

Learning through Experimentation

CS246: Mining Massive Datasets
Jure Leskovec, Stanford University
<http://cs246.stanford.edu>



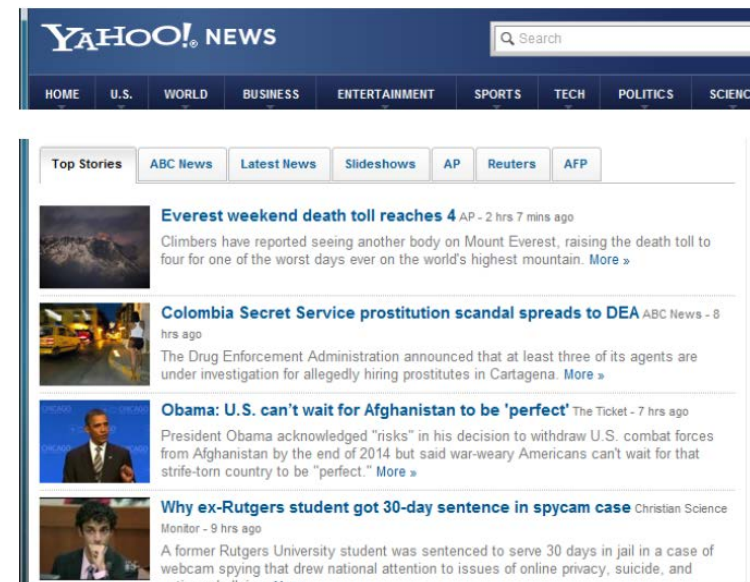
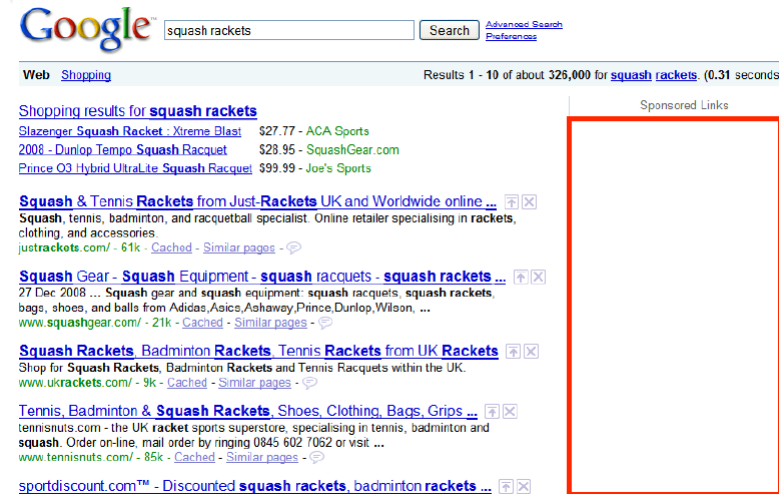
Learning through Experimentation

■ Web advertising

- We discussed how to match advertisers to queries in real-time
- But we did not discuss how to estimate the **CTR** (Click-Through Rate)

■ Recommendation engines

- We discussed how to build recommender systems
- But we did not discuss the **cold-start** problem



Learning through Experimentation

- What do **CTR** and **cold start** have in common?
- With every **ad we show/product we recommend** we gather more data about the **ad/product**
- Theme: Learning through experimentation

A screenshot of a Google search results page for the query "squash rackets". The search bar at the top shows the query and a "Search" button. Below the search bar, there are links for "Web" and "Shopping". The results section shows several organic search results, including links to "Squash & Tennis Rackets from Just-Rackets UK and Worldwide online...", "Squash Gear - Squash Equipment - squash racquets - squash rackets...", and "Squash Rackets, Badminton Rackets, Tennis Rackets from UK Rackets". To the right of the organic results, there is a red rectangular box labeled "Sponsored Links" which is currently empty.

