

Caliari, Daniele

**Working Paper**

## Rationality is not consistency

WZB Discussion Paper, No. SP II 2023-304

**Provided in Cooperation with:**

WZB Berlin Social Science Center

*Suggested Citation:* Caliari, Daniele (2023) : Rationality is not consistency, WZB Discussion Paper, No. SP II 2023-304, Wissenschaftszentrum Berlin für Sozialforschung (WZB), Berlin

This Version is available at:

<https://hdl.handle.net/10419/274074>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

# WZB

Wissenschaftszentrum Berlin  
für Sozialforschung



Daniele Caliori

## **Rationality is not Consistency**

**Discussion Paper**

SP II 2023–304

July 2023

Research Area

**Markets and Choice**

Research Unit

**Economics of Change**

Wissenschaftszentrum Berlin für Sozialforschung gGmbH  
Reichpietschufer 50  
10785 Berlin  
Germany  
[www.wzb.eu](http://www.wzb.eu)

Copyright remains with the authors.

Discussion papers of the WZB serve to disseminate the research results of work in progress prior to publication to encourage the exchange of ideas and academic debate. Inclusion of a paper in the discussion paper series does not constitute publication and should not limit publication in any other venue. The discussion papers published by the WZB represent the views of the respective author(s) and not of the institute as a whole.

Affiliation of the authors:

Daniele Caliori, WZB ([daniele.caliari@wzb.eu](mailto:daniele.caliari@wzb.eu))

Abstract

### ***Rationality is not Consistency\****

We challenge the standard definition of economic rationality as consistency by making use of a novel distinction between axioms of decision theory: consistency and preference axioms. We argue that this distinction has been overlooked by the literature and, as a result, evidence that consistency is a proxy of decision-making ability is often based on incorrect identification strategies. We conduct an experiment to investigate the factors that drive violations of consistency alone. While we find no evidence that consistency axioms are a proxy of decision-making ability, we provide suggestive evidence that some preference axioms are, confirming their potential role as confounding factors. Overall, our experimental evidence raises doubts about the choice of language that equates consistency with rationality in economics.

*Keywords:* Decision Theory, Experimental Design, Consistency, Rationality.

*JEL classification:* D00, D90, D91

---

\* I am indebted to Marco Mariotti and Christopher Tyson for their advice and guidance. I also thank Georgios Gerasimou, David Freeman, Aniol Llorente-Saguer, Ivan Soraperra, Maria Vittoria Levati, David Dillenberger, Asen Ivanov, Tomas Jagelka, Dorothea Kubler, Lorenzo Neri, Pietro Ortoleva, Ariel Rubinstein, Steffen Huck, Yiming Liu, Kai Barron, Valentino Dardanoni, Itzhak Gilboa, and the participants of ESEM-EEA European Summer Meeting 2019, EEA Virtual Meeting 2020, WZB MC-Reading Group, and the first MPI-WZB workshop, 2021. I also thank Queen Mary University of London for funding the experiment and University of St. Andrews for hosting it. The experiment was approved by the Queen Mary Ethics of Research Committee: ref. QMREC2102.

All discussion papers are downloadable:  
<http://www.wzb.eu/en/publications/discussion-papers/markets-and-choice>

# 1 Introduction

"Economists sometimes use the adjective rational in place of consistent, with the implied pejorative that choices that don't conform to their models are irrational. This is bad choice of language and is the source of all sorts of silly arguments with psychologists, sociologists, etc..." (Kreps, 2015).

"Rationality, they say, equals consistency... But this means only that those choices are consistent with one another *when viewed from the perspective of some theory* [italics in the original]" (Sugden, 1991).

In this paper, the theory under scrutiny is utility maximization. It is well known that a decision maker can be modelled as a utility maximizer if and only if she has transitive and complete preferences. This is the case if and only if her choices are consistent, i.e. they satisfy the Weak Axiom of Revealed Preference. Economists have attached the adjective "rational" to consistency requirements, implicitly assuming that it is *how* the decision maker chooses that defines her rationality level and not *what* she chooses.

Our first, and theoretical, contribution is to introduce two types of decision theory axioms: consistency axioms (henceforth **ConAx**)<sup>1</sup> and preference axioms (henceforth **PrAx**) which capture the ideas of *how* the decision maker ought to choose, and *what* she ought to choose, respectively. To clarify our intuition with well-known decision theory axioms, transitivity and completeness (**ConAx**) and monotonicity (**PrAx**); imagine a decision maker who has to choose between three alternatives: £5, £6, and £7. The set of transitive and complete preferences is neutral,<sup>2</sup> namely, it does not constrain the preferences of the decision-maker who could, for instance, prefer £5 to £7. On the other hand, the set of monotonic preferences - the singleton  $\text{£7} > \text{£6} > \text{£5}$  - is not neutral; the decision-maker ought to prefer £7.

---

<sup>1</sup>Notable examples of consistency axioms, sometimes referred to as revealed preference tests, are the Weak and Strong Axioms of Revealed Preference, Sen's Property  $\alpha$ ,  $\beta$  and  $\gamma$  (Sen, 1971) in deterministic choice or Independence from Irrelevant Alternatives, Stochastic Transitivity and Regularity in stochastic choice (Block & Marschak, 1960). A review of different notions of consistency and their violations can be found in Rieskamp et al. (2006).

<sup>2</sup>In Appendix A1 we formalize the notion of neutrality and provide a definition for **ConAx** and **PrAx**. In doing so, we define transitivity and completeness as structural axioms (**StAx**) because, differently from WARP, they are constraints on the primitives (preferences) and not on the observables (choices). Nonetheless, we exploit the equivalence between WARP and transitive & complete preferences to treat them as **ConAx**. A comprehensive analysis of the connection between **ConAx** and **StAx** can be found in Mahmoud (2017).

Contrary to consistency, rationality is generally a more ambiguous concept. Among the several attempts to provide a formal definition, Gilboa (2009) defines objective rationality as "modes of behaviour that can be explained to others so that these are convinced by them." It follows that to define objective rationality one needs to search for a widely accepted definition of good decision-making abilities. To convince us that consistency is a good measure of objective rationality, economists have shown correlations between consistency, real world economic outcomes, and individual characteristics (Choi et al. (2014), Andersson et al. (2016), Banks et al. (2019)). We argue that these findings are misleading because consistency has been tested together with other requirements (**PrAx**) that may have driven the results. Our motivation is to provide both theoretical and experimental evidence that, in the words of Kreps, "the use of the adjective rational in place of consistent is a bad choice of language". In the next two paragraphs, we give two examples from the recent literature in which consistency (**ConAx**) cannot be tested independently.

In their influential paper titled "Who is (more) rational?", Choi et al. (2014) find that consistency with GARP in a lab experiment is correlated with wealth. However, by Afriat's Theorem (Afriat, 1967), GARP is equivalent to the existence of a continuous, strictly increasing, and concave utility function that rationalizes the data and, therefore, it has a stronger content (e.g. monotonicity and concavity) than the mere existence of a utility function. We devote Section 3.5 to discussing how our paper relates to the literature on GARP.

Andersson et al. (2016) use a classical Multiple Price List design (Holt & Laury (2002), Andersen et al. (2008)) to study risk elicitation. They find a positive correlation between consistency and cognitive abilities. In this design, a failure of consistency is equated to multiple switches which also imply a violation of monotonicity. In fact, Andersson et al. (2016) write: "We define subjects as Consistent if their decisions are compatible with rational [transitive and complete] and monotonic preferences". Again, economic rationality, as consistency, is combined with monotonicity, making it unclear how to test the former without the latter.<sup>3</sup>

In this paper, we ask whether **ConAx** are, as the literature has suggested, *necessary* conditions for high decision-making ability and, therefore, whether they justify the use of the adjective rational. Firstly, **ConAx** are clearly not a *sufficient* condition for high

---

<sup>3</sup>Evidence of a misalignment between multiple switching behaviour (MSB) and cognitive abilities has been recently investigated by Yu et al. (2021) and Chew et al. (2022) through a novel distinction between regular and irregular MSB.

decision-making ability. Consider a decision maker who always chooses £5 over £6. **ConAx** suggest perfect consistency and a utility function can be constructed such that  $u(5) > u(6)$ . Nonetheless, it is hard to argue that she has good decision-making abilities since monotonicity (**PrAx**) is violated. However, we make the stronger argument against **ConAx** being *necessary* conditions for high-quality decision-making. There are many deterministic models - e.g. Manzini & Mariotti (2012), Manzini & Mariotti (2007), Masatlioglu et al. (2012) - and stochastic models - e.g. Machina (1985), Fudenberg et al. (2015), Cerreia-Vioglio et al. (2019) - that rely on a type of optimization process and that rationalize violations of consistency.<sup>4</sup> Particularly, we focus on the idea of deliberate randomization which has been theoretically studied by Cerreia-Vioglio et al. (2019) in the context of risk preferences, empirically validated in the same environment by Agranov & Ortoleva (2017), and investigated in a different domain (university applications) by Dwenger et al. (2018).

Our second contribution is empirical. Although novel in its formalization, ours is not the first theoretical criticism of the idea of economic rationality as consistency. However, we argue that this paper is the first attempt to provide a clear-cut empirical strategy for the test of **ConAx** independently from **PrAx** and to provide evidence about their relationship with proxies of decision-making abilities. We construct a choice elicitation experiment in which subjects are asked to make choices regarding delayed payment plans and gambles. Henceforth, we refer to the two environments as "Time" and "Risk". We investigate the factors that contribute to violations of **ConAx**, focusing on the number of violations of WARP. This measure, which will be our main dependent variable, can be simply defined as the sum - over all pairs of elements - of the product between the number of times an element  $\mathbf{x}$  is chosen when  $\mathbf{y}$  is available, and vice versa.<sup>5</sup> Our identification strategy relies on the fact that in some of our choice problems, which we call MAIN problems, only **ConAx** can be violated, while in the

<sup>4</sup>Inconsistent choices have been documented in both experimental and non-experimental settings for different reasons: mistakes (Cason & Plott, 2014), variation in tastes (Echenique et al., 2011), deliberate randomization (Agranov & Ortoleva (2017), Cerreia-Vioglio et al. (2019)), the attraction and the similarity effect (Tversky & Russo (1969), Natenzon (2019)), thinking costs (Caplin & Dean (2015), Fudenberg et al. (2015)), non-standard behavioural procedures (Caplin et al., 2011), and many more.

<sup>5</sup>Let  $C_{xy}$  be the number of times  $\mathbf{x}$  is chosen when  $\mathbf{y}$  is available, for each individual  $i$ :

$$WARP_i = \sum_{x,y} C_{xy} \cdot C_{yx}$$

As a robustness check, in the Online Appendix, we repeat our analysis using the Strong Axiom of Revealed Preference and confirm our results.



remaining ones violations of both **ConAx** and **PrAx** can arise. Importantly, since the MAIN problems are all the non-empty subsets of a set of four alternatives, our test of WARP (**ConAx**) is a test of the standard definition of economic rationality as utility maximization (Sen, 1971).

