

# Sequential decision making: modeling how we interact with the world

Alan Malek

April 9th, 2016

- Omer Atun: CEO of AgilOne Predictive Marketing
- Brienne Ghafouriar: CEO of Entefy
- Jeffrey Rothschild: founder Veritas Software and Mpath Interactive, Facebook VP of Infrastructure Software

- Omer Atun: CEO of AgilOne Predictive Marketing
- Brienne Ghafouriar: CEO of Entefy
- Jeffrey Rothschild: founder Veritas Software and Mpath Interactive, Facebook VP of Infrastructure Software
- Alan Malek: Graduate student

- Omer Atun: CEO of AgilOne Predictive Marketing
- Brienne Ghafouriar: CEO of Entefy
- Jeffrey Rothschild: founder Veritas Software and Mpath Interactive, Facebook VP of Infrastructure Software
- Alan Malek: Graduate student,



Harker alumnus '05

# Goals of this talk

- What is Grad school?
- Modern sequential decision making problems
- Whet your appetite with a cool problem
- Some advice: is grad school for you?

# What is a Grad school?

# What is a Grad school?

- A third the pay for a third the responsibilities

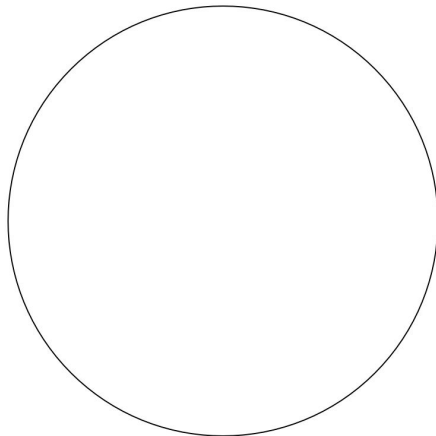
# What is a Grad school?

- A third the pay for a third the responsibilities
- Take classes for a few years
- Be a TA
- Do research
  - In theory: think about how to formalize problems, prove theorems
  - Or more applied: engineer solutions



# What is a PhD?

Imagine a circle that contains all of human knowledge:

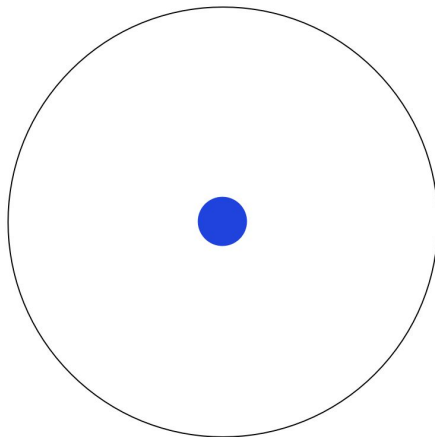


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

By the time you finish elementary school, you know a little:

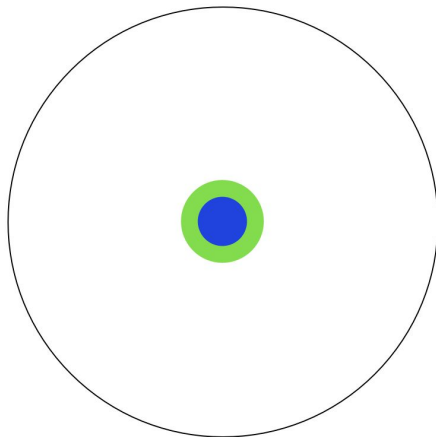


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

By the time you finish high school, you know a bit more:

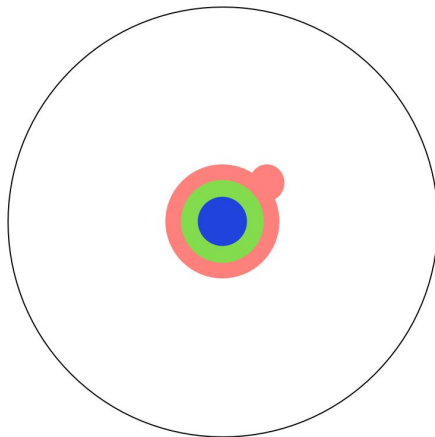


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

With a bachelor's degree, you gain a specialty:

