

Algorithmic Game Theory

Uri Feige

Robi Krauthgamer

Moni Naor

Lecture 9: Social Choice



Lecturer: Moni Naor

Social choice or Preference Aggregation

- Collectively choosing among outcomes
 - Elections,
 - Choice of Restaurant
 - Rating of movies
 - Who is assigned what job
 - Goods allocation
 - Should we build a bridge?
- Participants have **preferences** over outcomes
- **Social choice function** aggregates those preferences and **picks and outcome**

Voting

If there are **two** options and an odd number of voters

- Each having a clear preference between the options

Natural choice: **majority voting**

- Sincere/Truthful
- Order of queries has no significance
 - trivial

When there are more than two options:

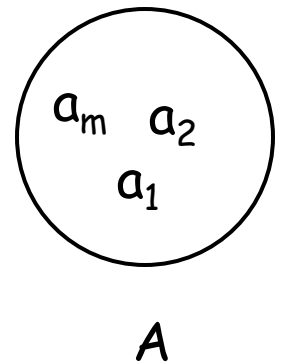
If we start pairing the alternatives:

- Order may matter

a_{10}, a_1, \dots, a_8

Assumption: n voters give their **complete** ranking on set A of alternatives

- L – the set of **linear orders** on A (permutation).
- Each voter i provides \prec_i in L
 - Input to the aggregator/voting rule is $(\prec_1, \prec_2, \dots, \prec_n)$



Goal

A function $f: L^n \mapsto A$ is called a **social choice function**

- Aggregates voters preferences and selects a **winner**

A function $W: L^n \mapsto L$, is called a **social welfare function**

- Aggregates voters preference into a **common order**

Example voting rules

Scoring rules: defined by a vector (a_1, a_2, \dots, a_m)

Being ranked i th in a vote gives the candidate a_i points

- **Plurality:** defined by $(1, 0, 0, \dots, 0)$
 - Winner is candidate that is **ranked first** most often
- **Veto:** is defined by $(1, 1, \dots, 1, 0)$
 - Winner is candidate that is **ranked last** the least often
- **Borda:** defined by $(m-1, m-2, \dots, 0)$



EUROVISION
SONG CONTEST

Jean-Charles de Borda 1770

Plurality with (2-candidate) runoff: top two candidates in terms of plurality score proceed to runoff.

Single Transferable Vote (STV, aka. Instant Runoff): candidate with lowest plurality score drops out; for voters who voted for that candidate: the vote is transferred to the next (live) candidate

Repeat until only one candidate remains

Marquis de Condorcet

Marie Jean Antoine Nicolas de Caritat,
marquis de Condorcet



1743-1794

- **There is something wrong with Borda! [1785]**

Condorcet criterion

- A candidate is the **Condorcet winner** if it wins all of its pairwise elections
- Does not always exist...

Condorcet paradox: there can be **cycles**

- Three voters and candidates:
 $a > b > c, b > c > a, c > a > b$
- a defeats b, b defeats c, c defeats a

Many rules do not satisfy the criterion

• For instance: **plurality:**

- $b > a > c > d$
- $c > a > b > d$
- $d > a > b > c$

• a is the Condorcet winner, but not the plurality winner

- Candidates a and b:
- Comparing how often a is ranked above b, to how often b is ranked above a

Also **Borda:**

- $a > b > c > d > e$
- $a > b > c > d > e$
- $c > b > d > e > a$

Even more voting rules...

- **Kemeny:**

- Consider all pairwise comparisons.
- Graph representation: edge from winner to loser
- Create an overall ranking of the candidates that has as few disagreements as possible with the pairwise comparisons.
 - Delete as few edges as possible so as to make the directed comparison graph acyclic

- 
- Honor societies
 - General Secretary of the UN

- **Approval** [not a ranking-based rule]: every voter labels each candidate as **approved** or **disapproved**. Candidate with the most approvals wins

How do we choose one rule from all of these rules?

- How do we know that there does not exist another, “perfect” rule?
- We will list some **criteria** that we would like our voting rule to satisfy

Arrow's Impossibility Theorem

Skip to the 20th Century

Kenneth Arrow, an economist. In his PhD thesis, 1950, he:

- Listed desirable properties of voting scheme
- Showed that no rule can satisfy all of them.