Our empirical strategy consists of three steps. First, we identify those subjects who use specific heuristics (simple rules) and those who deliberately randomize. The identification step relies, for the heuristics, on preference elicitation tools since, in our design, heuristics induce specific extreme preferences. Deliberate randomization is, instead, identified using the subjects' reported behaviour in the questionnaire as in Agranov & Ortoleva (2017). The questionnaire (see Appendix B) also allows us to corroborate our approach for the identification of heuristics, which is also confirmed by the response times, in line with Rubinstein (2013). The intuition behind inducing heuristics through the design and the consequent identification via preference elicitation tools were inspired by the following hypothesis (Rubinstein, 2013) in the scarce literature on the relation between **ConAx** and **PrAx**: "Consistency may reflect the use of a simple rule rather than greater sophistication." If this hypothesis was correct, not only subjects who follow simple rules should be more consistent; but, if the use of simple rules is correlated with lower sophistication or lower effort, we would empirically challenge the use of the adjective rational in place of consistent. Therefore, in the first step, we investigate the relationship between heuristics, **ConAx**, cognitive abilities, and response times.<sup>6</sup> In the second and third steps, we use cognitive abilities, response times, and the level of understanding of the experiment as proxies of decision-making abilities. First, we provide evidence that **PrAx** are potential confounding factors for the evaluation of consistency as a requirement for high decision-making ability. Second, we investigate if these factors, as well as violations of **PrAx**, are correlated with the violations of **ConAx** alone. What follows is the preview of our main results.

**Result 1.** In step 1, we confirm Rubinstein's hypothesis. We find a strong negative effect of the use of heuristics on the number of violations of WARP. This result is notably in contrast with Choi et al. (2014) who wrote: "Some subjects may therefore adopt simple decision rules, and this "simplification" may cause their choices to be inconsistent." The difference in the effect between Time and Risk is due to the predominant use of these rules in the former case (59%) versus the latter (20%). We also find that deliberate

---

<sup>6</sup>Response times are notoriously noisy. In Appendix C, we analyse the response times along different dimensions and show that, in our experiment, they contain important information regarding the complexity and the nature of the questions.

randomization is more widespread in Risk (19%) than in Time (5%). This behaviour is significantly positively correlated with violations of WARP. Together, heuristics and deliberate randomization explain the substantial difference in WARP violations between Time and Risk (Figure 2).

**Result 2.** In step 2, we provide suggestive evidence that, contrarily to **ConAx**, violations of some **PrAx** display characteristics that are naturally connected to the idea of bad decision-making abilities and therefore confirm their role of potential confounding factors. In Time, we find robust negative correlations with cognitive abilities, response times, and the level of understanding of the experiment. In Risk, these results are weaker and only partially replicated.

**Result 3.** In step 3, we find that violations of **ConAx** and **PrAx** are not significantly correlated in both Time and Risk, with the exception of a strong and positive correlation between violations of WARP and SOSD.<sup>7</sup> This finding, in line with previous experiments such as Sopher & Narramore (2000) and Agranov & Ortoleva (2017), is further evidence of the widespread use of deliberate randomization in Risk as modelled by Machina (1985) and Cerreia-Vioglio et al. (2019). To confirm our interpretation, we notice that whereas violations of SOSD are strongly correlated with WARP violations, violations of FOSD<sup>8</sup> are not. This finding is in line with Cerreia-Vioglio et al. (2019), whose model forbids deliberate randomization using lotteries that are FOSD, while it allows it for lotteries that are SOSD.

**Result 4.** We find an ambiguous correlation between cognitive abilities and violations of WARP. While it is negative in Time, it is positive in Risk. We conclude that, in our experiment, greater sophistication is connected to higher consistency in Time and lower consistency in Risk. This finding contrasts with those of Burks et al. (2009), Choi et al. (2014) and Andersson et al. (2016).

**Result 5.** We find a positive correlation between response times and violations of WARP in both Time and Risk. Hence, we confirm the findings of Rubinstein (2013): the use of quick simple rules generally increases consistency. If response times are considered as a measure of effort, this finding suggests that violations of WARP and effort are, at best, uncorrelated.

---

<sup>7</sup>We say  $x$  second-order stochastically dominates  $y$  if given  $F_x$  and  $F_y$  the respective cumulative distribution functions:  $\int_{-\infty}^a [F_y(t) - F_x(t)] dt \geq 0$  for all  $a$ .

<sup>8</sup>We say  $x$  first-order stochastically dominates  $y$  if given  $F_x$  and  $F_y$  the respective cumulative distribution functions:  $F_y(a) \geq F_x(a)$  for all  $a$ .

## 2 The Experiment

The experiment follows a choice elicitation design. In Part One and Two subjects were asked to choose from different sets of alternatives: delayed payment plans (Time) or lotteries (Risk). Each part had 25 choice problems that were designed to be non-trivial. Namely, unlike Tversky & Russo (1969), Manzini & Mariotti (2010), and McCausland et al. (2019), in none of the MAIN problems a dominant alternative was present.

Before the start of the experiment subjects received general instructions plus specific instructions about both parts.<sup>9</sup> Furthermore, at the beginning of Part One and Two, subjects answered three trial problems in order to make them familiar with the experiment's environment.

For both Time and Risk the alternatives were divided into two groups: Tables 1 and 2 describe the four MAIN alternatives (see the Online Appendix in Section 1 for a complete description of all the alternatives).

Table 1: LIST OF MAIN DELAYED PAYMENT PLANS

ALTERNATIVES	MONTHS				
	0	3	6	9	12
<b>One Shot (OS)</b>	160	0	0	0	0
<b>Decreasing (D)</b>	110	50	25	0	0
<b>Constant (K)</b>	50	50	50	50	0
<b>Increasing (I)</b>	0	15	40	170	0

Table 2: LIST OF MAIN LOTTERIES

ALTERNATIVES	TOKEN		PROBABILITIES		EV
<b>Degenerate (D)</b>	50	0	1	0	50
<b>Safe (S)</b>	65	25	0.8	0.2	57
<b>Fifty-Fifty (50)</b>	90	25	0.5	0.5	57.5
<b>Risky (R)</b>	300	5	0.2	0.8	64

NOTE -- The amounts are described in Token. The exchange rate was fixed at 20:1 pounds for Delayed Payment Plans and 10:1 pounds for Lotteries.

Each individual solved all the 11 choice problems involving the MAIN alterna-

<sup>9</sup>The instructions were available both on screen and on paper such that they could be consulted during the entire experiment. Note that the experiment has been designed to be paper-free.

tives.<sup>10</sup> The other problems were designed to obtain particular information about the following **PrAx**: (i) monotonicity implies that an individual should prefer more money than less. In choices among gambles, the individual never chooses first-order stochastically dominated gambles; (ii) impatience<sup>11</sup> implies that money has a decreasing value when moved ahead in time; (iii) risk aversion implies that individuals should never choose lotteries that are second-order stochastically dominated or, equivalently, individuals have a concave Bernoulli utility function.

The positions of the alternatives were randomized. The subjects could face two orders of problems and we also inverted Time and Risk elicitation such that we had a total of four treatments. A complete structure of the experiment as well as a description of the orders can be found in Section 2 of the Online Appendix.

The reward was measured in Token with an exchange rate of 1:10 for lotteries and 1:20 for delayed payment plans. The average reward was about £19 per subject and the experiment lasted on average 1:15 hours. The experiment took place at the University of St. Andrews between June and September 2019 and 145 subjects participated. The experiment has been performed using z-tree (Fischbacher, 2007).

### 3 Results

The results are structured as follows: we start, in Section 3.1, by presenting a stark difference in behaviour between Time and Risk. In Section 3.2, we introduce our identification strategy for subjects that use heuristics. In Time, we focus on two heuristics that resemble the behaviour of extremely patient and impatient subjects. In Risk, we focus on one heuristic based on the choice of the safest lottery (extreme risk aversion), and one based on the calculation of the expected value (risk neutrality). Then, we introduce the concept of deliberate randomization with specific reference (in Risk) to the work of Cerreia-Vioglio et al. (2019). Subsequently, in Section 3.3, we investigate the factors that are correlated with the violation of **PrAx**. Finally, in Section 3.4, we unify our results into a comprehensive analysis of violations of **ConAx**.

---

<sup>10</sup>The eleven sets of questions refer to all the non-empty subsets of the set of MAIN alternatives minus the singletons.

<sup>11</sup>By impatience we intend to refer to the violation of discounting models. The term "impatience" has been used by Fishburn & Rubinstein (1982) to denote Axiom A3.

### 3.1 Consistency in Time and Risk

Subjects' consistency differs significantly between Time and Risk. Figure 1 displays the number of violations of WARP in Risk on the x-axis and those in Time on the y-axis. The area of the circles is proportional to the number of subjects. The correlation is low and mostly driven by a small portion of consistent individuals in both Time and Risk. The magnitude of the WARP violations is similarly important. Figure 2 presents the distributions of violations in Time, Risk and those that would have been made by a random chooser. A Kolmogorov-Smirnov test reports a highly significant difference between the three distributions Time/Risk (p-value=0.000), Time/Random (p-value=0.000) and Risk/Random (p-value=0.000). To give an idea of the difference in magnitude, on average, subjects display 1.96 WARP violations in Time and 4.99 in Risk (t-test, p-value=0.000). This result is only the starting point of our analysis, nonetheless, it is novel and it suggests that individuals may behave very differently in Time and Risk.

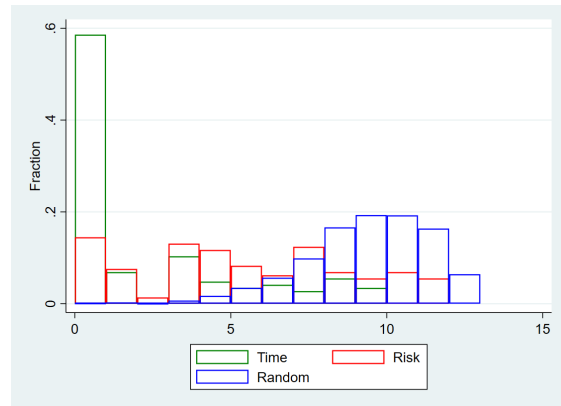
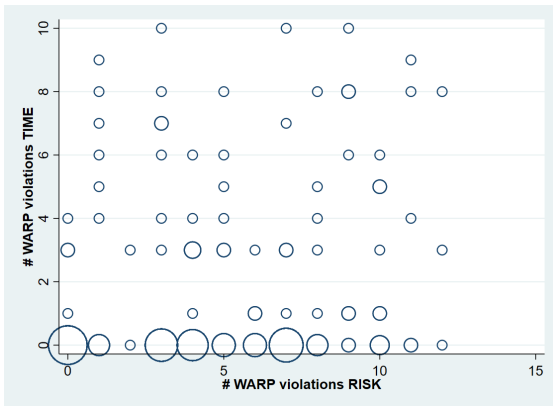


Figure 1: Violations of WARP in the entire dataset.

Figure 2: Distribution of the violations of WARP.

