

# Warm-up

Design an algorithm to determine the winner of three candidates  $a, b, c$  given the ranking provided by  $n$  individual voters, described by a  $3 \times n$  matrix  $M$

```
function voting( $M$ )
```

```
Input:  $M$  where  $M_{ij} \in \{a, b, c\}$  is the candidate at rank  $j$  for voter  $i$ 
```

```
Output:  $x \in \{a, b, c\}$  describes the winner
```

```
Return  $x$ 
```

Example Matrix  $M$

	Voter 1	Voter 2	Voter 3	Voter 4
Rank 1	a	c	b	a
Rank 2	b	b	c	b
Rank 3	c	a	a	c

# Announcements

Feedback (please don't forget!):

- [www.cmu.edu/hub/fce](http://www.cmu.edu/hub/fce)
- <https://www.egrad.cs.cmu.edu/ta/S25/feedback/>

Assignments:

- P5 due Thursday April 24
- HW11 due Friday April 25
- Get your questions ready for Wednesday's AMA! (poll)

Final Exam:

- All material is fair game, will focus disproportionately on material not yet covered on midterm exams
- Further details about final exam (including review) to follow shortly (Piazza)<sub>2</sub>

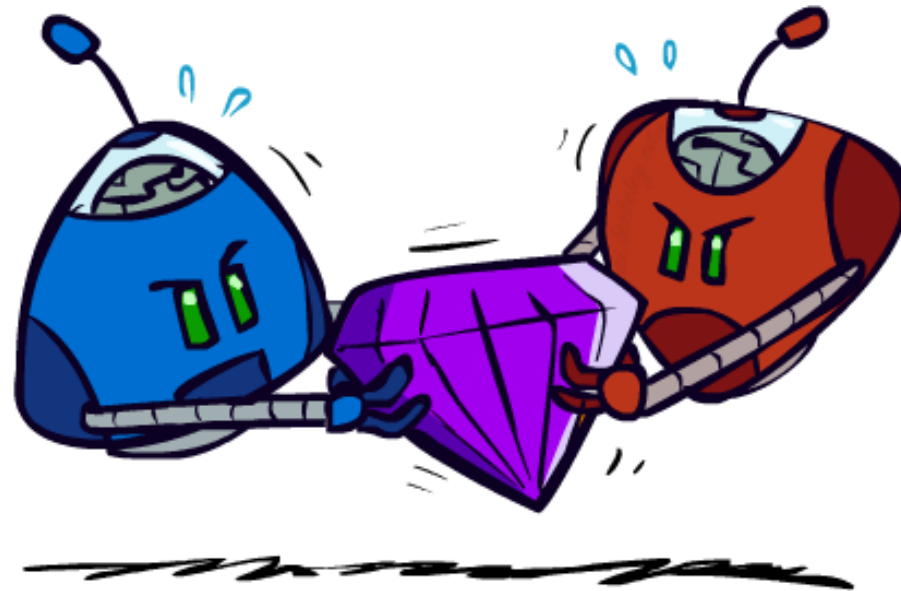
# BONUS assignment

Bonus assignment on voting (will be released later) so you get some practice

Relatively light assignment, instant feedback, but not due until right before final exam

# AI: Representation and Problem Solving

## Game Theory: Equilibrium (cont) & Social Choice



Instructors: Tuomas Sandholm and Vincent Conitzer

Slide credits: CMU AI, Fei Fang

# Normal-Form Games

A game in **normal form** consists of the following elements

- Set of players
- Set of actions for each player
- Payoffs / Utility functions
  - Determines the utility for each player given the actions chosen by all players (referred to as action profile)
- Bimatrix game is special case: two players, finite action sets

Players move **simultaneously** and the game ends immediately afterwards

# Practice for later

What is the Mixed Strategy Nash Equilibrium for this new problem?

		Berry	
		Football	Concert
Alex	Football	4,1	0,0
	Concert	0,0	3,3

# Solution Concepts in Games

How should one player play and what should we expect all the players to play?

- Dominant strategy and dominant strategy equilibrium
- Nash Equilibrium
- Minimax strategy
- Maximin strategy
- Stackelberg Equilibrium

# Power of Commitment

What are the PSNEs in this game, and the players' utilities?

