

# Notes on Decision Theory and Practice

02.229 - Decision Theory and Practice, 2019 Jan-April

Yustynn Panicker

February 6, 2019

## Contents

<b>1</b>	<b>Course Information</b>	<b>3</b>
1.1	Instructor . . . . .	3
1.2	Office Hours . . . . .	3
1.2.1	Zsombor's Classes, can find him after . . . . .	3
1.3	Textbook . . . . .	3
1.4	Grading . . . . .	3
1.4.1	Reading Summaries . . . . .	3
1.4.2	Case Study . . . . .	3
<b>2</b>	<b>TODO Current Questions to Ask</b>	<b>4</b>
2.1	<b>TODO</b> For Zsombor . . . . .	4
2.1.1	<b>TODO</b> Can you <b>trust</b> someone to act against your interest? . . . . .	4
2.2	<b>TODO</b> For Not Zsombor . . . . .	4
2.2.1	<b>TODO</b> What was Zsombor's Paradigm/Theory distinction thing? . . . . .	4
2.2.2	<b>TODO</b> What's completeness again? . . . . .	4
<b>3</b>	<b>Normative Decision Theory</b>	<b>4</b>
3.1	W1: Introduction. Decision trees. . . . .	4
3.1.1	Peterson CH1: Introduction . . . . .	4
3.1.1.1	Terms . . . . .	4
3.1.1.1.1	Self-Explanatory . . . . .	4
3.1.1.2	Right vs Rational . . . . .	4
3.1.1.2.1	Irrational Right Story: Battle of Narva . . . . .	5
3.1.1.3	Pragmatically, Normative Decision Theory > Descriptive Decision Theory . . . . .	5
3.1.1.3.1	My Problems with it . . . . .	5
3.1.1.4	Instrumental Rationality . . . . .	5
3.1.1.4.1	Is this aim always rational? . . . . .	6
3.1.1.5	Jean-Jacques Rousseau's <b>Stag Hunt</b> . . . . .	6
3.1.1.6	History of Decision Theory . . . . .	6
3.1.1.6.1	Period 1: Old Period (Ancient Greece) . . . . .	6
3.1.1.6.2	Period 2: Pioneering Period (>1650s) . . . . .	6
3.1.1.6.3	Period 3: Axiomatic Period (>1920s) . . . . .	6
3.1.2	Peterson CH2: The Decision Matrix . . . . .	6
3.1.2.1	Terms . . . . .	7
3.1.2.1.1	Scales . . . . .	7
3.1.2.2	3 Transformative Decision Rules . . . . .	7

3.1.3	Gilboa CH1: Feasibility and Desirability . . . . .	7
3.1.3.1	No direct causal link . . . . .	7
3.1.3.1.1	Zen and the Absurd (as in Camus Absurd) . . . . .	8
3.1.3.2	Uncertainty and Feasibility . . . . .	8
3.1.3.2.1	Example . . . . .	8
3.1.3.3	Link is mediated by information . . . . .	8
3.1.4	Class . . . . .	8
3.1.4.1	Decision Theory vs Social Choice Theory vs Game Theory (Zsombor's Distinctions) . . . . .	8
3.1.4.2	Distinctions . . . . .	8
3.1.4.2.1	For States . . . . .	8
3.1.4.2.2	For Acts . . . . .	8
3.2	W2: Decision rules under uncertainty. . . . .	9
3.2.1	Peterson CH3: Decisions Under Ignorance . . . . .	9
3.2.1.1	Effective Decision Rules . . . . .	9
3.2.1.1.1	Maximin and Maximax (ordinal) . . . . .	9
3.2.1.1.2	Further Constraints: Leximin and Leximin Modifications (ordinal)	9
3.2.1.1.3	Combination: Optimism-Pessimism Rule (a.k.a alpha-index rule) (cardinal) . . . . .	9
3.2.1.1.4	Problem: Relevance of non-extreme values [all] . . . . .	10
3.2.1.1.5	Problem: Unintuitive equivalence [Minimax, Maximax] . . . . .	10
3.2.1.1.6	Minimax Regret . . . . .	10
3.2.1.2	Transformative Decision Rules . . . . .	11
3.2.1.2.1	Principle of Insufficient Reason . . . . .	11
3.2.1.2.2	Randomized Acts . . . . .	12
3.2.1.3	Axiomatic Analysis of the Decision Rules . . . . .	12
3.2.1.3.1	Descriptions of Axioms . . . . .	12
3.2.1.3.2	Axiomatic Analysis . . . . .	13
3.2.2	Class . . . . .	14
3.2.2.1	Terms . . . . .	14
3.2.2.1.1	Distinctions . . . . .	14
3.2.2.2	Why necessary untruths are sometimes included as states . . . . .	14
3.2.2.3	Misc . . . . .	14
3.2.2.3.1	Some possible units for decision matrix . . . . .	14
3.3	W3: Probability theory and Bayes' Rule. Expected value maximization. . . . .	15
3.4	W4: Preferences and utility. Expected utility maximization. . . . .	15
<b>4</b>	<b>Behavioural Decision Theory</b>	<b>15</b>
4.1	W5: Discounting the future. The value of information. Biased choice. . . . .	15
4.2	W6: Paradoxes of choice. Prospect theory. Biases in probabilistic reasoning. . . . .	15
4.3	W8: Two system-theories. Behavioral design: nudging and fast-and-frugal heuristics. . . . .	15
4.4	W9: Social preferences and choice. . . . .	15
<b>5</b>	<b>Philosophy and Decision Theory</b>	<b>15</b>
5.1	W10: Interpretations of probability. The problem of induction. . . . .	15
5.2	W11: Causal, evidential and functional decision theory. . . . .	15
5.3	W12: Applications: I. Pascal's Wager. II. Discounting. III. The value of life. . . . .	15
5.4	W13: Superintelligence and the AI alignment problem. . . . .	15

