

Last time

- PAC learning
- Gradient descent & SGD

Today

- online learning
Multi-weight update alg
- applications including
zero-sum games

Figure credits

- Avrim Blum slides
- Emily Fox 4+6 slides
- book on Boosting by Schapire & Freund
- Nate Jensen, Gabrielle Cohen

$$A = \{P_1, P_2, \dots, P_m\}$$

\uparrow
 R_1^t (R_2^t) \dots R_m^t

Online learning.

A : action set $|A|=n$

reward vector
on day t

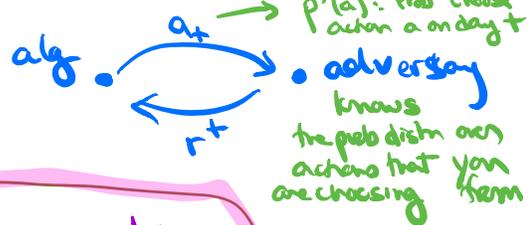
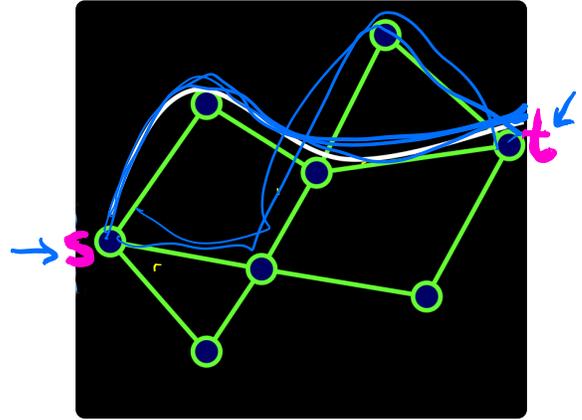
$$r^t: A \rightarrow [-1, 1]$$

$r^t(a)$: payoff on day t
if you choose action a

Alg: on day t , chooses $a_t \in A$

knows only r^1, r^2, \dots, r^{t-1}

$t=1, \dots, T$



What can you hope to achieve?

Best possible outcome

$$\sum_{t=1}^T \max_a r^t(a)$$

$$\sum_{a \in A} p^t(a) = 1$$

This benchmark is way too strong.

$$-p^t(1) + p^t(2) \leq 0$$

$$A = \{1, 2\}$$

$$\left. \begin{array}{l} \text{if } p^t(1) \geq \frac{1}{2} \\ \text{o.w. } p^t(2) \geq \frac{1}{2} \end{array} \right\} \left. \begin{array}{l} r^t(1) = -1 \\ r^t(2) = 1 \\ r^t(1) = 1 \\ r^t(2) = -1 \end{array} \right\}$$

benchmark T

combining expert advice
or doing as well as best
expert.



$$\text{Regret}(a^1, a^2, \dots, a^T) = \max_{a \in A} \sum_{t=1}^T r^t(a) - \sum_{t=1}^T r^t(a_t)$$

best reward if you use same action every day
alg total reward.

on day t , need to choose $a \in A$.

"Follow the leader"

Claim: without randomization, impossible to do well

on day t $r^t(a_t) = 0$ \implies alg total reward of 0.
 $r^t(a) = 1$ $a \neq a_t$

n actions $\sum_{t=1}^T \sum_{a \in A} r^t(a) = (n-1)T$

\exists action w/ reward at least $\frac{(n-1)T}{n} = (1 - \frac{1}{n})T$

$\geq \frac{T}{2}$

What is the best we can hope for?

$n=2$

toss a fair coin:



Expected reward of any online alg: 0

if pick action 1,

$\#H's - \#T's$

2,

$\#T's - \#H's.$

$E(|\#H's - \#T's|) = \Theta(\sqrt{T})$

central limit theorem.

Regret is $\Omega(\sqrt{T})$

$\frac{\text{regret}}{T} = O(\frac{1}{\sqrt{T}}) \rightarrow 0$

Thm: \exists online alg with exp regret $2\sqrt{T \ln n}$
 $\max_a \sum_t r_t(a) - E(\text{reward of online alg}) \leq 2\sqrt{T \ln n}$

$$\frac{1}{T} E(\text{reward of online alg}) \geq \frac{1}{T} \max_a \sum_t r_t(a) - \frac{2\sqrt{T \ln n}}{T}$$

total reward of best action

$-1 \leq r_t(a) \leq 1 \quad \forall a$

$2\sqrt{\frac{\ln n}{T}}$

MWU algorithm
 initialize $w^1(a) = 1 \quad \forall a$ initialize all weights to 1
 for $t=1$ to T
 pick action a^t with probability proportional to $w^t(a)$
 given r^t , update weights.
 $w^{t+1}(a) = w^t(a) \cdot (1 + \eta \cdot r^t(a))$ (*)
learning rate parameter $\eta \in (0, \frac{1}{2})$