What action should player 2 choose if player 1 commits to playing  $b$ ?  
What is player 1's utility?

What action should player 2 choose if player 1 commits to playing  $a$  and  $b$  uniformly randomly? What is player 1's expected utility?

		Player 2	
		c	d
Player 1	a	2,1	4,0
	b	1,0	3,2



# Stackelberg Equilibrium

## Stackelberg Game

- Leader commits to a strategy first
- Follower responds after observing the leader's strategy

## Stackelberg Equilibrium

- Follower best responds to leader's strategy
- Leader commits to a strategy that maximizes her utility assuming follower best responds

		Player 2	
		c	d
Player 1	a	2,1	4,0
	b	1,0	3,2

# Stackelberg Equilibrium

If the leader can only commit to a pure strategy, or you know that the leader's strategy in equilibrium is a pure strategy, the equilibrium can be found by enumerating the leader's pure strategies

If ties for the follower are broken by the follower such that the leader benefits, the leader can exploit this. This is the **strong Stackelberg equilibrium (SSE)**

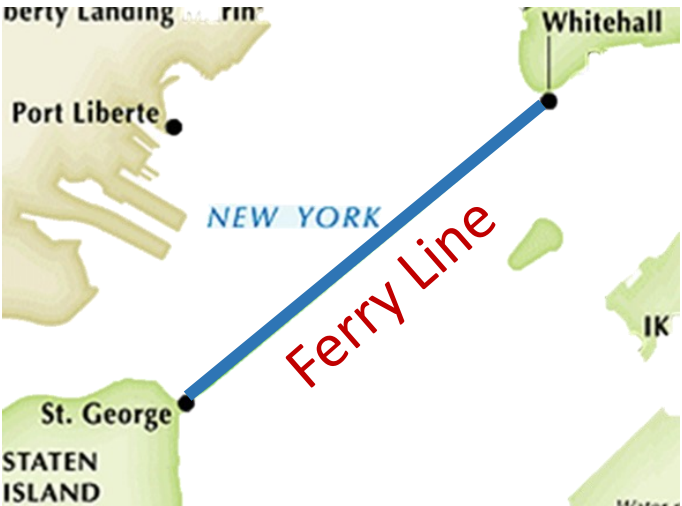
In general, the leader can commit to a mixed strategy and in that case, for the leader:  $u^{SSE} \geq u^{NE}$  (**first-mover advantage**)! & solvable by linear programming!

[Conitzer & Sandholm 2006, von Stengel and Zamir 2010; see also Prof. Fei Fang's work here at CMU (on the following slides)]