# 1 Course Information

## 1.1 Instructor

Zsombor Zoltan Meder

## 1.2 Office Hours

- Nothing fixed, just email him / go to his office
  - Note: He doesn't like email / respond fast to email

### 1.2.1 Zsombor's Classes, can find him after

Tue: 9-10 / 2.404???? Thu: 15-17 / 1.508 Fri: 11.30-13.30 / 1.508

## 1.3 Textbook

Loads, but main is Peterson, M.: An introduction to decision theory. Cambridge University Press, 2017.

## 1.4 Grading

Percentage	Component Name	Notes
35%	Midterm	W8; tested on W1-W6
35%	Case Study	Proposal 25/3, 2350; Due 28/4, 2359
20%	Reading Summaries	Weekly (12 total); lowest 3 dropped
10%	Short surprise quizzes	6 total; lowest 1 dropped

### 1.4.1 Reading Summaries

- Due after class every Monday
- Roughly 1 page (a little more or less is k)
- Grade  $\in \{ 0, 1, 2 \}$

### 1.4.2 Case Study

3k-4k words, max 8k

Can be based on fiction

Feedback till 31/3

## 2 TODO Current Questions to Ask

### 2.1 TODO For Zsombor

#### 2.1.1 TODO Can you trust someone to act against your interest?

Probs yes, but I wanna be sure

### 2.2 TODO For Not Zsombor

#### 2.2.1 TODO What was Zsombor's Paradigm/Theory distinction thing?

#### 2.2.2 TODO What's completeness again?

## 3 Normative Decision Theory

### 3.1 W1: Introduction. Decision trees.

#### 3.1.1 Peterson CH1: Introduction

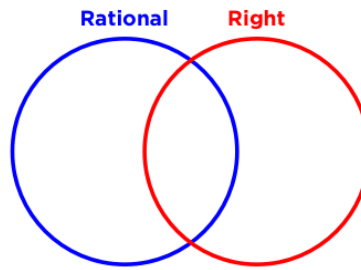
Term	Definition/Explanation	Notes
Risk	$\mathbb{P}(\text{outcomes})$ known	
Ignorance	$\mathbb{P}(\text{outcomes})$ unknown	
Uncertainty	$(\text{Ignorance}) \vee (\text{Risk} \cup \text{Ignorance})$	Context-dependent :(
Rational	Most reasonable outcome (ex ante)	I assume reasonability is based on available info
Right	Outcome is (at least weakly) pareto dominant. (ex post)	
Social Choice Theory	More than one decision maker	Some collective entities can be reduced to single decision makers ( $\therefore$ not SCT)

##### 3.1.1.1 Terms

##### 3.1.1.1.1 Self-Explanatory

- Decision-maker
- Set of alternatives
- True state of the world
- Outcome
- Principle of maximizing expected value