### 3.2 Heuristics and deliberate randomization

In our experiment, some natural heuristics are mapped to extreme preferences ordered by discount factor and risk aversion parameter. This observation inspires our identification strategy which consists of focusing on those subjects who chose in line with extreme preferences. Since choices are noisy, in order to identify preferences we use an optimal weighting algorithm, as described in Caliari (2023). As mentioned in the introduction, we provide several robustness checks to our approach. First, in Appendix B, we show that our identification through elicited preferences highly overlaps with the

direct revelation of the use of heuristics through the questionnaire.<sup>12</sup> Second, in Section 3.4 of the Online Appendix, we show that our main analysis on the violations of **ConAx** is confirmed when heuristics are identified using direct revelation instead of the elicited preferences. Finally, to show that our results do not rely on the specific optimal weighting algorithm, in Section 3.3 of the Online Appendix, we show that the correlations between the use of heuristics, violations of **ConAx**, and cognitive abilities are robust both to the use of reported preferences and the use of different algorithms such as the Minimum Swaps algorithm (Apestegua & Ballester, 2015) and the Sequential algorithm (Horan & Sprumont, 2016).

As presented in Tables 1 and 2, we have denoted the risk alternatives: Degenerate [D], Safe [S], Fifty-Fifty [50] and Risky [R]; and the time alternatives: One Shot Payment [OS], Decreasing [D], Constant [K], Increasing [I]. The choice of the Degenerate lottery (or One Shot Payment plan) involves particularly high (low) risk aversion parameter (discount factor).<sup>13</sup> When these alternatives are consistently chosen the existence of a heuristic is particularly probable.<sup>14</sup> Importantly, these heuristics do not require calculations. Highly risk-averse subjects simply search for the option with the highest probability to win the bigger amount within the gamble (represented by pies, see the Online Appendix, Section 6), while highly impatient ones search for the highest, and first, histogram in the graphical description of the payment plans (similarly, see the Online Appendix, Section 6). On the other side of the preferences, we identify two heuristics that involve calculations: the risk neutral subjects calculate the expected value of the lotteries, and the patient ones sum all the payments in the plan and choose the one with the highest value.

Overall, we identify four heuristics: we refer to subjects who choose according to the preference  $D > S > 50 > R$  as **Most Risk Averse**,  $R > 50 > S > D$  as **Risk Neutral**,  $OS > D > K > I$  as **Most Impatient** and  $I > K > D > OS$  as **Patient**. Crucially, these individuals are not, a priori, more consistent than others. For example, a subject that behaves according to a CRRA utility function with a risk parameter between 0.8 and 1.2 should consistently choose following the preference  $S > D > 50 > R$ . A complete breakdown of the consistency of the subjects by elicited preferences is presented in Section 3.1 of the Online Appendix.

Finally, we use the questionnaire to identify a behaviour denoted as **Deliberate**

<sup>12</sup>A full description of the questionnaire can be found in Section 4 of the Online Appendix.

<sup>13</sup>See the Online Appendix for more details regarding how the alternatives have been chosen.

<sup>14</sup>Iyengar & Kamenica (2010) show that, particularly in cases of choice overload, subjects have a preference for simplicity. In the case of gambles, they observe a preference for degenerate ones.

**Randomization**<sup>15</sup> in both Time and Risk, as in Agranov & Ortoleva (2017). Samples of the reported answers are presented in Appendix B.

### 3.2.1 Time preferences

We start by presenting the relation between heuristics, deliberate randomization and WARP violations, Raven’s scores, and Response Times. The box plots report both the median (red line) and the mean (blue star). Above each box plot, we report the probabilities related to two tests performed in the comparison with the group of remaining subjects, denoted as **Others**: unpaired t-test for different mean  $p_1$ , and Wilcoxon rank sum test for different median  $p_2$ .

Figure 3 confirms Rubinstein’s hypothesis. Heuristics (simple rules) imply a significantly higher level of consistency. These rules are widely adopted as 59% of the subjects are found to be either very impatient or perfectly patient. Figure 4 shows that this result is not driven by greater sophistication, Raven’s scores are not significantly different among groups. Finally, Figure 5 confirms the hypothesis that **Most Impatient** subjects follow a heuristic which does not require calculations and therefore answer questions more quickly than **Patient** subjects, who follow the summation rule, and **Other** subjects. **Deliberate Randomization** is very limited in Time as only 5% of the subjects have reported this behaviour in the questionnaire.<sup>16</sup> These subjects violate WARP more often than those who use heuristics but do not report significant differences in terms of Raven’s scores and Response Times.

---

<sup>15</sup>Examples of randomization in choices among lotteries are Hey & Carbone (1995), Ballinger & Wilcox (1997), Sopher & Narramore (2000), Hey (2001), Agranov & Ortoleva (2017), while Dwenger et al. (2018) shows similar evidence from university applications. Evidence of randomization between delayed payment plans can be found in studies involving Multiple Price Lists such as Andersen et al. (2008).

<sup>16</sup>Only two subjects are overlapping between the group of **Deliberate Randomization** and **Patient**, while the group **Most Impatient** is disjoint.

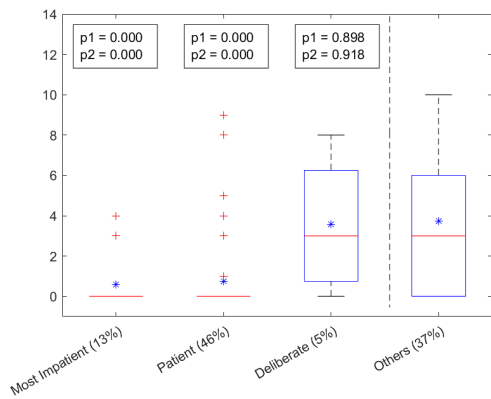


Figure 3: WARP violations by Group in Time.

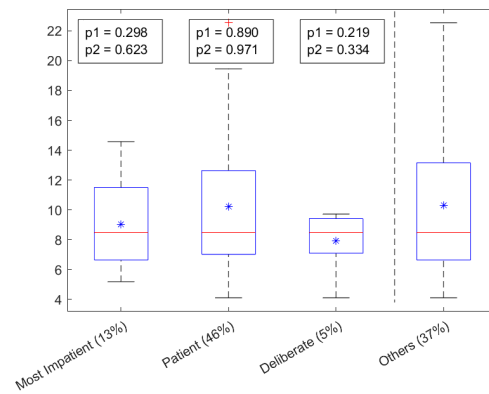


Figure 4: Raven's scores by Group in Time.

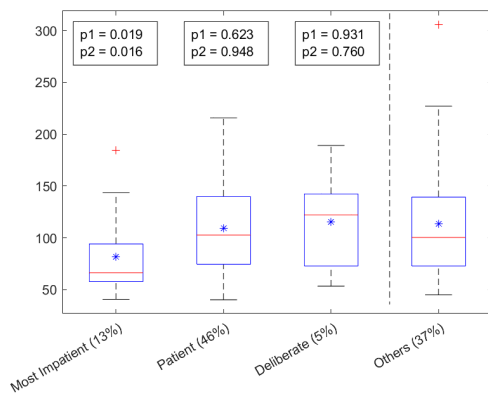


Figure 5: Response Times by Group in Time.

## NOTES:

The box plots contain results about WARP violations (Figure 4), Raven's score (Figure 5) and Response Times (Figure 6) divided into Groups: **Most Impatient**, **Patient**, **Deliberate**, and **Others**. We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first two box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

### 3.2.2 Risk preferences

First, we notice that the **Deliberate Randomization** behaviour is more widespread in Risk as 19% of the subjects reported it in the questionnaire.<sup>17</sup> Given that the construction of this group relies on direct reporting, it is possible that the 19% is underestimated. Different from Time, here we can also rely on the theoretical model of Cerreia-Vioglio et al. (2019) (see Appendix A2 for details) to identify subjects who deliberately randomize. This aspect is important and will be investigated in the next sections as it is based on the introduction of **PrAx**.

Analysing the figures below, Rubinstein's hypothesis is again confirmed. Figure 6 shows that subjects who use heuristics are significantly more consistent, however, contrarily to **Deliberate Randomization**, their use is far less common than in Time. Only

<sup>17</sup>Only one subject has reported to deliberately randomize and, at the same time, chose according to simple rules. This confirms how these behaviours are fundamentally different.



20% of the subjects have extreme preferences.<sup>18</sup> The use of heuristics and deliberate randomization is the main driver of the zero correlation observed in Figure 1 and the stark difference in WARP violations between Time and Risk as reported in Figure 2.

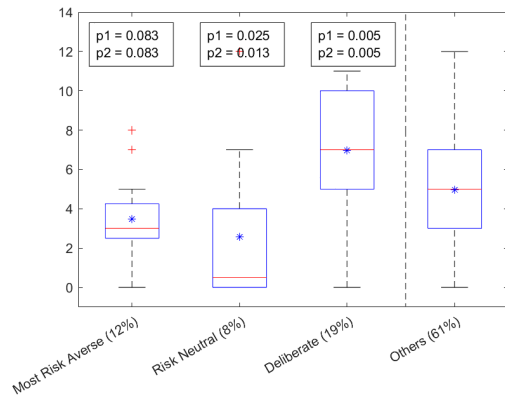


Figure 6: WARP violations by Group in Risk.

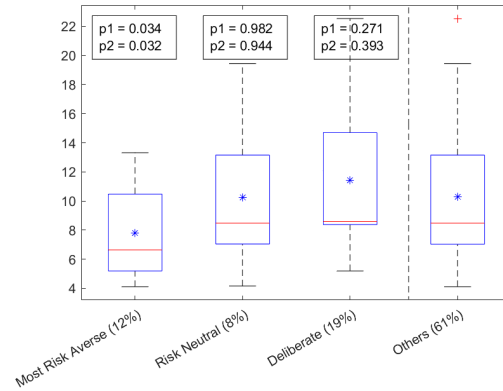


Figure 7: Raven's scores by Group in Risk.

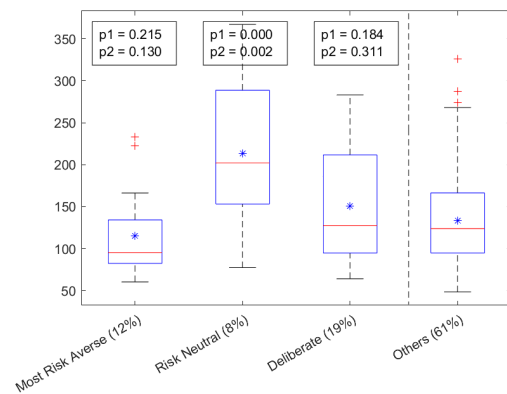


Figure 8: Response Times by Group in Risk.

#### NOTES:

The box plots contain results about WARP violations (Figure 7), Raven's score (Figure 8) and Response Times (Figure 9) divided into Groups: **Most Risk Averse**, **Risk Neutral**, **Deliberate**, and **Others**. We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first three box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

As expected, we find that subjects in the **Deliberate Randomization** group are significantly less consistent than **Others**. This is even more true if we compare them with the **Most Risk Averse** and **Risk Neutral** groups. This finding suggests that these subjects are fundamentally different in terms of behaviour, independent of their preferences.