		Berry	
		Football	Concert
Alex	Football	2,1	0,0
	Concert	0,0	1,2

		Player 2	
		c	d
Player 1	a	2,1	4,0
	b	1,0	3,2

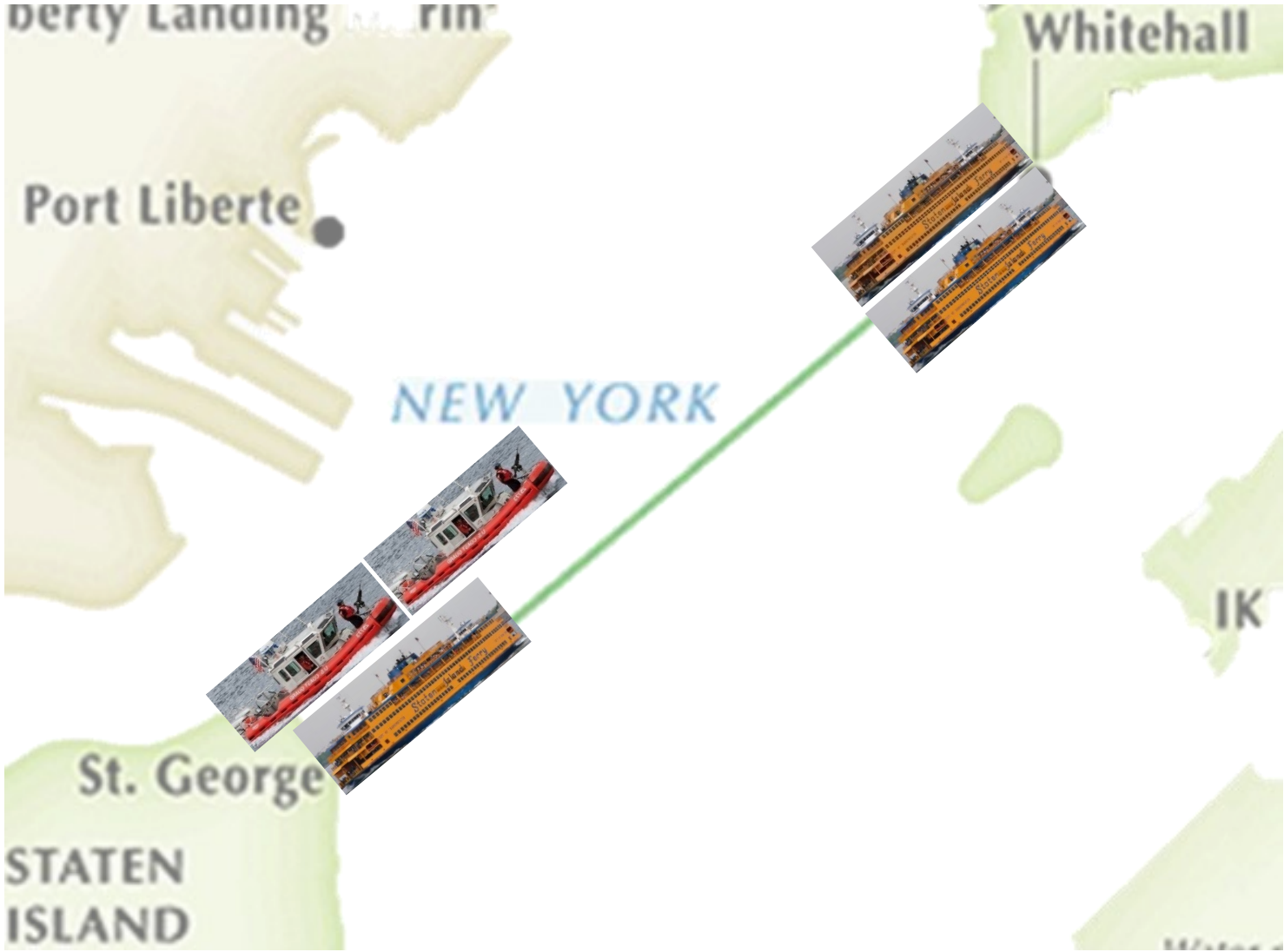
# Protecting Staten Island Ferry



# Protecting Staten Island Ferry



# Previous USCG Approach



# Problem



Optimal Patrol Strategy for Protecting Moving Targets with Multiple Mobile Resources  
Fei Fang, Albert Xin Jiang, Milind Tambe  
In AAMAS-13: The Twelfth International Conference on Autonomous Agents and Multiagent Systems, May 2013

# Game Model and Linear Programming-based Solution

Stackelberg game: Leader – Defender, Follower – Attacker

Attacker's payoff:  $u_i(t)$  if not protected, 0 otherwise

Zero-sum → Strong Stackelberg Equilibrium=Nash Equilibrium  
=Minimax (Minimize Attacker's Maximum Expected Utility)

$$\min_{p_r, v} v$$

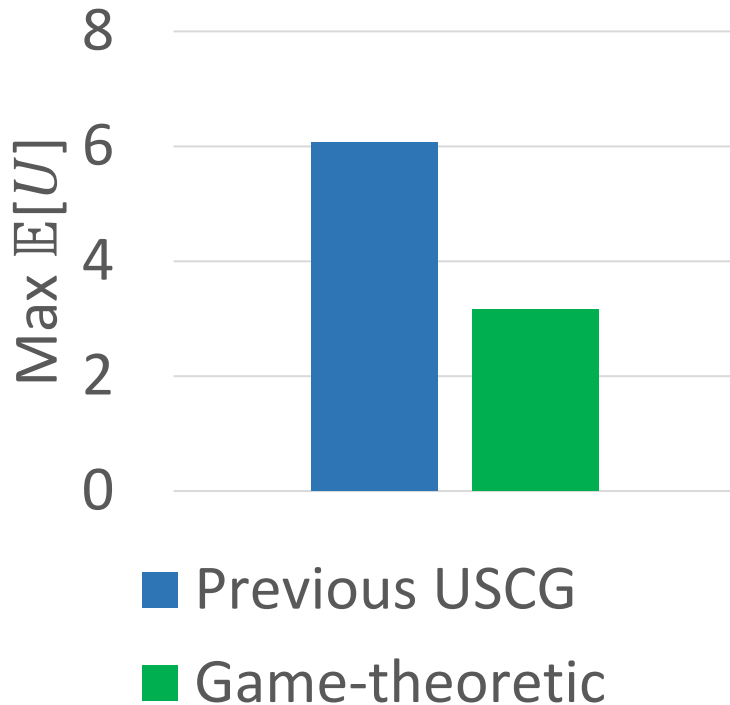
$$\text{s.t. } v \geq \mathbb{E}[U^{att}(i, t)] = u_i(t) \times \mathbb{P}[\text{unprotected}(i, t)], \forall i, t$$