**3.1.1.2 Right vs Rational** Due to imperfect information, rationality does not necessarily correspond to rightness.



### 3.1.1.2.1 Irrational Right Story: Battle of Narva

#### 1. Context

- 20 November 1700
- Sweden vs Russia, on border of Estonia
  - Note: Estonia did not exist back then
- King Carl of Sweden: 8,000 troops
- Tsar Peter the Great (Russia): ~80,000 troops
- No strategic reason for Sweden to attack (little to gain)

#### 2. Story

- Sweden attacked (irrationally)
- Unexpected blizzard blinded Russian army
- Swedes won
- Battle ended <2h
- Swedes lost 667 men
- Russians lost ~15,000 men

**3.1.1.3 Pragmatically, Normative Decision Theory > Descriptive Decision Theory** Author claims that people behave rationally most of the time as they have good lives. Possibly flimsy argument (I don't like it).

#### 3.1.1.3.1 My Problems with it

- **Rational  $\neq$  right**, as seen above. E.g. maybe good lives are due to instinct rather than rationally correct behaviour
- People aren't living close to their best possible lives. Most of their lives suck. IMO people operate on habit more than reason

**3.1.1.4 Instrumental Rationality** Presupposes an aim (which is external to decision theory)

#### 3.1.1.4.1 Is this aim always rational?

- Widely thought that single aims cannot be evaluated in terms of rationality (though sets of aims can be irrational, e.g. inconsistent)
- John Rawls argues some aims are irrational (e.g. counting blades of grass on a courthouse lawn is too unimportant to be rational)
  - IMO, problematic argument. Importance varies according to values, values vary between people and even within the same person they are temporally inconsistent

#### 3.1.1.5 Jean-Jacques Rousseau's Stag Hunt

	stag	rabbit
stag	5 , 5	0 , 3
rabbit	3 , 0	3 , 3

- Tension between risk minimization and outcome maximization
- Rational choice is solely and directly dependent on **trust**

#### 3.1.1.6 History of Decision Theory

##### 3.1.1.6.1 Period 1: Old Period (Ancient Greece)

- Normative decision examples instead of rules
- Followed by 1500 years of decision theory stagnation

##### 3.1.1.6.2 Period 2: Pioneering Period (>1650s)

- Probability theory developed (Pascal and Fermat through letter correspondence)
- Some resistance by Catholic Church in normative moral theory (of course)
- 1738, **moral value** (now known as utility) was coined

##### 3.1.1.6.3 Period 3: Axiomatic Period (>1920s)

- Attempt to make axioms from principles of rational decision making
- 1950s was a golden age for decision theory
  - Still highly relevant to today

#### 3.1.2 Peterson CH2: The Decision Matrix

Notation	Term	Definition/Explanation
(square in decision tree)	Choice Node	-
(circle in decision tree)	Chance Node	-
$\pi$	<i>Formal decision problem</i> , $\pi = \langle A, S, O \rangle$	$\langle A, S, O \rangle =$ Acts, States, Outcomes
$t(\pi) \succeq \pi$	-	$t(\pi)$ is <i>at least</i> as reasonable as $\pi$
$t(\pi) \sim \pi$	-	$t(\pi)$ is <i>exactly</i> as reasonable as $\pi$
$a \circ b$	-	$(a \circ b)(\pi) = b(a(\pi))$
-	Transformative decision rule	Decision rule that modifies formalization of a decision problem
-	Effective decision rule	Filter that singles out some acts to produce a set of recommended acts
-	Rival formalizations	$\geq 2$ formalizations of same problem that are both 1. equally reasonable and 2. strictly better than other formalizations

### 3.1.2.1 Terms

#### 3.1.2.1.1 Scales

Scale	Strictly increasing	Difference information	Ratio information	Allowed Transforms
Ordinal	Yes	No	No	+ve Monotone
Cardinal: Interval	Yes	Yes	No	+ve Linear
Cardinal: Ratio	Yes	Yes	Yes	+ve Scalar

By information we mean  $h(f(a), f(b)) = h(f(c), f(d))$ , where  $h$  is whatever information (e.g. difference)

1. **TODO** A scale is not a function, it s a collection of functions. [4/2/19] Zsombor sending definition soon... Need to read and understand

#### 3.1.2.2 3 Transformative Decision Rules

1. Order-Independence (OI) If OI-condition holds for all  $\pi \in \Pi$ :
  - $(u \circ t)(\pi) = (t \circ u)(\pi)$
2. The Principle of Insufficient Reason (ir) If state probabilities are unknwon,  $\pi$  may be transformed into  $\pi'$  in which *equal probabilities are assigned to all states*
3. Merger of states (ms) If  $\geq 2$  states yield identical outcomes under all acts, they can be collapsed into one (with probabilities summed up, if known)

### 3.1.3 Gilboa CH1: Feasibility and Desirability

Can (feasibility) vs want (desirability)

#### 3.1.3.1 No direct causal link Usually, straightforwardly independent

**3.1.3.1.1 Zen and the Absurd (as in Camus Absurd)** Under some cases (e.g. mathematicians who like challenges), feasibility itself appears to have direct negative causal link with desirability.