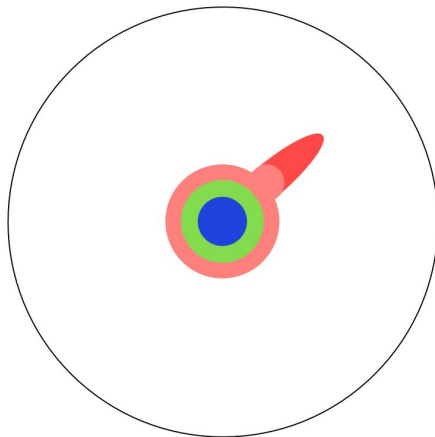


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

A master's degree deepens that specialty:

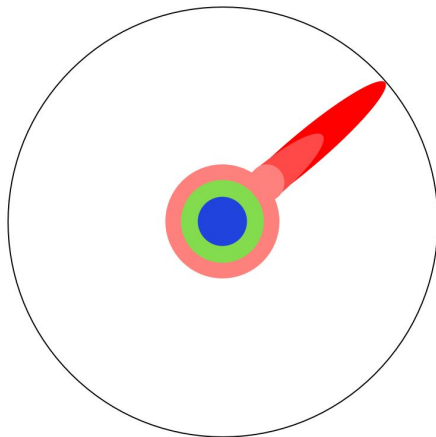


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

Reading research papers takes you to the edge of human knowledge:

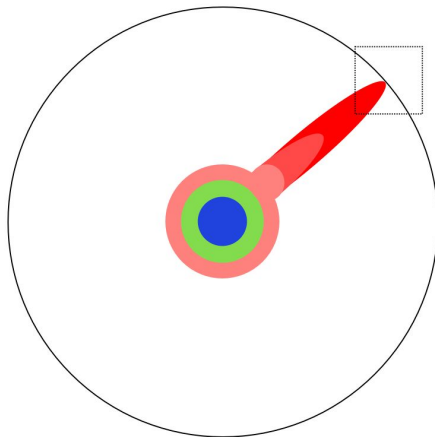


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

Once you're at the boundary, you focus:

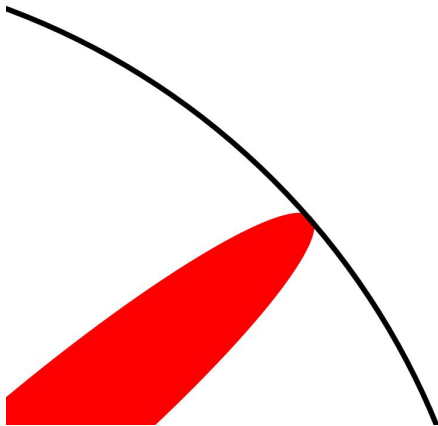


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

You push at the boundary for a few years:



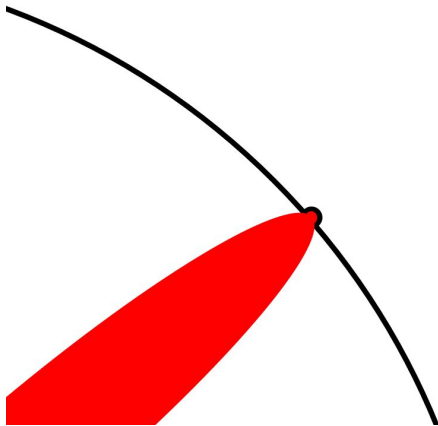
Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>



# What is a PhD?

Until one day, the boundary gives way:

