  $\mu$  into two groups:  $u$  such that  $u(q_2), u(q) > u(p) = 0$ , and  $u$  such that  $u(q_1) > u(q) > u(p) = 0 > u(q_2)$ . In the first case,

$$\begin{aligned} S_{(u,N)} (\alpha \delta_{\{p,q_1\}} + (1-\alpha) \delta_{\{p,q_2\}}) - S_{(u,N)} (\delta_{\{p,q\}}) \\ = (1-\alpha) s(q_2) [\nu(q_2 \notin W(p)) - \nu(q \notin W(p))] > 0. \end{aligned}$$

In the second case,

$$S_{(u,N)} (\alpha \delta_{\{p,q_1\}} + (1-\alpha) \delta_{\{p,q_2\}}) - S_{(u,N)} (\delta_{\{p,q\}}) = -(1-\alpha) s(q_2) \nu(q \notin W(p)) > 0.$$

Integrating over all  $u$  such that  $u(q) \geq u(p)$ , we conclude that the social planner is strictly better off imposing information than withholding it.

To show that  $S_{(u,N)} (\{\delta_{\{p,q\}}, \alpha \delta_{\{p,q_1\}} + (1-\alpha) \delta_{\{p,q_2\}}\})$  is strictly between these two extremes, it suffices to show that one positive-measure group of DMs voluntarily acquires information, and another positive-measure group avoids it. For the first part, consider a DM with beliefs  $\tilde{\nu}$  and utility  $u$  such that  $u(q_1) > u(q) > u(p) = 0 > u(q_2)$ . It is strictly optimal for him to acquire information if

$$\alpha u(q_1) [\tilde{\nu}(q_1 \notin W(p)) - \tilde{\nu}(q \notin W(p))] - (1-\alpha) u(q_2) \tilde{\nu}(q \notin W(p)) > 0. \quad (\text{A.17})$$

Recall that  $q_2 \in W_{\mathcal{M}}(p)$  and  $q \notin W_{\mathcal{M}}(p)$ . Since  $q$  is a convex combination of  $q_1$  and  $q_2$  and  $W_{\mathcal{M}}(p)$  is convex,  $q_1 \notin W_{\mathcal{M}}(p)$ , so  $\tilde{\nu}(q_1 \notin W(p)) \geq \tilde{\nu}(q \notin W(p))$ . Thus, (A.17) will hold if  $\tilde{\nu}(q \notin W(p)) > 0$ , i.e. if  $\tilde{\nu}(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \mathcal{M}\}) > 0$ . Let  $E$  be the set of  $\tilde{\nu} \in \text{supp}(N)$  that satisfy this



condition. By definition of self-knowledge,

$$\begin{aligned} \int_{\tilde{\nu}} \tilde{\nu}(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \mathcal{M}\}) dN(\tilde{\nu}) &= N(E) \int_{\tilde{\nu} \in E} \tilde{\nu}(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \mathcal{M}\}) dN(\tilde{\nu}|E) \\ &= \nu(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \mathcal{M}\}) \\ &> 0. \end{aligned}$$

This implies  $N(E) > 0$ , so (A.17) holds for a positive-measure group of DMs.

Now consider a DM with beliefs  $\tilde{\nu}$  and utility  $u$  such that  $u(q_2) > u(q) > 0 > u(q_1)$ . A DM in this set strictly prefers to avoid information if

$$(1 - \alpha)u(q_2) [\tilde{\nu}(q_2 \notin W(p)) - \tilde{\nu}(q \notin W(p))] - \alpha u(q_1) \tilde{\nu}(q \notin W(p)) < 0. \quad (\text{A.18})$$

For any  $\tilde{\nu}$  such that  $\tilde{\nu}(q_2 \notin W(p)) < \tilde{\nu}(q \notin W(p))$ , this condition will hold if  $u(q_2)$  is sufficiently large relative to  $-u(q_1)$ . We can use the same arguments as the previous step to show that  $N$  puts positive probability on  $\tilde{\nu}(q_2 \notin W(p)) < \tilde{\nu}(q \notin W(p))$ . Thus, (A.18) holds for a positive-measure group of DMs.

#### A.1.17 PROOF OF PROPOSITION 8

Suppose that, for some  $\mathcal{M} \in \text{supp}(\nu)$ , there exist  $m_1, m_2 \in \mathcal{M}$  such that  $m_1(q_1) \geq m_1(p)$  and  $m_2(q_2) \geq m_2(p)$ . Suppose further that  $m(p) > \min\{m(q_1), m(q_2)\}$  for all  $m \in \mathcal{M}$ . This implies

$$m(p) < \max \{m(p + \lambda(p - q_1)), m(p + \lambda(p - q_2))\} \quad (\text{A.19})$$

for all  $m \in \mathcal{M}$  and all  $\lambda > 0$  such that  $p + \lambda(p - q_1), p + \lambda(p - q_2) \in \Delta(Z)$ . (Since  $p$  is interior, some such  $\lambda$  must exist.)

We showed in the proof of Theorem 3 (specifically, in the necessity proof of Convexity) that (A.19) implies the existence of  $\alpha \in [0, 1]$  such that

$$m(p + \lambda(p - (\alpha q_1 + (1 - \alpha)q_2))) > m(p)$$

for all  $m \in \mathcal{M}$ . Rearranging, we have

$$m(p) > m(\alpha q_1 + (1 - \alpha)q_2)$$

for all  $M \in \mathcal{M}$ . That is,  $q \in W_{\mathcal{M}}(p)$ . Let  $q := \alpha q_1 + (1 - \alpha)q_2$ . Since  $m_1(q_1) \geq m_1(p)$  and  $m_2(q_2) \geq m_2(p)$ ,  $q \neq q_1$  and  $q \neq q_2$ .

Fix any utility  $u$  such that  $\min\{u(q_1), u(q_2)\} > u(p) = 0$  and any belief  $\tilde{v}$ . We have

$$\begin{aligned} & U_{(u, \tilde{v})}(\alpha \delta_{\{p, q_1\}} + (1 - \alpha) \delta_{\{p, q_2\}}) - U_{(u, \tilde{v})}(\delta_{\{p, q\}}) \\ &= \alpha u(q_1)[\tilde{v}(q_1 \notin W(p)) - \tilde{v}(q \notin W(p))] + (1 - \alpha)u(q_2)[\tilde{v}(q_2 \notin W(p)) - \tilde{v}(q \notin W(p))]. \end{aligned} \tag{A.20}$$

Let  $\underline{\mathcal{M}}$  be the member of  $\text{supp}(\nu)$  such that

$$\nu(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \underline{\mathcal{M}}\}) = \max\{\nu(q_1 \notin W(p)), \nu(q_2 \notin W(p))\},$$

and let  $\bar{\mathcal{M}}$  be the member of  $\text{supp}(\nu)$  such that

$$\nu(\{\tilde{\mathcal{M}} : \tilde{\mathcal{M}} \supseteq \bar{\mathcal{M}}\}) = \nu(q \notin W(p)).$$

Since we have  $\mathcal{M}$  such that  $q_1, q_2 \notin W_{\mathcal{M}}(p)$  but  $q \in W_{\mathcal{M}}(p)$ , we must have  $\underline{\mathcal{M}} \subset \bar{\mathcal{M}}$ . For any  $\tilde{v}$

such that

$$\tilde{v} \left( \left\{ \tilde{\mathcal{M}} : \bar{\mathcal{M}} \supset \tilde{\mathcal{M}} \supset \underline{\mathcal{M}} \right\} \right) > 0, \quad (\text{A.2I})$$

we have  $\max\{\tilde{v}(q_1 \notin W(p)), \tilde{v}(q_2 \notin W(p))\} > \tilde{v}(q \notin W(p))$ , so (A.2o) is strictly positive. The DM strictly prefers to acquire information. Otherwise, the DM is indifferent to information. By assumption, he acquires information in this case as well.

Now consider a social planner with  $s(p) = 0 > \max\{s(q_1), s(q_2)\}$  and a utility  $u$  such that  $\min\{u(q_1), u(q_2)\} > u(p) = 0$ . We have shown that a DM with utility  $u$  and beliefs  $\tilde{v}$  strictly prefers to acquire information if and only if (A.2I) holds. Fix any self-knowledge  $N$ , and let  $E$  be the set of  $\tilde{v} \in \text{supp}(N)$  that satisfy (A.2I). An equivalent definition is

$$E := \{\tilde{v} \in \text{supp}(N) : \max\{\tilde{v}(q_1 \notin W(p)), \tilde{v}(q_2 \notin W(p))\} > \tilde{v}(q \notin W(p))\}.$$

By definition of self-knowledge,

$$\begin{aligned} \int_{\tilde{v}} \tilde{v} \left( \left\{ \tilde{\mathcal{M}} : \bar{\mathcal{M}} \supset \tilde{\mathcal{M}} \supset \underline{\mathcal{M}} \right\} \right) dN(\tilde{v}) &= N(E) \int_{\tilde{v} \in E} \tilde{v} \left( \left\{ \tilde{\mathcal{M}} : \bar{\mathcal{M}} \supset \tilde{\mathcal{M}} \supset \underline{\mathcal{M}} \right\} \right) dN(\tilde{v}|E) \\ &= \nu \left( \left\{ \tilde{\mathcal{M}} : \bar{\mathcal{M}} \supset \tilde{\mathcal{M}} \supset \underline{\mathcal{M}} \right\} \right) \\ &> 0. \end{aligned}$$

Thus,  $N(E) > 0$ . If all the DMs with beliefs in  $E$  acquire information, the cost to the social planner is

$$\begin{aligned} N(E) \left[ \alpha[-s(q_1)] \int_E \tilde{v}(q_1 \notin W(p)) - \tilde{v}(q \notin W(p)) N(d\tilde{v}) \right. \\ \left. + (1 - \alpha)[-s(q_2)] \int_E \tilde{v}(q_2 \notin W(p)) - \tilde{v}(q \notin W(p)) N(d\tilde{v}) \right] > 0. \end{aligned}$$

Since information does not change the behavior of any DM with beliefs outside  $E$ , we conclude that

$$S_{(u,N)}(\delta_{\{p,q\}}) > S_{(u,N)}(\alpha\delta_{\{p,q_1\}} + (1-\alpha)\delta_{\{p,q_2\}}).$$

Fix any  $\mu$  such that  $\mu(\{u : \min\{u(q_1), u(q_2)\} > u(p)\}) = 1$ . Integrating over  $\text{supp}(\mu)$ , we can replace subscript  $u$  with subscript  $\mu$ .

Since the DMs in this society acquire information with probability 1, we also have

$$S_{(u,N)}(\alpha\delta_{\{p,q_1\}} + (1-\alpha)\delta_{\{p,q_2\}}) = S_{(u,N)}(\{\delta_{\{p,q\}}, \alpha\delta_{\{p,q_1\}} + (1-\alpha)\delta_{\{p,q_2\}}\}).$$

## A.2 MODEL VARIANTS

### A.2.1 CONTINUITY WITHOUT RECOVERABILITY

We slightly strengthen local non-satiation.

**Definition 59** (Locally non-satiated\*).  *$\mathcal{M}$  is locally non-satiated\* if, for any  $a \in \mathcal{A}$ , there exists  $Z \in \mathcal{F}(B_\varepsilon(a) \cap \mathcal{Z})$  such that  $Z$  is strictly preferred to  $a$  by all  $m \in \mathcal{M}$ , and  $a$  is strictly preferred to  $Z$  by  $u$ .*

To get a version of Theorem 2 with local non-satiation\* in place of non-satiation, we slightly strengthen Improvability.

**Axiom 27** (Improvability\*). *For any  $a \in \mathcal{A}$  and any  $\varepsilon > 0$ , there is some  $Z \in \mathcal{F}(B_\varepsilon(a) \cap \mathcal{Z})$  such that  $a \succ Z$  and  $a \notin c(\{a\} \cup Z)$ .*

**Proposition 22.** *Fix  $(\succsim, c)$ , and let*

$$\hat{c}(A) := \{a \in A : a \sim c(A)\}.$$

If  $(\succsim, c)$  has a representation  $(u, \mathcal{M})$  with  $\mathcal{M}$  closed and locally non-satiated\*, then  $(\succsim, \hat{c})$  has a recoverable representation  $(u, \hat{\mathcal{M}})$  with  $\hat{\mathcal{M}}$  closed and locally non-satiated\*.

*Proof.* Take any representation  $(u, \mathcal{M})$  for  $(\succsim, c)$ , where  $\mathcal{M}$  closed and locally non-satiated\*. Since  $\mathcal{M}$  is closed,

$$\bigcap_{m \in \mathcal{M}} \left\{ a \in \mathcal{A} : v(a) < \max_{b \in B} m(b) \right\}$$

is open for any  $B \in \mathcal{F}(\mathcal{A})$ . Enumerate the indifference classes of  $B$  from best to worst according to  $u$ :  $u(B_1) > \dots > u(B_n)$ . The set

$$\bigcup_{i=1}^n \left\{ a \in \mathcal{A} : u(a) \geq u(B_i) \text{ and } \forall m \in \mathcal{M} \ v(a) < \max_{b \in B_i \cup \dots \cup B_n} m(b) \right\}.$$

need not be open. Specifically, there may be  $a$  such that  $u(a) = u(B_i)$  and  $m(a) < \max_{b \in B_i \cup \dots \cup B_n} m(b)$  for all  $m \in \mathcal{M}$ . In choice terms, there may be  $a \sim B_i$  such that  $a \notin c(\{a\} \cup B_i \cup \dots \cup B_n)$ .  $\hat{c}$  eliminates this problem;  $a \in \hat{c}(\{a\} \cup B_i \cup \dots \cup B_n)$ .  $a \notin \hat{W}(B_i \cup \dots \cup B_n)$ , so  $a \notin \hat{W}(B)$ . We conclude that  $\hat{W}(B)$  is open, so  $\hat{c}$  satisfies Continuity.

It remains to show that  $\hat{c}$  satisfies C-IUA and Improvability. Consider C-IUA first. Suppose that  $d \in W(B)$  and  $B \subset \bar{W}(A)$ . We can assume  $d \succsim B$  and  $d \notin c(\{d\} \cup B)$ . We will have  $d \notin \hat{c}(\{d\} \cup B)$  only if  $d \succ c(\{d\} \cup B)$ ; assume this is the case. By IUA, we can eliminate any  $b \in B$  such that  $b \sim d$  without changing choice. Thus, we can assume  $d \succ B$ . We can also assume that there are sequences  $b_i \rightarrow b$  and  $A_i \rightarrow A$  such that  $b_i \in \hat{W}(A_i)$  for all  $i$ , and that  $b \succsim A$  for some  $b \in B$ . We conclude that  $d \succ A$ . By C-IUA, we have  $d \notin c(\{d\} \cup A)$ . Since  $d$  is not indifferent to any item in  $A$ ,  $d \notin \hat{c}(\{d\} \cup A)$  as well. A parallel argument goes through if  $d \in \bar{W}(B)$  and  $B \subset W(A)$ . Since  $\hat{c}$  does not add any new item to any  $W$  or  $\bar{W}$  this is all we need to check.

Now consider Improvability. Since  $\mathcal{M}$  is locally non-satiated\*, we can find  $Z \in \mathcal{F}(\mathcal{Z})$  arbitrarily close to  $a$  such that  $a \succ Z$  and  $a \notin c(\{a\} \cup Z)$ . This is still true for  $\hat{c}$ , so  $(\succsim, \hat{c})$  satisfies Improvabil-

ity\*.

By Theorem 2,  $\hat{c}$  has a recoverable representation  $(u, \hat{\mathcal{M}})$  with  $\hat{\mathcal{M}}$  closed and locally non-satiated\*.

□

### A.2.2 PREFERENCES FIXED, CONSTRAINTS RANDOM, $|\mathcal{A}| = 3$

Let  $\rho$  be a stochastic choice function on  $\mathcal{A}$ . Let  $\Pi$  be the set of strict preferences on  $\mathcal{A}$ .