Figure 7 reports two important findings. First, subjects who deliberately randomize do not present significant differences in Raven's scores when compared to **Others**, but they do have significantly higher scores than **Most Risk Averse** subjects. The result is somewhat expected because deliberate randomization involves an optimization problem

<sup>18</sup>Our data are in line with Andersson et al. (2016); in one of their treatment (statistics are similar in the other treatment), they find 14.8% of extremely risk-averse subjects. Our finding is 11.7%.

without attention or information costs. Second, more generally and in line with the literature,<sup>19</sup> **Most Risk Averse** subjects have significantly lower Raven's scores than the remaining subjects. Finally, Figure 8 confirms, on one hand, that the heuristic **Risk Neutral**, which requires the calculation of the expected value, implies significantly higher response times than the other groups; on the other hand, that the heuristic **Most Risk Averse** requires low response times.

### 3.3 Violations of preference axioms

In the previous subsection, we investigated the characteristics of subjects who used heuristics and those who deliberately randomized in our experiment. Now, we turn our attention to **PrAx** as we aim to investigate their potential role of confounding factors in determining the relationship between **ConAx** and measures of decision-making abilities. As mentioned in Section 2, we consider two **PrAx** in Time: monotonicity and impatience; and two in Risk: FOSD and SOSD (or equivalently, monotonicity and concavity of the Bernoulli utility).

#### 3.3.1 Time preferences

We construct dummies that take value 1 if a **PrAx** is violated and 0 otherwise.<sup>20</sup> In Time, violations of **PrAx** are relatively rare in our experiment as 20% of the subjects violated impatience and only 9% violated monotonicity.

First, we document a negative correlation between violations of impatience and monotonicity, and Raven's scores (resp. -0.276, -0.115), Response Times (resp. -0.259, -0.094), and Understanding (resp. -0.342, -0.143).<sup>21</sup> Importantly, the Response Times are calculated focusing only on problems in which **PrAx** could be violated, hence outside the MAIN ones. The variable "Understanding" is constructed by aggregating the first two questions of Questionnaire 1 presented in Section 4 of the Online Appendix.<sup>22</sup>

<sup>19</sup>See Dohmen et al. (2010) as a notable example and Andersson et al. (2016) for a review of the literature.

<sup>20</sup>The use of dummies does not influence our analysis as the distribution of violations of **PrAx** is skewed towards one given the relatively rare possibility of violating them in our experiment. In Time, 92.3% (resp. 62.0%) of the subjects who violated monotonicity (resp. impatience) did it only once while in Risk, these proportions for FOSD and SOSD are respectively 92.1% and 35.2%.

<sup>21</sup>All correlations regarding impatience are significantly different from zero even correcting for multiple hypothesis testing (Bonferroni correction), while those involving monotonicity are not significantly different from zero.

<sup>22</sup>In these questions, we asked if subjects had a good overall understanding of the experiment and if

We represent these results using box plots and report the differences in Raven's scores, Response Times, and Understanding between the group of subjects who violated the **PrAx**, and those who did not.

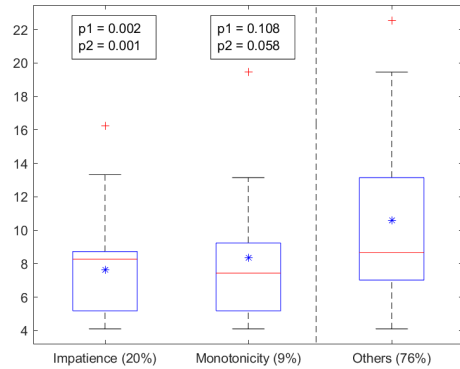


Figure 9: Raven's scores by **PrAx** in Time.

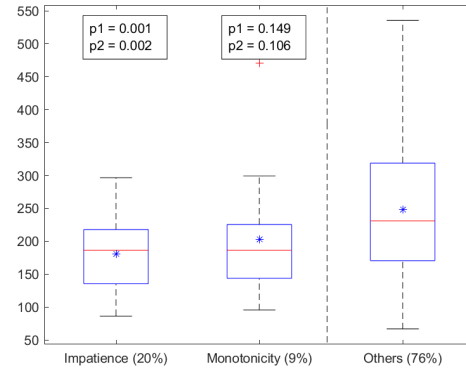


Figure 10: Response Times by **PrAx** in Time.

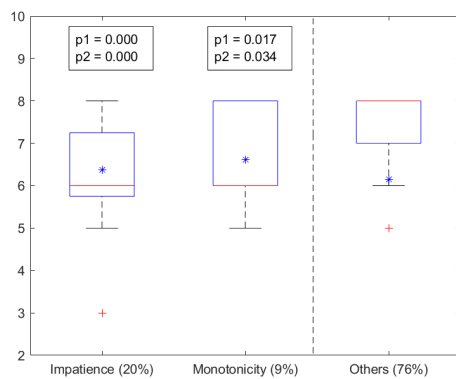


Figure 11: Understanding by **PrAx** in Time.

#### NOTES:

The box plots contain results about Raven's score (Figure 9), Response Times (Figure 10), and Understanding (Figure 11) divided into Groups: **Impatience** (subjects who violated impatience at least once), **Monotonicity** (subjects who violated monotonicity at least once), **Others** (remaining subjects). We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first three box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

Figures 9-11 show suggestive evidence that subjects who violated either impatience or monotonicity have lower decision-making abilities - or displayed lower effort - as they have lower Raven's scores, quicker response times, and a lower understanding of the experiment. This evidence strengthens the role of **PrAx** as potential confounding factors in the existing literature (Choi et al., 2014), (Andersson et al., 2016).

To provide more insights regarding violations of **PrAx**, in particular to control for the instructions provided were enough to understand the experiment's tasks. The answers were ordered from "strongly agree" to "strongly disagree", with values from 1 to 5. Overall, the subjects showed a good understanding. The two questions had a mean of 1.55 and 1.37 and a standard deviation of 0.6 and 0.62, respectively. We then take 10 minus the overall score to give a more immediate interpretation of the variable.

the use of heuristics, we perform a regression analysis that we report in Section 3.2 of the Online Appendix. We summarize our analysis in Figure 12.<sup>23</sup> In this figure and all the following ones, we report the 95% confidence interval.

First, heuristics have clear and easy-to-interpret effects. On one hand, **Most Impatient** subjects rarely violate impatience since the heuristic of choosing the first and highest histogram is a perfect remedy against violations of impatience. On the other hand, **Patient** subjects rarely violate monotonicity. In a similar fashion, the simple rule of choosing the plan that yields the highest summation of payments is a perfect remedy against violations of monotonicity.<sup>24</sup>

Even after controlling for the use of heuristics, we find that both Raven's scores, Response Times, and Understanding are negatively correlated with violations of impatience confirming our interpretation in line with lower decision-making abilities irrespective of the use of simple rules. Importantly, the coefficients of Response Times (Bonferroni-adjusted  $p < 0.001$ ) and Understanding ( $p = 0.0216$ ) are significant even after correcting the standard errors for multiple hypothesis testing. Non-surprisingly, Raven's scores and Understanding are significantly positively correlated (+0.225) and, if the variable Understanding is omitted (see the Online Appendix for details), the coefficient of Raven's scores becomes highly significant ( $p = 0.003$ ). This suggests that cognitive abilities affect violations of impatience partly through a lower understanding of the experiment.

We repeat the analysis for monotonicity. Even if the direction of the coefficients is confirmed, we do not find significant effects of Raven's scores, Response Times, and Understanding. Instead, we document that the correlation between violations of monotonicity and impatience is strong and significant both unconditionally (+0.266) and in our regression analysis (see the Online Appendix).

---

<sup>23</sup>In this figure, we report the regression coefficients of specification (4) in Tables A7a and A7b, Section 3.2 of the Online Appendix.

<sup>24</sup>Given the low proportion of subjects who reported **Deliberate Randomization**, we exclude them from Figure 12.

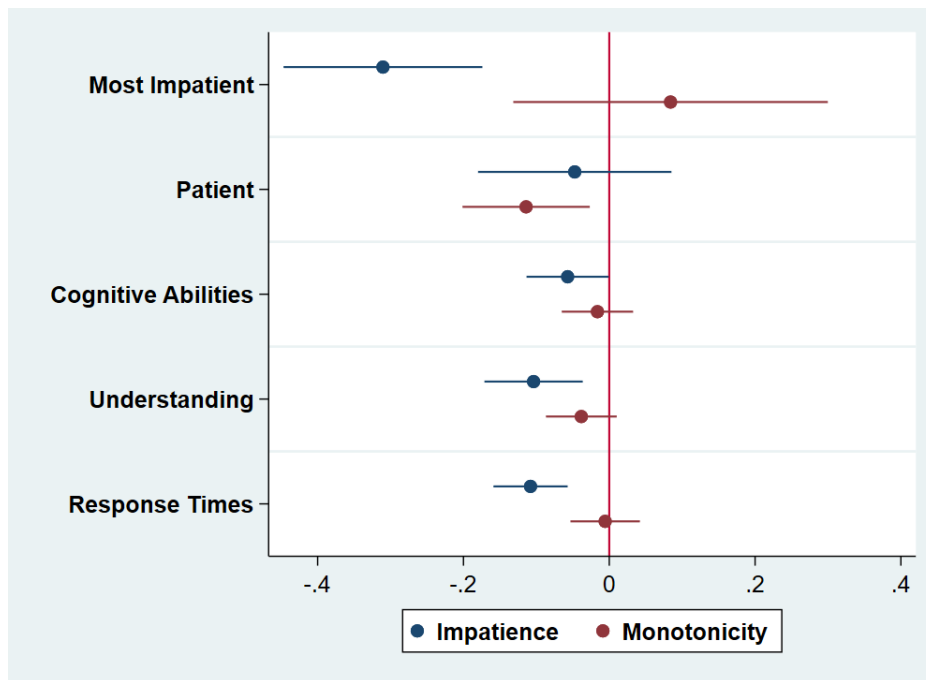


Figure 12: Violations of **PrAx** in Time

### 3.3.2 Risk preferences

In Risk, exploiting the results of Andersson et al. (2016) and Agranov & Ortoleva (2017), we first provide evidence of the role of **PrAx** as confounding factors by showing the difference between violations of **ConAx** with and without **PrAx**. As anticipated, the Multiple Price List design does not allow to test **ConAx** and **PrAx** separately. Andersson et al. (2016) documented that 14.8% and 30.48% of their subjects were inconsistent and non-monotonic, respectively, in their first and second Multiple Price Lists. We find that 26% of our subjects violate FOSD. This is in stark contrast with the number of subjects who violate only **ConAx** found both in the literature and in our experiment. Agranov & Ortoleva (2017) documented that 90% of their subjects were inconsistent while we find that 85.5% violated WARP only in our MAIN problems (Figure 2). These findings are already, per sé, strong evidence of the role of confounding factors that **PrAx** may play. We, now, investigate the factors that are correlated with the violations of **PrAx** in Risk.

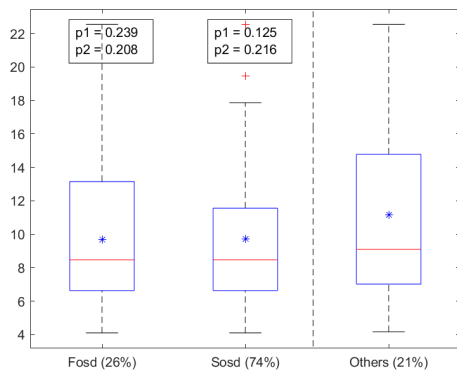


Figure 13: Raven's scores by **PrAx** in Risk.

