Introduction to Mechanism Design

Mehdi Dastani BBL-521 M.M.Dastani@cs.uu.nl

Thanks to Paul Harrenstein and Mathijs de Weerdt who provided me some of these slides.



What is Mechanism Design Trying to Accomplish?

- Takahashi wishes to sell a precious Ming vase to either Ann or Bob, depending on who values it most.
- Yet, Takahashi does not know the valuations of either Ann or Bob.
- Hence, the goal he sets himself is to sell the vase to Ann if she values it most, and to Bob otherwise. (Social choice function).
- Then, he tries to construct a protocol (mechanism) that guarantees the vase to be sold to the one who values it most, for all possible valuations Ann and Bob may have. (Implement the social choice function).

What is Mechanism Design Trying to Accomplish?

- Fix a set of possible outcomes O.
- Given a preference profile (≥1,...,≥n) certain outcomes O* ∈ O are more desirable than others from an outsider's point of view (e.g., assigning an object to the agent that values it most).
- However, preferences of the players are unknown to the designer, or hard to obtain.
- Issue: Design a game (mechanism) such that the outcome given a particular (fixed) solution concept of this game generates one of the desired outcomes in O* for all (relevant) preference profiles (≿1,...,≿n).

Strategic Behavior and the Importance of Truthfulness

Principles of voting make an election more of a game of skill than a real test of the wishes of the electors. My own opinion is that it is better for elections to be decided according to the wish of the majority than of those who happen to have most skill at the game.

(C.L. Dodgson)



Strategic Behavior and the Importance of Truthfulness

Consider the Borda rule.

| ≿1 | ≿2 | ≿₃ | ≿1 | ≿2 | ≿′ ₃ |
|----|----|----|----|----|-----------------|
| а | а | d | а | а | b |
| b | b | С | b | b | С |
| С | d | b | С | d | d |
| d | С | а | d | С | а |

- ▶ Borda winner given (\geq_1, \geq_2, \geq_3) is a with 6 points
- ▶ Borda winner given $(\succeq_1, \succeq_2, \succeq_3')$ is b with 7 points

Conclusion: If \geq_3 are player 3's true preferences, he had better lie about them!

- Preferences as strategies (lie or tell the truth)
- Combination of such strategies is a preference profile
- Outcomes are defined by a social choice function f
- Social choice functions reflect how to determine the desired (by the designer) outcomes relative to the agents' preferences. E.g., assigning the object to the agent that values it most
- The true preferences of the agents are assumed to be unknown to the designer
- ► Each possible preference profile ≥ yields such a game

| | ≿2 | ≿′2 |
|-----|---------------------------|------------------------------|
| ≿1 | $f(\geq_1,\geq_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succsim_1',\succsim_2')$ |

| (\succeq_1,\succeq_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|-------------------------|---------------------------|----------------------------|----------------------------|---------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| (≿'₁,≿2) | ≿2 | ≿′2 | (z'_1,z'_2) | ≿2 | ≿′2 |
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | | ≿1 | $f(\geq_1,\geq_2)$ | $f(\geq_1,\geq_2')$ |
| ≿′1 | $f(\succeq'_1,\succeq_2)$ | $f(\succeq'_1,\succeq'_2)$ | ≿′1 | $f(\succeq'_1,\succeq_2)$ | $f(\succeq_1',\succeq_2')$ |

| (\succsim_1,\succsim_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|--|---------------------------|--|----------------------------|--------------------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\gtrsim_1',\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | (\succeq_1',\succeq_2') | ≿2 | ≿′2 |
| (≿ ₁ ',≿ ₂) ≿1 | | \gtrsim_2' $f(\gtrsim_1,\gtrsim_2')$ | | \gtrsim_2 $f(\gtrsim_1,\gtrsim_2)$ | |