**Definition 60** (Fixed-preference representation). *A fixed-preference representation for  $\rho$  is  $(\succ, \nu) \in \Pi \times \Delta(\mathcal{F}(\Pi))$  such that*

$$\rho(a|\mathcal{A}) = \nu(\{\mathcal{M} \in \mathcal{F}(\Pi) : a = \arg \max(\mathcal{M}(A), \succ)\})$$

where

$$\mathcal{M}(A) := \bigcup_{\succ_m \in \mathcal{M}} \arg \max(A, \succ_m).$$

**Proposition 23.** *Suppose  $|\mathcal{A}| = 3$ . The following are equivalent:*

1.  $\rho$  has a fixed-preference representation.
2. There is at most one pair  $(x, y) \in \mathcal{A}^2$  such that

$$\rho(x|\{x, y\}) < \rho(x|\mathcal{A}).$$

*Proof.* Write  $\mathcal{A} = \{a, b, d\}$ . Suppose that

$$\rho(b|\{a, b\}) < \rho(b|\mathcal{A}).$$

(If there is no pair  $x, y$  such that  $\rho(x|\{x, y\}) < \rho(x|\mathcal{A})$ , then we have a standard Random Utility representation. We can define  $\succ$  as usual, and let  $\nu = \delta_{\succ}$ .)

We take  $\succ$  to be  $a \succ b \succ d$ . Now we show how to construct an appropriate  $\nu$ . It is helpful to divide  $\mathcal{F}(\Pi)$  into subsets (“states”), where each state is defined by the restrictions that prevent the DM from maximizing  $\succ$ . To formalize this, fix any  $\mathcal{M} \in \mathcal{F}(\Pi)$ , any  $X \in \mathcal{F}(\mathcal{A})$  and  $y \notin X$ . Write  $X \triangleright_{\mathcal{M}} y$  if the following three conditions are met: (1)  $y \succ X$ , (2) for each proper subset  $X'$  of  $X$ ,  $y \succ_m X'$  for some  $\succ_m \in \mathcal{M}$ , and (3)  $\neg(y \succ_m X)$  for all  $\succ_m \in \mathcal{M}$ . Finally, define the state  $X \triangleright y$  to be the set of  $\mathcal{M}$  such that  $X \triangleright_{\mathcal{M}} y$ , but  $\neg(Z \triangleright_{\mathcal{M}} w)$  for all  $(Z, w) \neq (X, x)$ .

When  $|\mathcal{A}| = 3$  and  $a \succ b \succ d$ , there are eight states:

- (1) :  $\{b, d\} \triangleright a$
- (2) :  $d \triangleright a$
- (3) :  $d \triangleright a$  and  $d \triangleright b$
- (4) :  $b \triangleright a$
- (5) :  $d \triangleright a$  and  $b \triangleright a$
- (6) :  $d \triangleright a$  and  $d \triangleright b$  and  $b \triangleright a$
- (7) :  $d \triangleright b$
- (8) : no constraints

The probabilities of states (7) and (8) are uniquely determined:

$$\Pr(7) = \rho(d|b, d) - \rho(d|\mathcal{A})$$

$$\Pr(8) = \rho(b|b, d) - \rho(b|\mathcal{A}).$$

The probabilities of the remaining states are not uniquely determined, but certain sums are:

$$\Pr(1) + \Pr(2) + \Pr(3) = \rho(a|a, b) - \rho(a|\mathcal{A})$$

$$\Pr(4) + \Pr(5) + \Pr(6) = \rho(b|a, b)$$

$$\Pr(1) + \Pr(4) = \rho(a|a, d) - \rho(a|\mathcal{A})$$

$$\Pr(2) + \Pr(5) = \rho(d|a, d) - \rho(d|\mathcal{A})$$

$$\Pr(3) + \Pr(6) = \rho(d|\mathcal{A}).$$

Notice that states (1), (2), (3) can be transformed into (4), (5), (6) respectively by imposing  $b \triangleright a$ . Similarly, states (1), (4) can be transformed into (2), (5) by adding  $d \triangleright a$ , and then states (2), (5) can be transformed into (3), (6) by adding  $d \triangleright b$ . We will assume that, conditional on states (1)-(6) (equivalently, conditional on  $a$  being ruled out by a subset of  $\{b, d\}$ ),  $b \triangleright a$  is independent of  $d \triangleright a$  and  $d \triangleright a, d \triangleright b$ . This gives

$$\Pr(1) = \frac{(\rho(a|a, b) - \rho(a|\mathcal{A})) (\rho(a|a, d) - \rho(a|\mathcal{A}))}{1 - \rho(a|\mathcal{A})}$$

$$\Pr(2) = \frac{(\rho(a|a, b) - \rho(a|\mathcal{A})) (\rho(d|a, d) - \rho(d|\mathcal{A}))}{1 - \rho(a|\mathcal{A})}$$

$$\Pr(3) = \frac{(\rho(a|a, b) - \rho(a|\mathcal{A})) \rho(d|\mathcal{A})}{1 - \rho(a|\mathcal{A})}$$

$$\Pr(4) = \frac{\rho(b|a, b) (\rho(a|a, d) - \rho(a|\mathcal{A}))}{1 - \rho(a|\mathcal{A})}$$

$$\Pr(5) = \frac{\rho(b|a, b) (\rho(d|a, d) - \rho(d|\mathcal{A}))}{1 - \rho(a|\mathcal{A})}$$

$$\Pr(6) = \frac{\rho(b|a, b) \rho(d|\mathcal{A})}{1 - \rho(a|\mathcal{A})}.$$

Given our assumptions on  $\rho$ , all these probabilities will be positive. Moreover, the probabilities of all the states will sum to 1. It is easy to check that this assignment delivers the right predictions on all



menus.

To complete the proof, we just need to pick a  $\mathcal{M}$  that corresponds to each “state.” For the state with no restrictions, we could pick  $\{\succ\}$ . For (7), we could pick  $\{\succ_1\}$  such that  $a \succ_1 d \succ_1 b$ . For (1), we could pick  $\{\succ_2, \succ_3\}$  such that  $b \succ_2 a \succ_2 d$  and  $d \succ_3 a \succ_3 b$ . Proceeding in this way, we get an appropriate  $\nu \in \Delta(2^\Pi \setminus \emptyset)$ .  $\square$

As suggested in the proof, the distribution over exclusion relationships is not unique, even for  $|\mathcal{A}| = 3$ . One distribution can be picked out by imposing a particular independence assumption. It says: given that  $a$  is excluded by  $\{b, d\}$  or a subset, the probability that  $b$  excludes  $a$  does not depend on whether  $d$  excludes  $a$  and  $b, d$  excludes  $a$  alone, or  $d$  excludes neither  $a$  nor  $b$ . Similarly, the probability that  $d$  excludes  $a$  (or that  $d$  excludes both  $a$  and  $b$ ) does not depend on whether  $b$  excludes  $a$ . It may be possible to extend this idea to  $|\mathcal{A}| > 3$ , but the number of states is large even for  $|\mathcal{A}| = 4$ , and it is not obvious how to extend the independence assumption.

### A.2.3 PREFERENCES RANDOM, CONSTRAINTS FIXED

Take  $|\mathcal{A}| < \infty$ . Take a stochastic choice function  $\rho$  on  $\mathcal{F}(\mathcal{A})$ . Let  $\Pi$  be the set of strict preferences on  $\mathcal{A}$ .

**Definition 61** (Fixed-constraint representation). *A fixed-constraint representation is  $\mathcal{M} \subseteq \Pi$  and  $\mu \in \text{int}(\Delta(\Pi))$  such that*

$$\rho(a|A) = \mu(\{\succ \in \Pi : a \succ \mathcal{M}(A)\})$$

where

$$\mathcal{M}(A) = \bigcup_{\succ_m \in \mathcal{M}} \arg \max(\succ_m, A).$$

For any  $A \subseteq \mathcal{A}$ , let

$$S(A) = \text{supp}(\rho(\cdot|A)).$$

**Axiom 28** (Plott). For any  $A, B \in \mathcal{F}(\mathcal{A})$ ,

$$S(A \cup B) = S(S(A) \cup B).$$

**Axiom 29** (Support Dependence). For any  $A, B \in \mathcal{F}(\mathcal{A})$ ,  $S(A) = S(B)$  implies  $\rho(\cdot|A) = \rho(\cdot|B)$ .

Let

$$S := \{(a, A) \in \mathcal{A} \times \mathcal{F}(\mathcal{A}) : a \in A \text{ and } A = S(A)\}.$$

**Axiom 30** (No Arbitrage). Fix  $\lambda \in \mathbb{R}^{|S|}$ . If

$$\sum_{(a,A) \in S} \lambda_i \mathbf{1}\{a \succsim A\} \geq 0, \quad (\text{A.22})$$

for all  $\succ \in \Pi$ , then

$$\sum_{(a,A) \in S} \lambda_i \rho(a|A) \geq 0 \quad (\text{A.23})$$

If in addition (A.22) holds with strict inequality for some  $\succ$ , then (A.23) holds with strict inequality.

**Proposition 24.**  $\rho$  has a fixed-constraint representation if and only if it satisfies Plott, Support Dependence and No Arbitrage.

*Proof.* **Plott (1973)** showed that the Plott axiom delivers  $\mathcal{M} \subseteq \Pi$  such that

$$S(A) = \bigcup_{\succ_m \in \mathcal{M}} \arg \max(\succ_m, A).$$

By Support Dependence, it suffices to show that some  $\mu \in \text{int}(\Delta(\Pi))$  explains choice on  $\{A \in \mathcal{F}(\mathcal{A}) : A = S(B) \text{ for some } B\}$ . Take any  $A, B$  such that  $A = S(B)$ . For each  $a \in A$ , there exists  $\succ_m \in \mathcal{M}$  such that  $a \succ_m B$ . Since  $A \subseteq B$ ,  $a \succ_m B$  implies  $a \succ_m A$ , so  $A = S(A)$ . Thus, it suffices to show that some  $\mu$  explains choice on  $\{A \in \mathcal{F}(\mathcal{A}) : A = S(A)\}$ .

This follows from No Arbitrage. To see why, index the preferences in  $\Pi$  from  $\succ_1$  to  $\succ_{|\Pi|}$ , and index the item-menu pairs in  $S$  from  $(a_1, A_1)$  to  $(a_{|S|}, A_{|S|})$ . Let  $X$  be the  $|S|$ -by- $|\Pi|$  matrix such that

$$X_{ij} = 1\{a_i \succ_j A_i\}$$

for each  $i, j$ . Let  $y$  be the  $|S|$ -length vector in which

$$y_i = \rho(a_i | A_i).$$

We need a strictly positive  $|\Pi|$ -length vector  $\mu$  such that  $X\mu = y$ . It is well known that such a  $\mu$  exists if and only if there is no vector  $\lambda$  such that  $X'\lambda \geq 0$  and  $y'\lambda \leq 0$  with at least one inequality strict. This is precisely the No Arbitrage condition.<sup>1</sup> Now we show that  $\mu$  has a unit sum. Fix any  $a \in \mathcal{A}$ . Since  $a \succ_i \{a\}$  for all  $i$  and since  $\{a\} = S(\{a\})$ , there is at least one row of  $X$  that consists entirely of ones. Since  $\rho(a | \{a\}) = 1$ , we have  $\sum_j \mu_j = 1$ .  $\square$

#### A.2.4 CONSTRAINED INFORMATION CHOICE

We extend the notion of a justification distribution to allow first-stage information choice to be less constrained than second-stage choice (without making it fully unconstrained). This is done by having each DM draw two different sets of justifications, one of which is weakly larger than the other. The larger set of justifications constrains information choice, while the smaller set constrains second-stage choice as usual.

**Definition 62** (Augmented justification distribution). *An augmented justification distribution is  $\nu \in \Delta(\mathfrak{U} \times \mathfrak{U})$  that satisfies the following conditions:*

1. *The marginal of  $\nu$  in its first dimension,  $\nu_1$ , is a justification distribution.*

---

<sup>1</sup>The application of No Arbitrage to random utility is not new; it was done by [Clark \(1996\)](#). Clark used a slightly different version of the axiom, which does not require  $\mu \in \text{int}(\Delta(\Pi))$ .

2. For each  $\mathcal{M}_1$  in  $\text{supp}(\nu_1)$ , the conditional of  $\nu$  on  $\mathcal{M}_1$ ,  $\nu_2(\cdot; \mathcal{M}_1)$ , satisfies

$$\nu_2(\{\mathcal{M}_2 : \mathcal{M}_2 \subseteq \mathcal{M}_1\}; \mathcal{M}_1) = 1.$$

The notion of self-knowledge can easily be extended to this setup by allowing private information about the second-stage justifications to depend on  $\mathcal{M}_1$ . Formally, a self-knowledge  $N$  must now satisfy

$$\int_{\tilde{v}} \tilde{v}(\cdot) N(d\tilde{v}; \mathcal{M}_1) = \nu_2(\cdot; \mathcal{M}_1).$$

Consider the following choice procedure, working backward from second-stage choice. Suppose the DM faces menu  $\{p, q\}$  at  $t = 2$  because he chose ignorance at  $t = 1$ . In that case, he maximizes his primary utility  $u$  over  $\mathcal{M}_2(\mathcal{M}_1^{\text{avoid}}(\{p, q\}))$ . Now suppose the DM faces menu  $\{p, q\}$  at  $t = 2$  because he was not offered information at  $t = 1$ . In that case, he maximizes  $u$  over  $\mathcal{M}_2(\{p, q\})$  as usual. Finally, suppose the DM faces menu  $\{p, q_i\}$  at  $t = 2$  because he chose information, or was compelled to receive information, at  $t = 1$ , and  $q_i$  was realized. Again, he maximizes  $u$  over  $\mathcal{M}_2(\{p, q_i\})$  as usual.

Now consider a DM who faces a choice between information and ignorance at  $t = 1$ . He knows  $\mathcal{M}_1$  and has self-knowledge  $N(\cdot; \mathcal{M}_1)$  about  $\mathcal{M}_2$ . If  $\mathcal{M}_1^{\text{avoid}}$  is empty, he acquires information. Otherwise, he maximizes his expected utility given  $N$  and the  $t = 2$  behavior spelled out above.

This version of the model connects the two extremes considered in Section 1.5.2. We recover the unconstrained-information case by setting  $\nu_1 := \delta_{\mathcal{U}}$ . We recover the fully-constrained information case by setting  $\nu_2(\cdot; \mathcal{M}_1) = \delta_{\mathcal{M}_1}$ .

Unlike the fully-constrained case, intermediate cases can predict information avoidance. To illustrate, consider a DM who knows both  $\mathcal{M}_1$  and  $\mathcal{M}_2$  at  $t = 1$ . Consider  $p, q_1, q_2$  such that  $q_1, q_2 \succ p$

and

$$\mathcal{M}_1^{\text{avoid}}(\{p, q\}) = \{p, q\}$$

$$\mathcal{M}_2(\{p, q\}) = \{p, q\}$$

$$\mathcal{M}_2(\{p, q_1\}) = \{p, q_1\}$$

$$\mathcal{M}_2(\{p, q_2\}) = \{p\}.$$

(This will happen if  $\mathcal{M}_1$ , but not  $\mathcal{M}_2$ , contains a preference such that  $q_1, q_2 \succsim_m p$ , while both  $\mathcal{M}_1$  and  $\mathcal{M}_2$  contain a preference such that  $q_1 \succ_m q \succ_m p \succ_m q_2$ .) If the DM avoids information at  $t = 1$ , he will get  $q$  for sure, since

$$\mathcal{M}_2(\mathcal{M}_1^{\text{avoid}}(\{p, q\})) = \{p, q\}.$$

If the DM acquires information at  $t = 1$ , he will have to choose  $p$  when  $q_2$  is realized. Thus, the DM will avoid information.