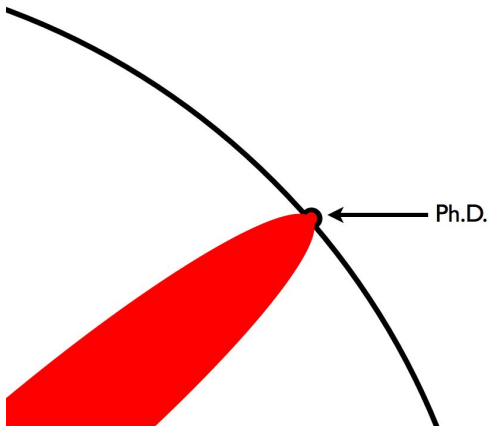


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

And, that dent you've made is called a Ph.D.:

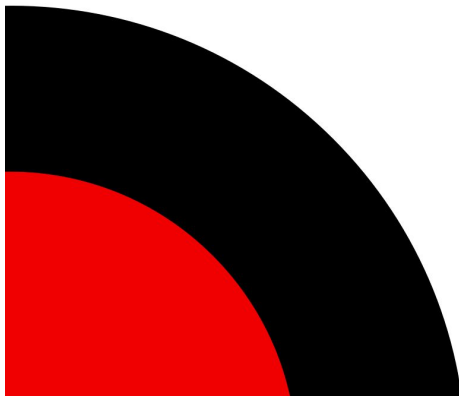


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

Of course, the world looks different to you now:

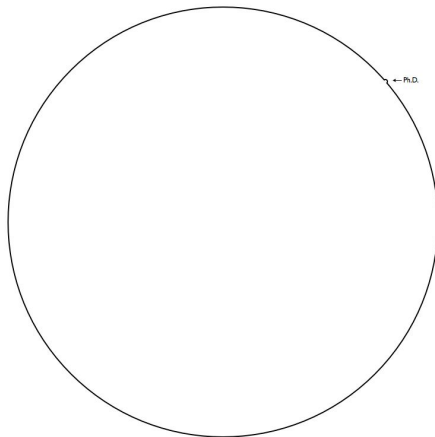


Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# What is a PhD?

So, don't forget the bigger picture:



Credit: Matt Might

<http://matt.might.net/articles/phd-school-in-pictures/>

# Goals of this talk

- What is Grad school?
- Modern sequential decision making problems
- Whet your appetite with a cool problem
- Give some advice: how to prepare for grad school

# Big Data

- Tons of data:  $2.5 \times 10^{18}$  bytes a day <sup>1</sup>
- 90% of the world's data created in the last 2 years <sup>1</sup>
- Google: 100 billion searches a month, half from mobile <sup>2</sup>

---

<sup>1</sup><http://www-01.ibm.com/software/data/bigdata/what-is-big-data.html>

<sup>2</sup>[http://blogs.wsj.com/digits/2015/10/08/](http://blogs.wsj.com/digits/2015/10/08/google-says-mobile-searches-surpass-those-on-pcs/)

# Big Data

- Tons of data:  $2.5 \times 10^{18}$  bytes a day <sup>1</sup>
- 90% of the world's data created in the last 2 years <sup>1</sup>
- Google: 100 billion searches a month, half from mobile <sup>2</sup>
- Personalization: Andrew Ng's Coursera example



<sup>1</sup><http://www-01.ibm.com/software/data/bigdata/what-is-big-data.html>






<sup>2</sup><http://blogs.wsj.com/digits/2015/10/08/>

[google-says-mobile-searches-surpass-those-on-pcs/](http://blogs.wsj.com/digits/2015/10/08/google-says-mobile-searches-surpass-those-on-pcs/)

## Other side: Interactions

- Ubiquity of devices

---

<sup>3</sup>Pouchter, Jacob “Smartphone Ownership and Internet Usage Continues to Climb in Emerging Economies, Pew Research Center”. Pewglobal.org. Retrieved 2016-02-23.     



