# Designing Matching Mechanism

In both a centralized and decentralized market, stated preferences of agents are needed to produce stable matchings. In fact, a matchmaker or algorithm depends on it. But what if agents don't report true preferences, but instead something else so as to get themselves a better match? Is it possible to to incentivize all agents to report preferences truthfully all of the time?

(I) Matching mechanisms in One-to-one MARKET

Define $\mathcal{E} = (T, B; (Z_i)_{i \in T \cup B})$. Assume a designer (match maker) knows agents in T and B but does not know their preferences.

Let $\mathcal{P}_i$ be the set of all complete and transitive preference relations.

Think of designer as asking each agent $i$ to report a preference relation.

Let $\hat{Z}_i$ be the reported preference which may or may not be true (i.e. the true pref $Z_i$ for agent $i$)

Denote set of matches as $\mathcal{M}$ and so $M \in \mathcal{M}$.

$M$ will depend on $T, B$.

Then $\mathcal{P}$ is set of all preferences and $\mathcal{M}$ set of all matches

## Definition MATCHING MECHANISM

A function, that, for each input, describes the matching output that the algorithm produces

$$m : \prod_{i \in T \cup B} \mathcal{P}_i \to M$$

Formally, a matching mechanism $m$ is a function $m$ whose range is the set of all possible inputs $(T, B; (\succeq_i)_{i \in T \cup B})$, or, rather set of all preferences for each agent, and whose output is a matching of $T$ and $B$.

Note: designer does not know true preference relations of the agents, he only knows what they report.

2

# Example : T-proposing Deferred Acceptance Mechanism.

Notice difference in title. We say Mechanism and not algorithm because we apply the algorithm to <u>reported</u> preferences which may or may not be true.

We will get out as output a match. The match will be stable relative to the <u>reported</u> preferences. Key Question: When is such a matching, then, stable relative to true preferences in spite of reported preferences?

<u>DEFINITION</u>: Call mechanism $m : \prod_{i \in T \cup B} P_i \to m$ a <u>stable</u> <u>mechanism</u> if it produces:

a stable matching relative to each reported preference.

<u>Example</u>: Let $T = \{t_1, t_2\}$ $B = \{b_1, b_2\}$

<u>True preferences</u>:

$t_1 : b_1 > b_2 > t_1$

$t_2 : b_2 > b_1 > t_2$

$b_1 : t_2 > t_1 > b_1$

$b_2 : t_1 > t_2 > b_2$

T-prop DAA would give:

$M_{T_0}(t_1) = b_1$

$M_{T_0}(t_2) = b_2$

Would there be an incentive to misreport preferences?

- Of course not for $t$ agents who get their best match — that is if all agents are reporting truthfully, then no $t$ agent has incentive to misrepresent preferences.

Suppose $t_1, t_2, b_2$ each tell the truth. Then $b_1$ has incentive to misrepresent.

If $b_1$ reports $\hat{\succsim}_{b_1}$ such that:

$$t_2 \; \hat{\succ}_{b_1} \; b_1 \; \hat{\succ}_{b_1} \; t_1$$

Then, T-proposing DA Algorithm will do the following:

Round 1 : "$\rightarrow$" proposes

$t_1 \rightarrow b_1$

$t_2 \rightarrow b_2$

$b_1$ unmatched

$b_2$ accepts $t_2$

Round 2

$t_1 \rightarrow b_2$

$t_2 \rightarrow b_2$

$b_2$ accepts $t_1$

$b_1$ unmatched

Round 3 :

$t_1 \rightarrow b_2$   $b_1$ matches with $t_2$

$t_2 \rightarrow b_1$   $b_2$ matches with $t_1$

When one $b$ misreports, the matching becomes B-optimal.

4

<u>Question</u>: Can we design mechanism such that no agent has incentive to misreport?

Need to be more precise...

- Want no incentive before the match at the reporting stage.

- Want ex post that no one has incentive to deviate from a match (need stability)

<u>DEFINITION</u>: A strategy for agent $i$: $S_i = \mathcal{P}_i \to \mathcal{P}_i$. For every true preference relation (input) the output is a reported preference relation.