Author argues that the act of challenge is sought rather than the state of infeasibility, and thus the causal link still does not exist.

Violates Occam's razor, but makes sense

**3.1.3.2 Uncertainty and Feasibility** Feasibility of states need not propagate to feasibility of states

**3.1.3.2.1 Example** You can certainly perform the act try to solve a math problem for 2h, without being certain about whether the state of having it solved is indeed achievable)

**3.1.3.3 Link is mediated by information**

1. Example

- Desire: Buy strawberries
- Situation: End of day; 1 box of strawberries left over (∴ feasible)
- Thought process: Why did no one buy that last box? Does it suck?
- Decision: Don't buy, even though feasible, as feasibility

**3.1.4 Class**

Less Wrong

**3.1.4.1 Decision Theory vs Social Choice Theory vs Game Theory (Zsombor's Distinctions)**  
**Decision Theory** is about single decision makers making decisions

**SCT** is about decision makers acting as a collective agent.

**Game Theory** is about players acting strategically (oppositionally almost)

**3.1.4.2 Distinctions**

**3.1.4.2.1 For States**

- Certainty / Uncertainty
- Possible / Impossible
- Desirable / Undesirable

**3.1.4.2.2 For Acts**

- Feasible / Infeasible



## 3.2 W2: Decision rules under uncertainty.

### 3.2.1 Peterson CH3: Decisions Under Ignorance

Term	Definition/Explanation	Notes
Dominance	(just pareto stuff)	Strong/weak dominance distinction

#### 3.2.1.1 Effective Decision Rules

##### 3.2.1.1.1 Maximin and Maximax (ordinal) They're the same, just at different extreme ends.

###### 1. Maximin

**Informal** Principle of choosing the act with the largest minimal outcome obtainable

**Formal**  $a_i \succeq a_j \iff \min(a_i) \geq \min(a_j)$

###### 2. Maximax

**Informal:** Principle of choosing the act with the largest maximal outcome obtainable

**Formal:**  $a_i \succeq a_j \iff \max(a_i) \geq \max(a_j)$

##### 3.2.1.1.2 Further Constraints: Leximin and Leximin Modifications (ordinal) Below is for leximin. Leximax is just the opposite

###### 1. Explanation

**Description:** Essentially a way to filter the dominance space from **maximin** (also an **effective decision rule**)

**Procedure:** Iteratively compare the next minimal outcomes under the states until you find a difference. Remove the act(s) with a lower outcome

**Equivalence:** The only remaining equivalent acts are acts which are equivalent in every state

###### 2. Formal Definition

$a_i \succ a_j \iff$  there exists some positive integer  $n$  such that  $\min^n(a_i) > \min^n(a_j)$  and  $\min^m(a_i) = \min^m(a_j)$  for all  $m < n$

##### 3.2.1.1.3 Combination: Optimism-Pessimism Rule (a.k.a alpha-index rule) (cardinal) **Informal:** Weighted combination of maximin and maximax. Weight parameter $\alpha$ reflects optimism

**Formal:**  $a_i \succeq a_j \iff \alpha \cdot \max(a_i) + (1 - \alpha) \cdot \min(a_i) \geq \alpha \cdot \max(a_j) + (1 - \alpha) \cdot \min(a_j)$

###### 1. Random Thought You can probably combine with leximin/leximax to some degree too, to get an even better framework

**3.2.1.1.4 Problem: Relevance of non-extreme values [all]** Particularly easy to see when the mins are close. E.g. below.

Applying the maximin principle selects option  $a_2$ , but intuitively  $a_1$  seems far better.

Note: Example formulated for **maximin/leximin**. Just flip the logic for **maximax**

	$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$	$s_7$
$a_1$	1	0.99	99999	99999	99999	99999	99999
$a_2$	1	1	1	1	1	1	1

Pedantic Note: you don't know the probability distribution of the state space. Maybe  $a_2$  is better after all...

