

Graduate AI

Lecture 22:

Game Theory IV

Teachers:

Zico Kolter

Ariel Procaccia (this time)

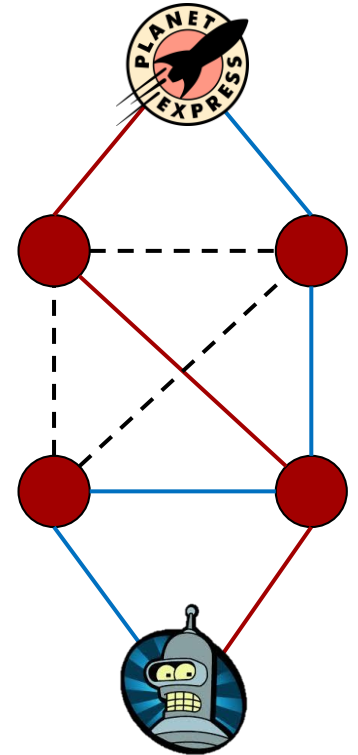
REMINDER: THE MINIMAX THEOREM

- Theorem [von Neumann, 1928]: Every 2-player zero-sum game has a unique value v such that:
 - Player 1 can guarantee value at least v
 - Player 2 can guarantee loss at most v
- We will prove the theorem via no-regret learning



HOW TO REACH YOUR SPACESHIP

- Each morning pick one of n possible routes
- Then find out how long each route took
- Is there a strategy for picking routes that does almost as well as the best fixed route **in hindsight**?



53 minutes

47 minutes

...

THE MODEL

- View as a matrix (maybe infinite #columns)

Adversary

Algorithm							

- Algorithm picks row, adversary column
- Alg pays cost of (row,column) and gets column as feedback
- Assume costs are in $[0,1]$

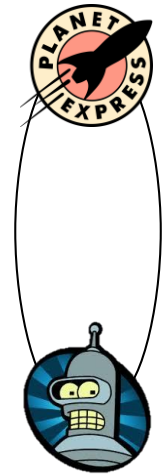
THE MODEL

- Define **average regret** in T time steps as (average per-day cost of alg) – (average per-day cost of best fixed row in hindsight)
- **No-regret algorithm**: $\text{regret} \rightarrow 0$ as $T \rightarrow \infty$
- Not competing with adaptive strategy, just the best **fixed** row



EXAMPLE

- **Algorithm 1:** Alternate between U and D
- **Poll 1:** What is algorithm 1's worst-case average regret?
 1. $\Theta(1/T)$
 2. $\Theta(1)$
 3. $\Theta(T)$
 4. ∞

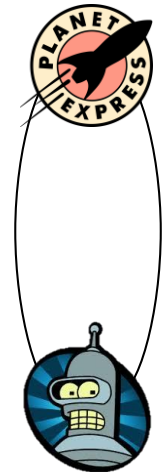


Adversary

	1	0
Algorithm	0	1

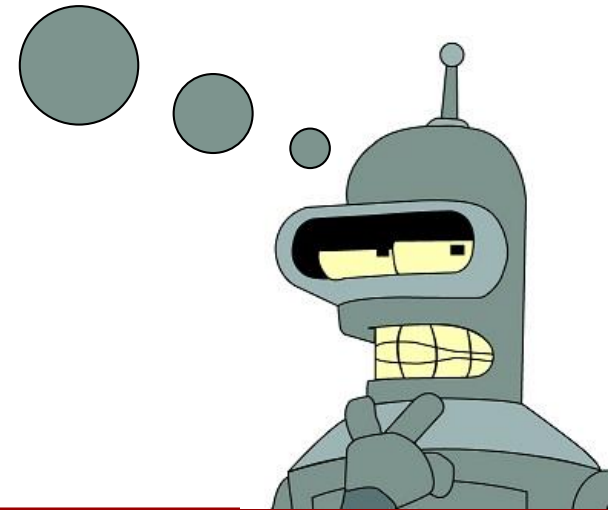
EXAMPLE

- **Algorithm 2:** Choose action that has lower cost so far
- **Poll 2:** What is algorithm 2's worst-case average regret?
 1. $\Theta(1/T)$
 2. $\Theta(1/\sqrt{T})$
 3. $\Theta(1/\log T)$
 4. $\Theta(1)$



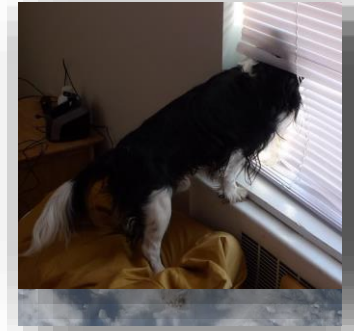
		Adversary	
Algorithm	1	1	0
	0	0	1

What can we say
more generally
about deterministic
algorithms?