This is not to say that the DM will avoid information just as often as he would in the fully unconstrained case. Consider  $r_1, r_2$  such that  $r_1, r_2 \succ p$  and

$$\mathcal{M}_1^{\text{avoid}}(\{p, r\}) = \{p\}$$

$$\mathcal{M}_2(\{p, r\}) = \{p, r\}$$

$$\mathcal{M}_2(\{p, r_1\}) = \{p, r_1\}$$

$$\mathcal{M}_2(\{p, r_2\}) = \{p\}.$$

(This will happen if neither  $\mathcal{M}_1$  nor  $\mathcal{M}_2$  contains any preference such that  $r_1, r_2 \succsim_m p$ , but both contain some preference such that  $r_1 \succ_m r \succ_m p \succ_m r_2$ .) If the DM avoids information at  $t = 1$ ,

he will have to choose  $p$ . If he acquires information, he can choose  $r_1$  when it is realized. Thus, the DM will not avoid information. By contrast, he would avoid information if information choice were fully unconstrained, i.e. if  $\mathcal{M}_1$  were replaced with  $\mathcal{U}$ . He would get  $r$  for sure by remaining ignorant, but would have to choose  $p$  if he observed  $r_2$ .

Proposition 8 holds without modification in the augmented model. This is because no DM has an incentive to avoid information, so the constraints on information choice are not relevant.

Now consider Proposition 7. It is clear that the social planner will still prefer forcibly informing everyone to withholding information, since information choice does not play a role in that result. Moreover, the planner will still prefer providing information freely to withholding it, since some DMs will voluntarily become informed even if information choice is unconstrained (and even more DMs will make that choice if information choice is constrained). The only question is whether the planner will prefer forcibly informing everyone to withholding information—equivalently, whether some DMs will avoid information. This will certainly not be the case if information choice is fully constrained, so we will need to place additional restrictions on the augmented justification distribution  $\nu$ . A natural restriction is, for each  $\mathcal{M}_1^* \in \text{supp}(\nu_1)$ ,

$$\nu_2(\cdot; \mathcal{M}_1) = \nu_1(\cdot | \mathcal{M}_1 \subseteq \mathcal{M}_1^*). \quad (\text{A.24})$$

This restriction can be interpreted as follows: the DM draws  $\mathcal{M}_1$  from a standard justification distribution, and then draws  $\mathcal{M}_2$  from the same justification distribution, conditional on  $\mathcal{M}_2 \subseteq \mathcal{M}_1$ . Notice that the subset of DMs who draw  $\mathcal{M}_1 = \mathcal{U}$  will behave exactly like the population of DMs in Proposition 7. Since  $\nu_1(\mathcal{U}) > 0$ , there is a positive mass of such DMs. We showed in the proof of Proposition 7 that a positive mass of these DMs will avoid information. Thus, Proposition 7 holds given (A.24).

# B

## Appendix to Chapter 2

### B.1 PROOFS OF RESULTS IN TEXT

#### B.1.1 PROOF OF THEOREM 6

Notation:  $f(x)(y)$  is the probability that lottery  $f(x) \in \Delta(X \cup \{\diamond\})$  places on prize  $y \in X \cup \{\diamond\}$ .

$y_x \diamond$  is an act that delivers  $y$  for sure if  $x$  happens, and  $\diamond$  otherwise.

- Step 1: By Independence and Mixture Continuity, each  $\succsim_D$  has a linear representation  $U_D$  :

$\mathcal{F} \rightarrow \mathbb{R}$ .

$$nU\left(\frac{1}{n}f + \frac{n-1}{n}\diamond\right) = U(f) + (n-1)U(\diamond).$$

Normalizing  $U(\diamond) = 0$ ,

$$\begin{aligned} U_D(f) &= nU_D\left(\frac{1}{n}\sum_{x \in \mathcal{X}}(f(x))_x \diamond\right) \\ &= \sum_{x \in \mathcal{X}} U_D((f(x))_x \diamond) \\ &= \sum_{x \in \mathcal{X}} U_D\left(\sum_{y \in \mathcal{X}} f(x)(y) y_x \diamond\right) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} f(x)(y) U_D(y_x \diamond). \end{aligned}$$

For any constant act  $p$ ,

$$\begin{aligned} U_D(p) &= \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} U_D(x_y \diamond) \\ &= n \sum_{x \in \mathcal{X}} p(x) U_D\left(\frac{1}{n}x + \frac{n-1}{n}\diamond\right) \\ &= \sum_{x \in \mathcal{X}} p(x) U_D(x) \\ &= \mathbb{E}_p[U_D(x)]. \end{aligned}$$

It is convenient to rescale  $U_D$  so that  $\sum_{x \in \mathcal{X}} U_D(x) = 1$ . Since  $U_D$  is non-constant and weakly positive, this is always possible.

- Step 2: Setting  $D = \emptyset$ , we have

$$U_\emptyset(f) = \sum_{x \in \mathcal{X}} U_\emptyset((f(x))_x \diamond).$$



By Dynamic Consistency,

$$U_\emptyset(p_x \diamond) \geq U_\emptyset(q_x \diamond) \iff U_x(p) \geq U_x(q).$$

By cardinal uniqueness of EU representations, and the normalization  $U_\emptyset(\diamond) = U_x(\diamond) = 0$ ,

$$U_\emptyset(p_x \diamond) = P(x)U_x(p)$$

for some  $P(x) > 0$ . We have

$$\begin{aligned} U_\emptyset(y_x \diamond) &= P(x)U_x(y) \\ \implies \sum_y U_\emptyset(y_x \diamond) &= P(x) \\ \implies \sum_x \sum_y U_\emptyset(y_x \diamond) &= \sum_x P(x) \\ \implies \sum_y U_\emptyset(y) &= \sum_x P(x) \\ \implies 1 &= \sum_x P(x). \end{aligned}$$

Putting everything together,

$$U_\emptyset(f) = \sum_{x \in \mathcal{X}} P(x)U_x(f(x)) = \sum_{x \in \mathcal{X}} P(x)\mathbb{E}_{f(x)}U_x(y).$$

- Step 3: In this step, we will treat  $D$  as an ordered vector  $(d_1 \cdots d_n)$  rather than a set. We will show in the next step that the order does not matter.

Suppose that we have strictly positive probabilities  $P(D)$ ,  $\{P(Dx)\}_{x \in \mathcal{X}}$  such that

$$\sum_x P(Dx) = P(D)$$

$$P(D)U_D(f) = \sum_{x \in \mathcal{X}} P(Dx)U_{Dx}(f(x)) = \sum_{x \in \mathcal{X}} P(Dx)\mathbb{E}_{f(x)}U_{Dx}(y).$$

We will find strictly positive probabilities  $\{P(Dxy)\}_{(x,y) \in \mathcal{X}^2}$  such that

$$P(Dxy) = P(Dyx) \tag{B.1}$$

$$\sum_y P(Dxy) = P(Dx) \tag{B.2}$$

$$P(Dx)U_{Dx}(f) = \sum_{y \in \mathcal{X}} P(Dxy)U_{Dxy}(f(y)) = \sum_{x \in \mathcal{X}} P(Dxy)\mathbb{E}_{f(y)}U_{Dxy}(z). \tag{B.3}$$

We are looking for a strictly positive solution  $p$  to  $b = Ap$ , where

$$\begin{aligned}
 b &= \begin{bmatrix} P(Dx_1)U_{Dx_1}(x_1) \\ \vdots \\ P(Dx_1)U_{Dx_1}(x_n) \\ \vdots \\ P(Dx_n)U_{Dx_n}(x_1) \\ \vdots \\ P(Dx_n)U_{Dx_n}(x_n) \end{bmatrix} \\
 A &= \begin{bmatrix} U_{Dx_1x_1}(x_1) & \cdots & U_{Dx_1x_n}(x_1) & 0 & \cdots & 0 & 0 & \cdots & 0 \\ & \ddots & & & \ddots & & & \ddots & \\ U_{Dx_1x_1}(x_n) & \cdots & U_{Dx_1x_n}(x_n) & 0 & \cdots & 0 & 0 & \cdots & 0 \\ & \ddots & & & \ddots & & & \ddots & \\ 0 & \cdots & 0 & 0 & \cdots & 0 & U_{Dx_nx_1}(x_1) & \cdots & U_{Dx_nx_n}(x_1) \\ & \ddots & & & \ddots & & & \ddots & \\ 0 & \cdots & 0 & 0 & \cdots & 0 & U_{Dx_nx_1}(x_n) & \cdots & U_{Dx_nx_n}(x_n) \end{bmatrix} \\
 p &= \begin{bmatrix} P(Dx_1x_1) \\ \vdots \\ P(Dx_1x_n) \\ \vdots \\ P(Dx_nx_1) \\ \vdots \\ P(Dx_nx_n) \end{bmatrix}
 \end{aligned}$$

Since the  $k$ -th  $n$  rows of  $b$  sum to  $P(Dx_k)$ , and since the  $k$ -th  $n$  rows of  $Ap$  sum to  $\sum_y P(Dx_ky)$ ,

(B.2) will hold for any solution  $p$ . To guarantee (B.1), we add

$$0 = P(Dx_i x_j) - P(Dx_j x_i)$$

to the matrix equality for each  $i > j$ . (This is the reason for handling all the  $\{Dxy\}_{(x,y) \in \mathcal{X}^2}$  together, rather than each  $\{Dxy\}_{y \in \mathcal{X}}$  separately.) A version of Farkas' Lemma guarantees a weakly positive solution provided there is no

$$\begin{bmatrix} v \\ w \end{bmatrix} = \begin{bmatrix} v_1(x_1) \\ \vdots \\ v_1(x_n) \\ \vdots \\ v_n(x_1) \\ \vdots \\ v_n(x_n) \\ w_{12} \\ \vdots \\ w_{1n} \\ \vdots \\ w_{(n-1)n} \end{bmatrix}$$

such that

$$\begin{aligned} \sum_i P(Dx_i) \sum_y v_i(y) U_{Dx_i}(y) &< 0 \\ \forall i \quad \sum_y v_i(y) U_{Dx_i x_i}(y) &\geq 0 \\ \forall i > j \quad \sum_y v_i(y) U_{Dx_i x_j}(y) + w_{ij} &\geq 0 \\ \forall j < i \quad \sum_y v_j(y) U_{Dx_j x_i}(y) - w_{ij} &\geq 0. \end{aligned}$$

We use  $Dx_i x_j = Dx_j x_i$  to rewrite the last three lines:

$$\forall i, j \quad \sum_y \left( \frac{1}{2} v_i(y) + \frac{1}{2} v_j(y) \right) U_{Dx_i x_j}(y) \geq 0.$$

This gets rid of  $w$ . Now suppose there is a  $v$  that satisfies the above conditions:

$$\begin{aligned} \sum_x P(Dx) \sum_{y: v_x(y) > 0} v_x(y) U_{Dx}(y) &< \sum_x P(Dx) \sum_{y: v_x(y) < 0} (-v_x(y)) U_{Dx}(y) \\ \forall x, y \quad \frac{1}{2} \sum_{z: v_x(z) > 0} v_x(z) U_{Dxy}(z) + \frac{1}{2} \sum_{z: v_y(z) > 0} v_y(z) U_{Dxy}(z) &\geq \\ \frac{1}{2} \sum_{z: v_x(z) < 0} (-v_x(z)) U_{Dxy}(z) + \frac{1}{2} \sum_{z: v_y(z) > 0} (-v_y(z)) U_{Dxy}(z) & \end{aligned}$$

We can always take  $\sum_{y: v_i(y) > 0} v_i(y) \leq 1$  and  $\sum_{y: v_i(y) \leq 0} (-v_i(y)) \leq 1$ . This allows us to replace the weighted sums above with valuations of lotteries. (If  $\sum_{y: v_i(y) > 0} v_i(y) < 1$ , the

lottery will need to place positive weight on  $\diamond$ .) There must be some  $f, g \in \mathcal{F}$  such that

$$\sum_x P(Dx) \mathbb{E}_{f(x)} U_{Dx}(y) < \sum_x P(Dx) \mathbb{E}_{g(x)} U_{Dx}(y)$$

$$\forall x, y \quad \mathbb{E}_{\frac{1}{2}f(x) + \frac{1}{2}f(y)} U_{Dxy}(z) \geq \mathbb{E}_{\frac{1}{2}g(x) + \frac{1}{2}g(y)} U_{Dxy}(z).$$

By the inductive hypothesis and linearity of  $U_{Dxy}$ , this is equivalent to

$$P(D)U_D(f) < P(D)U_D(g)$$

$$\forall x, y \quad U_{Dxy} \left( \frac{1}{2}f(x) + \frac{1}{2}f(y) \right) \geq U_{Dxy} \left( \frac{1}{2}g(x) + \frac{1}{2}g(y) \right).$$

In terms of preferences,

$$f \prec_D g$$

$$\forall x, y \quad \frac{1}{2}f(x) + \frac{1}{2}f(y) \succsim_{Dxy} \frac{1}{2}g(x) + \frac{1}{2}g(y)$$

which contradicts Forward Exchangeability. We conclude that there is a weakly positive solution  $p$  to the matrix equality. It remains to show that each  $P(Dxy) > 0$ . Suppose some  $P(Dxy) = 0$ . Since

$$U_{Dx}(z_y \diamond) = P(Dxy)U_{Dxy}(z) = 0,$$

we must have  $z_y \diamond \sim_{Dx} \diamond$  for all  $z \in \mathcal{X}$ . But since  $z \succ_{Dxy} \diamond$  for some  $z$ , Dynamic Consistency implies  $z_y \diamond \succ_{Dx} \diamond$  for some  $z$ —contradiction.

- Step 4: By the Daniell extension theorem (see [Kallenberg \(2006a\)](#) Theorem 5.14 and Corol-

lary 5.15), there is a unique probability  $P$  on  $\mathcal{X}^\infty$  that satisfies

$$P((x_1, \dots, x_m) = (d_1, \dots, d_n)) = P(d_1 \cdots d_n).$$

We show that  $P$  is exchangeable:

$$P(d_1 \cdots d_{m+2}) = P(d_{\pi(1)} \cdots d_{\pi(m+2)})$$

for any permutation  $\pi$  of  $\mathcal{X}$ . Suppose that

$$P(d_1 \cdots d_{m+1}) = P(d_{\pi(1)} \cdots d_{\pi(m+1)}).$$

Recall that  $P(d_1 \cdots d_{m+2})$  depends on  $P(d_1 \cdots d_{m+1})$ ,  $U_{d_1 \cdots d_{m+1}x}$ , and  $U_{d_1 \cdots d_{m+1}y}$ . None of these depends on the order of the first  $m + 1$  elements, so  $P(d_1 \cdots d_{m+2})$  doesn't either. If  $\pi(m + 2) = m + 2$ , the last element is already correct, so we are done. Otherwise, we can reorder the first  $m + 1$  elements so that  $d_{\pi(m+2)}$  comes second-to-last (just before  $d_{m+2}$ ). By construction of  $P$ , we can switch the last two elements. The last element is now correct. Permuting the remaining items once again, we get the correct ordering.