<u>DEFINITION</u>: A strategy is truthful if $S_i(z_i) = \succsim_i$ for each $z_i \in \mathcal{P}_i$ (or truthful reporting)

<u>DEFINITION</u>: A strategy $S_i$ is dominant for mechanism $m : \prod_{i \in \mathcal{TC}_B} \mathcal{P}_i \to m$ if for each preference relation $z_i \in \mathcal{P}_i$ and for each $\hat{\succsim}_{-i} \in \mathcal{P}_j$, $j \in \mathcal{TC}_B \setminus \{i\}$ such that

$$ m\left(\underbrace{S_i(z_i)}_{\substack{\text{my report}}}, \underbrace{\hat{\succsim}_{-i}}_{\substack{\text{others} \\ \text{report}}}\right)(i) \underset{\substack{\downarrow \\ i's\ pref \\ relation}}{\succsim_i} m\left(\underbrace{\hat{\succsim}_i}_{\substack{\text{my report}}}, \hat{\succsim}_{-i}\right)(i) $$

$\underbrace{\hphantom{m(S_i(z_i), \hat{\succsim}_{-i})(i)}}_{\substack{\text{who } i \text{ is matched with} \\ \text{when } i \text{ reports } S_i(z_i) \text{ and} \\ \text{everyone else reports } \hat{\succsim}_{-i}}}$
$\qquad\qquad\qquad\qquad$
$\underbrace{\hphantom{m(\hat{\succsim}_i, \hat{\succsim}_{-i})(i)}}_{\substack{\text{who } i \text{ is matched with when } i \text{ reports } \hat{\succsim}_i \\ \text{and everyone else reports } \hat{\succsim}_{-i}}}$

<u>Definition</u>: A mechanism is strategy proof if for each agent $i \in T \cup B$ the truthful strategy is dominant.

<u>Question</u>: Does there exist a stable matching that is strategy proof?

## <u>NO</u>

<u>Theorem 1</u>  If $\min \{|T|, |B|\} \geq 2$ There is no stable strategy proof mechanism (Impossibility Thm) $_{pg\,87}$

In other words, no stable matching mechanism exists for which stating true preferences is a dominant strategy for every agent.

<u>Proof</u>:

Let $T = \{t_1, \ldots, t_{|T|}\}$  $|T| \geq 2$
$B = \{b_1, \ldots, b_{|B|}\}$  $|B| \geq 2$

<u>Group Agents</u>

Suppose:

$A(t_1) = A(t_2) = \{b_1, b_2\}$
$A(b_1) = A(b_2) = \{t_1, t_2\}$
$A(t) = \phi$  $t \in T \setminus \{t_1, t_2\}$
$A(b) = \phi$  $b \in B \setminus \{b_1, b_2\}$

Recall $A(i) = \{j \in J, j \geq_i i\}$ or acceptable set. This proof says there is no stable strategy (for all strategy proof m's). So, contradiction based proof only needs one deviation ($\exists$)

↳ proof matching mech.

# Preferences Over Agents $\succsim^*$

$$b_1 \succ^*_{t_1} b_2 \qquad t_2 \succ^*_{b_1} t_1 \qquad \text{if } M \text{ is stable for } \succsim^*$$

$$b_2 \succ^*_{t_2} b_1 \qquad t_1 \succ^*_{b_2} t_2$$

$$M(t) = t \quad \text{if } t \in T \setminus \{t_1, t_2\}$$
$$M(b) = b \quad \text{if } b \in B \setminus \{b_1, b_2\}$$

$$\text{Stable} \atop \text{for } \succ^* \left\{ \begin{array}{l} M(t_1), M(t_2) \in \{b_1, b_2\} \\ M(b_1), M(b_2) \in \{t_1, t_2\} \end{array} \right.$$

So, if there were a stable strategy proof mechanism:

① For each $\succsim \in \prod_{i \in T \cup B} P_i : m(\succsim)$ is stable relative to $\succsim$.

② For every $i \in T \cup B$ and every $\hat{\succsim}_{-i} \in \prod_{j \neq i} P_j$

$$m(\succsim_i, \hat{\succsim}_{-i})(i) \succsim_i m(\hat{\succsim}_i, \hat{\succsim}_{-i})(i)$$

for all $\hat{\succsim}_i \in P_i$.

To show this is false: show there exists $\succsim^* \in \prod P_i$ such that we cannot have a stable match with ② satisfied with $\succsim^*$.