USING EXPERT ADVICE

- Want to predict the stock market
- Solicit advice from n experts
 - Expert = someone with an opinion



Day	Expert 1	Expert 2	Expert 3	Charlie	Truth
1	-	-	+	+	+
2	+	-	+	-	-
...

- Can we do as well as best in hindsight?



WEIGHTED MAJORITY

- **Idea:** Experts are penalized every time they make a mistake
- **Weighted Majority Algorithm:**
 - Start with all experts having weight 1
 - Predict based on weighted majority vote
 - Penalize mistakes by cutting weight in half

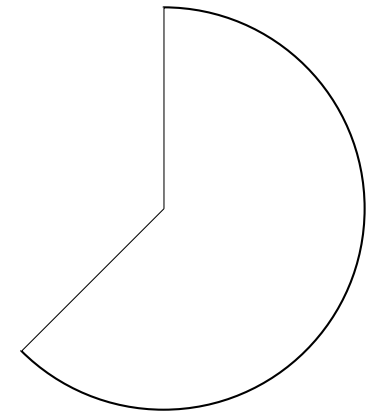
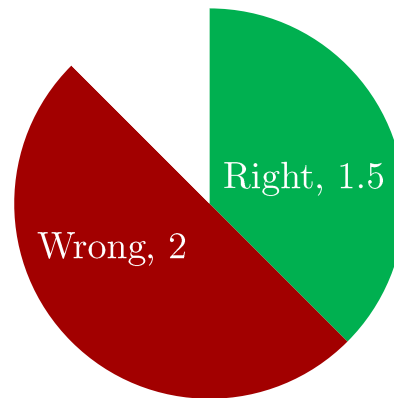
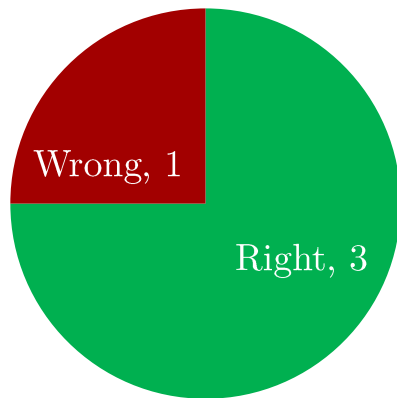


	Expert 1	Expert 2	Expert 3	Charlie
Weights	1	1	1	1
Prediction	-	+	+	+
Weights	0.5	1	1	1
Prediction	+	+	-	-
Weights	0.5	1	0.5	0.5

Alg	Truth
-----	-------

+	+
---	---

-	+
---	---



WEIGHTED MAJORITY: ANALYSIS

- $M = \#$ mistakes we've made so far
- $m = \#$ mistakes of best expert so far
- $W =$ total weight (starts at n)
- For each mistake, W drops by at least 25%
 \Rightarrow after M mistakes: $W \leq n(3/4)^M$
- Weight of best expert is $(1/2)^m$

$$\left(\frac{1}{2}\right)^m \leq n \left(\frac{3}{4}\right)^M \Rightarrow \left(\frac{4}{3}\right)^M \leq n2^m \Rightarrow M \leq 2.5(m + \lg n)$$



RANDOMIZED WEIGHTED MAJORITY

- Randomized Weighted Majority

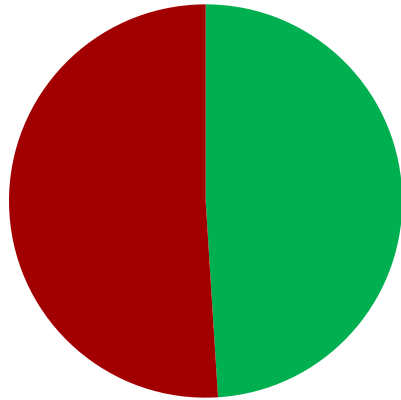
Algorithm:

