

Self-interested agents

- What does it mean to say that an agent is **self-interested**?
 - not that they want to harm other agents
 - not that they only care about things that benefit them
 - that the agent has its own description of states of the world that it likes, and that its actions are motivated by this description

Self-interested agents

- What does it mean to say that an agent is **self-interested**?
 - not that they want to harm other agents
 - not that they only care about things that benefit them
 - that the agent has its own description of states of the world that it likes, and that its actions are motivated by this description
- Utility theory:
 - **quantifies** degree of preference across alternatives
 - understand the impact of **uncertainty** on these preferences
 - **utility function**: a mapping from states of the world to real numbers, indicating the agent's level of happiness with that state of the world
 - **Decision-theoretic rationality**: take actions to maximize expected utility.

Example: friends and enemies

- Alice has three options: club (c), movie (m), watching a video at home (h)
- On her own, her utility for these three outcomes is 100 for c , 50 for m and 50 for h
- However, Alice also cares about Bob (who she hates) and Carol (who she likes)
 - Bob is at the club 60% of the time, and at the movies otherwise
 - Carol is at the movies 75% of the time, and at the club otherwise
- If Alice runs into Bob at the movies, she suffers disutility of 40; if she sees him at the club she suffers disutility of 90.
- If Alice sees Carol, she enjoys whatever activity she's doing 1.5 times as much as she would have enjoyed it otherwise (taking into account the possible disutility caused by Bob)

Example: friends and enemies

- Alice has three options: club (c), movie (m), watching a video at home (h)
- On her own, her utility for these three outcomes is 100 for c , 50 for m and 50 for h
- However, Alice also cares about Bob (who she hates) and Carol (who she likes)
 - Bob is at the club 60% of the time, and at the movies otherwise
 - Carol is at the movies 75% of the time, and at the club otherwise
- If Alice runs into Bob at the movies, she suffers disutility of 40; if she sees him at the club she suffers disutility of 90.
- If Alice sees Carol, she enjoys whatever activity she's doing 1.5 times as much as she would have enjoyed it otherwise (taking into account the possible disutility caused by Bob)
- What should Alice do (show of hands)?

What activity should Alice choose?

	$B = c$	$B = m$
$C = c$	15	150
$C = m$	10	100
	$A = c$	

	$B = c$	$B = m$
$C = c$	50	10
$C = m$	75	15
	$A = m$	

What activity should Alice choose?

	$B = c$	$B = m$
$C = c$	15	150
$C = m$	10	100
	$A = c$	

	$B = c$	$B = m$
$C = c$	50	10
$C = m$	75	15
	$A = m$	

- Alice's expected utility for c :

$$0.25(0.6 \cdot 15 + 0.4 \cdot 150) + 0.75(0.6 \cdot 10 + 0.4 \cdot 100) = 51.75.$$

- Alice's expected utility for m :

$$0.25(0.6 \cdot 50 + 0.4 \cdot 10) + 0.75(0.6(75) + 0.4(15)) = 46.75.$$

- Alice's expected utility for h : 50.

Alice prefers to go to the club (though Bob is often there and Carol rarely is), and prefers staying home to going to the movies (though Bob is usually not at the movies and Carol almost always is).

Why utility?

- Why would anyone argue with the idea that an agent's preferences could be described using a utility function as we just did?

Why utility?

- Why would anyone argue with the idea that an agent's preferences could be described using a utility function as we just did?
 - why should a single-dimensional function be enough to explain preferences over an arbitrarily complicated set of alternatives?
 - Why should an agent's response to uncertainty be captured purely by the *expected value* of his utility function?
- It turns out that the claim that an agent has a utility function is substantive.