- Step 5: We have constructed an exchangeable probability measure  $P$  on  $\mathcal{X}^\infty$ . The identity mapping defines a random sequence  $(\tilde{d}_1, \tilde{d}_2, \dots)$  in this space. We can define another random sequence as follows: for any  $(d_1, d_2, \dots) \in \mathcal{X}^\infty$ ,

$$(U_t(d_1, d_2, \dots))_{t=0}^\infty := (U_{\{d_1, \dots, d_t\}})_{t=0}^\infty.$$

Since

$$U_D = \sum_{x \in \mathcal{X}} \frac{P(Dx)}{P(D)} U_{Dx} = \mathbb{E}_P[U_{Dx}],$$

we have

$$U_t = \mathbb{E}_P[U_{t+1} | d_1, \dots, d_t].$$

This says  $U$  is a martingale with respect to the filtration induced by  $(\tilde{d}_1, \tilde{d}_2, \dots)$ . Since  $U_D(x) \in [0, 1]$  for all  $x \in \mathcal{X}$  and all  $D \in \mathcal{D}$ ,  $U$  is uniformly bounded, so uniformly integrable. By Doob's martingale convergence theorem (see [Kallenberg \(2006a\)](#) Theorem 6.21), there exists a random variable  $U_\infty$  such that  $U_t \rightarrow U_\infty$  almost surely, and

$$U_t = \mathbb{E}_P[U_\infty | d_1, \dots, d_t]. \tag{B.4}$$

- Step 6: We show that  $U_\infty$  is measurable with respect to the exchangeable  $\sigma$ -algebra, which contains precisely those sets invariant under finite permutations of  $\mathbb{N}$ . We have

$$\begin{aligned} U_t(d_1, \dots, d_t, d_{t+1}, \dots) &= U_{\{d_1, \dots, d_t\}} \\ &= U_{\{d_{\pi(1)}, \dots, d_{\pi(t)}\}} \\ &= U_t(d_{\pi(1)}, \dots, d_{\pi(t)}, d_{t+1}, \dots). \end{aligned}$$

The same argument can be used to show

$$U_s(d_1, \dots, d_t, d_{t+1}, \dots) = U_s(d_{\pi(1)}, \dots, d_{\pi(t)}, d_{t+1}, \dots)$$



for all  $s > t$ , so

$$\begin{aligned} U_\infty(d_1, \dots, d_t, d_{t+1}) &= \lim_t U_t(d_1, \dots, d_t, d_{t+1}, \dots) \\ &= \lim_t U_t(d_{\pi(1)}, \dots, d_{\pi(t)}, d_{t+1}, \dots) \\ &= U_\infty(d_{\pi(1)}, \dots, d_{\pi(t)}, d_{t+1}, \dots) \end{aligned}$$

as desired.

To complete the argument, we need de Finetti's theorem, stated in [Kallenberg \(2006b\)](#) Theorem 1.1 and Proposition 1.4. We also need [Kallenberg \(2006b\)](#) Corollary 1.6, which says that the  $\sigma$ -algebra generated by the directing random measure  $\mu$  of  $P$  is precisely the exchangeable sigma-algebra. Since  $U_\infty$  is measurable with respect to the latter, we can write  $U_\infty(\mu)$  instead of  $U_\infty(d_1, d_2, \dots)$ . Recall that  $U_\infty(\mu)$  is a vector of values, one for each  $x \in \mathcal{X}$ . We write  $u(x, \mu)$  to pick out the value corresponding to  $x$ .  $u(\cdot, \mu)$  might not be defined for some zero-measure subset of  $\mu$ , but we can fill it in arbitrarily to get a function  $u : \mathcal{X} \times \Delta(\mathcal{X}) \rightarrow \mathbb{R}$ . Substituting this into (B.4), we have

$$U_t(x) = \mathbb{E}_P [u(x, \mu) | d_1, \dots, d_t].$$

This can be rewritten as

$$U_D(x) = \mathbb{E}_P [u(x, \mu) | D].$$

Since  $U_D$  is linear,

$$U_D(p) = \mathbb{E}_P \left[ \sum_x p(x) u(x, \mu) | D \right].$$

We can let

$$u(p, \mu) := \sum_x p(x) u(x, \mu).$$

Finally,

$$\begin{aligned}
U_D(f) &= \sum_x P(x|D) U_{Dx}(f(x)) \\
&= \sum_x P(x|D) \mathbb{E}_P [u(f(x), \mu) | D] \\
&= \mathbb{E}_P [u(f(x), \mu) | D].
\end{aligned}$$

### B.1.2 PROOF OF COROLLARY 5

By cardinal uniqueness of Bernoulli utilities, we have

$$V_D(x) - V_D(\diamond) = A(D) (U_D(x) - U_D(\diamond)) \quad (\text{B.5})$$

for some  $A : \mathcal{D} \rightarrow \mathbb{R}_{++}$ . This implies

$$A(D) = \frac{v(p, \mu) - v(\diamond, \mu)}{u(p, \mu) - u(\diamond, \mu)}$$

where  $p$  is the constant act that delivers an equally weighted lottery over  $X$  with certainty. The numerator is bounded above because  $v$  is. For the denominator, note that  $u(p, \mu) > u(\diamond, \mu)$  for all  $\mu$  because  $u(\cdot, \mu)$  is non-constant and minimized at  $\diamond$ . Since  $u$  is continuous in  $\mu$  and  $\Delta(X)$  is compact,  $u(p, \mu) - u(\diamond, \mu)$  must have a strictly positive maximum. Thus,  $A(D)$  is bounded. A similar argument establishes that  $1/A(D)$  is bounded.

Recall the construction of the  $P(D)$  in the proof of Theorem 6. Given Full Rank, the  $P(D)$  and  $Q(D)$  are unique, and are related by

$$\frac{A(\emptyset)}{A(D)} P(D) = Q(D). \quad (\text{B.6})$$

Plugging this into

$$Q(D) = \sum_x Q(Dx),$$

we have

$$\frac{1}{A(D)} = \sum_x \frac{1}{A(Dx)} \frac{P(Dx)}{P(D)};$$

$\frac{1}{A}$  is a martingale under  $P$ . A similar argument establishes that  $A$  is a martingale under  $Q$ . Since  $1/A$  and  $A$  are bounded, we can apply martingale convergence:

$$\begin{aligned} A(D) &= \mathbb{E}_Q[A_\infty | D] \\ \frac{1}{A(D)} &= \mathbb{E}_P \left[ \frac{1}{A_\infty} \middle| D \right]. \end{aligned}$$

Just as in the proof of Theorem 6, we can pass from limiting values to functions of  $\mu$ :

$$\begin{aligned} A(D) &= \mathbb{E}_Q [A(\mu) | D] \\ \frac{1}{A(D)} &= \mathbb{E}_P \left[ \frac{1}{A(\mu)} \middle| D \right] \\ &= \int \frac{1}{A(\mu)} \frac{\mu(D)}{P(D)} dP(\mu) \end{aligned}$$

Substituting (B.6) into the last equality,

$$Q(D) = \int \mu(D) \frac{A(\emptyset)}{A(\mu)} dP(\mu). \tag{B.7}$$

The de Finetti decomposition of  $Q$  is

$$Q(D) = \int \mu(D) dQ(\mu).$$

By [Kallenberg \(2006b\)](#) Proposition 1.4, the distribution over  $\mu$  that generates the  $Q(D)$  is unique.

Comparing the de Finetti decomposition with (B.7), we have

$$dQ(\mu) = \frac{A(\emptyset)}{A(\mu)} dP(\mu).$$

We can use the martingale property of  $A$  to replace  $A(\emptyset)$  with an expectation.

To complete the argument, let  $\tilde{U}_D(x) = U_D(x) - U_D(\diamond)$ . We have

$$\tilde{U}_D(x) = \mathbb{E}_P[\tilde{u}(x, \mu) | D] = \int \tilde{u}(\mu) \frac{\mu(D)}{P(D)} dP(\mu).$$

where  $\tilde{u}(x, \mu) = u(x, \mu) - u(\diamond, \mu)$ . Substituting  $dQ$  for  $dP$ ,

$$\tilde{U}_D = \frac{1}{P(D)} \int \tilde{u}(\mu) \mu(D) \frac{A(\mu)}{A(\emptyset)} dQ(\mu).$$

Substituting  $\tilde{V}_D$  for  $\tilde{U}_D$  and  $Q$  for  $P$ ,

$$\tilde{V}_D = \int A(\mu) \tilde{u}(\mu) \frac{\mu(D)}{Q(D)} dQ(\mu) = \mathbb{E}_Q[A(\mu) \tilde{u}(\mu) | D].$$

$\tilde{V}_D$  must also satisfy

$$\tilde{V}_D = \mathbb{E}_Q[\tilde{v}(\mu) | D].$$

By uniqueness of the limiting values,

$$\tilde{v}(\mu) = A(\mu) \tilde{u}(\mu).$$

### B.1.3 PROOF OF PROPOSITION 9

If  $\Delta u(x, \mu)$  is increasing in  $\mu$ , then  $u$  exhibits decreasing differences. We show that  $\mathbb{E}_P[u(x, \mu)]$  inherits this property. Fix  $y > x$  and  $\nu >_{FOSD} \mu$ . Let  $f_P(x)$  denote the distribution at quantile  $x$  of  $P$ , and likewise for  $Q$ . We have

$$\begin{aligned} u(y, \nu) - u(x, \nu) &\geq u(y, \mu) - u(x, \mu) \\ \forall z \in [0, 1] \quad u(y, f_P(z)) - u(x, f_P(z)) &\geq u(y, f_Q(z)) - u(x, f_Q(z)) \\ \int_0^1 u(y, f_P(z)) - u(x, f_P(z)) dx &\geq \int_0^1 u(y, f_Q(z)) - u(x, f_Q(z)) dz \\ \mathbb{E}_P[u(y, \nu) - u(x, \nu)] &\geq \mathbb{E}_Q[u(y, \nu) - u(x, \nu)]. \end{aligned}$$

The result follows from Topkis's monotonicity theorem.

### B.1.4 PROOF OF PROPOSITION 10

Let

$$U(x, Q) = \mathbb{E}_Q[u(x, \mu)].$$

For brevity, let

$$P = \sum_i \lambda_i P_i.$$

We show

$$\Delta U \left( \sum_i \lambda_i x_i^*, P \right) \geq \sum_i \lambda_i \Delta U(x_i^*, P_i). \quad (\text{B.8})$$

Notice that

$$\begin{aligned}
U(x, P) &= \mathbb{E}_P[u(x, \mu)] \\
&= \sum_i \lambda_i \mathbb{E}_{P_i}[u(x, \mu)] \\
&= \sum_i \lambda_i U(x, P_i)
\end{aligned}$$

so it suffices to show

$$\sum_i \lambda_i \left( \Delta U \left( \sum_j \lambda_j x_j^* \right) - \Delta U(x_i^*, P_i) \right) \geq 0.$$

Since  $\Delta u(x, \mu)$  is concave in  $x$ , so is  $\Delta U(x, P)$ . The left-hand-side is weakly greater than

$$\begin{aligned}
&\sum_i \lambda_i \left[ \sum_j \lambda_j \Delta U(x_j^*, M_i) - U(x_i^*, M_i) \right] \\
&= \sum_i \lambda_i \left[ \sum_{j \neq i} \lambda_j \left( \Delta U(x_j^*, M_i) - \Delta U(x_i^*, M_i) \right) \right] \\
&= \sum_i \lambda_i \sum_{j < i} \lambda_j \left[ \Delta U(x_i^*, M_j) + \Delta U(x_j^*, M_i) - \Delta U(x_i^*, M_i) - \Delta U(x_j^*, M_j) \right].
\end{aligned}$$

By Proposition 9,  $M_i >_{FOSD} M_j$  implies  $x_i^* \geq x_j^*$ . Since  $\Delta U(x, P)$  inherits decreasing differences from  $\Delta u(x, \mu)$ ,

$$\Delta U(x_i^*, M_j) + \Delta U(x_j^*, M_i) - \Delta U(x_i^*, M_i) - \Delta U(x_j^*, M_j) \geq 0.$$

Plugging this in delivers (B.8).

Since  $x_i^*$  is optimal for  $M_i$ ,

$$\Delta U(x_i^*, M_i) > 0$$

for all  $\Delta > 0$ . Plugging this into (B.8), we have

$$\Delta U \left( \sum_i \lambda_i x_i^*, M \right) > 0$$

for all  $\Delta > 0$ . Given beliefs  $M$ ,  $\sum_i \lambda_i x_i^*$  is strictly better than any  $x < \sum_i \lambda_i x_i^*$ . We conclude that  $x^* \geq \sum_i \lambda_i x_i^*$ .

#### B.1.5 PROOF OF COROLLARY 6

For the first part: the MLRP implies  $P|Dx >_{FOSD} P|Dy$  if  $x > y$ . The result follows from Proposition 9.

The second part is a direct implication of Proposition 10; just replace  $P_i$  with  $P|Dx$  and  $\lambda_i$  with  $P(x|D)$ .

#### B.1.6 PROOF OF PROPOSITION 11

We show how to transform a social model into an observationally equivalent pure learning model.

For each  $(T, u)$ , let

$$\hat{u}_{T,u}(x, \theta) := u(x, \theta, T \circ (x^*)^{-1}).$$

Clearly,  $\hat{x}^*(\theta; \hat{u}_{T,u}) = x^*(\theta; T, u)$ . Moreover,  $\hat{u}$  inherits the required properties in Assumption 2.

Now let

$$\hat{P}((T, \hat{u}_{T,u}) | (T, u) \in B) := P((T, u) \in B).$$

We have

$$\int_{T, \hat{u}} \hat{u}(x, \theta) \prod_{d \in D} \frac{d}{dd} T \circ (\hat{x}^*)^{-1}(d; \hat{u}) d\hat{P}(T, \hat{u}) = \int_{T, u} u(x, \theta, T \circ (x^*)^{-1}) \prod_{d \in D} \frac{d}{dd} T \circ (x^*)^{-1}(d; T, u) dP(T, u)$$

so preferences are the same across the two models.

### B.1.7 PROOF OF PROPOSITION 12

First, consider reversals  $R_p$ .

Since  $U_D$  depends on beliefs about  $T$ , we write  $U_D^p$  for clarity. No index is needed for  $U_\emptyset$ , though.

In the pure learning model,  $T$  is only used for interpreting the evidence about  $u$  contained in  $D$ .

Uncertainty about  $T$  is therefore irrelevant when  $D = \emptyset$ .

$P(D)\Delta U_D$  is linear in beliefs about  $T$ :

$$\begin{aligned} P(D)\Delta U_D^p &= \int_{T, u} \Delta u f(D|T, u) dP(T, u) \\ &= \int_T \left[ \int_u \Delta u f(D|T, u) P(dU|T) \right] P(dT) \\ &= \sum_i \lambda_i \int_T \left[ \int_u \Delta u f(D|T, u) P(dU|T) \right] P_i(dT) \\ &= \sum_i \lambda_i \int_{T, u} \Delta u f(D|T, u) dP_i(u, T) \\ &= \sum_i \lambda_i \left[ P_i(D)\Delta U_D^{p_i} \right]. \end{aligned}$$

Since  $1_{(z \leq 0)}|z|$  is convex and  $\Delta U_D \Delta U_\emptyset$  is linear,

$$1_{P(D)\Delta_D^p \Delta_\emptyset \leq 0} |P(D)\Delta_D^p \Delta_\emptyset|$$



is convex. We can rewrite this function as

$$P(D)1_{\Delta_D^p \Delta_\emptyset \leq 0} |\Delta_D^p \Delta_\emptyset|.$$

Summing over all  $D$  of size  $n$  gives the desired result.

Now consider movement  $M_p$ . We can show linearity as above:

$$\begin{aligned} P(D) (U_D^p(x) - U_\emptyset(x)) &= \sum_i \lambda_i P_i(D) U_D^{p_i}(x) - \sum_i \lambda_i P_i(D) U_\emptyset \\ &= \sum_i \lambda_i P_i(D) (U_D^{p_i}(x) - U_\emptyset(x)). \end{aligned}$$

By convexity of the absolute value function,

$$P(D) |U_D^p(x) - U_\emptyset(x)| \leq \sum_i \lambda_i P_i(D) |U_D^{p_i}(x) - U_\emptyset(x)|.$$

Again, summing over all  $D$  of size  $n$  gives the result.