$p^t(a) = \frac{w^t(a)}{\sum_{a \in A} w^t(a)} \leftarrow p^t$ sum of all wts at time t

$E[\text{reward of alg at time } t]$
 $= \sum_a p^t(a) r^t(a)$
 $= \frac{1}{p^t} \sum_a w^t(a) r^t(a)$ $= \sum_a p^t(a) r^t(a)$

$p^{t+1} = \sum_a w^t(a) (1 + \eta r^t(a))$
 $= p^t + \eta \sum_a w^t(a) r^t(a)$
 $p^t \sigma_t$

Suppose a^* is best action in hindsight

$p^{T+1} \geq w^{T+1}(a^*)$
 $= w^1(a^*) \prod_{t=1}^T (1 + \eta r^t(a^*))$
 $e^{\sum_{t=1}^T \ln(1 + \eta r^t(a^*))}$

$p^{t+1} \leq p^t (1 + \eta \sigma_t)$
 $p^{t+1} \leq p^t e^{\eta \sigma_t^2} \quad \forall t$ $1+x \leq e^x$

$p^T \leq p^1 e^{\sum_{t=1}^T \eta \sigma_t^2}$
 $p^T \leq \frac{1}{n} e^{\sum_{t=1}^T \eta \sigma_t^2}$ (b)
#actions

$\ln(1+x) \geq x - x^2$
 $\ln(1+x) = \frac{x - x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} \dots$

$\geq e^{\sum_t \ln(1 + \eta r^t(a^*))}$
 $\geq e^{\sum_t [\eta r^t(a^*) - \eta^2 r^t(a^*)^2]}$
 $\geq e^{\left[\sum_t \eta r^t(a^*) \right] - \eta^2 T}$ (a)

(a) \leq (b) $\ln(a) \leq \ln(b)$

$$\frac{\sum_{t=1}^T \eta r_t(a^*) - \eta^2 T}{\eta} \leq \frac{\ln n}{\eta} + \eta \sum_{t=1}^T \sigma_t^2$$

$\uparrow \ln(a)$
 \uparrow
 $E[\text{reward of MWU}]$

$E[\text{reward of MWU}] \geq \text{total reward of } a^*$

$$-\eta T - \frac{\ln n}{\eta}$$

$$\eta T = \frac{\ln n}{\eta}$$

$$\eta^2 = \frac{\ln n}{T}$$

$$\eta = \sqrt{\frac{\ln n}{T}}$$

$$\sqrt{T \ln n}$$

with $\eta = \sqrt{\frac{\ln n}{T}}$

$$\frac{1}{T} E[\text{reward of MWU}] \geq \frac{1}{T} \max_a \sum_{t=1}^T r_t(a) - 2 \sqrt{\frac{T \ln n}{T}}$$

$$2 \sqrt{\frac{\ln n}{T}}$$

per-step regret $\rightarrow 0$ as $T \rightarrow \infty$

2-player Zero-sum games



Penalty Kicks.

zero-sum means that one player's gain is other player's loss.

col player: goalkeeper

	L	R
row player: <u>kicker</u>	L	0 1
	R	1 0

payoff to row player

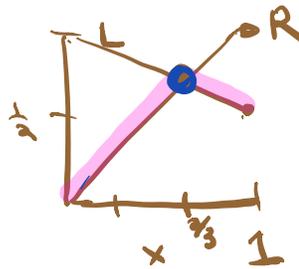
rows (cols) are called "pure" strategies.

mixed strategy is prob distn over pure strategies.

if $x = \frac{2}{3}$

		goalkeeper:	
		L	R
kicker	L	x $\frac{1}{2}$ 1	
	R	$1-x$ 1 0	

What is opt mixed strategy for kicker if he has to go first.