A screenshot of a Yahoo! News page. The header shows the "YAHOO! NEWS" logo and a search bar. Below the header, there are navigation links for "HOME", "U.S.", "WORLD", "BUSINESS", "ENTERTAINMENT", "SPORTS", "TECH", "POLITICS", and "SCIENCE". The main content area is titled "Top Stories" and features four news items, each with a small image and a headline. The first item is "Everest weekend death toll reaches 4 AP - 2 hrs 7 mins ago". The second is "Colombia Secret Service prostitution scandal spreads to DEA ABC News - 8 hrs ago". The third is "Obama: U.S. can't wait for Afghanistan to be 'perfect' The Ticket - 7 hrs ago". The fourth is "Why ex-Rutgers student got 30-day sentence in spycam case Christian Science Monitor - 9 hrs ago".

Example: Web Advertising

- **Google's goal: Maximize revenue**
- **The old way: Pay by impression (CPM)**
 - **Best strategy: Go with the highest bidder**
 - But this ignores “effectiveness” of an ad
- **The new way: Pay per click! (CPC)**
 - **Best strategy: Go with expected revenue**
 - What's the expected revenue of ad a for query q ?
 - $E[\text{revenue}_{a,q}] = P(\text{click}_a \mid q) * \text{amount}_{a,q}$

Prob. user will click on ad a given
that she issues query q
(Unknown! Need to gather information)

Bid amount for
ad a on query q
(Known)