Properties

- Unanimity
- Independence of irrelevant alternatives
- Not Dictatorial

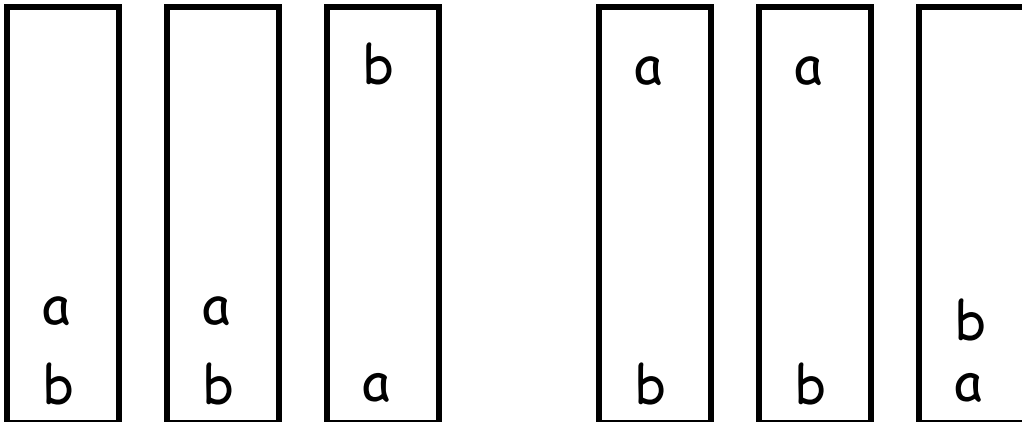


Kenneth Arrow

1921-

Independence of irrelevant alternatives

- Independence of irrelevant alternatives criterion: if
 - the rule ranks **a** above **b** for the current votes,
 - we then change the votes but do not change which is ahead between **a** and **b** in each votethen **a** should still be ranked ahead of **b**.
- None of our rules satisfy this property
 - Should they?



Arrow's Impossibility Theorem

Every **Social Welfare Function** W over a set A of at least 3 candidates:

- If it satisfies

- **Unanimity** (if all voters agree on \langle on the result is \langle)

$$W(\langle, \langle, \dots, \langle) = \langle$$

for all \langle in L

- **Independence of irrelevant alternatives**

Then it is **dictatorial** : there exists a voter i where

$$W(\langle_1, \langle_2, \dots, \langle_n) = \langle_i$$

for all $\langle_1, \langle_2, \dots, \langle_n$ in L

Is there hope for the truth?

- At the very least would like our voting system to encourage voters to tell there true preferences

Strategic Manipulations

- A **social choice function** f can be **manipulated** by voter i if for some $\langle_1, \langle_2, \dots, \langle_n$ and \langle'_i and we have $\mathbf{a} = f(\langle_1, \dots, \langle_i, \dots, \langle_n)$ and $\mathbf{a}' = f(\langle_1, \dots, \langle'_i, \dots, \langle_n)$ but $\mathbf{a} \prec_i \mathbf{a}'$

voter i prefers \mathbf{a}' over \mathbf{a} and can get it by changing his vote

f is called **incentive compatible** if it cannot be manipulated

Gibbard-Satterthwaite Impossibility Theorem

- Suppose there are at least 3 alternatives
- There exists no **social choice function** f that is simultaneously:
 - **Onto**
 - for every candidate, there are some votes that make the candidate win
 - **Nondictatorial**
 - **Incentive compatible**

Implication of Gibbard-Satterthwaite Impossibility Theorem

- All mechanism design problems can be modeled as a social choice problem.
- This theorem seems to quash any hope for designing incentive compatible social choice functions.
- The whole field of Mechanism Design is trying to escape from this impossibility results.
- Introducing “money” is one way to achieve this.

Proof of Arrow's Impossibility Theorem

Claim(Pairwise Unanimity): Every **Social Welfare Function** W over a set A of at least 3 candidates

- If it satisfies

- **Unanimity** (if all voters agree on \prec on the result is \prec)

$$W(\prec, \prec, \dots, \prec) = \prec$$

for all \prec in L

- **Independence of irrelevant alternatives**

Then it is **Pareto efficient**

If $W(\prec_1, \prec_2, \dots, \prec_n) = \prec$ and for all i $a \prec_i b$ then $a \prec b$

Proof of Arrow's Theorem

Claim (Neutrality): let

- $\langle_1, \langle_2, \dots, \langle_n$ and $\langle'_1, \langle'_2, \dots, \langle'_n$ be two profiles
- $\langle = W(\langle_1, \langle_2, \dots, \langle_n)$ and $\langle' = W(\langle'_1, \langle'_2, \dots, \langle'_n)$
- and where for all i

$$a \langle_i b \Leftrightarrow c \langle'_i d$$

Then $a \langle b \Leftrightarrow c \langle' d$

Proof: suppose $a \langle b$ and $c \neq b$

Create a single preference π_i from \langle_i and \langle'_i : where c is just below a and d just above b .