- Start with all experts having weight 1
- Predict **proportionally** to weights: the total weight of + is w_+ and the total weight of - is w_- , predict + with probability $\frac{w_+}{w_+ + w_-}$ and - with probability $\frac{w_-}{w_+ + w_-}$
- Penalize mistakes by removing ϵ **fraction** of weight

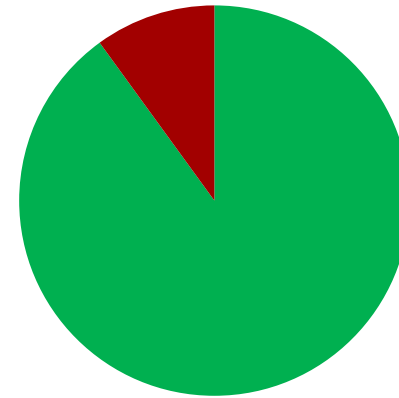


RANDOMIZED WEIGHTED MAJORITY

Idea: smooth out the worst case



The worst-case is
~50-50: now we have
a 50% chance of
getting it right



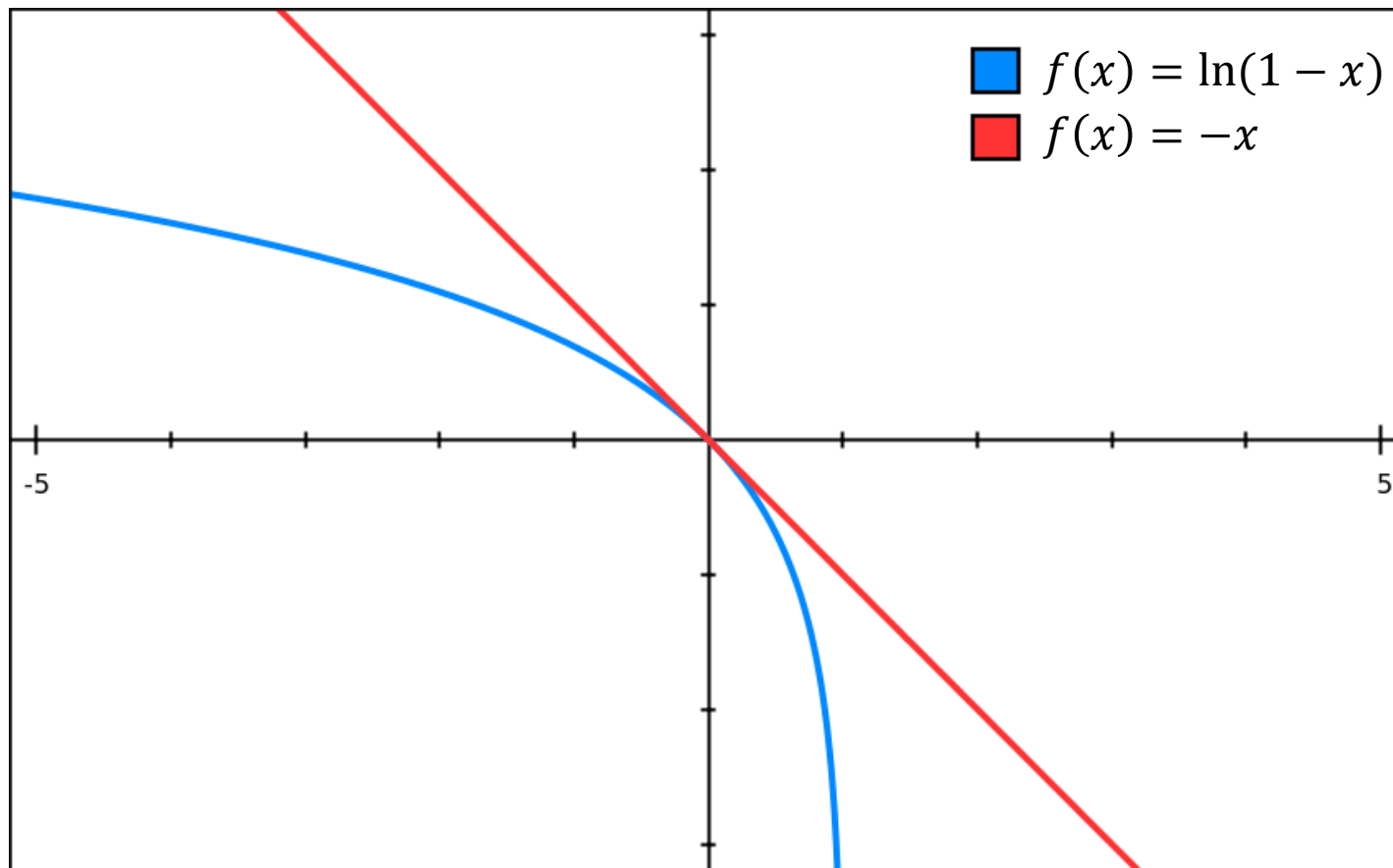
What about 90-10?
We're very likely to
agree with the
majority