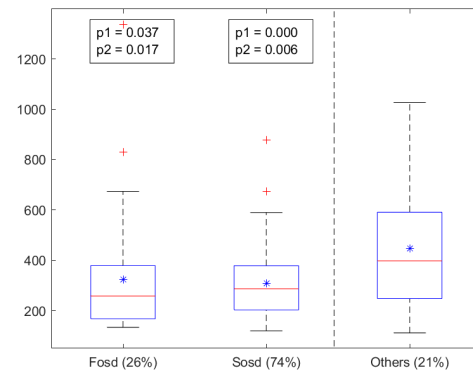


Figure 14: Response Times by **PrAx** in Risk.

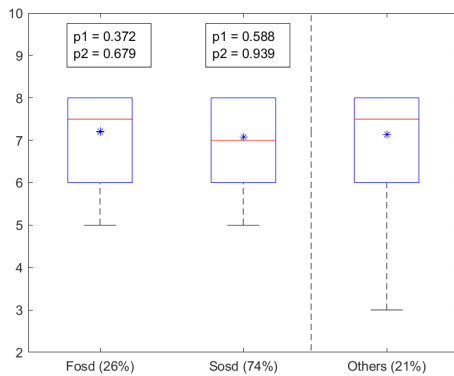


Figure 15: Understanding by **PrAx** in Risk.

#### NOTES:

The box plots contain results about Raven's score (Figure 13), Response Times (Figure 14), and Understanding (Figure 15) divided into Groups: **Fosd** (subjects who violated FOSD at least once), **Sosd** (subjects who violated SOSD at least once), **Others** (remaining subjects). We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first three box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

Even if violations of **PrAx** are more common in Risk than in Time, with 26% of subjects violating FOSD and 74% violating SOSD, we only partially replicate the previous results regarding impatience and monotonicity. We find weak correlations between violations of FOSD and SOSD, and Raven's scores (resp. -0.040, -0.101), Response Times (resp. -0.062, -0.315), and Understanding (resp. +0.074, +0.001). These results are equivalently reported in Figures 13-15 which, except for Response Times, do not display significant differences between subjects who did and did not violate **PrAx**.

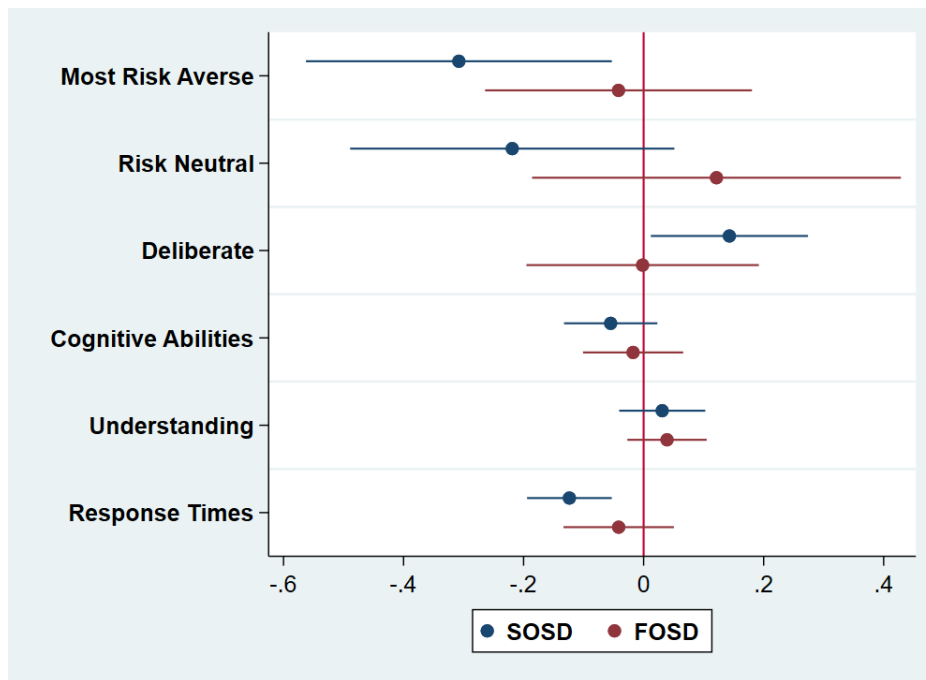


Figure 16: Violations of **PrAx** in Risk

So, why do subjects violate **PrAx** in Risk in our experiment? We start by describing the observable implications of deliberate randomization. As described in Appendix A2, Cerreia-Vioglio et al. (2019) show that subjects who deliberately randomize should not choose alternatives that are first-order stochastically dominated. Hence, if chosen by mistake outside the MAIN problems they should not predict inconsistencies inside it. Conversely, they do not constrain the choice of second-order stochastically dominated ones. Therefore, subjects who are described as **Deliberate** should violate SOSD more often.

More generally, to take into account the role of heuristics, we repeat the regression analysis of the previous section. Figure 16 summarizes our results.<sup>25</sup> First, **Deliberate Randomization** subjects report a significantly higher number of violations of SOSD, while there is no difference for FOSD. Other heuristics, on the other hand, reduce the violations of SOSD with a greater effect found for **Most Risk Averse** subjects. This is expected since highly risk averse individuals are unlikely to choose lotteries that would be preferred by risk loving ones. Even in this case, no significant effect is found for FOSD. Second, the correlations between violations of **PrAx**, Raven's scores, and Understanding are negative but we find no significant results. The coefficient of Response

<sup>25</sup>In this figure, we report the regression coefficients of specification (4) in Tables A8a and A8b, Section 3.2 of the Online Appendix.

times, instead, is significantly negative (Bonferroni-adjusted  $p=0.004$ ).

To summarize, on one hand, we find strong evidence of how heuristics affect violations of SOSD. These factors account for the majority of the explained variation. On the other hand, we do not find factors that can explain violations of FOSD; the amount of the explained variation is surprisingly low ( $R^2=0.019$ ). Furthermore, we detect an important difference between Time and Risk. Unlike violations of impatience and monotonicity, violations of FOSD and SOSD are not significantly correlated (+0.097). Furthermore, even if monotonicity and FOSD represent the violation of the same assumption on the utility function, the context matters as they are also not significantly correlated (+0.033).

### 3.4 Main results: violations of consistency axioms

We begin presenting the correlation between violations of **ConAx**, Raven's scores, Response Times, and Understanding and compare it to the correlation with violations of **PrAx**.

In Time, violations of WARP are - not significantly - negatively correlated with Raven's scores (-0.066) and Understanding (-0.115), and significantly positively correlated with Response Times (+0.204). In Risk, violations of WARP are not significantly correlated with Raven's scores (+0.018), Response Times (-0.005), and Understanding (+0.051). In both Time and Risk, these correlations are weaker, or even opposite sign in the case of Response Times, than those involving **PrAx** and presented in previous subsections. We take this as suggestive evidence that **ConAx** cannot be immediately connected to decision-making abilities.

We then ask, do subjects who violate **PrAx** also violate **ConAx**? Our findings indicate only weak correlations in Time: +0.109 and +0.08 respectively with impatience and monotonicity. In Risk, instead, the correlation between violations of FOSD and WARP is weak (+0.053), while it is positive and significant between violations of SOSD and WARP (+0.304).

As we have seen in Section 3.3, this latter result is driven by **Deliberate Randomization** subjects, but more generally, both the correlation between violations of **ConAx** and **PrAx**, and the relationship between violations of **ConAx**, Raven's scores, Response Times, and Understanding are heavily influenced by heuristics. Therefore, we perform a regression analysis that we summarize in Figure 17 for both Time and Risk.<sup>26</sup> The complete analysis, which involves also robustness checks for the measure-

<sup>26</sup>In both Time and Risk, the coefficients regard the specifications (7) of the regression analyses pre-



ment of heuristics, can be found in Section 3.4 of the Online Appendix. The dependent variable in the represented regressions is the number of violations of WARP.

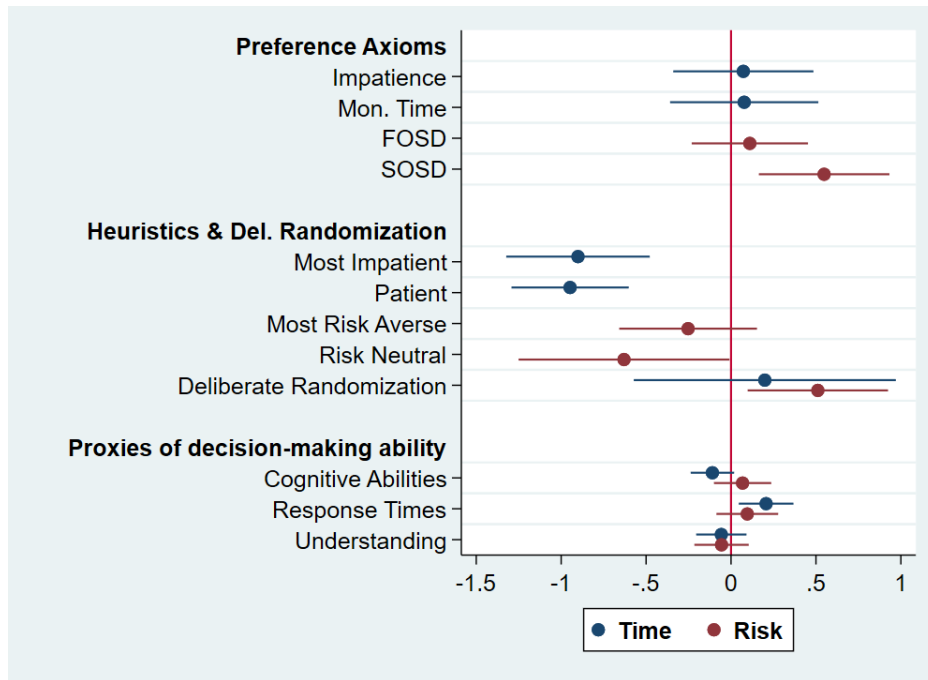


Figure 17: Factors correlated with Violations of **ConAx**.

First, we find no correlation between violations of **ConAx** and **PrAx** - except SOSD as expected. This result is relatively weak given the size of the standard errors of the coefficients of **PrAx**. However, these correlations, which are unconditionally positive and weak as shown above, disappear (see the Online Appendix for details) when we control for heuristics. This confirms that heuristics are remedies against both the violation of **PrAx** and **ConAx** and, outside their use, the connection between **PrAx** and **ConAx** seems to be weak. Non-surprisingly, instead, we find a significant positive correlation between WARP violations and SOSD. This result is perfectly in line with the predictions of the model of deliberate randomization as described in Cerreia-Vioglio et al. (2019).

Second, heuristics strongly reduce the number of WARP violations. Deliberate randomization, on the other hand, strongly increases the number of WARP violations but only in Risk, as in Time only a small proportion of the subjects (5%) reported this behaviour. Taken together, heuristics and deliberate randomization are the main factors driving WARP violations in both Time and Risk.

sented in Section 3.4 of the Online Appendix.