Other Applications

- **Clinical trials:**

- Investigate effects of different treatments while minimizing patient losses

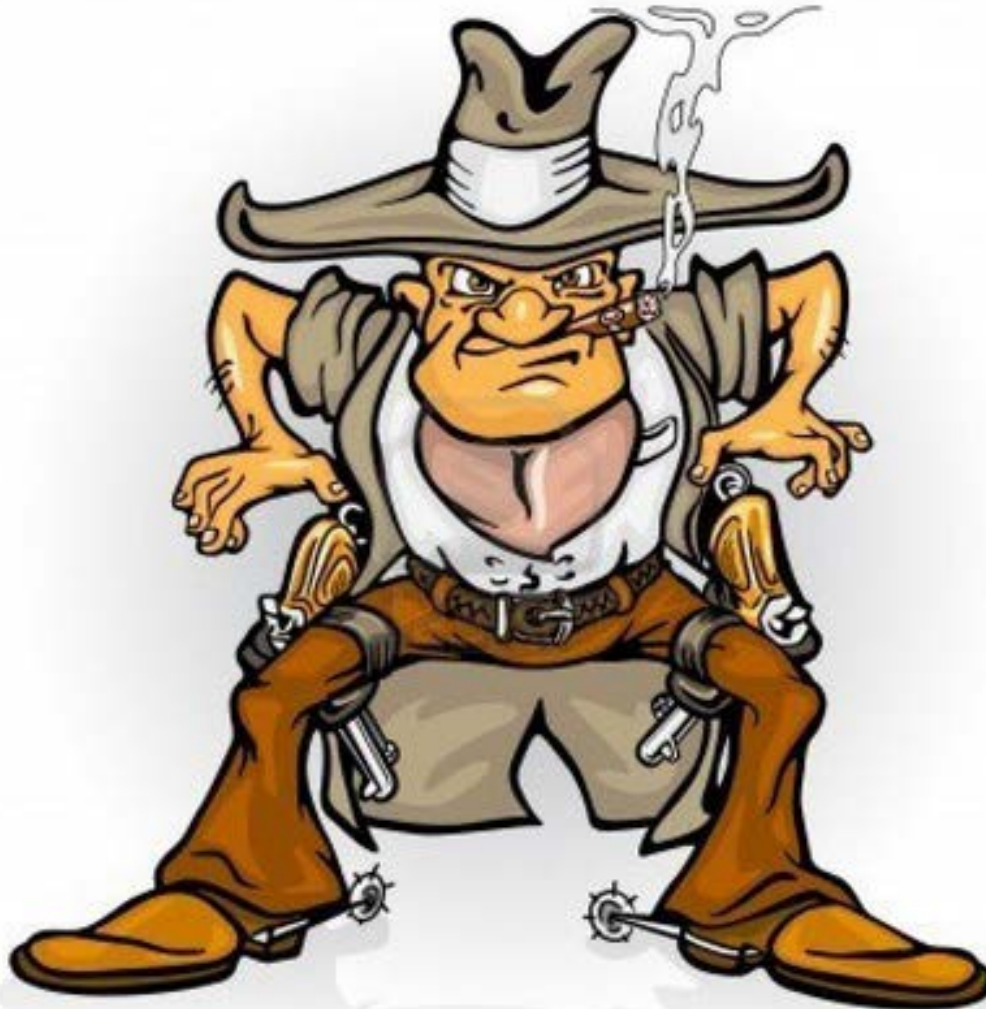
- **Adaptive routing:**

- Minimize delay in the network by investigating different routes

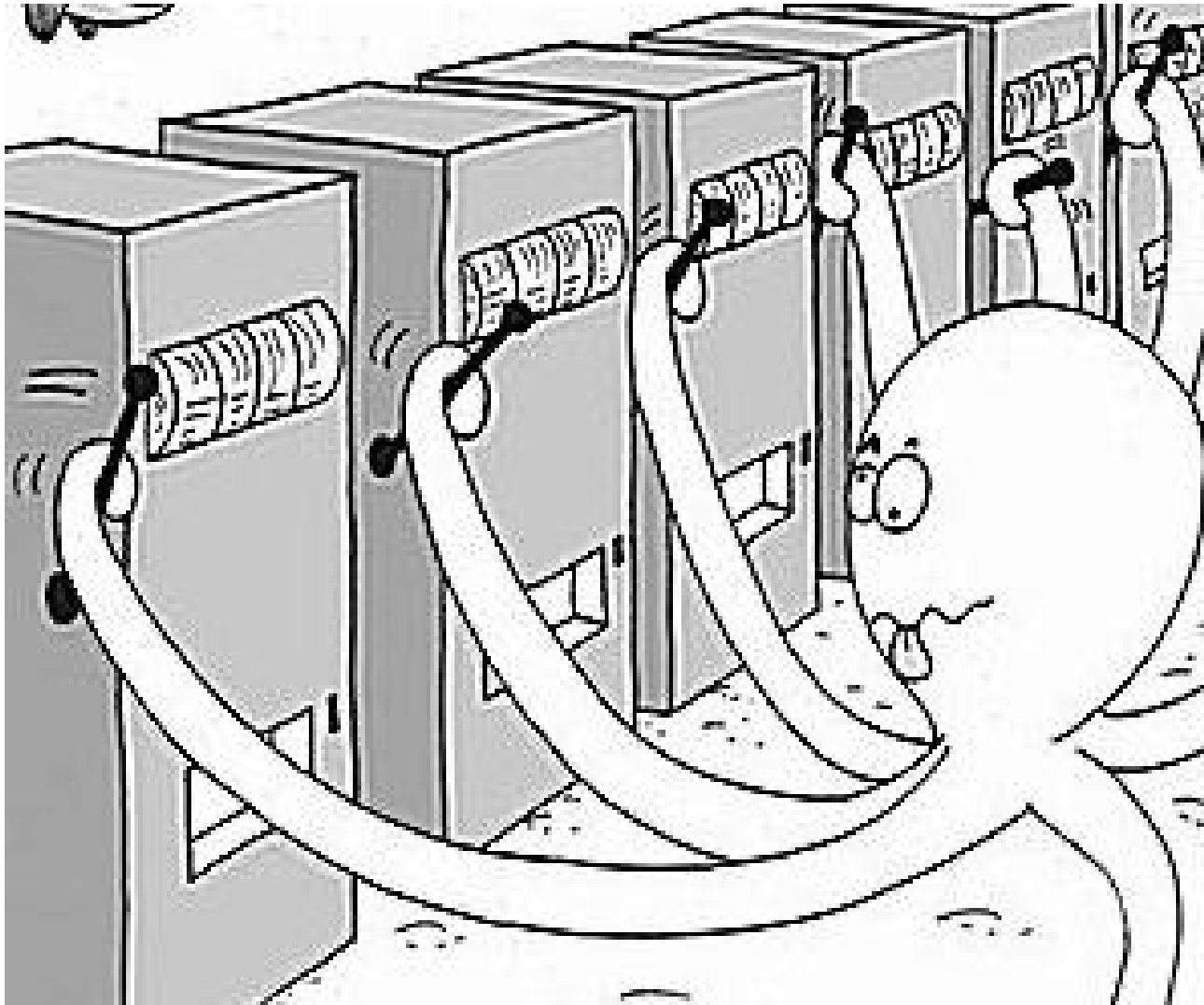
- **Asset pricing:**

- Figure out product prices while trying to make most money