		<b>Adversary</b>			
		10:00:00 AM Target 1	10:00:01 AM Target 1	...	10:30:00 AM Target 3
<b>Defender</b>	$p_r$ ↓ 30%	Purple Route	-5, 5	-4, 4	0, 0
	40%	Orange Route			
	20%	Blue Route			
	.....				

$$\sum_r p_r \leq 1$$

# Evaluation: Simulation & Real-World Feedback

Reduce potential risk by 50%



Deployed by US Coast Guard

USCG evaluation

- Point defense to zone defense
- Increased randomness

Professional mariners:

- Apparent increase in Coast Guard patrols



# Game Theory: Social Choice

arXiv > cs > arXiv:2404.10271

Search... All fields Search

Help | Advanced Search

Computer Science > Machine Learning

[Submitted on 16 Apr 2024 (v1), last revised 4 Jun 2024 (this version, v2)]

## Social Choice Should Guide AI Alignment in Dealing with Diverse Human Feedback

Vincent Conitzer, Rachel Freedman, Jobst Heitzig, Wesley H. Holliday, Bob M. Jacobs, Nathan Lambert, Milan Mossé, Eric Pacuit, Stuart Russell, Hailey Schoelkopf, Emanuel Tewelde, William S. Zwicker

Foundation models such as GPT-4 are fine-tuned to avoid unsafe or otherwise problematic behavior, such as helping to commit crimes or producing racist text. One approach to fine-tuning, called reinforcement learning from human feedback, learns from humans' expressed preferences over multiple outputs. Another approach is constitutional AI, in which the input from humans is a list of high-level principles. But how do we deal with potentially diverging input from humans? How can we aggregate the input into consistent data about "collective" preferences or otherwise use it to make collective choices about model behavior? In this paper, we argue that the field of social choice is well positioned to address these questions, and we discuss ways forward for this agenda, drawing on discussions in a recent workshop on Social Choice for AI Ethics and Safety held in Berkeley, CA, USA in December 2023.

Comments: 15 pages, 4 figures

Subjects: **Machine Learning (cs.LG)**; Artificial Intelligence (cs.AI); Computation and Language (cs.CL); Computers and Society (cs.CY); Computer Science and Game Theory (cs.GT)

MSC classes: 68T01, 68T50, 91B14

### Access Paper:



cs.GT

### References & Citations

- [NASA ADS](#)
- [Google Scholar](#)
- [Semantic Scholar](#)

[Export BibTeX Citation](#)

### Bookmark



Social Choice for AI Ethics and Safety  
**SC4AI25**  
AAMAS 2025 Workshop — Detroit

## HANDBOOK of COMPUTATIONAL SOCIAL CHOICE

Felix Brandt • Vincent Conitzer • Ulle Endriss  
Jerome Lang • Ariel Procaccia



# Warm-up

Design an algorithm to determine the winner of three candidates  $a, b, c$  given the ranking provided by  $n$  individual voters, described by a  $3 \times n$  matrix  $M$

**function voting( $M$ )**

Input:  $M$  where  $M_{ij} \in \{a, b, c\}$  is the candidate at rank  $j$  for voter  $i$

Output:  $x \in \{a, b, c\}$  describes the winner

Return  $x$

Example Matrix  $M$

	Voter 1	Voter 2	Voter 3	Voter 4
Rank 1	a	c	b	a
Rank 2	b	b	c	b
Rank 3	c	a	a	c

# Social Choice Theory

A mathematical theory that deal with aggregation of individual preferences

Wide applications in economics, public policy, etc.