| $\left(\succsim_1,\succsim_2\right)$ | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|--------------------------------------|---------------------------|--|----------------------------|--------------------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | (\succeq_1',\succeq_2') | ≿2 | ≿′2 |
| (\gtrsim'_1,\gtrsim_2) \gtrsim_1 | | $ \gtrsim_2' $ $ f(\gtrsim_1,\gtrsim_2') $ | _ | \gtrsim_2 $f(\gtrsim_1,\gtrsim_2)$ | |

| (\succeq_1,\succeq_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|--|---|---|----------------------------|--|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\gtrsim_1',\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | (\succeq_1',\succeq_2') | ≿2 | ≿′2 |
| (≿' ₁ ,≿ ₂) ≿1 | $\underset{(\succeq 1, \succeq 2)}{\succsim_2}$ | $ \succsim_2' $ $ f(\succsim_1,\succsim_2') $ |] | $ \gtrsim_2 $ $ f(\gtrsim_1,\gtrsim_2) $ | |

| (\succeq_1,\succeq_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|-------------------------|---------------------------|----------------------------|----------------------------|---------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| (≿'₁,≿2) | ≿2 | ≿′2 | (\succeq_1',\succeq_2') | ≿2 | ≿′2 |
| ≿1 | $f(\geq_1,\geq_2)$ | | ≿1 | $f(\geq_1,\geq_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′₁ | $f(\succeq'_1,\succeq_2)$ | f(~! ~!) | <u>≻′</u> | $f(\gtrsim'_1,\gtrsim_2)$ | f(>', >',) |

He sent for a sword, and when it was brought, he said, "Cut the living child in two and give each woman half of it". The real mother, her heart full of love for her son, said to the king, "Please, Your Majesty, don't kill the child! Give it to her!" But the other woman said, "Don't give it to either of us; go on and cut it in two". Then Solomon said, "Don't kill the child! Give it to the first woman, she is its real mother."

(1 Kings 3: 16-28)



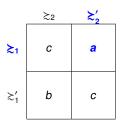
- Three outcomes:
 - First woman gets the baby.
 - Second woman gets the baby.
 - c: The baby is bisected.
- For each woman, two possible types (good mother, bad mother):

$$\geq_1$$
: $a > b > c$

$$\gtrsim_1$$
: $a > b > c$ \gtrsim_2 : $b > a > c$ \gtrsim_1 : $a > c > b$ \gtrsim_2 : $b > c > a$

$$\gtrsim_1'$$
: $a > c > b$

$$\gtrsim_2'$$
: $b > c > a$



| | ≿2 | ≿′2 |
|-----|----|-----|
| ≿1 | С | а |
| ≿′1 | b | С |

| | ≿2 | ≿₂′ |
|-----|----|-----|
| ≿1 | С | а |
| ≿′1 | b | С |

$$\begin{array}{ccc} & \gtrsim_1 & \gtrsim_2' \\ \hline a & b \\ b & c \\ c & a \end{array}$$

Observation: Be the true types given by (\geq_1, \geq'_2) , revealing her true preferences is **not** a dominant strategy for the second woman (the bad mother)!

| | ≿2 | ≿₂′ |
|-----|----|-----|
| ≿1 | С | а |
| ≿′1 | b | С |

Observation: Be the true types given by (\geq_1, \geq'_2) , revealing her true preferences is **not** a dominant strategy for the second woman (the bad mother)!

| | ≿2 | ≿₂′ |
|-----|----|-----|
| ≿1 | С | а |
| ≿′1 | b | С |

Observation: Be the true types given by (\geq_1, \geq'_2) , revealing her true preferences is **not** a dominant strategy for the second woman (the bad mother)!