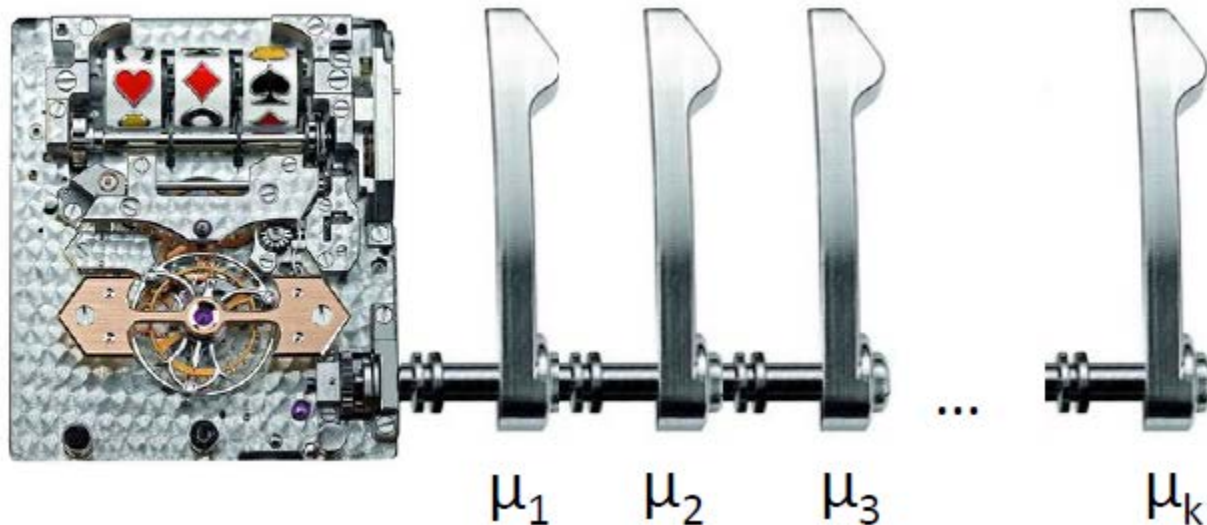
Approach: Bandits



Approach: Multiarmed Bandits

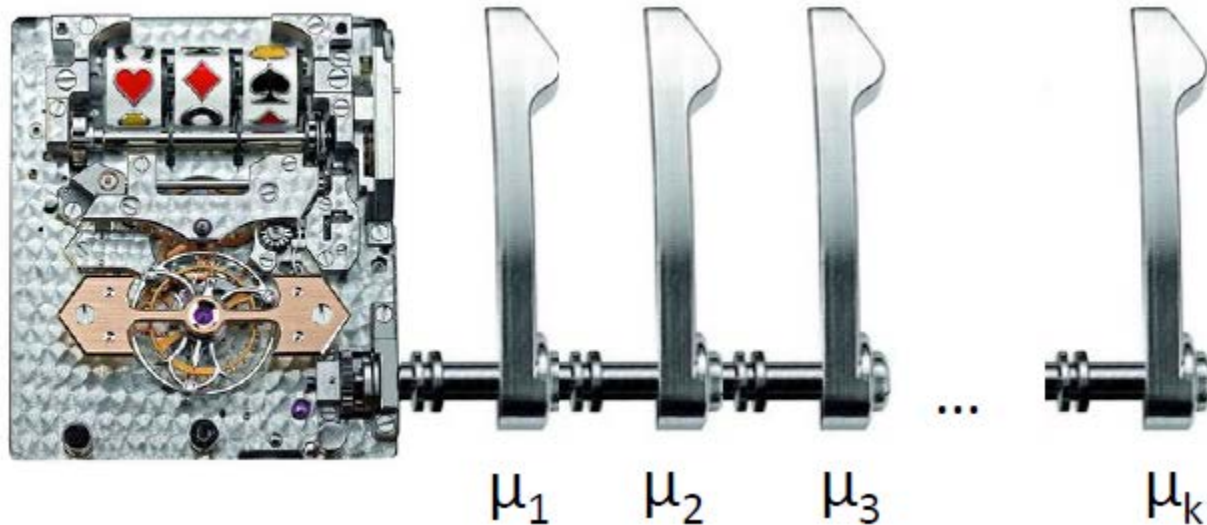


k-Armed Bandit



- **Each arm a**
 - **Wins** (reward=1) with fixed (unknown) prob. μ_a
 - **Loses** (reward=0) with fixed (unknown) prob. $1-\mu_a$
- All draws are independent given $\mu_1 \dots \mu_k$
- **How to pull arms to maximize total reward?**