Origins in Ancient Greece

18th century

- Formal foundations by Condorcet and Borda

19th Century

- Charles Dodgson

20th Century

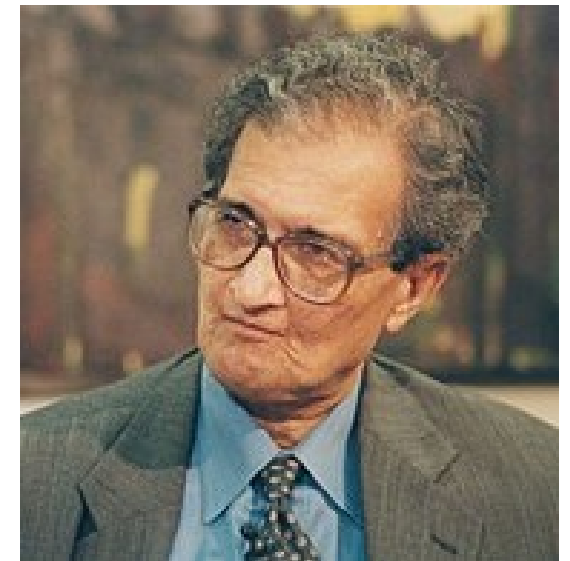
- Nobel Prize in Economics

20th Century – Winners of Nobel Memorial Prize in Economic Sciences

Kenneth Arrow



Amartya Kumar Sen



# Voting Model

## Model

- Set of voters  $N = \{1..n\}$
- Set of alternatives  $A$  ( $|A| = m$ )
- Each voter has a ranking over the alternatives
- Preference profile: collection of all voters' rankings

Voter ID	1	2	3	4
Ranking	a	c	b	a
	b	b	c	b
	c	a	a	c

# Voting Rules

Voting rule: function that maps preference profiles to alternatives that specifies the winner of the election

**function voting( $M$ )**

Input:  $M$  where  $M_{ij} \in \{a, b, c\}$  is the candidate at rank  $j$  for voter  $i$

Output:  $x \in \{a, b, c\}$  describes the winner

Example Matrix  $M$

a	c	b	a
b	b	c	b
c	a	a	c

Return  $x$

# Voting Rules

## Plurality (used in many political elections)

- Each voter gives one point to top alternative
- Alternative with most points wins

Who's the winner? a

Voter ID	1	2	3	4
Ranking	a	c	b	a
	b	b	c	b
	c	a	a	c

# Voting Rules

## Borda count (used for national election in Slovenia)

- Each voter awards  $m - k$  points to alternative ranked  $k^{th}$
- Alternative with most points wins

Who's the winner? b

Voter ID	1	2	3	4
Ranking	a	c	b	a
	b	b	c	b
	c	a	a	c

# Voting Rules

## Borda count (used for national election in Slovenia)

- Each voter awards  $m - k$  points to alternative ranked  $k^{th}$
- Alternative with most points wins

Who's the winner? b

Voter ID	1	2	3	4	$m - k$
Ranking	a	c	b	a	2
	b	b	c	b	1
	c	a	a	c	0

a:  $2+0+0+2=4$ ; b:  $1+1+2+1=5$ ; c:  $0+2+1+0=3$



# Pairwise Election

Alternative  $x$  beats  $y$  in pairwise election if majority of voters prefer  $x$  to  $y$

Who beats who in pairwise election?

b beats c

Voter ID	1	2	3	4
Ranking	a	c	b	a
	b	b	c	b
	c	a	a	c

# Voting Rules

## Plurality with runoff

- First round: two alternatives with highest plurality scores survive
- Second round: pairwise election between the two

$x$  beats  $y$  if majority of voters prefer  $x$  to  $y$

## Who's the winner?

Voter ID	1	2	3	4	5
Ranking	a	c	b	a	c
	b	b	c	b	b
	c	a	a	c	a

# Voting Rules

## Plurality with runoff

- First round: two alternatives with highest plurality scores survive
- Second round: pairwise election between the two

$x$  beats  $y$  if majority of voters prefer  $x$  to  $y$

Who's the winner?      a and c survive, and then c beats a

Voter ID	1	2	3	4	5
Ranking	a	c	b	a	c
	b	b	c	b	b
	c	a	a	c	a

# Voting Rules

## Single Transferable Vote (STV)

- (used in Ireland, Australia, New Zealand, Maine, San Francisco, Cambridge)
- $m - 1$  rounds: In each round, alternative with least plurality votes is eliminated
- Alternative left is the winner

Who's the winner?