## Other side: Interactions

- Ubiquity of devices
- Smartphone adoption rates:<sup>3</sup>

Rank	Country	% adoption
1	South Korea	88
2	Australia	77
3	Israel	74
4	United States	72
5	Spain	71

---

<sup>3</sup>Pouchter, Jacob "Smartphone Ownership and Internet Usage Continues to Climb in Emerging Economies, Pew Research Center". Pewglobal.org. Retrieved 2016-02-23

## Other side: Interactions

- Ubiquity of devices
- Smartphone adoption rates:<sup>3</sup>

Rank	Country	% adoption
1	South Korea	88
2	Australia	77
3	Israel	74
4	United States	72
5	Spain	71

- Sequential Interaction

---

<sup>3</sup>Pouchter, Jacob "Smartphone Ownership and Internet Usage Continues to Climb in Emerging Economies, Pew Research Center". Pewglobal.org. Retrieved 2016-02-23

## Other side: Interactions

- Ubiquity of devices
- Smartphone adoption rates:<sup>3</sup>

Rank	Country	% adoption
1	South Korea	88
2	Australia	77
3	Israel	74
4	United States	72
5	Spain	71

- Sequential Interaction
- Thought experiment: how did people arrange to meet before cell phones?

---

<sup>3</sup>Pouchter, Jacob "Smartphone Ownership and Internet Usage Continues to Climb in Emerging Economies, Pew Research Center". Pewglobal.org. Retrieved 2016-02-23

# Sequential Decision Making - running a newspaper

# Sequential Decision Making - running a newspaper

- First problem: which headline to choose?



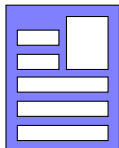
# Sequential Decision Making - running a newspaper

- First problem: which headline to choose?



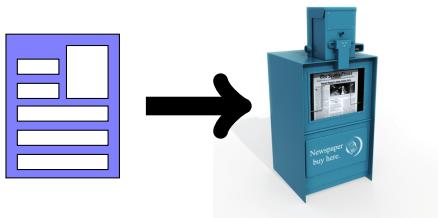
- Goal: pick the headline that sells the best

# Print Newspapers



1. Choose headline

# Print Newspapers



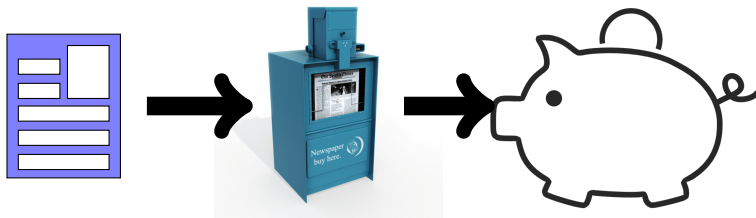
1. Choose headline

2. Sell paper

Too late: by the time feedback comes, your headlines are stale.



# Print Newspapers



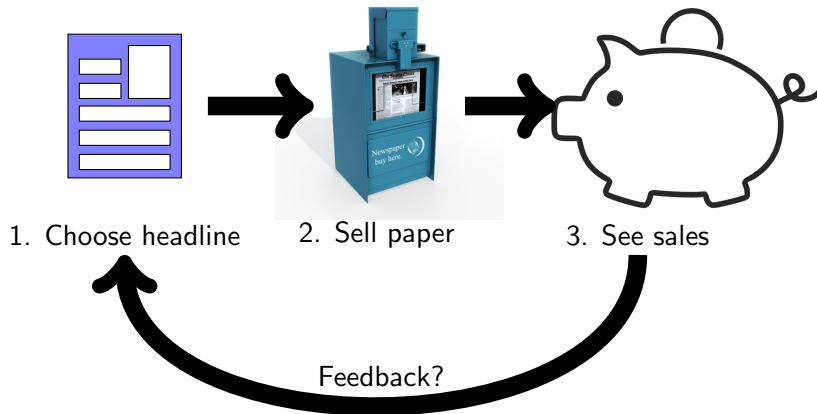
1. Choose headline

2. Sell paper

3. See sales

Too late: by the time feedback comes, your headlines are stale.