k-Armed Bandit



- How does this map to our setting?
- Each **query** is a **bandit**
- Each **ad** is an **arm**
- We want to estimate the arm's probability of winning μ_a (i.e., ad's the CTR μ_a)
- Every time we pull an arm we do an 'experiment'

Stochastic k-Armed Bandit

The setting:

- Set of k choices (arms)
- Each choice a is associated with unknown probability distribution P_a supported in $[0,1]$
- We play the game for T rounds
- In each round t :
 - (1) We pick some arm j
 - (2) We obtain random sample X_t from P_j
 - Note reward is independent of previous draws
- Our goal is to maximize $\sum_{t=1}^T X_t$
- But we don't know μ_a ! But every time we pull some arm a we get to learn a bit about μ_a

Online Optimization

■ Online optimization with limited feedback

Choices	X_1	X_2	X_3	X_4	X_5	X_6	...
a_1					1	1	
a_2	0		1	0			
...							
a_k		0					

Time →

■ Like in online algorithms:

- Have to make a choice each time
- But we only receive information about the chosen action

Solving the Bandit Problem

- **Policy:** a strategy/rule that in each iteration tells me which arm to pull
 - Hopefully policy depends on the history of rewards
- **How to quantify performance of the algorithm? Regret!**

Performance Metric: Regret

- Let be μ_a the mean of P_a
- Payoff/reward of **best arm**: $\mu^* = \max_a \mu_a$
- Let $i_1, i_2 \dots i_T$ be the sequence of arms pulled
- Instantaneous **regret** at time t : $r_t = \mu^* - \mu_{a_t}$
- **Total regret**:

$$R_T = \sum_{t=1}^T r_t$$

- Typical goal: **Want a policy (arm allocation strategy) that guarantees: $\frac{R_T}{T} \rightarrow 0$ as $T \rightarrow \infty$**
 - Note: Ensuring $R_T/T \rightarrow 0$ is stronger than maximizing payoffs (minimizing regret), as it means that in the limit we discover the true best hand.

Allocation Strategies

- If we knew the payoffs, which arm would we pull?

Pick $\arg \max_a \mu_a$

- What if we only care about estimating payoffs μ_a ?

- Pick each of k arms equally often: $\frac{T}{k}$

- Estimate: $\widehat{\mu}_a = \frac{k}{T} \sum_{j=1}^{T/k} X_{a,j}$

- Regret: $R_T = \frac{T}{k} \sum_a^k (\mu^* - \mu_a)$

$X_{a,j} \dots$ payoff received
when pulling arm a for
 j -th time

Bandit Algorithm: First try

- Regret is defined in terms of average reward
- So, if we can estimate avg. reward we can minimize regret
- Consider algorithm: *Greedy*
Take the action with the highest avg. reward
 - **Example:** Consider 2 actions
 - **A1** reward 1 with prob. 0.3
 - **A2** has reward 1 with prob. 0.7
 - Play **A1**, get reward 1
 - Play **A2**, get reward 0
 - Now avg. reward of **A1** will never drop to 0, and we will never play action **A2**

Exploration vs. Exploitation

- The example illustrates a classic problem in **decision making**:
 - We need to trade off **exploration** (gathering data about arm payoffs) and **exploitation** (making decisions based on data already gathered)
- **The Greedy does not explore sufficiently**
 - **Exploration**: Pull an arm we never pulled before
 - **Exploitation**: Pull an arm a for which we currently have the highest estimate of μ_a

Optimism

- The problem with our **Greedy** algorithm is that it is **too certain** in the estimate of μ_a
 - When we have seen a single reward of 0 we shouldn't conclude the average reward is 0
- **Greedy does not explore sufficiently!**