### T-Proposing DA :

$$M^+_{TD}(t_1) = b_1$$
$$M^+_{TD}(t_2) = b_2$$
$$M^*_{TD}(b_1) = t_1$$
$$M^*_{TD}(b_2) = t_2$$

### B-proposing :

$$M^*_{Bb}(t_1) = b_2$$
$$M^*_{BD}(t_2) = b_1$$
$$M^*_{BD}(b_1) = t_2$$
$$M^+_{Db}(b_2) = t_1$$

Since $m(z^*) \in \{M^*_{TO}, M^*_{BO}\}$ is stable, we will need to go through both cases and show someone always has incentive to misrepresent.

Case A: $m(z^*) = M^*_{TO}$

Let $b_1$ report $\hat{z}_{b_1}: t_2 \hat{\succ}_{b_1} b_1$ and $b_1 \hat{\succ}_b t$ for all $t \in T \setminus \{t_2\}$.
That is only $t_2$ is acceptable to $b_1$ and no other $t$ is.
Assume all other agents report truthfully. We want to know: Does $b_1$ prefer $m(\hat{z})$ to $M^*_{TO}$

Now, we do not know what $m(\hat{z})$ is. We do know $M(\hat{z})$ is stable relative to $\hat{z}$.

- We will show that there is only one stable match for preference profile $\hat{z}$, and it corresponds to T-proposing DA applied to $\hat{z}$.

To show that T-prop $^{DAA}$ applied to $\hat{z}$ is only stable match, it suffices to show that when we apply T-prop DA to $\hat{z}$ and B-prop DA to $\hat{z}$ we get the same match.

## T-prop DA on $\hat{z}$

### Round 1

$t_1 \rightarrow b_1$      $b_1$ rejects $t_1$

$t_2 \rightarrow b_2$      $b_2$ keeps $t_2$

### Round 2

$t_1 \rightarrow b_2$      $b_1$ unmatched

$t_2 \rightarrow b_2$      $b_2$ keeps $t_1$

### Round 3

$t_1 \rightarrow b_2$      $b_2$ accepts $t_1$

$t_2 \rightarrow b_1$      $b_1$ accepts $t_2$

## B-prop DA on $\hat{z}$

$b_1 \rightarrow t_2$      $t_1$ accepts $b_2$      $\Rightarrow$ Gives exact same match

$b_2 \rightarrow t_1$      $t_2$ accepts $b_1$      as above

### Conclude:

$$m(\hat{z}) = \hat{M}_{TD}$$

$$m(\hat{z}_{b_1}, z^*_{-b_1})(b_1) = t_2 \succ_{b_1} m(z^*_{b_1}, z^*_{-b_1}) = t_1$$

Then our assumption of truth telling being dominant is contradicted. That is, we assumed mechanism gave us something stable for all preferences. Then, we assumed that a truth telling preference generated matching would be preferred by all agents. Here we have our "counterexample"

But we are not done.

### Case B: $m(z^*) = M^*_{BO}$

- let $t_1$ report $\hat{z}_{t_1}$ such that $b_1 \hat{\succ}_{t_1} t_1$ and $t_1 \succ_{t_1} b$ for all $b \neq 1$

- Let all other agents report truthfully.

$$\hat{z}_i = \succ^*_i$$

Need to know what is $m(\hat{\succsim})$. We know that matchings generated from its preferences will be stable. Will show it is T-prop DAA applied to $\hat{\succsim}$, and that is the same as B-prop DA algorithm applied to $\hat{\succsim}$.

$\underline{T\text{-prop PA Algo applied to } \hat{\succsim}:}$

$t_1 \to b_1$    $b_1$ accepts $t_1$

$t_2 \to b_2$    $b_2$ accepts $t_2$    done

$\underline{B\text{-prop UA Algo}}$

$R_1 \begin{bmatrix} b_1 \to t_2 & b_2 \text{ accepts } b_1 \\ b_2 \to t_1 & t_1 \text{ rejects } b_2 \end{bmatrix} \Bigg|\ R_2 \begin{bmatrix} b_1 \to t_2 & t_2 \text{ keeps } b_2 \\ b_2 \to t_2 & b_1 \text{ no offer} \end{bmatrix}$

$R_3 \begin{bmatrix} b_1 \to t_1 & t_1 \text{ accepts } b_1 \\ b_2 \to t_2 & t_2 \text{ retains } b_2 \end{bmatrix}$

Then, both are the same! Hence: $\mu_{TB} = m(\hat{\succsim})$

$$m(\hat{\succsim}_{t_1}, \succsim^+_{-t_1}) = b_1 \succ_{t_1} m(\succsim^+_{t_1}, \succsim^*_{-t_1}) = b_2$$

Again, we have a contradiction, for the same reason as case A.