Preferences Over Outcomes

If o_1 and o_2 are outcomes

- $o_1 \succeq o_2$ means o_1 is at least as desirable as o_2 .
 - read this as “the agent **weakly prefers** o_1 to o_2 ”
- $o_1 \sim o_2$ means $o_1 \succeq o_2$ and $o_2 \succeq o_1$.
 - read this as “the agent is **indifferent** between o_1 and o_2 .”
- $o_1 \succ o_2$ means $o_1 \succeq o_2$ and $o_2 \not\succeq o_1$
 - read this as “the agent **strictly prefers** o_1 to o_2 ”

Lotteries

- An agent may not know the outcomes of his actions, but may instead only have a probability distribution over the outcomes.

Definition (lottery)

A **lottery** is a probability distribution over outcomes. It is written

$$[p_1 : o_1, p_2 : o_2, \dots, p_k : o_k]$$

where the o_i are outcomes and $p_i > 0$ such that

$$\sum_i p_i = 1$$

- The lottery specifies that outcome o_i occurs with probability p_i .
- We will consider lotteries to be outcomes.

Preference Axioms: Completeness

Definition (Completeness)

A preference relationship must be defined between every pair of outcomes:

$$\forall o_1 \forall o_2 \quad o_1 \succeq o_2 \text{ or } o_2 \succeq o_1$$

Preference Axioms: Transitivity

Definition (Transitivity)

Preferences must be transitive:

$$\text{if } o_1 \succeq o_2 \text{ and } o_2 \succeq o_3 \text{ then } o_1 \succeq o_3$$

- This makes good sense: otherwise $o_1 \succeq o_2$ and $o_2 \succeq o_3$ and $o_3 \succ o_1$.
- An agent should be prepared to pay some amount to swap between an outcome they prefer less and an outcome they prefer more
- Intransitive preferences mean we can construct a “money pump”!

Preference Axioms

Definition (Monotonicity)

An agent prefers a larger chance of getting a better outcome to a smaller chance:

- If $o_1 \succ o_2$ and $p > q$ then

$$[p : o_1, 1 - p : o_2] \succ [q : o_1, 1 - q : o_2]$$

Preference Axioms

Let $P_\ell(o_i)$ denote the probability that outcome o_i is selected by lottery ℓ . For example, if $\ell = [0.3 : o_1; 0.7 : [0.8 : o_2; 0.2 : o_1]]$ then $P_\ell(o_1) = 0.44$ and $P_\ell(o_3) = 0$.

Definition (Decomposability (“no fun in gambling”))

If $\forall o_i \in O, P_{\ell_1}(o_i) = P_{\ell_2}(o_i)$ then $\ell_1 \sim \ell_2$.

Preference Axioms

Definition (Substitutability)

If $o_1 \sim o_2$ then for all sequences of one or more outcomes o_3, \dots, o_k and sets of probabilities p, p_3, \dots, p_k for which $p + \sum_{i=3}^k p_i = 1$,
 $[p : o_1, p_3 : o_3, \dots, p_k : o_k] \sim [p : o_2, p_3 : o_3, \dots, p_k : o_k]$.

Preference Axioms

Definition (Continuity)

Suppose $o_1 \succ o_2$ and $o_2 \succ o_3$, then there exists a $p \in [0, 1]$ such that $o_2 \sim [p : o_1, 1 - p : o_3]$.

Preferences and utility functions

Theorem (von Neumann and Morgenstern, 1944)

If an agent's preference relation satisfies the axioms Completeness, Transitivity, Decomposability, Substitutability, Monotonicity and Continuity then there exists a function $u : O \rightarrow [0, 1]$ with the properties that:

- ❶ $u(o_1) \geq u(o_2)$ iff the agent prefers o_1 to o_2 ; and
- ❷ when faced about uncertainty about which outcomes he will receive, the agent prefers outcomes that maximize the expected value of u .

Proof idea:

- define the utility of the best outcome $u(\bar{o}) = 1$ and of the worst $u(\underline{o}) = 0$
- now define the utility of each other outcome o as the p for which $o \sim [p : \bar{o}; (1 - p) : \underline{o}]$.