ANALYSIS

- At time t we have a fraction F_t of weight on experts that made a mistake
- Prob. F_t of making a mistake, remove ϵF_t fraction of total weight
- $W_{final} = n \prod_t (1 - \epsilon F_t)$
- $\ln W_{final} = \ln n + \sum_t \ln(1 - \epsilon F_t)$
 $\leq \ln n - \epsilon \sum_t F_t = \ln n - \epsilon M$

↑
 $\ln(1 - x) \leq -x$
(next slide)

ANALYSIS



ANALYSIS

- Weight of best expert is $W_{best} = (1 - \epsilon)^m$
- $\ln n - \epsilon M \geq \ln W_{final} \geq \ln W_{best} = m \ln(1 - \epsilon)$
- By setting $\epsilon = \sqrt{\ln n / m}$ and solving, we get
$$M \leq m + 2\sqrt{m \ln n}$$
- Since $m \leq T$, $M \leq m + 2\sqrt{T \ln n}$
- Average regret is $(2\sqrt{T \ln n})/T \rightarrow 0$ ■



MORE GENERALLY

- Each **expert** is an **action** with cost in $[0,1]$
- Run Randomized Weighted Majority
 - Choose expert i with probability w_i/W
 - Update weights: $w_i \leftarrow w_i(1 - c_i\epsilon)$
- Same analysis applies:
 - Our expected cost: $\sum_j c_j w_j / W$
 - Fraction of weight removed: $\epsilon \sum_j c_j w_j / W$
 - So, fraction removed = $\epsilon \cdot$ (our cost)

PROOF OF THE MINIMAX THEOREM

- In a zero-sum game G , denote:
 - V_C is the smallest reward the column player can guarantee if he commits first
 - V_R is the largest reward the row player can guarantee if he commits first
- Obviously $V_C \geq V_R$, and the theorem says equality holds
- Assume for contradiction that $V_C > V_R$
- Scale matrix so that payoffs to row player are in $[-1,0]$, and let $V_C = V_R + \delta$



PROOF OF THE MINIMAX THEOREM

- Suppose the game is played repeatedly; in each round the row player commits, and the column player responds
- Let the row player play RWM, and let the column player respond optimally to current mixed strategy
- After T steps
 - $\text{ALG} \geq \text{best row in hindsight} - 2\sqrt{T \log n}$
 - $\text{ALG} \leq T \cdot V_R$



PROOF OF THE MINIMAX THEOREM

- **Claim:** Best row in hindsight $\geq T \cdot V_C$
 - Suppose the column player played s_t in round t
 - Define a mixed strategy y that plays each s_t with probability $1/T$ (multiplicities possible)
 - Let x be row's best response to y
 - $V_C \leq u_1(x, y) = \frac{1}{T} u_1(x, s_1) + \dots + \frac{1}{T} u_1(x, s_T)$
 - $u_1(x, s_1) + \dots + u_1(x, s_T) \leq \text{best row in hindsight}$ ■
- It follows that $T \cdot V_R \geq T \cdot V_C - 2\sqrt{T \log n}$
- $\delta T \leq 2\sqrt{T \log n}$ – contradiction for large T ■



SUMMARY

- Terminology:
 - Regret
 - No-regret learning
- Algorithms:
 - Randomized weighted majority
- Big ideas:
 - It is possible to achieve no-regret learning guarantees!
 - Connections between game theory and learning theory