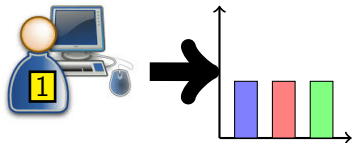
# Print Newspapers



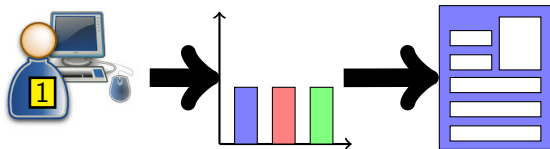
Too late: by the time feedback comes, your headlines are stale.



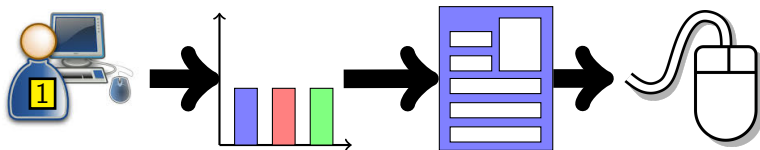
1. User arrives



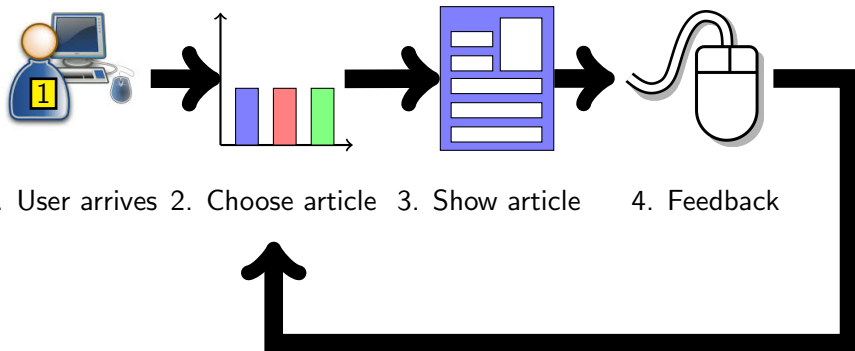
1. User arrives 2. Choose article

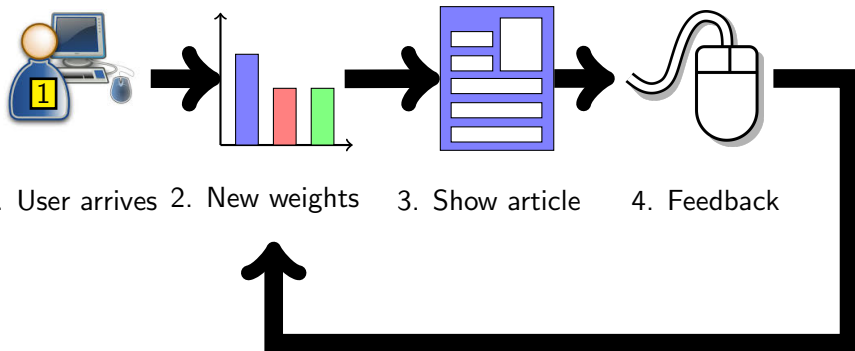


1. User arrives 2. Choose article 3. Show article

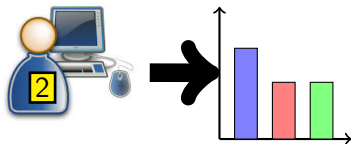


1. User arrives   2. Choose article   3. Show article   4. Feedback

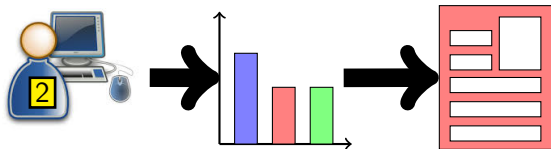




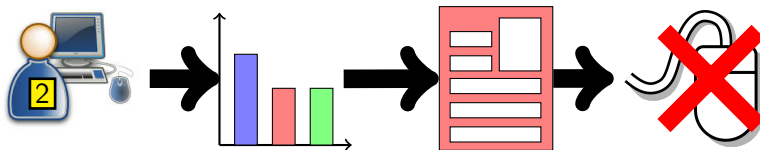




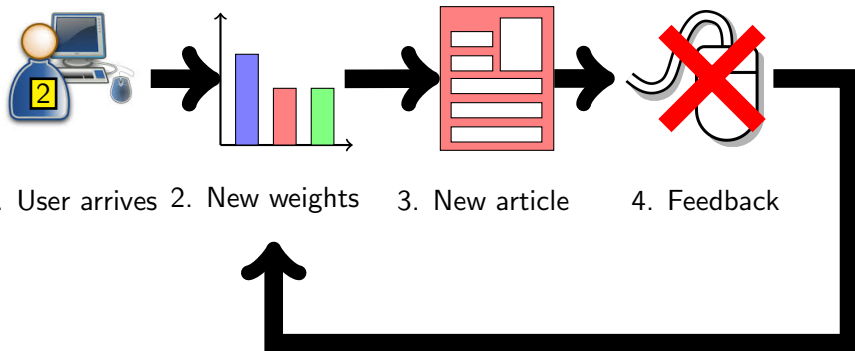
1. User arrives 2. New weights

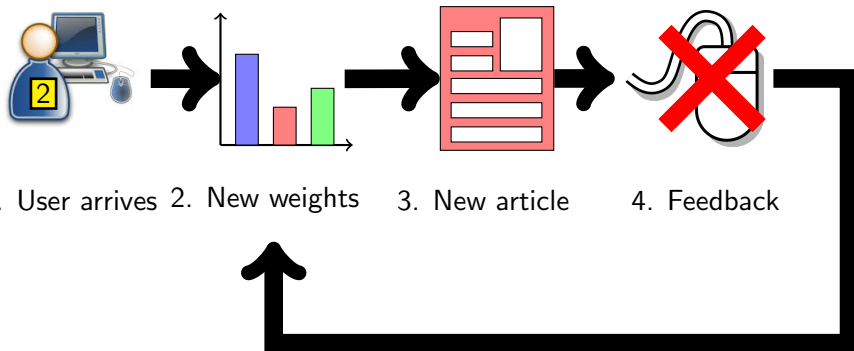


1. User arrives    2. New weights    3. New article



1. User arrives    2. New weights    3. New article    4. Feedback





# Sequential problem

- Unlike newspaper stand, we get repeated feedback
- Can modify our choices in *real time*

**exploration**

vs.

**exploitation**

# Sequential problem

- Unlike newspaper stand, we get repeated feedback
- Can modify our choices in *real time*

**exploration**

vs.

**exploitation**

Continue to explore headlines until we are very sure we haven't missed the best one.

# Sequential problem

- Unlike newspaper stand, we get repeated feedback
- Can modify our choices in *real time*

## **exploration**

Continue to explore headlines until we are very sure we haven't missed the best one.

vs.

## **exploitation**

Use the headline that has been the best so far.



# Sequential problem

- Unlike newspaper stand, we get repeated feedback
- Can modify our choices in *real time*

## exploration

Continue to explore headlines until we are very sure we haven't missed the best one.

vs.

## exploitation

Use the headline that has been the best so far.

- How can we formalize this?