Again, to prove this (that there is no stable strategy proof match) we assumed (i) all preferences in $\succsim_i$ give stable match in $M$. And, it is better for every agent to tell truth. Then, we presented case where no matter what, there was incentive to lie.

## Theorem 2

Fix an environment $\mathcal{E} = (T; B, (z_i)_{i \in T \cup B})$ with strict preferences that has at least two stable matches. Then, for any stable matching mechanism, there exists some agent $j \in T \cup B$ and some report $\hat{z}_j$:

$$m(\hat{z}_j, z_{-j}^*)(j) \succ_j^* m(z_j^*, z_{-j}^*)(j)$$

Moreover, $\hat{z}_j$ can be chosen such that

$$m(\hat{z}_j, z_{-j}^*)(j) \text{ is } \quad j\text{'s most preferred}$$

achievable outcome.

$$\{k : \mu(i) = k \text{ for stable } \mu\}$$

In English (RS 88): When any stable mechanism is applied to a marriage market in which preferences are strict and there is more than one stable matching, then at least one agent can profitably misrepresent his preferences, assuming others tell the truth. This agent can misrepresent in such a way as to be matched to his or her most preferred achievable mate under the true preferences at every stable matching under the misstated preferences.

<u>Proof</u>: Fix environment $\varepsilon = (T, B; (z))$ and assume preferences are strict and that there are at least two stable matches

- Then, b/c preferences are strict: $M^*_{TD} \neq M^*_{OD}$, as applied to $z^*$

- $M(z^*)$ $\quad \to$ don't know what it is but it is stable

  there exists some $I \in \{T \cup B\}$ such that
  $$M^*_{TD} = M(z^*),$$
  WLOG, take $I = T$, choose $t \in T$ such that:
  $$M^*_{TD}(t) \neq M(z^*)(t).$$

  We know such a $t$ exists. If there were no such $t$, we could conclude $M^*_{TD} = M(z^*)$

- Because $M(z^*)$ is stable and $M^*_{TD}$ is $t$-optimal (since prefs are strict):
  $$M^*_{TD}(t) \succsim_t M(z^*)(t)$$

- For this $t \in T$: $M^*_{TD}(t) \succsim_t M(z^*)(t)$

Want to show: If everyone reports truthfully, there is some misreport $t$ can make that would get him $M^*_{TD}$

Let $\hat{z}_t: M_{TO}(t) \succeq_t t$ and $t \succeq_t b \ \forall \ b \neq M_{TO}^*(t)$

$\underbrace{\phantom{M_{TO}(t) \succeq_t t}}_{M_{TO}^*(t) \neq t}$

To show: for any match $\hat{M}$ that is stable on $(\hat{z}_t, z_{-t})$

$$\hat{M}(t) = M_{TO}^*(t)$$

(a) Note: $M_{TO}$ is still stable for $(\hat{z}_t, z_{-t}^*)$

(b) Give match $\hat{M}$ that is stable for misreported preferences Profile when everyone is telling the truth.

By IR, $\hat{M}(t) \in \{M_{TO}^*(t), t\}$ if $\hat{M}(t) \neq M_{TO}(t)$ then $\hat{M}(t) = t$

Stable + strict pref. $\overbrace{\phantom{xxxxxx}}$ $M_{TO}^*(t) \in B$
No in RHT

By rural hospital theorem, (When preferences are strict and any hospital does not fill its quota at some stable matching is assigned precisely the same set of students at every stable matching) there cannot be a situation where $t$ is matched under DAA and not matched;

$$\Rightarrow \quad \hat{M}(t) = M_{TO}^*(t).$$

**Step 1**   there is some $t$ s.t. $M_{TD}^{*}(t) \geq_t m(z^*)(t)$

**Step 2**   there is report of $t$ $\bar{\sum}_t$ s.t. $m(\hat{z}_t, z^*_{-t})(t) = M_{TD}(t)$