# C

## Appendix to Chapter 3

### C.1 DOMAIN OF CHOICE

The domain of choice is similar to the space of consumption trees in [Chew and Epstein \(1991\)](#) (CE). I briefly outline their construction below. For a metric space  $Y$ , let  $\Delta(Y)$  be the space of Borel probability measures on  $Y$  endowed with the weak convergence topology. Let the space of payoffs  $C$  be an interval of  $\mathbb{R}$  (open, closed, or neither) with  $0 \in C$ . (CE allow for any separable metric space,

but we will not need the additional generality.) Let  $\bar{D}_{-1} = \{(c_0, c_1, \dots) : \forall t \ c_t \in C\}$ , and then for  $t \in \{0, 1, \dots\}$ ,

$$\bar{D}_t = \Delta(C \times \bar{D}_{t-1}).$$

Thus,  $\bar{D}_0$  is a measure over deterministic consumption streams,  $\bar{D}_1$  is a measure over a time-0 consumption and an element of  $\bar{D}_0$ , and so on.  $\bigcup_0^\infty \bar{D}_t$  contains consumption trees for which all uncertainty is resolved by some finite date. To add consumption trees for which uncertainty is never resolved, CE proceed as follows. For any tree  $d^t \in \bar{D}_t$ , it's possible to construct  $d^{t-1} \in \bar{D}_{t-1}$  by resolving at  $t - 1$  in  $d^{t-1}$  all the uncertainty resolved at  $t$  in  $d^t$ . (CE make this formal.) CE define an element of their domain  $\bar{D}$  as a sequence  $(d^0, d^1, \dots)$  where each  $d^t \in \bar{D}_t$ , and each  $d^{t-1}$  can be constructed from  $d^t$  as mentioned above. CE show that  $\bar{D}$  is a separable metric space such that  $\bar{D}$  is homeomorphic to  $\Delta(C \times \bar{D})$ .

For this paper, the domain will be  $D \subseteq \bar{D}$ . I do not require  $D = \bar{D}$  because this would force the Bernoulli utility function that appears in the representation to be bounded. However, I do require  $D$  to satisfy several conditions.

Explaining these conditions requires some additional notation. First, we need a convenient way of referring to certain deterministic consumption streams. For any  $a, b \in C$ , denote a tree that delivers  $a$  for the first  $t$  periods and  $b$  forever after as  $(a^t, b^\infty)$ . Denote a tree that delivers  $a$  in every period as  $a^\infty$ . Second, we need the notion of a subtree. A  $t$ -subtree  $d_t$  of  $d$  is a consumption tree that can be obtained by conditioning on a particular path through  $d$  from period 0 to  $t - 1$ . Said another way,  $d_t$  is a tree an agent could face at  $t$  if he started with  $d$  at 0. Third, we need a convenient way of referring to subtrees that consist of measures over constant consumption streams. There is a clear homeomorphism between the set of these subtrees and the set of measures over payoffs. We will call any such subtree a "lottery," and we will identify it with the cdf  $F$  of the corresponding element of  $\Delta(C)$ . It is a "simple lottery" if  $F$  has finite range.

Now we can state the conditions on  $D$ . The first pertains to deterministic consumption streams: (1)  $D$  must contain all trees of the form  $(0^t, a, b^\infty)$ . The next three conditions pertain to lotteries resolved at  $t = 0$ . They are necessary to obtain a standard expected utility representation over such lotteries. (2)  $D$  must contain every simple lottery  $F$  such that  $\text{supp}(F) \subset \text{int}(C)$ . (3) If  $D$  contains the lotteries  $F$  and  $F'$ , it must also contain  $\lambda F + (1 - \lambda)F'$  for all  $\lambda \in (0, 1)$ . (4) If  $D$  contains the lottery  $F$ , then it must also contain  $F^a$  and  $F_a$  for all  $a \in C$ , where  $F^a$  ( $F_a$ ) is constructed from  $F$  by reassigning all mass above (below)  $a$  to  $a$ .

The next two conditions pertain to lotteries resolved at  $t > 0$ . They are necessary to obtain an expected utility representation with distorted probabilities. (5)  $D$  must contain every tree of the form  $(0^t, F)$  where  $F$  is a binary lottery  $\text{supp}(F) \subset \text{int}(C)$ . (6) if  $D$  contains  $(0^t, F)$ , it must also contain  $F_{\leq F(c)}$  (the lottery  $F$  conditional on getting no more than  $c$ ) and  $F_{> F(c)}$  (the lottery  $F$  conditional on getting more than  $c$ ) for all  $c$  such that  $F(c), 1 - F(c) > 0$ .

The final two conditions pertain to subtrees that are not necessarily lotteries. They are necessary to obtain a recursive representation. Specifically, they allow us to arrive at the certainty equivalent of a tree by working backwards from any period  $t$  (without worrying that the intermediate trees constructed during this process will fall outside  $D$ ). (7) If  $d_t$  is a  $t$ -subtree of some  $d \in D$ , then  $(0^t, d_t) \in D$ . (8)  $D$  must have a recursive structure (defined below).

**Definition 63.**  $D$  has a recursive structure if the following conditions hold.

1. If for some  $t$ , for each  $t$ -subtree  $d_t$  of  $d$ , there exists  $c$  such that

$$(0^t, c^\infty) \sim (0^t, d_t),$$

then  $d' \in D$ , where  $d'$  is constructed from  $d$  by replacing each  $d_t$  with the corresponding  $c^\infty$ .

2. Take  $d = (0^t, d_t)$  where  $d_t$  is a measure over deterministic sequences of the form  $(c_t, c_{t+1}^\infty)$ . If

for each  $(c_t, c_{t+1}^\infty) \in \text{supp}(d_t)$  there exists  $c$  such that

$$(0^t, c_t, c_{t+1}^\infty) \sim (0^t, c^\infty),$$

then  $d' \in D$ , where  $d'$  is constructed from  $d$  by replacing each  $(c_t, c_{t+1}^\infty)$  with the corresponding  $c^\infty$ .

## C.2 PROBABILITY DISTORTIONS

Lemma 25 provides a necessary and sufficient condition for  $\pi_t$  to be a probability distortion function for  $u$ .

**Lemma 25.** *Suppose  $\pi_t$  is continuous and satisfies  $\pi_t(0, c_L, c_H) = 0$ ,  $\pi_t(1, c_L, c_H) = 1$ , and  $\forall q \in (0, 1)$   $\pi_t(q, c_L, c_H) \in (0, 1)$ . Then the following condition is necessary and sufficient for  $\pi_t$  to be a probability distortion function for  $u$ .*

$$\begin{aligned} \forall \Delta \in [0, 1 - q], c \in (c_L, c_H) \\ \pi_t \left( q + \Delta, u^{-1} \left( \frac{q}{q + \Delta} u(c_L) + \frac{\Delta}{q + \Delta} u(c) \right), c_H \right) \\ \geq \pi_t \left( q, c_L, u^{-1} \left( \frac{\Delta}{1 - q} u(c) + \frac{1 - q - \Delta}{1 - q} u(c_H) \right) \right). \end{aligned}$$

This condition is still somewhat difficult to interpret. By imposing a differentiability assumption, we arrive at something more tractable. The condition in Proposition 25 was used to generate the parametric special cases of  $\pi_t$  in Section 3.2.2.

**Proposition 25.** *Suppose  $\pi_t$  is continuous and satisfies  $\pi_t(0, c_L, c_H) = 0$  and  $\pi_t(1, c_L, c_H) = 1$ . Suppose further that  $\pi_t$  and  $u$  are differentiable. Then the following condition is necessary and sufficient*

for  $\pi_t$  to be a probability distortion function for  $u$ :

$$\pi_{t,1}(q, c_L, c_H) + \min \left\{ \frac{\pi_{t,2}(q, c_L, c_H)}{u'(c_L)q}, \frac{\pi_{t,3}(q, c_L, c_H)}{u'(c_H)(1-q)} \right\} (u(c_H) - u(c_L)) \geq 0.$$

If  $\pi_t$  takes the form

$$\frac{\pi_t(q, c_L, c_H)}{1 - \pi_t(q, c_L, c_H)} = g(u(c_L), u(c_H)) \frac{q}{1 - q}$$

then the condition reduces to

$$g(u_L, u_H) + \min \{0, g_1(u_L, u_H), g_2(u_L, u_H)\} (u(c_H) - u(c_L)) \geq 0.$$

### C.3 PROOFS OF RESULTS IN TEXT

#### C.3.1 PROOF OF LEMMA 25

Necessity: Define  $F$  as follows:

$$F(x) \begin{cases} 0 & x < c_L \\ q & x \in [c_L, c) \\ q + \Delta & x \in [c, c_H) \\ 1 & x \geq c_H \end{cases}.$$

If the inequality in the lemma is reversed, we have  $\hat{F}_t(F)(c) < \hat{F}_t(F)(c_L)$  even though  $c > c_L$ . Thus,  $\pi_t$  is not a probability distortion function.

Sufficiency: Take  $z' > z$ , and suppose that  $\hat{F}_t(F)(z') < \hat{F}_t(F)(z)$ . We have

$$\begin{aligned} & \pi_t \left( F(z'), u^{-1} \left( \frac{1}{F(z')} \int_{\leq z'} u(x) dF(x) \right), u^{-1} \left( \frac{1}{1 - F(z')} \int_{> z'} u(x) dF(x) \right) \right) \\ & < \pi_t \left( F(z), u^{-1} \left( \frac{1}{F(z)} \int_{\leq z} u(x) dF(x) \right), u^{-1} \left( \frac{1}{1 - F(z)} \int_{> z} u(x) dF(x) \right) \right). \end{aligned}$$

Letting  $q = F(z)$ ,  $\Delta = F(z') - F(z) \geq 0$ ,

$$\begin{aligned} c_L &= u^{-1} \left( \frac{1}{F(z)} \int_{\leq z} u(x) dF(x) \right), \\ c_H &= u^{-1} \left( \frac{1}{1 - F(z')} \int_{> z'} u(x) dF(x) \right), \\ \text{and } c &= u^{-1} \left( \frac{1}{F(z') - F(z)} \int_{z < x \leq z'} u(x) dF(x) \right), \end{aligned}$$

we have a violation of the inequality in the lemma.

### C.3.2 PROOF OF PROPOSITION 20

**Definition 64** (Recursive certainty equivalent functions). *A function  $I_t$  from the space of  $t$ -subtrees to  $C$  is a recursive certainty equivalent function of  $\succeq$  if, for any tree  $d \in D$ ,  $d' \sim d$ , where  $d'$  is constructed from  $d$  by replacing each  $t$ -subtree  $d_t$  with  $I_t(d_t)^\infty$ .*

First, we show that the axioms imply the intermediate representation

$$\begin{aligned} I_t(d) &= \mu_t(F_t(d)) \\ F_t(d)(B) &= d\{(c, d') : W(c, I_{t+1}(d')) \in B\}. \end{aligned}$$

By Recursive Certainty Equivalence and History/Counterfactual Irrelevance for trees,  $\succeq$  has recursive certainty equivalent functions  $I_0, I_1, \dots$ . By Monotonicity, each  $I_t$  is unique.

Define  $W : C \times C \rightarrow C$  by  $W(a, b) = I_0(a, b^\infty)$ . For any  $t$ , for any  $d'$  such that  $I_{t+1}(d') = z$ ,  $W(c, z) = I_t(c, d')$ . To see why:

$$\begin{aligned} W(c, I_{t+1}(d')) &= I_0(c, I_{t+1}(d')^\infty) \text{ by definition of } W \\ &= I_t(c, I_{t+1}(d')^\infty) \text{ by No Savoring} \\ &= I_t(c, d') \text{ by definition of } I_{t+1}. \end{aligned}$$

$W$  is increasing in both arguments because of Monotonicity. For continuity of  $W$ , we need  $\{(a, b) : W(a, b) \in B\}$  to be open for any open  $B$ . Equivalently, we need  $\{(a, b) : I_0(a, b^\infty) \in B\}$  to be open. Rewriting once more,  $\{(a, b) : \exists c \in B (a, b^\infty) \sim c\}$  must be open. This is precisely what Deterministic Continuity says.

We can use  $\{I_0, I_1, \dots\}$  to define  $F_t(d)$  by

$$F_t(d)(B) = d\{(c, d') : I_t(c, d') \in B\}.$$

We showed above that  $I_t(c, d') = W(c, I_{t+1}(d'))$ . This gives the desired definition of  $F_t$ .

Now we define  $\mu_t$ . If  $F = F_t(d_t)$  for some lottery  $d_t$  such that  $(0^t, d_t) \in D$ , let  $\mu_t(F) = I_t(d_t)$ . (There cannot be two distinct lotteries  $d$  and  $d'$  with  $F_t(d) = F_t(d')$ , so  $\mu_t$  is well defined.) Now take some  $F = F_t(f_t)$  where  $f_t$  is not a lottery, but  $(0^t, f_t) \in D$ . Notice that there exists a lottery  $d_t$  such that  $F = F_t(d_t)$  and  $(0^t, d_t) \in D$ . (To construct it, replace each  $(c, f_{t+1})$  in the support of  $f$  with  $I_t(c, f_{t+1})^\infty$ .) We now show  $I_t(d_t) = I_t(f_t)$ . Since  $F_t(d_t) = F_t(f_t)$ , we know that for all  $B$ ,

$$d_t\{c^\infty : c \in B\} = f_t\{(c, f_{t+1}) : W(c, I_{t+1}(f_{t+1})) \in B\}.$$

Construct  $f'_t$  from  $f_t$  by replacing each  $(c, f_{t+1})$  in the support of  $f_t$  with  $(c, I_{t+1}(f_{t+1})^\infty)$ . By defini-



tion of  $I_{t+1}$ ,  $(0^t, f'_t) \sim (0^t, f_t)$ . We have

$$d_t\{c^\infty : c \in B\} = f'_t\{(c, z^\infty) : W(c, z) \in B.\}.$$

Now construct  $f''_t$  from  $f'_t$  by replacing each  $(c, z^\infty)$  in the support of  $f'_t$  with  $W(c, z)^\infty$ . By definition of  $W$  and Counterfactual Irrelevance for streams,  $(0^t, f'_t) \sim (0^t, f''_t)$ . We have

$$d_t\{c^\infty : c \in B\} = f''_t\{c^\infty : c \in B\}.$$

This says  $d_t = f''_t$ , so of course  $(0^t, d_t) \sim (0^t, f''_t) \sim (0^t, f_t)$ . This implies  $(0^t, I_t(d_t)^\infty) \sim (0^t, I_t(f_t)^\infty)$ . The desired result  $I_t(d_t) = I_t(f_t)$  follows from Monotonicity. We conclude that  $\mu_t(F) = I_t(d_t)$  for any  $F$  such that  $F = F_t(d_t)$  for some  $d_t$  with  $(0^t, d_t) \in D$  (equivalently, for some  $t$ -subtree  $d_t$  of any  $d \in D$ ).

We now show that  $\mu_0$  is expected utility by showing that the agent's preference over lotteries satisfies the six conditions in Theorem 3.6 of [Wakker \(1993\)](#). (Specifically, we use the version with a rich domain closed under convex combinations. This version requires independence for all distributions, not just simple ones, and imposes a weaker conditional monotonicity condition.) Four of these conditions don't require much discussion. The restriction of  $\succeq$  to lotteries is clearly (1) a weak order. Initial Independence implies (2) independence. Recursive Certainty Equivalence implies (5) step equivalence, and Truncation Continuity implies Wakker's (6) truncation continuity in the presence of RCE.