Third, cognitive abilities are correlated with violations of WARP with opposite signs in Time and Risk. This is a novel result since the literature has normally found negative correlations. In Time, cognitive abilities can be interpreted traditionally, for instance, as being connected to preference stability. In Risk, instead, subjects who choose to deliberately randomize act as if they diversify risk. This behaviour is correlated to higher cognitive abilities, while the risk averse subjects who do not randomize are normally correlated to lower cognitive abilities (Dohmen et al., 2010).

Fourth, response times are positively correlated with WARP violations. Therefore, we confirm previous findings by Rubinstein (2013). The interpretation is again related to the adoption of simple rules that guarantee high consistency and low response times.

Overall, our experiment shows that the relationship between **ConAx** and decision-making abilities is, at best, unclear. We document that violations of WARP display: (1) ambiguous correlations (positive in Risk and negative in Time) with cognitive abilities, (2) relatively robust positive correlations with Response Times, and (3) non-significant correlations with Understanding. These findings are substantially different from those reported in relation to **PrAx**.

### 3.5 Discussion and further research

Before concluding our paper, we discuss the relationship between our study and the literature on GARP. As mentioned in the introduction, **ConAx** are often considered as necessary conditions for high quality behaviour. For instance, Kariv & Silverman (2013), Choi et al. (2014), and Carvalho & Silverman (2019) argue that GARP is a necessary but not sufficient condition for high-quality decision-making. Choi et al. (2014) write: "if decisions are high-quality then there exists a utility function the choices maximize". The authors (see Section II.B in Choi et al. (2014)) stress the fact that GARP does not impose normative restrictions on the utility function. They propose a strong argument showing that GARP allows for violations of FOSD. We believe this point is worth a discussion.

Choi et al. (2007), Choi et al. (2014), and Dembo et al. (2022) find that the vast majority of subjects violate GARP, even if mildly. The result may be expected given the high number of choices subjects encountered in their experiment (see Figure 4, in Choi et al. (2007)). However, this may not be the reason why such a high number of violations arise. First, Dembo et al. (2022) show that most of the subjects who violate GARP also violate FOSD, hence, violations of GARP and FOSD seem not to differ substantially. Further, the high number of GARP violations in these experiments, as well as the over-

lapping with violations of FOSD, seems to be a strikingly different result from many others observed in the literature. For instance, on one hand, in Agranov & Ortoleva (2017), as well as in our experiment, we observe a big difference between violations of WARP and FOSD in choices among gambles. On the other hand, in Manzini & Mariotti (2010), as well as in our experiment, we observe high consistency in Time. Finally, in several experiments involving MPLs (among many, see Andersen et al. (2008), Andersson et al. (2016)) FOSD seems to be violated by a low percentage of subjects ranging from 10% to 30%. On the contrary, Dembo et al. (2022) observe no subjects satisfying FOSD in a 3-dimensional case, while only a handful satisfy it in the 2-dimensional case (see the Appendix in Dembo et al. (2022)) equivalent to Choi et al. (2007) and Choi et al. (2014).

We propose the following explanation. FOSD in the context of budget sets is a complicated theoretical concept as it involves the interaction between states of the world. Monotonicity, in the sense of FOSD, is required on the vNM utility function (hence, between states of the world) as shown by Dembo et al. (2022). GARP, on the other hand, poses constraints (monotonicity and concavity) on the Bernoulli utility function, namely only within states (see Hansen et al. (1978) for a discussion of state-by-state stochastic dominance). In view of these considerations, to reconcile the listed findings, firstly note that the high number of violations of GARP is in line with the high number of violations of WARP observed in choices among gambles in many experiments, including ours. Secondly, the low number of violations of FOSD in experiments involving MPLs can be rationalized by the fact that, in this context, states of the world are unlabeled, hence FOSD is equivalent to monotonicity within states of the world as in Hansen et al. (1978). We argue that subjects may be able to apply stochastic dominance when this task is relatively simple as in MPLs or in our experiment, while they fail when the task becomes hard. To conclude, we believe further research is needed to have a deeper understanding of why decision-makers may violate (or satisfy) different stochastic dominance properties, as well as other **PrAx**, and how these properties are related to **ConAx**.

## 4 Conclusion

We challenge the definition of rationality as consistency using the novel concepts of consistency and preference axioms. We aim to show that an idea of rationality connected to good decision-making abilities is not correlated with consistency axioms, and

highly depends on the context (time and risk preferences). We design an experiment that allows us to test consistency alone. In our experiment, subjects answer a series of questions regarding time and risk outcomes. Subjects can violate consistency axioms only in a subgroup of questions while, in the remaining part, they can violate also preference axioms.

We find substantial differences in behaviour between Time and Risk. We break down the analysis by heuristics, preferences, cognitive abilities, response times, and level of understanding of the experiment. The main result is that the idea that consistency, or standard economic rationality, is a good proxy for decision-making ability is often misleading. We find no correlation between violations of preference and consistency axioms. Testing consistency alone, we find that, in Time, inconsistencies are strongly affected by heuristics. Subjects who use heuristics are far more consistent. In Risk, we confirm the results that connect heuristics to consistency, as well as strong use of deliberate randomization behaviour, confirming the findings of Sopher & Naramore (2000) and Agranov & Ortoleva (2017). Finally, in our experiment, we do not find evidence that justifies the use of consistency as a measure of decision-making ability measured by cognitive abilities, response times, and level of understanding of the experiment.

## Appendix A1 - Structural, Preference, and Consistency Axioms

In this appendix, we present the framework in which we build our distinction between **ConAx** and **PrAx**. Let  $X$  be a set of alternatives and  $\geq \in X \times X$  be a binary relation.  $\mathcal{X}$  is the set of non-empty subsets of  $X$ . The set of all binary relations is denoted as  $\mathcal{P}$ . A choice correspondence is a mapping  $c : \mathcal{X} \rightarrow \mathcal{X}$  with  $c(A) \in A$  for all  $A \in \mathcal{X}$ . The set of all choice correspondences is denoted as  $\mathcal{C}$ . Given a generic set  $S$ , an axiom **Ax** is a constraint  $\mathbf{Ax}(S) \subseteq S$ . Here we analyse axioms on the sets  $\mathcal{P}$  and  $\mathcal{C}$ . We use the term structure to define mathematical objects that we endow on the set  $X$ . For instance, a topological structure is a couple  $(X, \tau)$  where  $\tau$  is a topology. An order structure is a couple  $(X, \geq)$  such as a poset. What follows is the description of three families of axioms: structural axioms  $\mathbf{StAx}(\mathcal{P}) \subseteq \mathcal{P}$ , preference axioms  $\mathbf{PrAx}(\mathcal{P}) \subseteq \mathcal{P}$ , and consistency axioms  $\mathbf{ConAx}(\mathcal{C}) \subseteq \mathcal{C}$ .

### Structural Axioms and Preference Axioms

The difference between **StAx** and **PrAx** lies in the symmetric nature of their constrained sets. **StAx** do not discriminate alternatives by their labels while **PrAx** do. Let  $\pi : X \rightarrow X$  be a permutation over  $X$ . For all  $\pi$  and for all **StAx** we have that  $\mathbf{StAx}(\mathcal{P}) = \mathbf{StAx}(\pi(\mathcal{P}))$ . On the contrary, for all **PrAx** there are some  $\pi$  such that  $\mathbf{PrAx}(\mathcal{P}) \neq \mathbf{PrAx}(\pi(\mathcal{P}))$ . In the paper, we refer to this property as neutrality.

But why does this difference arise? The reason has to be searched in the structures that endow the set of alternatives. When the structure is neutral like a topology or a metric then only **StAx** can arise, for instance, continuity or local non-satiation. If no structure exists again we have only **StAx** such as completeness and transitivity. Instead when the structure is non-neutral like an order structure then **PrAx** arise. For instance, let  $X = \mathfrak{R}$  and  $\geq$  be the decreasing order. We can define monotonicity (**Mon**) as: if  $x \geq y$  then  $x \geq y$ . Note how **Mon** is not neutral. Let  $x = 5$  and  $y = 6$ , we have  $\mathbf{Mon}(\mathcal{P}) = \{(x, y)\}$  and  $\mathbf{Mon}(\pi(\mathcal{P})) = \{(y, x)\}$ .

### Structural and Consistency Axioms

There is a tight connection between **StAx** and **ConAx** that is guaranteed by the adoption of binary relations as primitives of some model  $\sigma$ . An example is the maximization model, which is the object of study in this paper. More formally, for all  $A \in \mathcal{X}$  we have

$c(A) = \mathbf{Max}(A, \geq)$ . This model creates a mapping  $\sigma : \mathbf{StAx} \rightarrow \mathbf{ConAx}$ ,<sup>27</sup> which can be described as follows:  $\sigma[\mathbf{Com}(\mathbf{Tr}(\mathcal{P}))] = \mathbf{WARP}(C)$ , namely  $c$  satisfies **WARP** if and only if there exists a transitive and complete preference relation  $\geq$  that rationalizes  $c$ . Therefore, in the paper, our analysis of **ConAx** can be seen as an indirect analysis of **StAx**. Note that **ConAx** are more general than **StAx**. This is because there exist many more models than **Max** that produce choice functions. However, as for **StAx**, also **ConAx** satisfy the neutrality property on the set  $C$ .

In one of our motivating examples, we argue that a test of **GARP** implies more than a simple test of **ConAx**. In this case, the problem arises from the particular convex structure of the budget sets on which **GARP** is defined. This additional structure on  $X$  implies that a **ConAx**, such as **GARP**, may induce not only the existence of a transitive and complete preference relation  $[\mathbf{StAx}(\mathcal{P})]$ , but also a particular subset of these preference relations with respect to monotonicity and concavity  $[\mathbf{PrAx}(\mathcal{P})]$ .

## Appendix A2 - Deliberate Randomization

The idea of Deliberate Randomization has been modelled by Machina (1985) and more recently by Cerreia-Vioglio et al. (2019). It describes subjects that have a deterministic preference over the convex hull of a set of lotteries and randomize to obtain the optimum. We use the framework of Cerreia-Vioglio et al. (2019) to formalize this idea. Let  $[w, b] \subseteq \mathfrak{X}$  be an interval on monetary prizes and  $\Delta$  be the set of lotteries over  $[w, b]$ . Denote  $\mathcal{A}$  as the collection of all finite, non-empty subsets of  $\Delta$ ; and  $co(A)$  as the convex hull of  $A \in \mathcal{A}$ . A stochastic choice function  $\rho$  is a map that assigns to each  $A$  a probability distribution  $\rho(A)$ . Particularly,  $\rho(A)$  is a compound lottery and  $\overline{\rho(A)}$  is the induced lottery over monetary outcomes:

$$\overline{\rho(A)} = \sum_{q \in A} \rho(A)(q)q$$

A stochastic choice function  $\rho$  has a Deliberate Stochastic Choice representation if there exists a complete preorder  $\geq$  over  $\Delta$  such that:

1. For every  $A \in \mathcal{A}$ :

$$\overline{\rho(A)} \geq q \quad \forall \quad q \in co(A)$$

2. For all  $p, q \in \Delta$ ,  $p \succ_{FOSD} q$  implies  $p \succ q$ .

---

<sup>27</sup>The general properties of this mapping are analysed by Mahmoud (2017).