**3.2.1.1.5 Problem: Unintuitive equivalence [Minimax, Maximax]** Note: Example formulated for **maximin**. Just flip the logic for **maximax**

	$s_1$	$s_2$
$a_1$	1	99999
$a_2$	1	1

Under vanilla **maximin**, both acts are equally reasonable. Obviously, this is weird

### 3.2.1.1.6 Minimax Regret

1. Explanation

- Essentially an attempt to formalize the concept of **regret**

2. Procedure (won't formally describe)

<b>Before</b>	$a_1$	12	8	20	<b>20</b>
	$a_2$	10	<b>15</b>	16	8
	$a_3$	<b>30</b>	6	25	14
	$a_4$	20	4	<b>30</b>	10
<div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> <div style="color: blue; font-weight: bold; font-size: 1.2em;">-30</div> <div style="color: blue; font-size: 2em;">↓</div> </div> <div style="text-align: center;"> <div style="color: blue; font-weight: bold; font-size: 1.2em;">-15</div> <div style="color: blue; font-size: 2em;">↓</div> </div> <div style="text-align: center;"> <div style="color: blue; font-weight: bold; font-size: 1.2em;">-30</div> <div style="color: blue; font-size: 2em;">↓</div> </div> <div style="text-align: center;"> <div style="color: blue; font-weight: bold; font-size: 1.2em;">-20</div> <div style="color: blue; font-size: 2em;">↓</div> </div> </div>					
<b>After</b>	$a_1$	<b>-18</b>	-7	-10	0
	$a_2$	<b>-20</b>	0	-14	-12
	$a_3$	0	<b>-9</b>	-5	-6
	$a_4$	-10	<b>-11</b>	0	-10

3. Note It's not globally accepted that this concept is relevant to rational decision making. But a substantial number of theorists think it is.

4. Problem: Argument from irrelevant alternatives

- Ranking can be altered by adding a non-optimal alternative
- Breaks intuition about how a "normatively plausible decision rule must not be sensitive to the addition of irrelevant alternatives"

(a) Example: Addition of  $a_5$

Table 3.14

$a_1$	12	8	20	20
$a_2$	10	<b>15</b>	16	8
$a_3$	<b>30</b>	6	25	14
$a_4$	20	4	<b>30</b>	10
$a_5$	-10	10	10	<b>39</b>

Table 3.15

$a_1$	-18	-7	-10	<b>-19</b>
$a_2$	-20	0	-14	<b>-31</b>
$a_2$	0	-9	-5	<b>-25</b>
$a_3$	-10	-11	0	<b>-29</b>
$a_5$	<b>-40</b>	-5	-20	0

(b) Counter

- Prima facie intuition is wrong. It's "rational to compare alternatives with the entire set of alternatives".

### 3.2.1.2 Transformative Decision Rules

#### 3.2.1.2.1 Principle of Insufficient Reason Pro: Decision under ignorance $\rightarrow$ Decision under risk

1. Problem: Which states should be considered? Modeling problem

- More states means lower probability for each state (direct influence on choice strategy)
- Choosing relevant states is often not easy
- Traditional argument for ir is from symmetric states (e.g. dice sides). Many problems have no such symmetry

2. Problem: Uniform probability assumption seems arbitrary Under ignorance, any probability distribution *seems to be* equally justifiable as any other. Assumption of equality seems arbitrary

(a) Counter: Symmetry

- Assume every probability distribution is equally justifiable
- Use lens of toy 2-state case  $S = \{s_1, s_2\}$
- Every probability distribution has a symmetric partner (e.g.  $\{p_{s_1} = 0.6, p_{s_2} = 0.4\}$  has  $\{p_{s_1} = 0.4, p_{s_2} = 0.6\}$ 
  - Exception: Uniform distribution. Suggests uniform distribution is a collapsed state of (in this case) 2 identical probability distributions. Making it multiplicatively more reasonable as any other case (in this case, 2x more reasonable)

i. Problem

- Beautiful argument, but the assumption of the uniform distribution being an additively collapsed one is a bit dubious imo
3. Problem: Practically, it's often not complete ignorance You generally know some things or at least have a sense of ordinal ranking for the probabilities of some of the states. Why not use it?