# Multi-Armed Bandit

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$  :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$

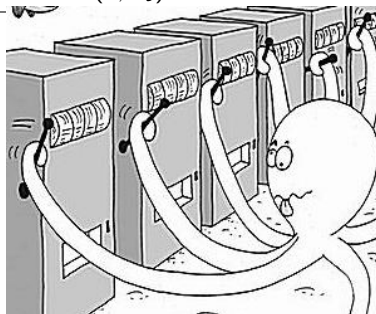
Figure: game protocol

# Multi-Armed Bandit

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$ :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$



# Multi-Armed Bandit

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$  :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$

Figure: game protocol

Learner does not  
see rewards for  
other actions

# Multi-Armed Bandit

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$  :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$

Learner needs to  
randomize

Figure: game protocol

Learner does not  
see rewards for  
other actions

# Multi-Armed Bandit

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$ :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$

Learner needs to  
randomize

Figure: game protocol

Learner does not  
see rewards for  
other actions

$$\text{Regret}_T = \underbrace{\max_{k'} \sum_{t=1}^T r(t, k')}_{\text{reward of best arm}} - \underbrace{\sum_{t=1}^T r(t, k_t)}_{\text{Learner's reward}}$$

# Simple problem, but already interesting

Given: game length  $T$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots, T$  :

- Adversary chooses rewards  $r(t, k) \in [0, 1]$
- Learner chooses an arm  $k_t$
- Learner gets reward  $r(t, k_t)$

Figure: game protocol

- Exploitation vs. exploration
- Stochastic vs. Adversarial data

# Stochastic vs. Adversarial data

- Stochastic

- Most of statistics, machine learning
- Assume that data are generated from some random variable
- E.g. arm  $k$  has mean  $\mu_k$  and

$$r(t, k) = \begin{cases} 1 & \text{with probability } \mu_k \\ 0 & \text{with probability } 1 - \mu_k. \end{cases}$$

- Past data and future data are similar



# Stochastic vs. Adversarial data

- Stochastic

- Most of statistics, machine learning
- Assume that data are generated from some random variable
- E.g. arm  $k$  has mean  $\mu_k$  and

$$r(t, k) = \begin{cases} 1 & \text{with probability } \mu_k \\ 0 & \text{with probability } 1 - \mu_k. \end{cases}$$

- Past data and future data are similar

- Adversarial

- No assumptions on  $r(t, k)$ ; much harder
- Adversary could choose  $r(t, k)$  based on your choices for time  $1, \dots, t - 1$  to make you do badly
- Need regret; minimizing loss is hopeless

$$\text{Regret}_T = \underbrace{\max_{k'} \sum_{t=1}^T r(t, k')}_{\text{reward of best arm}} - \underbrace{\sum_{t=1}^T r(t, k_t)}_{\text{Learner's reward}}$$

# Solutions to the bandit problem

- $\epsilon$ -greedy (warm-up)
- EXP3

Given: parameter  $\epsilon$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots$ :

- With probability  $\epsilon$ , pick arm  $k_t$  uniform at random
- Otherwise, pick  $k_t = \operatorname{argmax}_k \sum_{s=1}^{t-1} r(s, k)$

Given: parameter  $\epsilon$ , arms  $1, \dots, K$

For  $t = 1, 2, \dots$ :

- With probability  $\epsilon$ , pick arm  $k_t$  uniform at random
- Otherwise, pick  $k_t = \operatorname{argmax}_k \sum_{s=1}^{t-1} r(s, k)$

- Easy to defeat with adversarial data
- Not even optimal for stochastic data

# EXP3 algorithm

Given: parameter  $\eta$ , arms  $1, \dots, K$

Set:  $w_i(1) = 1$  for  $i = 1, \dots, K$

For  $t = 1, 2, \dots$ :

- Set  $p_i(t) = \frac{w_i(t)}{\sum_{j=1}^K w_j(t)}$   $i = 1, \dots, K$
- Draw  $k_t$  randomly proportional to  $p_1(t), \dots, p_K(t)$
- Get reward  $r(t, k_t)$
- For  $i = 1, \dots, K$  set

$$\hat{r}(t, i) = \begin{cases} r(t, j)/p_j(t) & \text{if } i = k_t \\ 0 & \text{otherwise,} \end{cases}$$

$$w_j(t+1) = w_j(t) \exp(\eta \hat{r}(t, i)) = \exp\left(\eta \sum_{s=1}^t \hat{r}(s, i)\right)$$