| (\succsim_1,\succsim_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|---------------------------|---------------------------|--|----------------------------|---|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | (\succeq_1',\succeq_2') | ≿2 | ≿′2 |
| | | $\underset{(\succeq 1, \succeq'_2)}{\succeq'_2}$ | | $ \begin{array}{c c} $ | _ |

| (\gtrsim_1,\gtrsim_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|--|--|--|-----------------------------|--------------------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | $(\succsim_1',\succsim_2')$ | ≿2 | ≿′2 |
| (≿ ₁ ',≿ ₂) ≿1 | $ \gtrsim_2 $ $ f(\gtrsim_1,\gtrsim_2) $ | $\underset{\sim}{\succsim_2'} f(\succsim_1,\succsim_2')$ | _ | \gtrsim_2 $f(\gtrsim_1,\gtrsim_2)$ | |

| (\succeq_1,\succeq_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|-------------------------|---------------------------|----------------------------|----------------------------|---------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| (≿'₁,≿2) | ≿2 | ≿′2 | (z_1',z_2') | ≿2 | ≿′2 |
| ≿1 | $f(\geq_1,\geq_2)$ | | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| | | | | | |

| (\succsim_1,\succsim_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|---------------------------|---------------------------|----------------------------|----------------------------|---------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| (≿'₁,≿₂) | ≿2 | ≿′2 | (z'_1,z'_2) | ≿2 | ≿′2 |
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | $f(\gtrsim_1,\gtrsim_2')$ |
| | | | | | |

| (\succsim_1,\succsim_2) | ≿2 | ≿′2 | (\succsim_1,\succsim_2') | ≿2 | ≿′2 |
|---------------------------|--|----------------------------|----------------------------|--------------------------------------|----------------------------|
| ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | | ≿1 | $f(\gtrsim_1,\gtrsim_2)$ | |
| ≿′1 | $f(\succeq_1',\succeq_2)$ | $f(\succeq_1',\succeq_2')$ | ≿′1 | $f(\gtrsim'_1,\gtrsim_2)$ | $f(\succeq_1',\succeq_2')$ |
| | | | | | |
| (\succeq_1',\succeq_2) | ≿2 | ≿′2 | (\succeq'_1,\succeq'_2) | ≿2 | ≿′2 |
| (z'_1, z_2) z_1 | $ \gtrsim_2 $ $ f(\gtrsim_1,\gtrsim_2) $ | _ | · - | \gtrsim_2 $f(\gtrsim_1,\gtrsim_2)$ | |

Definition: A social choice function f is dominant strategy incentive compatible w.r.t. a set A of type profiles, if for all $\geq^* \in A$, all $\geq, \geq' \in A$ and all $i \in N$:

$$f(\geq_1,\ldots,\geq_i^*,\ldots,\geq_n) \geq_i^* f(\geq_1,\ldots,\geq_i',\ldots,\geq_n)$$

Intuition: A social choice function f is truthfully implementable if for no player there are situations in which telling the truth can hurt.

Remark: Truthful implementation is important, for otherwise the mechanism does not guarantee that the goals of the designer, which are defined relative to the preferences of individuals, are reached.

Gibbard-Satterthwaite Theorem

Definition: An scf f is *dictatorial* if there is a player i with $f(\geq) \geq_i o$ for all \geq and all o in O.

Intuition: For every preference profile \gtrsim , the outcome of the social choice function, i.e. $f(\gtrsim)$, is among the dictator's most preferred outcomes.





Gibbard-Satterthwaite Theorem

Definition: An scf f is *dictatorial* if there is a player i with $f(\geq) \geq_i o$ for all \geq and all o in O.

Intuition: For every preference profile \gtrsim , the outcome of the social choice function, i.e. $f(\gtrsim)$, is among the dictator's most preferred outcomes





Theorem (Gibbard, 1973 and Satterthwaite, 1975): If |O| = 3, every incentive compatible social choice function onto O is dictatorial.

Intuition: For all non-trivial (i.e., non-dictatorial) social choice functions there are circumstances (i.e., type profiles) in which it is profitable for some player to lie about his true preferences.

Implementation and Solution Concepts

Let *S* be a solution concept and $\phi: A \to 2^{\mathcal{O}}$ a social choice correspondence.

- A mechanism is a game form G = (N, A, O, g)
- ▶ For each $\geq \in A$, (G, \geq) is a game
- ▶ For each $\geq \in A$, S selects a subset of strategy profiles of the game (G, \geq)
- ▶ Via g, for each $\geq \in A$, S selects a subset of *outcomes* in O
- ▶ For each $\geq \in A$, ϕ also selects a subset of *outcomes* in O

Implementation and Solution Concepts

Let *S* be a solution concept and $\phi \colon A \to 2^O$ a social choice correspondence.