## Appendix B - Heuristics in the Questionnaire

In the last part of the experiment, we collected information through a questionnaire. In particular, we used open questions to ask which model of behaviour did subjects adopt in Time and Risk. Hence, we use the answers to identify those who used the heuristics identified by our extreme preferences, as well as deliberate randomization.

Below, we provide examples of the reported modes of behaviour that we identify as heuristics:<sup>28</sup>

### Patient

- "Highest summation every time."
- "Added the tokens up, and chose the highest."
- "always went for the choice that over the long run gave the greatest income"

### Most Impatient

- "I strongly preferred payment plans that paid most or all of the total amount today, even if the total amount was smaller than other payment plans that involved long delays."
- "highest payment in time 0."
- "It's better to get money now since money tends to lose its value within time. In UK after BREXIT there is a high risk of increasing inflation, therefore money I get in 12 months may cost nothing."

### Deliberate Randomization (Time)

- "I did usually calculate the summation of the plan, but sometimes chose a plan which would be paid more quickly."
- "I considered the total payment as a factor but if a high fraction of the largest total payment option was available very quickly then I was tempted to switch."
- "Oftentimes used highest summation criterion, but also sometimes went for the instant gratification option due to the relatively low amount of money at stake."

---

<sup>28</sup>The dataset with the complete list of answers, as well as dataset and codes, are available upon request.

**Most Risk Averse**

- "least risky."
- "I chose the options that were the safest but still offered a decent amount of tokens such as 100% probability for 50 tokens since this guaranteed tokens."
- "most likely to win something. less risk."
- "Highest probability to win."

**Risk Neutral**

- "highest expected value."
- "Calculate expected value, note down. If a question only contains seen lotteries, look up previous values. I may have chosen some 300|15 when I thought they were 300|20."
- "Disciplined myself into being risk-neutral and chose highest expected value lotteries even though it was hard to choose a less secure option"

**Deliberate Randomization (Risk)**

- "I like to take risk to get the best payment, but sometimes I will take a middle one."
- "Highest likelihood of a reasonable gain. Once or twice I went big with high-risk, high-reward options."
- "in some questions, I make decisions based on expected value while in some other questions, I prefer the certain gain."



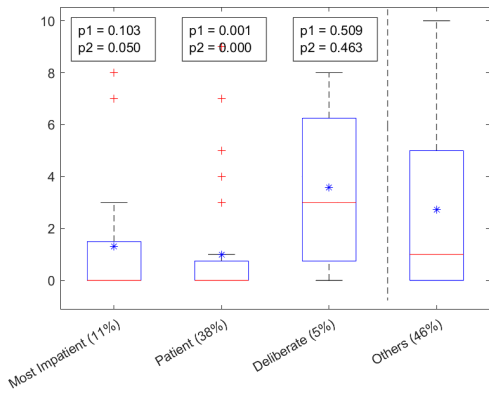


Figure 18: WARP violations by Group in Time.

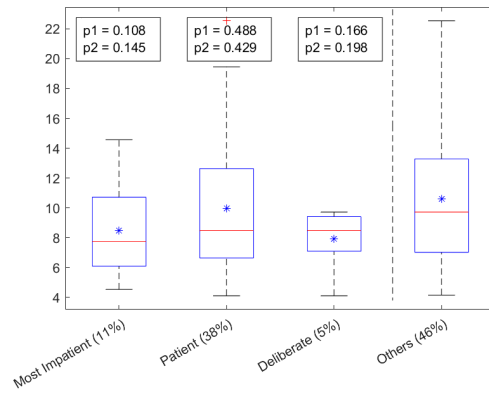


Figure 19: Raven's scores by Group in Time.

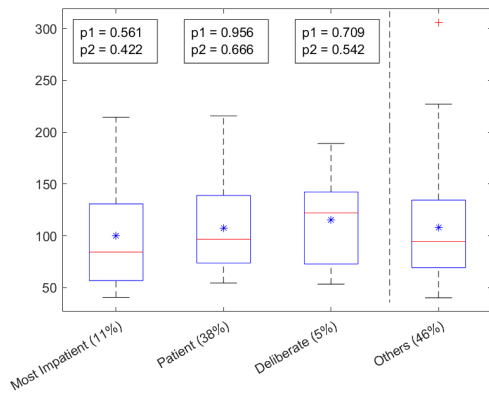


Figure 20: Response Times by Group in Time.

NOTES:

The box plots contain results about WARP violations (Figure 4), Raven's score (Figure 5) and Response Times (Figure 6) divided into Groups: **Most Impatient**, **Patient**, and **Others**. We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first two box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

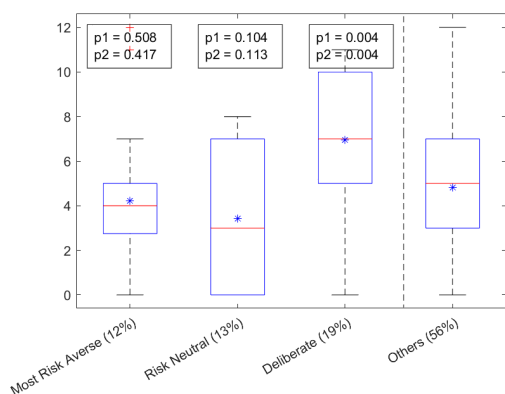


Figure 21: WARP violations by Group in Risk.

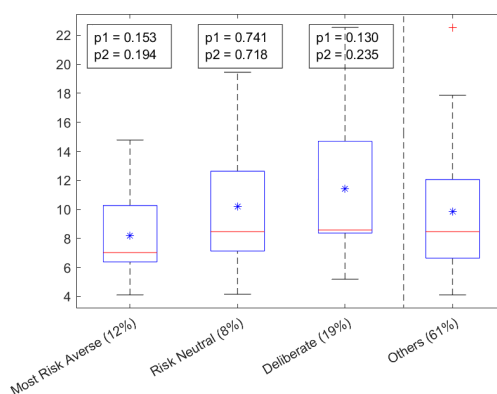


Figure 22: Raven's scores by Group in Risk.

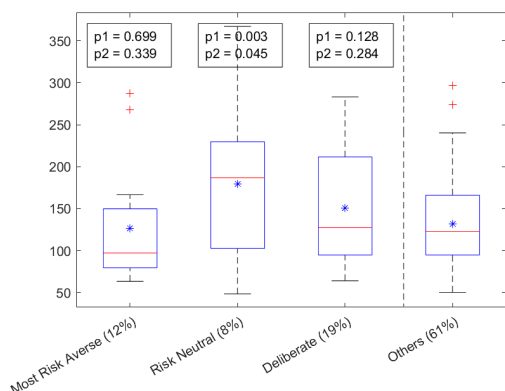


Figure 23: Response Times by Group in Risk.

## NOTES:

The box plots contain results about WARP violations (Figure 7), Raven's score (Figure 8) and Response Times (Figure 9) divided into Groups: **Most Risk Averse**, **Risk Neutral**, **Deliberate**, and **Others**. We report the median (red line), mean (blue star), 25<sup>th</sup> and 75<sup>th</sup> percentile and outliers (over 1.5 times the interquartile range above the median). Above the first three box plots we report two statistics that compare these Groups with **Others**: unpaired t-test for equal mean (p1) and Wilcoxon rank sum test for equal median (p2).

First, our reported measures are highly correlated with the heuristics identified through the elicitation of extreme preferences. In particular, for **Patient** subjects the correlation is 0.56 (p-value=0.000), for **Most Impatient** it is 0.65 (p-value=0.000), for **Risk Neutral** it is 0.40 (p-value=0.000), and for **Most Risk Averse** it is 0.13 (p-value=0.11). In this latter case, it is likely that the subjects had a clear idea of their best option but struggled to choose among the remaining ones. The correlation between **Most Risk Averse** subjects identified using the questionnaire and those whose elicited best element is the degenerate lottery is 0.34 (p-value=0.000).

As the reader can see comparing Figures 3-5 and Figures 18-20, the results are very similar in Time, with subjects who use heuristics being significantly more consistent. Similarly, comparing Figures 6-8 and Figures 21-23, we can see that our results are confirmed also in Risk. As a further robustness check, we have replicated our regression analysis of Section 3.4 using the reported heuristics and the results are confirmed (see specification 4 in Table A10 and A11 in the Online Appendix).

## Appendix C - Response Times

Response times are notoriously noisy and our experiment is not an exception. However, we present some indisputable evidence regarding their importance to explain several features of our experiment. Furthermore, we consider more generally the results that follow as sanity checks for our experiment. We represent our results using tables and we report only the Wilcoxon ranksum test and not the t-test to ease the exposition. The interpretation of the results is not affected.

TABLE 3: RESPONSE TIMES - TIME vs RISK

	<b>TIME</b>	<b>RISK</b>	<b>Wilcoxon</b>	
<b>Binary</b>	7.60	9.64	-4.88	0.000
<b>Ternary</b>	8.71	11.21	-3.73	0.000
<b>Quaternary</b>	9.73	13.44	-3.47	0.001
<b>Asymmetric Dominance</b>	14.07	16.23	-1.54	0.123
<b>Big Sets</b>	30.98	39.04	-3.78	0.000

NOTE -- The table reports the median response times in each category of questions. We also report the Wilcoxon ranksum test for equal median with the corresponding p-value. See the Online Appendix for a detailed description of each category.

In Table 3 we report the response times in the same domains in Time and Risk. The differences are highly significant in any domain. We reinforce the idea that choices among gambles are generally more complicated than choices among delayed payment plans. Table 4 reports, within Time and Risk, the z-statistics for the hypothesis of equal median between response times in different domains. The evidence shows that the higher cardinality of the set, the higher the response times. Furthermore, the presence of behavioural effects, particularly asymmetric dominance, increases the response time. This is somewhat surprising since one would expect the dominant alternative to be chosen quickly. Instead, we report higher response times in AD sets than in both ternary and quaternary sets. We infer that the presence of a new alternative, even if dominated, creates a learning process.

TABLE 4: RESPONSE TIMES BY CARDINALITY

	TIME					RISK					
	Bin	Ter	Quat	AD	Big	Bin	Ter	Quat	AD	Big	
Bin	-	-3.72***	-3.78***	-10.08***	-14.21***	Bin	-	-3.08***	-3.77***	-7.04***	-13.97***
Ter	-	-	-0.63	-6.48***	-13.23***	Ter <td>-</td> <td>-</td> <td>-1.40</td> <td>-4.14***</td> <td>-13.13***</td>	-	-	-1.40	-4.14***	-13.13***
Quat	-	-	-	-5.10***	-12.07***	Quat <td>-</td> <td>-</td> <td>-</td> <td>-2.08**</td> <td>-11.02***</td>	-	-	-	-2.08**	-11.02***
AD	-	-	-	-	-10.98***	AD <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>-11.22***</td>	-	-	-	-	-11.22***
Big	-	-	-	-	-	Big <td>-</td> <td>-</td> <td>-</td> <td>-</td> <td>-</td>	-	-	-	-	-

NOTE -- The table reports the z-values associated with the Wilcoxon ranksum test of the hypothesis of equal median between the mean of response times in the category of questions on the columns against those on the rows. Notice how all reported values are negative, meaning that the response times in the of rows is always smaller than the category on the columns. We add \*\*\* for significance at 1%, \*\* for significance at 5%, \* for significance at 10% and no stars otherwise.