(3) Step-vNM-continuity requires

$$F \succ F' \succ F'' \Rightarrow \exists p, q \in (0, 1) pF + (1-p)F'' \succ F' \succ qF + (1-q)F''.$$

By Recursive Certainty Equivalence, we can always find  $c, c'$  such that  $c \sim F$  and  $c' \sim F'$ . By Initial

Independence (applied twice),

$$\forall \alpha \in (0, 1) \alpha c \oplus (1 - \alpha)c'' \sim \alpha F + (1 - \alpha)F''.$$

Thus, it suffices to show

$$c \succ c' \succ c'' \Rightarrow \exists p, q \in (0, 1) pc \oplus (1 - p)c'' \succ c' \succ qc \oplus (1 - q)c''.$$

Let  $B = \{c : c \succ c'\}$ . By Monotonicity, this is just  $\{c : c > c'\}$ , so it is an open set. Binary Lottery Continuity then implies that  $\{pa \oplus (1 - p)b : \exists z \succ c' pa \oplus (1 - p)b \sim z\}$  is open. Since  $\delta_c$  is in this set,  $pc \oplus (1 - p)c''$  will also be in the set for  $p$  sufficiently close to 1. An analogous argument establishes the second half of the desired condition.

(4) Conditional monotonicity requires, for any lottery  $F$  and any simple lottery  $F'$ ,

$$\Pr_F(\{z : z \succeq F'\}) = 1 \Rightarrow F \succeq F'$$

$$\Pr_F(\{z : z \preceq F'\}) = 1 \Rightarrow F \preceq F'.$$

In the presence of RCE, it suffices to consider  $F' = \delta_{c'}$ . Given  $c \succeq c' \Leftrightarrow c \geq c'$  (one implication of Monotonicity), the above can be rewritten

$$\forall c' < c F(c') = 0 \Rightarrow F \succeq c$$

$$F(c) = 1 \Rightarrow F \preceq c.$$

This is another implication of Monotonicity.

Wakker's Theorem 3.6 delivers an EU representation with  $u$  cardinally unique and with all integrals finite. Monotonicity implies  $u$  strictly increasing. Moreover,  $u$  must be continuous. Suppose it

is not, and has  $u(a) > \lim u(a^n)$  for some increasing  $a^n \rightarrow a$ . By RCE, we can always find  $c^n$  such that  $c^n \sim (1/2)a \oplus (1/2)a^n$ . This is equivalent to  $u(c^n) = (1/2)u(a) + (1/2)u(a^n)$ , which implies  $\lim u(c^n) = (1/2)u(a) + (1/2)\lim u(a^n)$ . For  $m$  large, then, we must have  $u(c^m) > \lim u(a^n)$ . Fix some  $m$  satisfying this condition. We know  $c^m < a$ , so we must have  $a^n > c^m$  for  $n$  large enough. This implies  $\lim u(a^n) \geq u(c^m)$ , which generates a contradiction. A parallel argument rules out  $u(a) < \lim u(a^n)$  for some decreasing  $a^n \rightarrow a$ .

We now show that  $\mu_t$  is expected utility with appropriately distorted probabilities and the same Bernoulli utility as period 0. First, we define the probability distortion function:

$$\pi_t(q, c_H, c_L) = \frac{u(c_H) - u(\mu_t(qc_L \oplus (1-q)c_H))}{u(c_H) - u(c_L)}.$$

Clearly,  $\pi_t$  satisfies  $\pi_t(0, c_L, c_H) = 0$  and  $\pi_t(1, c_L, c_H) = 1$ . Monotonicity implies  $c_L < \mu_t(qc_L \oplus (1-q)c_H) < c_H$  for  $q \in (0, 1)$ , so  $\pi_t(q, c_L, c_H) \in (0, 1)$  if  $q \in (0, 1)$ . For continuity of  $\pi_t$ , continuity of  $u$  (already established) and continuity of  $\mu_t(qc_L \oplus (1-q)c_H)$  as a function of  $(q, c_L, c_H)$  will suffice. The second condition follows from Binary Lottery Continuity, which implies that  $\{qc_L \oplus (1-q)c_H : \mu_t(qc_L \oplus (1-q)c_H) > c\}$  and  $\{qc_L \oplus (1-q)c_H : \mu_t(qc_L \oplus (1-q)c_H) < c\}$  are open sets.

Now, define  $\hat{F}_t(F)$  as follows:

$$\begin{aligned} \hat{F}_t(F)(c) &= \pi_t(F(c), I_0(F_{\leq c}), I_0(F_{> c})) \\ &= \pi_t\left(F(c), u^{-1}\left(\frac{1}{F(c)} \int_c^c u(x)dF(x)\right), u^{-1}\left(\frac{1}{1-F(c)} \int_c^{\bar{c}} u(x)dF(x)\right)\right). \end{aligned}$$

$\hat{F}_t$  will satisfy

$$I_0\left(\hat{F}_t(c)I_0(F_{\leq c}) \oplus (1 - \hat{F}_t(c))I_0(F_{> c})\right) = I_t(F(c)I_0(F_{\leq c}) + (1 - F(c))I_0(F_{> c})),$$

so Binary Distortions ensures that it is a cdf. This (together with earlier observations about  $\pi_t$ )

means that  $\pi_t$  is a probability distortion function. Binary Distortions also says

$$\mu_t(F) = \mu_0(\hat{F}_t) = u^{-1} \left( \int u(c) d\hat{F}_t(c) \right).$$

### C.3.3 PROOF OF COROLLARY 8

The usual uniqueness result for Bernoulli utility  $u$  applies here. Uniqueness of  $\mathcal{W}$  is clear from its definition in the proof. The  $\pi_t$  associated with a given  $u$  must satisfy

$$u^{-1} (\pi_t(q, c_L, c_H)u(c_L) + (1 - \pi_t(q, c_L, c_H))u(c_H)) = I_t(qc_L \oplus (1 - q)c_H).$$

This equation pins down  $\pi_t$  given  $u$ . Moreover, the implied value of  $\pi_t$  does not change if  $u$  is replaced with a positive affine transformation. Thus,  $\pi_t$  is unique.

### C.3.4 PROOF OF PROPOSITION 21

First part: The axiom says

$$\begin{aligned} qu(c'_L) + (1 - q)u(c'_H) &\geq \pi_t(p, c'_L, c'_H)u(c'_L) + (1 - \pi_t(p, c'_L, c'_H))u(c'_H) \\ \Rightarrow qu(c_L) + (1 - q)u(c_H) &\geq \pi_t(p, c_L, c_H)u(c_L) + (1 - \pi_t(p, c_L, c_H))u(c_H). \end{aligned}$$

This is equivalent to

$$q \leq \pi_t(p, c'_L, c'_H) \Rightarrow q \leq \pi_t(p, c_L, c_H)$$

which is equivalent to

$$\pi_t(p, c'_L, c'_H) \leq \pi_t(p, c_L, c_H).$$

Second part: It will be helpful here to consider a different way of doing the probability distort-

tions. Fix  $q$  and define  $F_{>q}$  and  $F_{<q}$  as follows:

$$F_{>q}(c) = \begin{cases} 0 & \text{if } F(c) \leq q \\ \frac{F(c)-q}{1-q} & \text{if } F(c) > q \end{cases} \quad F_{<q}(c) = \begin{cases} \frac{F(c)}{q} & \text{if } F(c) < q \\ 1 & \text{if } F(c) \geq q. \end{cases}$$

Let  $Q$  be the quantile function of  $F$ , and define  $\hat{Q}(Q)$  as follows:

$$Q(q) = c \Rightarrow \hat{Q}(Q) (\pi_t(q, I_0(F_{<q}), I_0(F_{>q}))) = c.$$

Then let

$$\hat{F}_t(F)(c) = \sup\{q \in [0, 1] : \hat{Q}(q) \leq c\}$$

. We show that this is the same  $\hat{F}_t(F)$  from the main representation. Fix  $c$ . We want to show that  $\hat{F}_t(F)(c)$  defined above equals

$$\pi_t \left( F(c), \frac{1}{F(c)} \int_{\leq c} u(x) dF(x), \frac{1}{1-F(c)} \int_{>c} u(x) dF(x) \right).$$

This is the same as

$$\pi_t (q, I_0(F_{<q}), I_0(F_{>q}))$$

where  $q = F(c)$ . Clearly,  $\hat{Q}(Q) (\pi_t(q, I_0(F_{<q}), I_0(F_{>q}))) = c$ . We want to show that, for any  $q'$ ,

$$\begin{aligned} \hat{Q}(Q) (\pi_t(q', I_0(F_{<q'}), I_0(F_{>q'}))) &\leq c \\ \Rightarrow \pi_t(q', I_0(F_{<q'}), I_0(F_{>q'})) &\leq \pi_t(q, I_0(F_{<q}), I_0(F_{>q})). \end{aligned}$$

Suppose this condition doesn't hold for some  $q'$ . If there exists  $c'$  such that  $F(c') = q'$ , then

$$\pi_t(F(c'), I_0(F_{<F(c')}), I_0(F_{>F(c')})) > \pi_t(F(c), I_0(F_{<F(c)}), I_0(F_{>F(c)})).$$

Since  $\pi_t$  is a probability distortion function, this implies  $c' > c$ . But then

$$\hat{Q}(Q)(\pi_t(q', I_0(F_{<q'}), I_0(F_{>q'}))) = c' > c,$$

a contradiction. So there must not exist  $c'$  such that  $F(c') = q'$ . Instead, let  $c' = Q(q') = \inf\{c : F(c) \geq q'\}$ . By assumption,

$$\pi_t(F(c'), I_0(F_{<F(c')}), I_0(F_{>F(c')})) \leq \pi_t(F(c), I_0(F_{<F(c)}), I_0(F_{>F(c)})).$$

(Otherwise, we are back in the previous case.) Since  $\pi_t$  is a probability distortion function, this implies  $c' \leq c$ . If we can show that

$$\pi_t(F(c'), I_0(F_{<F(c')}), I_0(F_{>F(c')})) \geq \pi_t(q', I_0(F_{<q'}), I_0(F_{>q'})),$$

we are done. This can be rewritten

$$\begin{aligned} & \pi_t\left(F(c'), u^{-1}\left(\frac{q_L}{F(c')}u(\tilde{c}_L) + \frac{F(c') - q_L}{F(c')}u(c')\right), c_H\right) \\ & \geq \pi_t\left(q', u^{-1}\left(\frac{q' - q_L}{q'}u(c') + \frac{q_L}{q'}u(\tilde{c}_L)\right), u^{-1}\left(\frac{F(c') - q'}{1 - q'}u(c') + \frac{1 - F(c')}{1 - q'}u(c_H)\right)\right) \end{aligned}$$

where  $c_H = I_0(F_{>F(c')})$ ,  $\tilde{c}_L = I_0(F_{<q_L})$  and  $q_L = \lim F(c'_-)$ . After some substitution (in which

$F(c')$  becomes  $q + \delta$ ,  $q'$  becomes  $q$ ,  $c'$  becomes  $c$ , and

$$u^{-1} \left( \frac{q' - q_L}{q'} u(c') + \frac{q_L}{q'} u(\tilde{c}_L) \right)$$

becomes  $c_L$ ), we recognize this as the condition imposed on  $\pi_t$  (that makes it a probability distortion function).

Now take  $F, G$  such that  $F \succeq_{FOSD} G$ . We saw above that we can construct  $\hat{F}_t(F)$  as follows. Consider the quantile function of  $F$  plotted in payoff-probability space. Map  $(c, q)$  into  $(c, \pi_t(q, I_0(F_{<q}), I_0(F_{>q})))$  to obtain a new quantile function.  $\hat{F}_t(F)$  is the cdf associated with this quantile function.

If  $F \succeq_{FOSD} G$ , then the original plot for  $F$  will lie (weakly) to the right of the plot for  $G$ . That is, if we have  $(c_F, q)$  on the  $F$  plot and  $(c_G, q)$  on the  $G$  plot,  $c_F \geq c_G$ . The former will get mapped into  $(c_F, \pi_t(q, I_0(F_{<q}), I_0(F_{>q})))$ , and the latter into  $(c_G, \pi_t(q, I_0(G_{<q}), I_0(G_{>q})))$ . Dominance implies  $I_0(F_{<q}) \geq I_0(G_{<q})$  and  $I_0(F_{>q}) \geq I_0(G_{>q})$ , so the axiom gives  $\pi_t(q, I_0(F_{<q}), I_0(F_{>q})) \leq \pi_t(q, I_0(G_{<q}), I_0(G_{>q}))$ . This ensures that the new plot for  $F$  (the quantile function of  $\hat{F}_t(F)$ ) still lies to the right of the new plot for  $G$  (the quantile function of  $\hat{F}_t(G)$ ). Thus,  $\hat{F}_t(F) \succeq_{FOSD} \hat{F}_t(G)$  as desired.

If  $F >_{FOSD} G$ , the original plot for  $F$  will lie strictly to the right of the plot for  $G$  at some  $q$ . By the above argument, this will still be true of the new plots. (In particular, the  $F$  plot will lie strictly to the right of the  $G$  plot at  $\pi(q, I_0(G_{<q}), I_0(G_{>q}))$ .) Thus,  $\hat{F}_t(F) >_{FOSD} \hat{F}_t(G)$ .

### C.3.5 PROOF OF PROPOSITION 16

“ $\succeq$  fantasizes about  $t$ ” is equivalent to “ $\succeq$  fantasizes more than  $\succeq^B$  about  $t$ ”, where  $\succeq^B$  has expected-utility preferences at all dates, with the same Bernoulli utility as  $\succeq$ . Thus, the proof of the next proposition suffices to establish this one.

### C.3.6 PROOF OF PROPOSITION 17

Only if: Since  $I_0^A(F) = I_0^B(F)$  for all  $F$ ,  $A$  and  $B$  have the same induced preferences over time-0 lotteries. Cardinal uniqueness of Bernoulli utility functions then implies  $u_A = \alpha u_B + \beta$  (with  $\alpha > 0$ ).

Specializing  $I_t^A(F) \geq I_t^B(F)$  to an arbitrary binary lottery  $q c_L \oplus (1 - q) c_H$ ,

$$\begin{aligned} & u_A^{-1}(\pi_t^A(q, c_L, c_H)u_A(c_L) + (1 - \pi_t^A(q, c_L, c_H))u_A(c_H)) \\ & \geq u_B^{-1}(\pi_t^B(q, c_L, c_H)u_B(c_L) + (1 - \pi_t^B(q, c_L, c_H))u_B(c_H)). \end{aligned}$$

Eliminating  $u_B$ , this becomes

$$\begin{aligned} & u_A^{-1}(\pi_t^A(q, c_L, c_H)u_A(c_L) + (1 - \pi_t^A(q, c_L, c_H))u_A(c_H)) \\ & \geq u_A^{-1}(\pi_t^B(q, c_L, c_H)u_A(c_L) + (1 - \pi_t^B(q, c_L, c_H))u_A(c_H)) \end{aligned}$$

which is equivalent to

$$\pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, c_L, c_H).$$

If: It suffices to show that  $\hat{I}_t^A(F) \geq_{FOSD} \hat{I}_t^B(F)$  for all  $F$ . (If this condition holds, then  $I_0^A(\hat{F}_t^A) \geq I_0^B(\hat{F}_t^B)$  since  $I_0$  is just expected utility and  $u_A$  and  $u_B$  are the same up to a positive affine transformation. This implies  $I_t^A(F) \geq I_t^B(F)$  as desired.) We need

$$\pi_t^A(q, I_0^A(F_{\leq q}), I_0^A(F_{> q})) \leq \pi_t^B(q, I_0^B(F_{\leq q}), I_0^B(F_{> q})).$$

Since  $I_0^A = I_0^B$ , this follows directly from

$$\pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, c_L, c_H).$$



### C.3.7 PROOF OF PROPOSITION 18

Suppose that  $\succeq^A$  and  $\succeq^B$  have distorted-probability representations with  $u^A$  and  $u^B$  bounded and  $\text{range}(u_A) = \text{range}(u_B)$ . We show that  $\varphi = u_B^{-1}(u_A)$  is the unique function that gives

$$\forall F \quad I_0^A(F) = \varphi^{-1}(I_0^B(\varphi(F))) \text{ and } I_0^B(F) = \varphi(I_0^A(\varphi^{-1}(F))).$$

Recall that  $I_0$  is expected utility. By cardinal uniqueness of Bernoulli utility functions, any  $\varphi$  must satisfy

$$\exists \alpha > 0, \beta \forall c \in C \quad u_B(\varphi(c)) = \alpha u_A(c) + \beta.$$

This gives

$$\varphi(c) = u_B^{-1}(\alpha u_A(c) + \beta).$$

Let  $\inf(u_A) = \inf(u_B) = \underline{u}$  and  $\sup(u_A) = \sup(u_B) = \bar{u}$ . For  $\varphi$  to be well defined for all  $c \in C$ , we need

$$[\alpha \underline{u} + \beta, \alpha \bar{u} + \beta] \subseteq [\underline{u}, \bar{u}].$$

$\varphi^{-1}$  must also be well defined for all  $c \in C$ . Since

$$\varphi^{-1}(c) = u_A^{-1}\left(\frac{u_B(c) - \beta}{\alpha}\right),$$

we need

$$\left[\frac{\underline{u} - \beta}{\alpha}, \frac{\bar{u} - \beta}{\alpha}\right] \subseteq [\underline{u}, \bar{u}].$$

Combining these two requirements gives  $\alpha = 1$  and  $\beta = 0$ , so  $\varphi = u_B^{-1}u_A$  is the only candidate to satisfy the desired condition. It is easy to verify that it does.