Voter ID	1	2	3	4	5
Ranking	a	d	b	a	b
	b	b	c	b	d
	d	c	a	d	a
	c	a	d	c	c

# Voting Rules

## Single Transferable Vote (STV)

- (used in Ireland, Australia, New Zealand, Maine, San Francisco, Cambridge)
- $m - 1$  rounds: In each round, alternative with least plurality votes is eliminated
- Alternative left is the winner

Who's the winner?      c is eliminated, then d, then a, leaving b as the winner.

Voter ID	1	2	3	4	5
Ranking	a	d	b	a	b
	b	b	c	b	d
	d	c	a	d	a
	c	a	d	c	c

Note: When d is eliminated, the vote from voter 2 is effectively transferred to b

# Representation of Preference Profile

Identity of voters does not matter

Only record how many voters have a preference

33 voters	16 voters	3 voters	8 voters	18 voters	22 voters
a	b	c	c	d	e
b	d	d	e	e	c
c	c	b	b	c	b
d	e	a	d	b	d
e	a	e	a	a	a

# Tie Breaking

Commonly used tie breaking rules include

- Borda count
- Having the most votes in the first round
- ...

# Social Choice Axioms

How do we choose among different voting rules? What are the desirable properties?



# Majority consistency

**Majority consistency:** Given a voting rule that satisfies Majority Consistency, if a majority of voters ( $> 50\%$  of voters) rank alternative  $x$  first, then  $x$  should be the final winner.

# Poll 1

Which rules are NOT majority consistent?

- A. Plurality: Each voter give one point to top alternative
- B. Borda count: Each voter awards  $m - k$  points to alternative ranked  $k^{th}$
- C. Plurality with runoff: Pairwise election between two alternatives with highest plurality scores
- D. STV: In each round, alternative with least plurality votes is eliminated
- E. None

# Condorcet Consistency

Recall:  $x$  beats  $y$  in a pairwise election if majority of voters prefer  $x$  to  $y$

**Condorcet winner** is an alternative that beats every other alternative in pairwise election

Does a Condorcet winner always exist?

Condorcet paradox = cycle in majority preferences

Voter ID	1	2	3
Ranking over alternatives (first row is the most preferred)	a	c	b
	b	a	c
	c	b	a

Condorcet Consistency: a Condorcet winner should always win

If a rule satisfies majority consistency, does it satisfy Condorcet consistency? Vice versa?

Which of the introduced voting rules (Plurality, Borda count, Plurality with runoff, STV) are Condorcet consistent?

# Poll 2

## Which rules ARE Condorcet consistent?

- A. Plurality: Each voter give one point to top alternative
- B. Borda count: Each voter awards  $m - k$  points to alternative ranked  $k^{th}$
- C. Plurality with runoff: Pairwise election between two alternatives with highest plurality scores
- D. STV: In each round, alternative with least plurality votes is eliminated
- E. None

# Condorcet Consistency

## Winner under different voting rules in this example

- Plurality:
- Borda:
- Plurality with runoff:
- STV:
- Condorcet winner:

33 voters	16 voters	3 voter	8 voters	18 voters	22 voters
a	b	c	c	d	e
b	d	d	e	e	c
c	c	b	b	c	b
d	e	a	d	b	d
e	a	e	a	a	a

# Condorcet Consistency

## Winner under different voting rules in this example

- Plurality: a
- Borda: b
- Plurality with runoff: e
- STV: d
- Condorcet winner: c

33 voters	16 voters	3 voter	8 voters	18 voters	22 voters
a	b	c	c	d	e
b	d	d	e	e	c
c	c	b	b	c	b
d	e	a	d	b	d
e	a	e	a	a	a

# Strategy-Proofness

## Using Borda Count

Who is the winner?