Let $\langle_\pi = W(\pi_1, \pi_2, \dots, \pi_n)$

We must have: (i) $a \langle_\pi b$ (ii) $c \langle_\pi a$ and (iii) $b \langle_\pi d$

And therefore $c \langle_\pi d$ and $c \langle' d$

Preserve the order!

Proof of Arrow's Theorem: Find the Dictator

Claim: For any a, b in A consider sets of profiles

Voters

1	ab	ba	ba	...	ba
2	ab	ab	ba	...	ba
...	ab	ab	ab	...	ba
...
n	ab	ab	ab	...	ba
	0	1	2		n

$a < b$

Profiles

$b < a$

Hybrid argument

Change must happen at some profile i^*

- Where voter i^* changed his opinion

Claim: this i^* is the dictator!

Proof of Arrow's Theorem: i^* is the dictator

Claim: for any $\prec_1, \prec_2, \dots, \prec_n$ and $\prec = W(\prec_1, \prec_2, \dots, \prec_n)$ and c, d in A . If $c \prec_{i^*} d$ then $c \prec d$.

Proof: take $e \neq c, d$ and

- for $i < i^*$ move e to the bottom of \prec_i
- for $i > i^*$ move e to the top of \prec_i
- for i^* put e between c and d

For resulting preferences:

- Preferences of e and c like a and b in profile $i^* - 1$.
- Preferences of e and d like a and b in profile i^* .

$c \prec e$

Therefore $c \prec d$

$e \prec d$

Gibbard-Satterthwaite Impossibility Theorem

- Suppose there are at least 3 alternatives
- There exists no **social choice function** f that is simultaneously:
 - Onto
 - for every candidate, there are some votes that make the candidate win
 - Nondictatorial
 - Incentive compatible

Proof of the Gibbard-Satterthwaite Theorem

Construct a Social Welfare function W_f based on f .

$W_f(\langle_1, \dots, \langle_n) = a$ where $a < b$ iff

$$f(\langle_1^{\{a,b\}}, \dots, \langle_n^{\{a,b\}}) = b$$

Keep everything in order but
move a and b to top

Lemma: if f is an **incentive compatible** social choice function which is onto A , then W_f is a **social welfare function**

- If f is non dictatorial, then W_f also satisfies **Unanimity** and **Independence of irrelevant alternatives**

Proof of the Gibbard-Satterthwaite Theorem

Claim: for all $\langle_1, \dots, \langle_n$ and any subset S of A we have $f(\langle_1^S, \dots, \langle_n^S) \in S$

Keep everything in order but move elements of S to top

Take $a \in S$. There is some $\langle'_1, \langle'_2, \dots, \langle'_n$ where

$$f(\langle'_1, \langle'_2, \dots, \langle'_n) = a.$$

Sequentially change \langle'_i to \langle^S_i

- At no point does f output b not in S .
- Due to the incentive compatibility

Proof of Well Form Lemma

- Antisymmetry: implied by claim for $S=\{a,b\}$
- Transitivity: Suppose we obtained contradicting cycle $a < b < c < a$

take $S=\{a,b,c\}$ and suppose $a = f(\prec_1^S, \dots, \prec_n^S)$

Sequentially change \prec_i^S to $\prec_i^{\{a,b\}}$

Non manipulability implies that

$f(\prec_1^{\{a,b\}}, \dots, \prec_n^{\{a,b\}}) = a$ and $b < a$.

- Unanimity: if for all i , $b \prec_i a$ then

$(\prec_1^{\{a,b\}})^{\{a\}} = \prec_1^{\{a,b\}}$ and $f(\prec_1^{\{a,b\}}, \dots, \prec_n^{\{a,b\}}) = a$

Will repeatedly use the claim to show properties

Proof of Well Form Lemma

- Independence of irrelevant alternatives: if there are two profiles $\langle_1, \langle_2, \dots, \langle_n$ and $\langle'_1, \langle'_2, \dots, \langle'_n$ where for all i $b \langle_i a$ iff $b \langle'_i a$, then

$$f(\langle_1^{\{a,b\}}, \dots, \langle_n^{\{a,b\}}) = f(\langle'_1^{\{a,b\}}, \dots, \langle'_n^{\{a,b\}})$$

by sequentially flipping from $\langle_i^{\{a,b\}}$ to $\langle'_i^{\{a,b\}}$

- Non dictator: preserved