New Algorithm: Epsilon-Greedy

Algorithm: Epsilon-Greedy

- **For $t=1:T$**

- Set $\varepsilon_t = O(1/t)$
- **With prob. ε_t : Explore** by picking an arm chosen uniformly at random
- **With prob. $1 - \varepsilon_t$: Exploit** by picking an arm with highest empirical mean payoff

- **Theorem [Auer et al. '02]**

For suitable choice of ε_t it holds that

$$R_T = O(k \log T) \Rightarrow \frac{R_T}{T} = O\left(\frac{k \log T}{T}\right) \rightarrow 0$$

k ...number
of arms

Issues with Epsilon Greedy

- What are some issues with **Epsilon Greedy**?
 - “**Not elegant**”: Algorithm explicitly distinguishes between exploration and exploitation
 - **More importantly**: Exploration makes **suboptimal choices** (since it picks any arm equally likely)
- **Idea**: When exploring/exploiting we need to **compare** arms

Comparing Arms

- **Suppose we have done experiments:**
 - Arm 1: 1 0 0 1 1 0 0 1 0 1
 - Arm 2: 1
 - Arm 3: 1 1 0 1 1 1 0 1 1 1
- **Mean arm values:**
 - Arm 1: 5/10, Arm 2: 1, Arm 3: 8/10
- **Which arm would you pick next?**
- **Idea: Don't just look at the mean (that is, expected payoff) but also the confidence!**

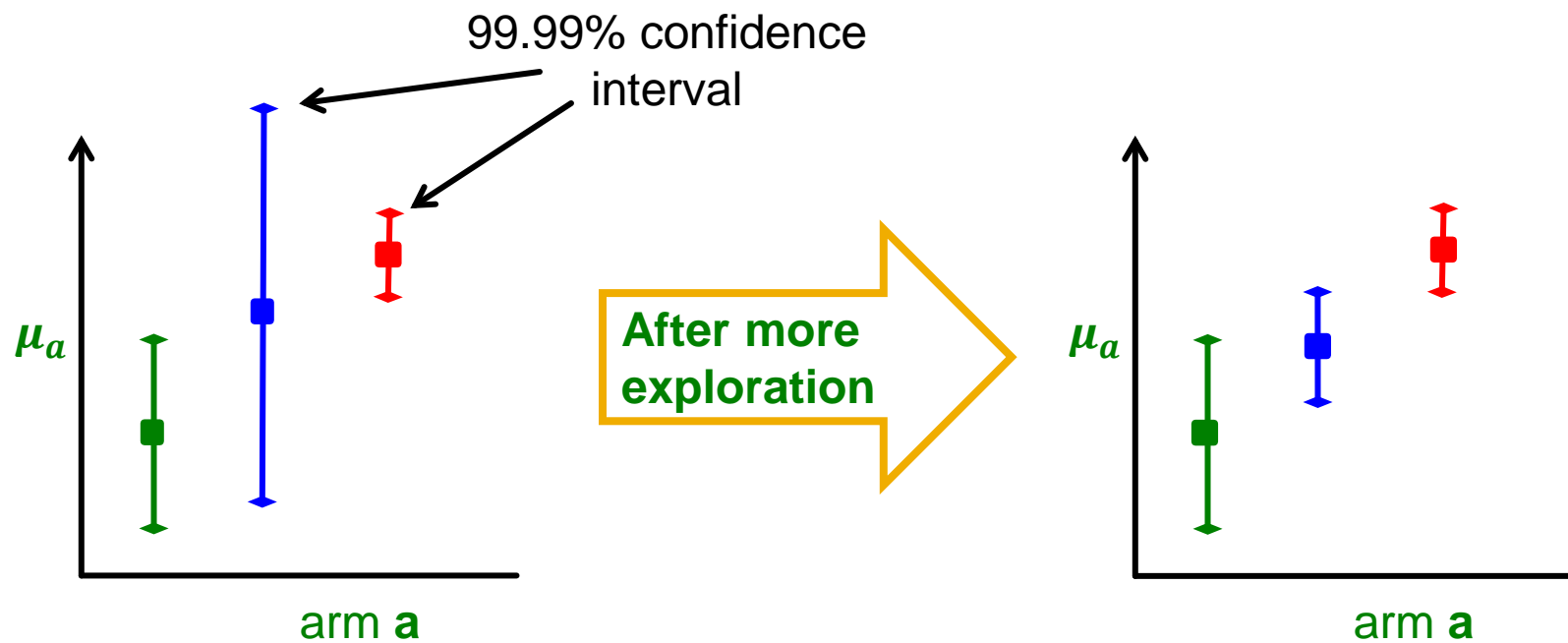
Confidence Intervals (1)