$$1 - .5 \cdot \frac{2}{3} = \frac{2}{3}$$

if goalkeeper goes L,
exp loss = $.5x + 1 - x$
= $1 - .5x$

if goalkeeper goes R
exp loss = x

$$1 - .5x = x$$

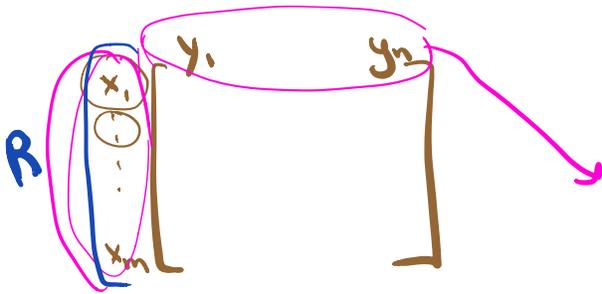
if kicker goes first maximize loss of goalkeeper

max x gain of kicker when goalkeeper best responds to x

$x = \frac{2}{3}$

minimax opt stratgy for row player
 zero sum game defined $m \times n$ matrix A

where a_{ij} is payoff to row player
 when he plays pure strategy i
 & col player plays pure strategy j



mixed strategy for col player
 $y_j = \text{Pr}(\text{play col } j)$

mixed strategy for row player
 $x_i \geq 0$
 $\sum_{i=1}^m x_i = 1$

$$E[\text{payoff to row player}] = \sum_{i=1}^m \sum_{j=1}^n \text{Pr}(\text{Row } i, \text{Col } j) a_{ij}$$

$x_i \cdot y_j$

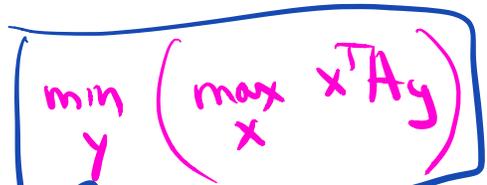
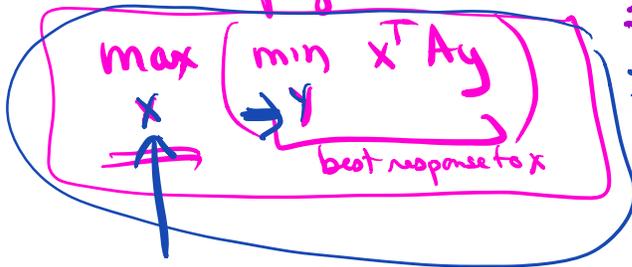
$$= \sum_i \sum_j x_i y_j a_{ij}$$

$$= x^T A y$$

if row player goes first & chooses \vec{x}

if col player first & chooses y

\Rightarrow col player will choose



x row player mixed str
 y col player mixed str.

"
 valued game.

Portfolio Selection.

Stock	P_i	$r_i = \frac{\text{closing price}}{\text{opening price}}$	r_i^2
1	0.5	1.5	0.9
2	0	1	0.1
3	0.5	1	0
4	0	0.5	0

0.3
0.1
0.4
0.2

\$1
 $\$W_0$

$$0.5 \cdot 1.5 + 0.5 \cdot 1 = \$1.25$$

$1.25 W_0$

p_i^t : fraction of your wealth that you put into stock i on day t .
 r_i^t : $\frac{\text{closing price of stock } i \text{ on day } t}{\text{opening price}}$

$$\frac{W^T}{W_0} = r_1^1 r_2^1 r_3^1 \dots r_1^T$$

$(W_0 r_1^2)$

$$\log\left(\frac{W^T}{W_0}\right) = \sum_{t=1}^T \log r_i^t$$

Lets pretend reward of stock i on day t is $\log r_i^t$

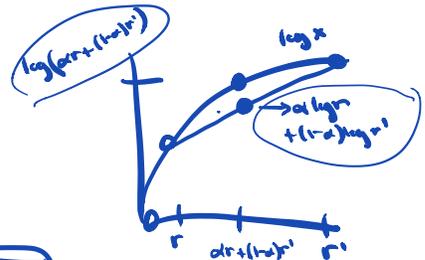
$$\sum_{t=1}^T \sum_i p_i^t \log r_i^t \geq \sum_{t=1}^T \log r_j^t - \epsilon$$

$\forall j$

What we really care about:

$$\frac{W^T}{W_0} = \prod_{t=1}^T \left(\sum_i p_i^t r_i^t \right)$$

$$\log\left(\frac{W^T}{W_0}\right) = \sum_{t=1}^T \log\left(\sum_i p_i^t r_i^t\right)$$



$$\sum_{t=1}^T \log\left(\sum_i p_i^t r_i^t\right) \geq \sum_{t=1}^T \sum_i p_i^t \log r_i^t \geq \sum_{t=1}^T \log r_j^t - \epsilon$$

$$\log\left(\frac{W^T}{W_0}\right) \geq \log \prod_{t=1}^T r_j^t - \epsilon$$

$$\frac{W_T(\omega_0)}{W_0} \approx \prod_{t=1}^T \left(1 + \frac{\epsilon}{\omega_0} \right) \approx 1 + \epsilon$$

$\forall \text{ stack } j$

CRB