Only if: We show that

$$\pi^A(q, c_L, c_H) \leq \pi^B(q, \varphi(c_L), \varphi(c_H))$$

where  $\varphi = u_B^{-1}(u_A)$ . Specializing

$$I_t^A(F) \geq \varphi^{-1}(I_t^B(\varphi(F)))$$

to an arbitrary binary lottery, we get

$$\begin{aligned} & u_A^{-1}(\pi_t^A(p, c_L, c_H)u_A(c_L) + (1 - \pi_t^A(p, c_L, c_H))u_A(c_H)) \\ & \geq \varphi^{-1}(u_B^{-1}(\pi_t^B(p, \varphi(c_L), \varphi(c_H))u_B(\varphi(c_L)) + (1 - \pi_t^B(p, \varphi(c_L), \varphi(c_H)))u_B(\varphi(c_H)))) \end{aligned}$$

Using the definition of  $\varphi$ , the right-hand side simplifies to

$$u_A^{-1}(\pi_t^B(p, \varphi(c_L), \varphi(c_H))u_A(c_L) + (1 - \pi_t^B(p, \varphi(c_L), \varphi(c_H)))u_A(c_H)).$$

The resulting inequality simplifies to

$$\pi_t^A(p, c_L, c_H) \leq \pi_t^B(p, \varphi(c_L), \varphi(c_H)).$$

If: Given  $\varphi = u_B^{-1}(u_A)$ , we show that

$$\forall F \quad I_t^A(F) \geq \varphi^{-1}(I_t^B(\varphi(F))).$$

It suffices to show that

$$\pi_t^A(q, I_0^A(F_{\leq q}), I_0^A(F_{> q})) \leq \pi_t^B(q, I_0^B(\varphi(F)_{\leq q}), I_0^B(\varphi(F)_{> q})).$$

Since  $\varphi$  is a strictly increasing function, it doesn't matter whether  $F$  is truncated at  $q$  before or after transforming the payoffs with  $\varphi$ . Thus, the right-hand side is the same as

$$\pi_t^B(q, I_0^B(\varphi(F_{\leq q})), I_0^B(\varphi(F_{> q}))).$$

Using the definition of  $\varphi$ , this becomes

$$\pi_t^B(q, \varphi(I_0^A(F_{\leq q})), \varphi(I_0^A(F_{> q}))).$$

Clearly, the desired inequality follows from

$$\forall q, c_L, c_H \quad \pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, \varphi(c_L), \varphi(c_H)).$$

### C.3.8 PROOF OF PROPOSITION 19

First part: We proceed in two steps. First, we show that the desired relation between  $\pi_t^A$  and  $\pi_t^B$  will hold for all  $\varphi$  if it holds for  $\varphi = u_B^{-1}u_A$ . This is done by verifying that

$$\pi_t^B(q, \varphi(c_L), \varphi(c_H)) = \pi_t^B(q, u_B^{-1}(u_A(c_L)), u_B^{-1}(u_A(c_H))).$$

Recall that any  $A$ -to- $B$  comparison function must satisfy

$$\exists \alpha > 0, \beta \quad \varphi(c) = u_B^{-1}(\alpha u_A(c) + \beta).$$

Assume that  $u_A$  and  $u_B$  are bounded below by 0. For  $\varphi$  to be well defined, we must then have  $\beta \geq 0$ . For  $\varphi^{-1}(u_A)$  to be well defined, we must also have  $\beta \leq 0$ , so  $\beta = 0$ . A parallel argument (with the signs reversed) works if  $u_A$  and  $u_B$  are bounded above by 0. We conclude that the set of  $\varphi$  is given

by  $u_B^{-1}(\alpha u_A)$  for  $\alpha > 0$ . Plugging this into the left-hand side above, we see that the desired result follows from the restriction on  $\pi_t^B$ .

Second, we show that the desired relation between  $I_t^A$  and  $I_t^B$  will hold for all  $\varphi$  if it holds for  $\varphi = u_B^{-1}u_A$ . This is done by verifying that

$$\begin{aligned}\varphi^{-1}(I_t^B(\varphi(F))) &= u_A^{-1}u_B(I_t^B(u_B^{-1}u_A(F))) \\ \varphi(I_t^A(\varphi^{-1}(F))) &= u_B^{-1}u_A(I_t^A(u_A^{-1}u_B(F))).\end{aligned}$$

We will discuss the first line below; a parallel argument establishes the second. Using the fact that  $I_t$  is expected utility with distorted probabilities and  $u_B(\varphi) = \alpha u_A$ , the desired equality becomes

$$u_A^{-1} \sum u_A(c_i) \hat{p}_t^B(\varphi(F))(\varphi(c_i)) = u_A^{-1} \sum u_A(c_i) \hat{p}_t^B(u_B^{-1}u_A(F))(u_B^{-1}u_A(c_i)).$$

Thus, it suffices to show that  $\hat{p}_t^B(\varphi(F))(\varphi(c_i)) = \hat{p}_t^B(u_B^{-1}u_A(F))(u_B^{-1}u_A(c_i))$ . For this, we need

$$\begin{aligned}\pi_t^B \left( q, u_B^{-1} \left( \sum_{c_i \leq c_q} u_B(\varphi(c_i)) p_i \right), u_B^{-1} \left( \sum_{c_i > c_q} u_B(\varphi(c_i)) p_i \right) \right) \\ = \pi_t^B \left( q, u_B^{-1} \left( \sum_{c_i \leq c_q} u_A(c_i) p_i \right), u_B^{-1} \left( \sum_{c_i > c_q} u_A(c_i) p_i \right) \right).\end{aligned}$$

Once we plug in  $u_B(\varphi) = \alpha u_A$ , it's clear that this too follows from the restriction on  $\pi_t^B$ .

### C.3.9 PROOF OF PROPOSITION 13

We want to find a probability distortion function  $\pi_t$  that satisfies

$$\pi_t(q, c_L, c_H) = \Phi \left( \frac{\ln(c_q) - \hat{\mu}(\mu, \sigma)}{\sigma} \right)$$

$$\text{where } c_q \text{ solves } q = \Phi \left( \frac{\ln(c_q) - \mu}{\sigma} \right)$$

and  $(\mu, \sigma)$  solves

$$qu(c_L) + (1 - q)u(c_H) = \frac{1}{1 - \gamma} \exp \left( (1 - \gamma)\mu + \frac{1}{2}(1 - \gamma)^2\sigma^2 \right)$$

$$qu(c_L) = \frac{1}{1 - \gamma} \exp \left( (1 - \gamma)\mu + \frac{1}{2}(1 - \gamma)^2\sigma^2 \right) \Phi \left( \Phi^{-1}(q) - (1 - \gamma)\sigma \right)$$

for some continuous  $\hat{\mu}$ . (The second two equations come from taking expectations of a lognormal.)

Solving for  $(\mu, \sigma)$  delivers the expressions in the proposition.

We need

$$\begin{aligned} \frac{d\pi_t}{dq} &\geq 0 \\ \frac{d}{dq} \left( \Phi^{-1}(q) - \frac{\hat{\mu}(\mu, \sigma) - \mu}{\sigma} \right) &\geq 0 \\ \forall c \in [c_L, c_H] \quad \frac{1}{\varphi(\Phi^{-1}(q))}x + \frac{1}{\varphi\left(\Phi^{-1}\left(\frac{qu(c_L)}{u(\bar{c})}\right)\right)}\frac{u(c)}{u(\bar{c})}(1 - x) &\geq 0 \\ \text{where } x &= \hat{\mu}_1 - \frac{\hat{\mu}_2}{(1 - \gamma)\sigma} + \frac{\hat{\mu} - \mu}{(1 - \gamma)\sigma^2} \end{aligned}$$

This condition will be satisfied if  $x \in [0, 1]$ .

C.3.10 PROOF OF COROLLARY 7

We know that  $u_A, u_B$  are either both bounded below by 0 or both bounded above by 0. Moreover,

$$\forall \alpha > 0 \quad \sigma^J(q, c_L, c_H) = \sigma^J(q, u_I^{-1}(\alpha u_I(c_L)), u_I^{-1}(\alpha u_I(c_H))).$$

Thus, we need to show that

$$\pi_t^A(q, c_L, c_H) \leq \pi_t^B(q, u_B^{-1}(u_A(c_L)), u_B^{-1}(u_A(c_H))).$$

Here, this reduces to

$$f^A(\sigma^A(q, c_L, c_H)) \geq f^B(\sigma^B(q, u_B^{-1}u_A(c_L), u_B^{-1}u_A(c_H))).$$

Plugging in the expressions for  $\sigma^A$  and  $\sigma^B$ , we get

$$f^A\left(\frac{1}{1-\gamma^A}x\right) \geq f^B\left(\frac{1}{1-\gamma^B}x\right)$$

where  $x$  can take any positive value if  $\gamma_A < 1$  and any negative value if  $\gamma_A > 1$ . This condition is equivalent to the one stated in the proposition.

### C.3.II PROOF OF PROPOSITION 15

We want to find a probability distortion function  $\pi_t$  that satisfies

$$\hat{\pi}_t(q, c_H) = 1 - \exp\left(\hat{\lambda}(\lambda)c_q\right)$$

$$\text{where } c_q \text{ solves } q = 1 - \exp(\lambda c_q)$$

$$\text{and } \lambda \text{ solves } (1 - q)u(c_H) = \int_q^1 u\left(-\frac{\ln(1-x)}{\lambda}\right) dx.$$

Solving for  $\lambda$  delivers the expression in the proposition.

We need

$$\begin{aligned} \frac{d\pi_t}{dq} &\geq 0 \\ \frac{d}{dq} \left( \frac{\hat{\lambda}}{\lambda} (-\ln(1-q)) \right) &\geq 0 \end{aligned}$$

$$\begin{aligned} \forall c \in [c_L, c_H] \quad &\frac{u(c_H) - u(c)}{(1-\gamma)u(c_H)} (-\ln(1-q)) \left( \frac{\hat{\lambda}}{\lambda} - \hat{\lambda}' \right) \\ &+ \frac{\ln(1-q)}{1-\gamma} \left( \frac{(-\ln(1-q))^{1-\gamma}}{\frac{1}{1-q} \int_q^1 (-\ln(1-x))^{1-\gamma} dx} - 1 \right) \hat{\lambda}' \\ &+ \frac{\hat{\lambda}}{\lambda} \left( 1 + \frac{\ln(1-q)}{1-\gamma} \left( 1 - \frac{(-\ln(1-q))^{1-\gamma}}{\frac{1}{1-q} \int_q^1 (-\ln(1-x))^{1-\gamma} dx} \right) \right) \geq 0. \end{aligned}$$

All three terms in the above expression will be positive if  $\hat{\lambda} \geq \lambda \hat{\lambda}'$  and  $\hat{\lambda}' \geq 0$ .

## References

- Abdellaoui, Mohammed**, “A Genuine Rank-Dependent Generalization of the Von Neumann-Morgenstern Expected Utility Theorem,” *Econometrica*, 2002, 70 (2), 717–736.
- Aizerman, Mark and Andrew Malishevski**, “General theory of best variants choice: Some aspects,” *IEEE Transactions on Automatic Control*, 1981, 26 (5), 1030–1040.
- Allcott, Hunt**, “Social norms and energy conservation,” *Journal of public Economics*, 2011, 95 (9–10), 1082–1095.
- Andreoni, James, Justin M Rao, and Hannah Trachtman**, “Avoiding the ask: A field experiment on altruism, empathy, and charitable giving,” *Journal of Political Economy*, 2017, 125 (3), 625–653.
- Banerjee, Abhijit V**, “A simple model of herd behavior,” *The quarterly journal of economics*, 1992, 107 (3), 797–817.
- Berger, Paul D and Gerald E Smith**, “The effect of direct mail framing strategies and segmentation variables on university fundraising performance,” *Journal of Direct Marketing*, 1997, 11 (1), 30–43.
- Bernheim, B Douglas**, “A theory of conformity,” *Journal of political Economy*, 1994, 102 (5), 841–877.
- Beshears, John, James J Choi, David Laibson, Brigitte C Madrian, and Katherine L Milkman**, “The effect of providing peer information on retirement savings decisions,” *The Journal of finance*, 2015, 70 (3), 1161–1201.
- Bicchieri, Cristina and Erte Xiao**, “Do the right thing: but only if others do so,” *Journal of Behavioral Decision Making*, 2009, 22 (2), 191–208.
- Bikhchandani, Sushil, David Hirshleifer, and Ivo Welch**, “A theory of fads, fashion, custom, and cultural change as informational cascades,” *Journal of political Economy*, 1992, 100 (5), 992–1026.



- Billot, Antoine, Itzhak Gilboa, Dov Samet, and David Schmeidler**, “Probabilities as similarity-weighted frequencies,” *Econometrica*, 2005, 73 (4), 1125–1136.
- Blanken, Irene, Niels van de Ven, and Marcel Zeelenberg**, “A meta-analytic review of moral licensing,” *Personality and Social Psychology Bulletin*, 2015, 41 (4), 540–558.
- Brady, Richard L and John Rehbeck**, “Menu-dependent stochastic feasibility,” *Econometrica*, 2016, 84 (3), 1203–1223.
- Brown, Jeffrey R, Zoran Ivković, Paul A Smith, and Scott Weisbenner**, “Neighbors matter: Causal community effects and stock market participation,” *The Journal of Finance*, 2008, 63 (3), 1509–1531.
- Bursztyn, Leonardo, Florian Ederer, Bruno Ferman, and Noam Yuchtman**, “Understanding mechanisms underlying peer effects: Evidence from a field experiment on financial decisions,” *Econometrica*, 2014, 82 (4), 1273–1301.
- , **Ingar K Haaland, Aakaash Rao, and Christopher P Roth**, “Disguising Prejudice: Popular Rationales as Excuses for Intolerant Expression,” Technical Report 2020.
- Cai, Hongbin, Yuyu Chen, and Hanming Fang**, “Observational learning: Evidence from a randomized natural field experiment,” *American Economic Review*, 2009, 99 (3), 864–82.
- Caplin, Andrew and John Leahy**, “Psychological expected utility theory and anticipatory feelings,” *The Quarterly Journal of Economics*, 2001, 116 (1), 55–79.
- Card, David, Alexandre Mas, Enrico Moretti, and Emmanuel Saez**, “Inequality at work: The effect of peer salaries on job satisfaction,” *American Economic Review*, 2012, 102 (6), 2981–3003.
- Cattaneo, Matias D, Xinwei Ma, Yusufcan Masatlioglu, and Elchin Suleymanov**, “A random attention model,” *Journal of Political Economy*, 2020, 128 (7), 2796–2836.
- Chambers, Christopher P, Tugce Cuhadaroglu, and Yusufcan Masatlioglu**, “Behavioral influence,” Technical Report, Mimeo 2019.
- Charness, Gary and Uri Gneezy**, “What’s in a name? Anonymity and social distance in dictator and ultimatum games,” *Journal of Economic Behavior & Organization*, 2008, 68 (1), 29–35.
- Chen, Yan, F Maxwell Harper, Joseph Konstan, and Sherry Xin Li**, “Social comparisons and contributions to online communities: A field experiment on movielens,” *American Economic Review*, 2010, 100 (4), 1358–98.
- Cherepanov, Vadim, Tim Feddersen, and Alvaro Sandroni**, “Revealed preferences and aspirations in warm glow theory,” *Economic Theory*, 2013, 54 (3), 501–535.