### 3.2.1.2.2 Randomized Acts

#### 1. Procedure

- Create a new act with expected values as outcomes
- If your decision making strategy selects random act, then randomly choose one of those initial acts
- Note: Choosing the random act is not a choice on its own, but a procedure to arrive at an actual choice

#### (a) Example

Introduce random act  $a_3$

Table 3.18

$a_1$	1	0
$a_2$	0	1

Table 3.19

$a_1$	1	0
$a_2$	0	1
$a_3$	1/2	1/2

#### 2. Potential Problem

I'm assuming the random choice doesn't have to be uniformly distributed. But this opens up a whole can of worms by allowing you to tweak the probability distribution of the random function to bias it towards whatever choice you irrationally want.

**3.2.1.3 Axiomatic Analysis of the Decision Rules** Taken directly and shamelessly from the textbook

#### 3.2.1.3.1 Descriptions of Axioms

1. **Ordering:**  $\succeq$  is transitive and complete. (See Chapter 5.)
2. **Symmetry:** The ordering imposed by  $\succeq$  is independent of the labeling of acts and states, so any two rows or columns in the decision matrix could be swapped.
3. **Strict Dominance:** If the outcome of one act is strictly better than the outcome of another under every state, then the former act is ranked above the latter.
4. **Continuity:** If one act weakly dominates another in a sequence of decision problems under ignorance, then this holds true also in the limit decision problem under ignorance.
5. **Interval scale:** The ordering imposed by  $\succeq$  remains unaffected by a positive linear transformation of the values assigned to outcomes.
6. **Irrelevant alternatives:** The ordering between old alternatives does not change if new alternatives are added to the decision problem.
7. **Column linearity:** The ordering imposed by  $\succeq$  does not change if a constant is added to a column.
8. **Column duplication:** The ordering imposed by  $\succeq$  does not change if an identical state (column) is added.
9. **Randomisation:** If two acts are equally valuable, then every randomisation between the two acts is also equally valuable.
10. **Special row adjunction:** Adding a weakly dominated act does not change the ordering of old acts.

### 3.2.1.3.2 Axiomatic Analysis

	Maximin	Optimism– pessimism	Minimax regret	Insufficient reason
1. Ordering	⊗	⊗	⊗	⊗
2. Symmetry	⊗	⊗	⊗	⊗
3. Strict dominance	⊗	⊗	⊗	⊗
4. Continuity	⊗	⊗	⊗	⊗
5. Interval scale	×	⊗	×	×
6. Irrelevant alternatives	⊗	⊗	–	⊗
7. Column linearity	–	–	⊗	⊗
8. Column duplication	⊗	⊗	⊗	–
9. Randomisation	⊗	–	⊗	×
10. Special row adjunction	×	×	⊗	×

Symbol	Meaning
–	Incompatible with decision rule
×	Compatible with decision rule
⊗	Necessary and sufficient for decision rule

### 3.2.2 Class

Term	Definition/Explanation	Notes
State	Event outside of DM's Control	Causally independent from acts
Act	A mapping of a state to an outcome	Under this definition, you may need a sequence of choices to constitute the act

#### 3.2.2.1 Terms

##### 3.2.2.1.1 Distinctions

- Preference (states)
- Dominance (acts)

Both use the  $\succeq$  signs

**3.2.2.2 Why necessary untruths are sometimes included as states** E.g. having two states: one for  $2^3 = 8$  and one for  $2^3 = 9$ .

It's necessary sometimes because the actor may not have information about whether it's true or false

##### 3.2.2.3 Misc

- Four Color Theorem: you can color any plane (e.g. a map) into any number of contiguous regions that

##### 3.2.2.3.1 Some possible units for decision matrix

- \$
- Utility
- Value

3.3 W3: Probability theory and Bayes' Rule. Expected value maximization.

3.4 W4: Preferences and utility. Expected utility maximization.

## 4 Behavioural Decision Theory

4.1 W5: Discounting the future. The value of information. Biased choice.

4.2 W6: Paradoxes of choice. Prospect theory. Biases in probabilistic reasoning.

4.3 W8: Two system-theories. Behavioral design: nudging and fast-and-frugal heuristics.

4.4 W9: Social preferences and choice.

## 5 Philosophy and Decision Theory

5.1 W10: Interpretations of probability. The problem of induction.

5.2 W11: Causal, evidential and functional decision theory.

5.3 W12: Applications: I. Pascal's Wager. II. Discounting. III. The value of life.

5.4 W13: Superintelligence and the AI alignment problem.