Voter ID	1	2	3	$m - k$
Ranking over alternatives (first row is the most preferred)	b	b	a	3
	a	a	b	2
	c	c	c	1
	d	d	d	0

Who is the winner now?

Voter ID	1	2	3	$m - k$
Ranking over alternatives (first row is the most preferred)	b	b	a	3
	a	a	c	2
	c	c	d	1
	d	d	b	0



# Strategy-Proofness

A single voter can manipulate the outcome!

Voter ID	1	2	3	$m - k$
Ranking over alternatives (first row is the most preferred)	b	b	a	3
	a	a	b	2
	c	c	c	1
	d	d	d	0

$$b: 2*3+1*2=8$$

$$a: 2*2+1*3=7$$

b is the winner

Voter ID	1	2	3	$m - k$
Ranking over alternatives (first row is the most preferred)	b	b	a	3
	a	a	c	2
	c	c	d	1
	d	d	b	0

$$b: 2*3+1*0=6$$

$$a: 2*2+1*3=7$$

a is the winner

# Strategy-Proofness

A voting rule is **strategyproof (SP)** if a voter can never **benefit** from lying about his preferences (regardless of what other voters do)

Do not lie: b is the winner

Voter ID	1	2	3
Ranking	b	b	a
	a	a	b
	c	c	c
	d	d	d

Lie: a is the winner

Voter ID	1	2	3
Ranking	b	b	a
	a	a	c
	c	c	d
	d	d	b

If a voter's preference is  $a > b > c$ , c will be selected w/o lying, and b will be selected w/ lying, then the voter still benefits

# Poll 3

Which of the introduced voting rules are strategyproof?

- A. Plurality: Each voter give one point to top alternative
- B. Borda count: Each voter awards  $m - k$  points to alternative ranked  $k^{th}$
- C. Plurality with runoff: Pairwise election between two alternatives with highest plurality scores
- D. STV: In each round, alternative with least plurality votes is eliminated
- E. None

# Greedy Algorithm for $f$ – Manipulation

Given voting rule  $f$  and preference profile of  $n - 1$  voters, how can the last voter report preference to let a specific alternative  $y$  *uniquely* win (no tie breaking)?

## Greedy algorithm for $f$ – Manipulation

Rank  $y$  in the first place

While there are unranked alternatives

    If  $\exists x$  that can be placed in the next spot without preventing  $y$  from winning

        place this alternative in the next spot

    else

        return false

return true (with final ranking)

Correctness proved under conditions (Bartholdi et al., 1989)

# Greedy Algorithm for $f$ – Manipulation

Example with Borda count voting rule

Voter ID	1	2	3
Ranking over alternatives (first row is the most preferred)	b	b	a
	a	a	
	c	c	
	d	d	

# Other Properties

A voting rule is **dictatorial** if there is a voter who always gets their most preferred alternative

A voting rule is **constant** if the same alternative is always chosen (regardless of the stated preferences)

A voting rule is **onto** if any alternative can win, for some set of stated preferences

Which of the introduced voting rules (Plurality, Borda count, Plurality with runoff, STV) are dictatorial, constant or onto?

# Results in Social Choice Theory

Constant functions and dictatorships are SP Why?

Theorem (Gibbard-Satterthwaite): If  $m \geq 3$ , then any voting rule that is SP and onto is dictatorial

- Any voting rule that is onto and nondictatorial is manipulable (=not strategy-proof)
- It is **impossible** to have a voting rule that is strategyproof, onto, and nondictatorial

# Other CS courses at CMU on game theory ( / social choice theory)

15-326 Computational Microeconomics (Vince)

15-784 Foundations of Cooperative AI (Vince)

15-888 Computational Game Solving (Tuomas) OFFERED F25

...



Activity: Favorite topics of 15281 (by approval voting)

# Learning Objectives

Understand the voting model

Find the winner under the following voting rules

- Plurality, Borda count, Plurality with runoff, Single Transferable Vote

Describe the following concepts, axioms, and properties of voting rules

- Pairwise election, Condorcet winner
- Majority consistency, Condorcet consistency, Strategyproof
- Dictatorial, constant, onto

Understand the possibility of satisfying multiple properties

Describe the greedy algorithm for voting rule manipulation