Finally, Table 5 presents an analysis of response times by parts. We exploit the fact that half the participants answered the Time part first while the other half the Risk part. Furthermore, we divide each part (Time and Risk) into two halves of 13 questions: we insert in both half the quaternary set that was always asked as the 13th question; hence the two parts are numerically and qualitatively symmetric. Note that, the use of different orders of questions allows us to reduce the probability that particular easy/difficult questions drive the result.

TABLE 5: RESPONSE TIMES BY PARTS

	First Part				Second Part				
	TIME	RISK	z	p	TIME	RISK	z	p	
First Half	289.84	367.49	-2.05	0.04	First Half	253.30	339.98	2.76	0.00
Second Half	186.12	238.70	-2.35	0.02	Second Half	164.68	196.08	1.57	0.12
z	4.99	5.00			z	5.53	5.50		
p	0.00	0.00			p	0.00	0.00		

NOTE -- The table reports the sum of response times of questions regarding the first and second halves of questions in Time and Risk when these were asked either as first or second part. We report the Wilcoxon ranksum test for equal median between response times on rows and columns. In the text you can find the results for the same test between tables.

The higher difficulty of Risk questions is confirmed at any level of the experiment. The rows show the z-statistics at each quarter of the experiment. A new insight is present: individuals became quicker in the second half of the 25 questions in both parts of the experiment and both in Time and Risk. This evidence is confirmed by the z-statistics in the columns. Finally, there are signs of institutional learning as defined by Day et al. (1987). Namely, individuals learnt the design of the experiment independently from the types of questions asked. This is confirmed by z-statistics between the tables. Since a three-dimensional table is not representable we limit ourselves to listing them. In the First Half, between the First and the Second Part, the answers were quicker in this latter both in Time and Risk; the p-values related to the Wilcoxon ranksum test are respectively 0.085 and 0.442 (hence not significant in this last case) in Time and Risk. In the Second Half, the p-values are respectively 0.082 and 0.024 in Time and

Risk confirming the natural hypothesis of institutional learning.

## References

- Afriat, S. N. (1967). The construction of utility functions from expenditure data. *International Economic Review*, 8(1), 67–77.
- Agranov, M., & Ortoleva, P. (2017). Stochastic choice and preferences for randomization. *Journal of Political Economy*, 125(1), 40–68.
- Andersen, S., Harrison, G. W., Lau, M. I., & Rutstrom, E. E. (2008). Eliciting time and risk preferences. *Econometrica*, 76(3), 583–618.
- Andersson, O., Holm, H. J., Tyran, J.-R., & Wengstrom, E. (2016). Risk aversion relates to cognitive ability: preferences or noise? *Journal of European Economic Association*, 14(5), 1129–1154.
- Apestequia, J., & Ballester, M. A. (2015). A measure of rationality and welfare. *Journal of Political Economy*, 6(123), 1278–1310.
- Apestequia, J., Ballester, M. A., & Lu, J. (2017). Single-crossing random utility models. *Econometrica*, 85(2), 661–674.
- Ballinger, T. P., & Wilcox, N. T. (1997). Decisions, error and heterogeneity. *The Economic Journal*, 107(443), 1090–1105.
- Banks, J., Carvalho, L. S., & Perez-Arce, F. (2019). Education, decision making, and economic rationality. *The Review of Economics and Statistics*, 101(3), 428–441.
- Block, H., & Marschak, J. (1960). Random orderings and stochastic theories of responses. *Contributions to Probability and Statistics*, (Stanford University Press).
- Burks, S. V., Carpenter, J. P., Goette, L., & Rustichini, A. (2009). Cognitive skills affect economic preferences, strategic behavior, and job attachment. *Proceedings of the National Academy of Sciences*, (106), 7745–7750.
- Caliari, D. (2023). Behavioral welfare analysis and revealed preference: theory and experimental evidence. *Working paper*.
- Caplin, A., & Dean, M. (2015). Revealed preference, rational inattention, and costly information acquisition. *American Economic Review*, 105(7), 2183–2203.

- Caplin, A., Dean, M., & Martin, D. (2011). Search and satisficing. *American Economic Review*, *101*(7), 2899–2922.
- Carvalho, L., & Silverman, D. (2019). Complexity and sophistication. *NBER Working Paper 26036*.
- Cason, T. N., & Plott, C. R. (2014). Misconceptions and game form recognition: challenges to theories of revealed preference and framing. *Journal of Political Economy*, *122*(6), 1235–1270.
- Cerreia-Vioglio, S., Dillenberger, D., Ortoleva, P., & Riella, G. (2019). Deliberately stochastic. *American Economic Review*, *7*(109), 2425–2445.
- Chew, S. H., Miao, B., Shen, Q., & Zhong, S. (2022). Multiple-switching behavior in choice-list elicitation of risk preference. *Journal of Economic Theory*, *204*.
- Choi, S., Fisman, R., Gale, D., & Kariv, S. (2007). Consistency and heterogeneity of individual behavior under uncertainty. *American Economic Review*, *97*(5), 1921–1938.
- Choi, S., Kariv, S., Müller, W., & Silverman, D. (2014). Who is (more) rational? *American Economic Review*, *104*(6), 1518–50.
- Day, B., Bateman, I., Carson, R., Dupont, D., Louviere, J., Morimoto, S., Scarpa, R., & Wang, P. (1987). Ordering effects and choice set awareness in repeat-response stated preference studies. *Journal of Environmental Economics and Management*, *1*(63), 73–91.
- Dembo, A., Kariv, S., Polisson, M., & Quah, J. K. (2022). Ever since allais. *Working paper*.
- Dohmen, T., Falk, A., Huffman, D., & Sunde, U. (2010). Are risk aversion and impatience related to cognitive ability? *American Economic Review*, *100*(3), 1238–1260.
- Dwenger, N., Kubler, D., & Weizsacker, G. (2018). Flipping a coin: evidence from university applications. *Journal of Public Economics*, (167), 240–250.
- Echenique, F., Lee, S., & Shum, M. (2011). The money pump as a measure of revealed preference violations. *Journal of Political Economy*, *119*(6), 1201–1223.
- Fischbacher, U. (2007). z-tree: Zurich toolbox for ready-made economic experiments. *Experimental Economics*, (10), 171–178.

- Fishburn, P. C., & Rubinstein, A. (1982). Time preference. *International economic review*, 23(3), 677–694.
- Fudenberg, D., Iijima, R., & Strzalecky, T. (2015). Stochastic choice and revealed perturbed utility. *Econometrica*, 83(6), 2371–2409.
- Gilboa, I. (2009). Theory of decision under uncertainty. *Cambridge university press*, (Vol. 45).
- Hansen, L., Holt, C., & Peled, D. (1978). A note on first degree stochastic dominance. *Economics Letters*, 1(4), 315–319.
- Hey, J. (2001). Does repetition improve consistency? *Experimental economics*, 4(1), 5–54.
- Hey, J., & Carbone, E. (1995). Stochastic choice with deterministic preferences: an experimental investigation. *Economics Letters*, 47(2), 161–167.
- Holt, C. A., & Laury, S. K. (2002). Risk aversion and incentive effects. *American Economic Review*, 92(5), 1644–1655.
- Horan, S., & Sprumont, Y. (2016). Welfare criteria from choice: An axiomatic analysis. *Games and Economic Behavior*, (99), 56–70.
- Iyengar, S. S., & Kamenica, E. (2010). Choice proliferation, simplicity seeking, and asset allocation. *Journal of Public Economics*, (94), 530–539.
- Kariv, S., & Silverman, D. (2013). An old measure of decision-making quality sheds new light on paternalism. *Journal of Institutional and Theoretical Economics*, 169(1), 29–44.
- Kreps, D. M. (2015). Choice, dynamic choice, and behavioral economics. *Stanford Graduate School of Business, lecture*, (<http://stanford.io/1GxjZfg>).
- Machina, M. (1985). Stochastic choice functions generated from deterministic preferences over lotteries. *The Economic Journal*, (95), 575–594.
- Mahmoud, O. (2017). On the consistency of choice. *Theory and Decision*, (83), 547–572.
- Manzini, P., & Mariotti, M. (2007). Sequentially rationalizable choice. *American Economic Review*, 97(5), 1824–1839.

- Manzini, P., & Mariotti, M. (2010). Revealed preference and boundedly rational choice procedures: an experiment. *Unpublished*.
- Manzini, P., & Mariotti, M. (2012). Choice by lexicographic semiorders. *Theoretical economics*, 7(1), 1–23.
- Masatlioglu, Y., Nakajima, D., & Ozbay, E. Y. (2012). Revealed attention. *American Economic Review*, 102(5), 2183–2205.
- McCausland, W. J., Davis-Stober, C., Marley, A., Park, S., & Brown, N. (2019). Testing the random utility hypothesis directly. *The Economic Journal*, (130), 183–207.
- Natenzon, P. (2019). Random choice and learning. *Journal of Political Economy*, 127(1), 419–457.
- Rieskamp, J., Busemeyer, J. R., & Mellers, B. A. (2006). Extending the bounds of rationality: Evidence and theories of preferential choice. *Journal of Economic Literature*, 44(3), 631–661.
- Rubinstein, A. (2013). Response time and decision making: an experimental study. *Judgement and Decision Making*, 8(5), 540–551.
- Sen, A. (1971). Choice functions and revealed preference. *The Review of Economic Studies*, 38(3), 307–317.
- Sopher, B., & Narramore, M. J. (2000). Stochastic choice and consistency in decision making under risk: an experimental study. *Theory and Decision*, 48(4), 323–350.
- Sugden, R. (1991). Rational choice: a survey of contributions from economics and philosophy. *The Economic Journal*, 101(407), 751–785.
- Tversky, A., & Russo, J. E. (1969). Substitutability and similarity in binary choices. *Journal of Mathematical Psychology*, 6(1), 1–12.
- Yu, C. W., Zhang, Y. J., & Zuo, S. X. (2021). Multiple switching and data quality in the multiple price list. *The Review of Economics and Statistics*, 103(1), 136–150.



## Discussion Papers of the Research Area Markets and Choice 2023

### Research Unit: **Market Behavior**

- Levent Neyse, Frank M. Fossen, Magnus Johannesson, Anna Dreber** SP II 2023-201  
Cognitive reflection and 2D:4D: Evidence from a large population sample
- Christian Basteck, Lars Ehlers** SP II 2023-202  
On the constrained efficiency of strategy-proof random assignment
- Hande Erkut, Ernesto Reuben** SP II 2023-203  
Social networks and organizational helping behavior: Experimental evidence from the helping game

### Research Unit: **Economics of Change**

- Kai Barron, Tilman Fries** SP II 2023-301  
Narrative persuasion
- Christina Timko, Maja Adena** SP II 2023-302  
Transparent app design reduces excessive usage time and increases willingness to pay compared to common behavioral design—a framed field experiment
- Daniele Caliari** SP II 2023-303  
Behavioural welfare analysis and revealed preference: Theory and experimental evidence
- Daniele Caliari** SP II 2023-304  
Rationality is not consistency