- , **Timothy Feddersen, and Alvaro Sandroni**, “Rationalization,” *Theoretical Economics*, 2013, 8 (3), 775–800.
- Chew, Soo H and Larry G Epstein**, “Recursive utility under uncertainty,” in “Equilibrium theory in infinite dimensional spaces,” Springer, 1991, pp. 352–369.
- Clark, Stephen A**, “The random utility model with an infinite choice space,” *Economic Theory*, 1996, 7 (1), 179–189.
- Coffman, Lucas C, Clayton R Featherstone, and Judd B Kessler**, “Can Social Information Affect What Job You Choose and Keep?,” *American Economic Journal: Applied Economics*, 2017, 9 (1), 96–117.
- Cripps, Martin W**, “Divisible updating,” Technical Report, mimeo 2018.
- Cunningham, Tom and Jonathan de Quidt**, “Implicit preferences inferred from choice,” *Available at SSRN 2709914*, 2015.
- d’Adda, Giovanna, Yu Gao, Russell Golman, and Massimo Tavoni**, “It’s so Hot in Here: Information Avoidance, Moral Wiggle Room, and High Air Conditioning Usage,” Technical Report, Fondazione Eni Enrico Mattei 2018.
- Dana, Jason, Daylian M Cain, and Robyn M Dawes**, “What you don’t know won’t hurt me: Costly (but quiet) exit in dictator games,” *Organizational Behavior and human decision Processes*, 2006, 100 (2), 193–201.
- , **Roberto A Weber, and Jason Xi Kuang**, “Exploiting moral wiggle room: experiments demonstrating an illusory preference for fairness,” *Economic Theory*, 2007, 33 (1), 67–80.
- Dillenberger, David**, “Preferences for One-Shot Resolution of Uncertainty and Allais-Type Behavior,” *Econometrica*, 2010, 78 (6), 1973–2004.
- **and Collin Raymond**, “Additive-Belief-Based Preferences,” Technical Report 2018.
- **and Philipp Sadowski**, “Ashamed to be selfish,” *Theoretical Economics*, 2012, 7 (1), 99–124.
- , **Andrew Postlewaite, and Kareen Rozen**, “Optimism and pessimism with expected utility,” *Journal of the European Economic Association*, 2017, 15 (5), 1158–1175.
- Duflo, Esther and Emmanuel Saez**, “The role of information and social interactions in retirement plan decisions: Evidence from a randomized experiment,” *The Quarterly journal of economics*, 2003, 118 (3), 815–842.
- Ehrich, Kristine R and Julie R Irwin**, “Willful ignorance in the request for product attribute information,” *Journal of Marketing Research*, 2005, 42 (3), 266–277.

- Epstein, Larry G**, “Living with risk,” *The Review of Economic Studies*, 2008, 75 (4), 1121–1141.
- Epstein, LR and Stanley E Zin**, “VSubstitution,” *Risk Aversion, and the Temporal Behavior of Consumption and Asset Returns: A Theoretical Framework*, *V Econometrica*, 1989, 57.
- Exley, Christine L**, “Excusing selfishness in charitable giving: The role of risk,” *The Review of Economic Studies*, 2016, 83 (2), 587–628.
- Falk, Armin**, “Facing Yourself-A Note on Self-image,” Technical Report, CESifo Working Paper Series 2017.
- Fong, Christina M and Felix Oberholzer-Gee**, “Truth in giving: Experimental evidence on the welfare effects of informed giving to the poor,” *Journal of Public Economics*, 2011, 95 (5-6), 436–444.
- Franzen, Axel and Sonja Pointner**, “Anonymity in the dictator game revisited,” *Journal of Economic Behavior & Organization*, 2012, 81 (1), 74–81.
- Freddi, Eleonora**, “Do people avoid morally relevant information? Evidence from the refugee crisis,” Technical Report, CentER 2017.
- Frey, Bruno S and Stephan Meier**, “Social comparisons and pro-social behavior: Testing” conditional cooperation” in a field experiment,” *American Economic Review*, 2004, 94 (5), 1717–1722.
- Frick, Mira, Ryota Iijima, and Tomasz Strzalecki**, “Dynamic random utility,” *Econometrica*, 2019, 87 (6), 1941–2002.
- Gerber, Alan S and Todd Rogers**, “Descriptive social norms and motivation to vote: Everybody’s voting and so should you,” *The Journal of Politics*, 2009, 71 (1), 178–191.
- Giarlotta, Alfio and Salvatore Greco**, “Necessary and possible preference structures,” *Journal of Mathematical Economics*, 2013, 49 (2), 163–172.
- Gneezy, Uri, Silvia Saccardo, and Roel Van Veldhuizen**, “Bribery: Behavioral drivers of distorted decisions,” *Journal of the European Economic Association*, 2019, 17 (3), 917–946.
- , —, —, **Marta Serra-Garcia, and Roel van Veldhuizen**, “Motivated self-deception, identity and unethical behavior,” in “Working paper” 2016.
- , —, —, —, **and —**, “Bribing the self,” *Games and Economic Behavior*, 2020, 120, 311–324.
- Godlonton, Susan and Rebecca Thornton**, “Peer effects in learning HIV results,” *Journal of Development Economics*, 2012, 97 (1), 118–129.

- Goldstein, Noah J, Robert B Cialdini, and Vldas Griskevicius**, “A room with a viewpoint: Using social norms to motivate environmental conservation in hotels,” *Journal of consumer Research*, 2008, 35 (3), 472–482.
- Goulas, Sofoklis and Rigissa Megalokonomou**, “Knowing who you are: The effect of feedback information on exam placement,” *University of Warwick, mimeo*, 2015.
- Grant, Simon, Atsushi Kajii, and Ben Polak**, “Intrinsic preference for information,” *Journal of Economic Theory*, 1998, 83 (2), 233–259.
- , —, and —, “Temporal Resolution of Uncertainty and Recursive Non-expected Utility Models,” *Econometrica*, 2000, 68 (2), 425–434.
- Green, Jerry R and Bruno Jullien**, “Ordinal independence in nonlinear utility theory,” *Journal of risk and uncertainty*, 1988, 1 (4), 355–387.
- Grossman, Zachary and Joel J Van Der Weele**, “Self-image and willful ignorance in social decisions,” *Journal of the European Economic Association*, 2017, 15 (1), 173–217.
- Gul, Faruk and Wolfgang Pesendorfer**, “The revealed preference theory of changing tastes,” *The Review of Economic Studies*, 2005, 72 (2), 429–448.
- and —, “Random expected utility,” *Econometrica*, 2006, 74 (1), 121–146.
- , **Paulo Natenzon, and Wolfgang Pesendorfer**, “Random evolving lotteries and intrinsic preference for information,” Technical Report, Working paper 2016.
- Haisley, Emily C and Roberto A Weber**, “Self-serving interpretations of ambiguity in other-regarding behavior,” *Games and economic behavior*, 2010, 68 (2), 614–625.
- Haley, Kevin J and Daniel MT Fessler**, “Nobody’s watching?: Subtle cues affect generosity in an anonymous economic game,” *Evolution and Human behavior*, 2005, 26 (3), 245–256.
- Hamman, John R, George Loewenstein, and Roberto A Weber**, “Self-interest through delegation: An additional rationale for the principal-agent relationship,” *American Economic Review*, 2010, 100 (4), 1826–46.
- Herden, Gerhard and Andreas Pallack**, “On the continuous analogue of the Szpilrajn Theorem I,” *Mathematical social sciences*, 2002, 43 (2), 115–134.
- Holt, Charles A**, “Preference reversals and the independence axiom,” *The American Economic Review*, 1986, 76 (3), 508–515.

- Johnson, Jennifer Wiggins and Annie Peng Cui**, “To influence or not to influence: External reference price strategies in pay-what-you-want pricing,” *Journal of Business Research*, 2013, 66 (2), 275–281.
- Kahneman, Daniel and Amos Tversky**, “Prospect Theory: An Analysis of Decision under Risk,” *Econometrica*, 1979, 47 (2), 263–292.
- Kajackaite, Agne**, “If I close my eyes, nobody will get hurt: The effect of ignorance on performance in a real-effort experiment,” *Journal of Economic Behavior & Organization*, 2015, 116, 518–524.
- Kalai, Gil, Ariel Rubinstein, and Ran Spiegler**, “Rationalizing choice functions by multiple rationales,” *Econometrica*, 2002, 70 (6), 2481–2488.
- Kallenberg, Olav**, *Foundations of modern probability*, Springer Science & Business Media, 2006.
- , *Probabilistic symmetries and invariance principles*, Springer Science & Business Media, 2006.
- Kamenica, Emir**, “Contextual inference in markets: On the informational content of product lines,” *American Economic Review*, 2008, 98 (5), 2127–49.
- Karni, Edi and Zvi Safra**, ““Preference reversal” and the observability of preferences by experimental methods,” *Econometrica: Journal of the Econometric Society*, 1987, pp. 675–685.
- Kreps, David M and Evan L Porteus**, “Temporal resolution of uncertainty and dynamic choice theory,” *Econometrica: journal of the Econometric Society*, 1978, pp. 185–200.
- Lehrer, Ehud and Roe Teper**, “Justifiable preferences,” *Journal of Economic Theory*, 2011, 146 (2), 762–774.
- Loewenstein, George, Samuel Issacharoff, Colin Camerer, and Linda Babcock**, “Self-serving assessments of fairness and pretrial bargaining,” *The Journal of Legal Studies*, 1993, 22 (1), 135–159.
- Majumdar, Dipjyoti**, “An axiomatic characterization of Bayes’ Rule,” *Mathematical social sciences*, 2004, 47 (3), 261–273.
- Manski, Charles F**, “Economic analysis of social interactions,” *Journal of economic perspectives*, 2000, 14 (3), 115–136.
- Manzini, Paola and Marco Mariotti**, “Sequentially rationalizable choice,” *American Economic Review*, 2007, 97 (5), 1824–1839.
- and —, “Stochastic choice and consideration sets,” *Econometrica*, 2014, 82 (3), 1153–1176.

- Masatlioglu, Yusufcan, Daisuke Nakajima, and Emre Ozdenoren**, “Willpower and compromise effect,” *Theoretical Economics*, 2020, 15 (1), 279–317.
- , —, and **Erkut Y Ozbay**, “Revealed attention,” *American Economic Review*, 2012, 102 (5), 2183–2205.
- Moulin, Hervé**, “Choice functions over a finite set: a summary,” *Social Choice and Welfare*, 1985, 2 (2), 147–160.
- Natenzon, Paulo**, “Random choice and learning,” *Journal of Political Economy*, 2019, 127 (1), 419–457.
- Norton, Michael I, Joseph A Vandello, and John M Darley**, “Casuistry and social category bias,” *Journal of personality and social psychology*, 2004, 87 (6), 817.
- Quimet, Paige and Geoffrey Tate**, “Learning from coworkers: Peer effects on individual investment decisions,” *The Journal of Finance*, 2020, 75 (1), 133–172.
- Plott, Charles R**, “Path independence, rationality, and social choice,” *Econometrica: Journal of the Econometric Society*, 1973, pp. 1075–1091.
- Quiggin, John**, “A theory of anticipated utility,” *Journal of Economic Behavior & Organization*, 1982, 3 (4), 323–343.
- Quinn, David M**, “Experimental evidence on teachers’ racial bias in student evaluation: The role of grading scales,” *Educational Evaluation and Policy Analysis*, 2020, 42 (3), 375–392.
- Riener, Gerhard and Christian Traxler**, “Norms, moods, and free lunch: Longitudinal evidence on payments from a Pay-What-You-Want restaurant,” *The Journal of Socio-Economics*, 2012, 41 (4), 476–483.
- Rodriguez-Lara, Ismael and Luis Moreno-Garrido**, “Self-interest and fairness: self-serving choices of justice principles,” *Experimental Economics*, 2012, 15 (1), 158–175.
- Rogers, Todd and Avi Feller**, “Discouraged by peer excellence: Exposure to exemplary peer performance causes quitting,” *Psychological science*, 2016, 27 (3), 365–374.
- Sacerdote, Bruce**, “Experimental and quasi-experimental analysis of peer effects: two steps forward?,” *Annu. Rev. Econ.*, 2014, 6 (1), 253–272.
- Safonov, Evgenii**, “Random choice with framing effects: a Bayesian model,” *Princeton University, mimeo*, 2017.
- , “Random choice with framing effects: a Bayesian model,” Technical Report 2018.

- Schons, Laura Marie, Mario Rese, Jan Wieseke, Wiebke Rasmussen, Daniel Weber, and Wolf-Christian Strotmann**, “There is nothing permanent except change—analyzing individual price dynamics in “pay-what-you-want” situations,” *Marketing Letters*, 2014, 25 (1), 25–36.
- Serra-Garcia, Marta and Nora Szech**, “The (in) elasticity of moral ignorance,” Technical Report, KIT Working Paper Series in Economics 2019.
- Shang, Jen and Rachel Croson**, “A field experiment in charitable contribution: The impact of social information on the voluntary provision of public goods,” *The economic journal*, 2009, 119 (540), 1422–1439.
- Snyder, Melvin L, Robert E Kleck, Angelo Strenta, and Steven J Mentzer**, “Avoidance of the handicapped: an attributional ambiguity analysis,” *Journal of personality and social psychology*, 1979, 37 (12), 2297.
- Strotz, Robert Henry**, “Myopia and inconsistency in dynamic utility maximization,” *The review of economic studies*, 1955, 23 (3), 165–180.
- Tversky, Amos and Daniel Kahneman**, “Advances in prospect theory: Cumulative representation of uncertainty,” *Journal of Risk and uncertainty*, 1992, 5 (4), 297–323.
- Wakker, Peter**, “Unbounded utility for Savage’s “Foundations of statistics,” and other models,” *Mathematics of Operations Research*, 1993, 18 (2), 446–485.
- , “Separating marginal utility and probabilistic risk aversion,” *Theory and decision*, 1994, 36 (1), 1–44.
- Woolley, Kaitlin and Jane L Risen**, “Closing your eyes to follow your heart: Avoiding information to protect a strong intuitive preference.” *Journal of personality and social psychology*, 2018, 114 (2), 230.
- Yaari, Menahem E**, “The dual theory of choice under risk,” *Econometrica: Journal of the Econometric Society*, 1987, pp. 95–115.

**T**HIS THESIS WAS TYPESET using L<sup>A</sup>T<sub>E</sub>X, originally developed by Leslie Lamport and based on Donald Knuth's T<sub>E</sub>X. The body text is set in 11 point Egenolff-Berner Garamond, a revival of Claude Garamont's humanist typeface. A template that can be used to format a PhD thesis with this look and feel has been released under the permissive MIT (X11) license, and can be found online at [github.com/suchow/Dissertate](https://github.com/suchow/Dissertate) or from its author, Jordan Suchow, at [suchow@post.harvard.edu](mailto:suchow@post.harvard.edu).