## EXP3 (weak) regret theorem

### Theorem

*The EXP3 algorithm has the following bound:*

$$\begin{aligned}\mathbb{E}[\text{Regret}_T] &= \max_{k'} \sum_{t=1}^T r(t, k') - \sum_{t=1}^T \mathbb{E}[r(t, k_t)] \\ &\leq \frac{\eta TK}{2} + \frac{\log(K)}{\eta}.\end{aligned}$$

*If we tune  $\gamma$  for  $T$ ,*

$$\mathbb{E}[\text{Regret}_T] \leq \sqrt{2KT \log(K)}.$$

## EXP3 (weak) regret theorem

### Theorem

*The EXP3 algorithm has the following bound:*

$$\begin{aligned}\mathbb{E}[\text{Regret}_T] &= \max_{k'} \sum_{t=1}^T r(t, k') - \sum_{t=1}^T \mathbb{E}[r(t, k_t)] \\ &\leq \frac{\eta TK}{2} + \frac{\log(K)}{\eta}.\end{aligned}$$

*If we tune  $\gamma$  for  $T$ ,*

*Regret per round  $\rightarrow 0$*

$$\mathbb{E}[\text{Regret}_T] \leq \sqrt{2KT \log(K)}.$$

## EXP3 (weak) regret theorem

### Theorem

The EXP3 algorithm has the following bound:

$$\begin{aligned}\mathbb{E}[\text{Regret}_T] &= \max_{k'} \sum_{t=1}^T r(t, k') - \sum_{t=1}^T \mathbb{E}[r(t, k_t)] \\ &\leq \frac{\eta TK}{2} + \frac{\log(K)}{\eta}.\end{aligned}$$

If we tune  $\gamma$  for  $T$ ,

$$\mathbb{E}[\text{Regret}_T] \leq \sqrt{2KT \log(K)}.$$

Regret per round  $\rightarrow 0$

lower bound of  $\Omega(\sqrt{KT})$



# Recap

- Started with real world problem
- Abstracted into Multi-Armed Bandit framework
- Proposed algorithms
- Proved upper bounds on their regret
- Compared to lower bounds

# Goals of this talk

- What is Grad school?
- Modern sequential decision making problems
- Whet your appetite with a cool problem
- Give some advice: how to prepare for grad school

# Things I wish I knew in college

# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads

# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads
- Go to their office hours, group meetings, ask about problems

# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads
- Go to their office hours, group meetings, ask about problems
- If research might be for you, get involved early (sophomore)

# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads
- Go to their office hours, group meetings, ask about problems
- If research might be for you, get involved early (sophomore)
- You will need three good letters to get a good grad school

# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads
- Go to their office hours, group meetings, ask about problems
- If research might be for you, get involved early (sophomore)
- You will need three good letters to get a good grad school
- Sometimes you will learn more from a research project than an extra class



# Things I wish I knew in college

- Your professors want to talk to you and meet undergrads
- Go to their office hours, group meetings, ask about problems
- If research might be for you, get involved early (sophomore)
- You will need three good letters to get a good grad school
- Sometimes you will learn more from a research project than an extra class
- *Failing is part of it!*

# Is a PhD for you?

Pros:

- Research is also rewarding; occasionally fun
- Very independent
- Become an expert in something!
- Exciting problems

# Is a PhD for you?

## Pros:

- Research is also rewarding; occasionally fun
- Very independent
- Become an expert in something!
- Exciting problems

## Cons:

- Research is frustrating; many more failed attempts
- Long hours, little pay
- Few jobs require a PhD
- Don't escape politics

Thank you!