- ▶ A mechanism is a game form G = (N, A, O, g)
- ▶ For each $\geq \in A$, (G, \geq) is a game
- ▶ For each $\geq \in A$, S selects a subset of strategy profiles of the game (G, \geq)
- ▶ Via g, for each $\geq \in A$, S selects a subset of *outcomes* in O
- ▶ For each $\geq \in A$, ϕ also selects a subset of *outcomes* in O

Definition An $\operatorname{scc} \phi \colon A \to 2^O$ is *S-implementable* if there is some *game form* G = (N, A, O, g) such that for all preference profiles \gtrsim :

$$\phi(\gtrsim) = \{g(s) \in A : s \in S(G, \gtrsim)\}$$

Intuition: The game form G implements ϕ in S if for each \gtrsim the sets of outcomes selected by S and ϕ coincide.

Implementation and Solution Concepts

| (≿1,≿2) | L | R |
|---------|---|---|
| Т | а | b |
| В | С | d |

| ,≿′2) | L | R |
|-------|---|---|
| Т | а | b |
| В | С | d |

(≿1

$$\begin{array}{c|cccc}
(z_1', z_2) & L & R \\
T & a & b \\
B & c & d
\end{array}$$

$$\begin{array}{c|cccc} (z_1',z_2') & L & R \\ \hline T & a & b \\ \hline B & c & d \\ \end{array}$$

| | $\phi(\succsim)$ | $S(G(\gtrsim))$ |
|--|------------------|-----------------|
| (≿1,≿2) | а | а |
| (\succsim_1,\succsim_2') | C | С |
| (\succeq_1',\succeq_2) | d | d |
| $\left(\succsim_1',\succsim_2'\right)$ | С | С |

Implementation in Dominant Strategies

Definition: A strategy profile s^* is a *dominant strategy equilibrium* in a game G whenever for all $i \in N$ and all $s, s' \in S$:

$$g(s_1,\ldots,s_i^*,\ldots,s_n) \gtrsim_i g(s_1,\ldots,s_i',\ldots,s_n)$$

- Very robust in that it assumes very little about the agents
- Dominant strategy equilibrium does not allow for much flexibility

Dominant Strategy Implementation with two Alternatives

| | okay | veto |
|------|------|------|
| okay | а | b |
| veto | b | b |

$$\begin{array}{cccc}
 & z^a & z^b \\
\hline
 & a & b \\
 & b & a
\end{array}$$

$$f^{\text{veto}}(\gtrsim_1,\gtrsim_2) = \begin{cases} a & \text{if } \gtrsim_1=\gtrsim^a \text{ and } \gtrsim_2=\gtrsim^a \\ b & \text{otherwise} \end{cases}$$

Dominant Strategy Implementation with two Alternatives

| (\gtrsim^a,\gtrsim^a) | okay | veto |
|-------------------------|------|------|
| okay | а | b |
| veto | b | b |

| (\gtrsim^a,\gtrsim^b) | okay | veto |
|-------------------------|------|------|
| okay | а | b |
| veto | b | b |

| ≿ª | ≿' |
|----|----|
| а | b |
| b | а |
| | |

| (\gtrsim^b,\gtrsim^a) | okay | veto |
|-------------------------|------|------|
| okay | а | b |
| veto | b | b |

$$(\gtrsim^b,\gtrsim^b)$$
 okay veto

okay a b

veto b b

The Revelation Principle

Lemma (*Revelation Principle*): If there exists any mechanism that implements an scf f in dominant strategies then there exists a direct mechanism that implements f in dominant strategies, i.e., an scf f is implementable in dominant strategies if and only if f is incentive compatible.

Intuition: For implementation in dominant strategy equilibria we can restrict ourselves to direct mechanisms.