- A confidence interval is a range of values within which we are sure the mean lies with a certain probability
 - We could believe μ_a is within $[0.2, 0.5]$ with probability 0.95
 - If we would have tried an action less often, our estimated reward is less accurate so the confidence interval is larger
 - Interval shrinks as we get more information (try the action more often)

Confidence Intervals (2)

- Assuming we know the confidence intervals
- Then, instead of **trying the action with the highest mean** we can **try the action with the highest upper bound on its confidence interval**
- This is called an **optimistic policy**
 - We believe an action is as good as possible given the available evidence

Confidence Based Selection



Calculating Confidence Bounds

Suppose we fix arm a :

- Let $Y_{a,1} \dots Y_{a,m}$ be the payoffs of arm a in the first m trials
 - So, $Y_{a,1} \dots Y_{a,m}$ are i.i.d. rnd. vars. taking values in $[0,1]$
- Mean payoff of arm a : $\mu_a = E[Y_{a,m}]$
- Our estimate: $\widehat{\mu}_{a,m} = \frac{1}{m} \sum_{\ell=1}^m Y_{a,\ell}$
- Want to find b such that with high probability $|\mu_a - \widehat{\mu}_{a,m}| \leq b$
 - Also want b to be as small as possible (why?)
- Goal: Want to bound $P(|\mu_a - \widehat{\mu}_{a,m}| \leq b)$

Hoeffding's Inequality

- **Hoeffding's inequality bounds** $\mathbf{P}(|\mu_a - \widehat{\mu}_{a,m}| \leq \mathbf{b})$
 - Let $X_1 \dots X_m$ be **i.i.d.** rnd. vars. taking values in **[0,1]**
 - Let $\mu = E[X]$ and $\widehat{\mu}_m = \frac{1}{m} \sum_{\ell=1}^m X_\ell$
 - **Then:** $\mathbf{P}(|\mu - \widehat{\mu}_m| \geq \mathbf{b}) \leq 2 \exp(-2b^2m) = \delta$
- **To find out the confidence interval \mathbf{b} (for a given confidence level δ) we solve:**
 - $2e^{-2b^2m} \leq \delta$ **then** $-2b^2m \leq \ln(\delta/2)$
 - **So:** $\mathbf{b} \geq \sqrt{\frac{\ln(\frac{2}{\delta})}{2m}}$

UCB₁ Algorithm

■ UCB₁ (Upper confidence sampling) algorithm

- Set: $\widehat{\mu}_1 = \dots = \widehat{\mu}_k = 0$ and $m_1 = \dots = m_k = 0$

- $\widehat{\mu}_a$ is our estimate of payoff of arm i
- m_a is the number of pulls of arm i so far

- For $t = 1:T$

- For each arm a calculate: $UCB(a) = \widehat{\mu}_a + \alpha \sqrt{\frac{2 \ln t}{m_a}}$
- Pick arm $j = \arg \max_a UCB(a)$
- Pull arm j and observe y_t
- Set: $m_j \leftarrow m_j + 1$ and $\widehat{\mu}_j \leftarrow \frac{1}{m_j} (y_t + (m_j - 1) \widehat{\mu}_j)$

Upper confidence
interval (Hoeffding's
inequality)



UCB₁: Discussion

■ $UCB(a) = \widehat{\mu}_a + \alpha \sqrt{\frac{2 \ln t}{m_a}}$

$$b \geq \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2m}}$$

- Confidence interval **grows** with the total number of actions t we have taken
- But **shrinks** with the number of times m_a we have tried arm a
- This ensures each arm is tried infinitely often but still balances exploration and exploitation

- α plays the role of δ : $\alpha = f\left(\frac{2}{\delta}\right)$

$$\alpha = 1 + \sqrt{\ln(2/\delta)/2}$$

“Optimism in face of uncertainty”:

The algorithm believes that it can obtain extra rewards by reaching the unexplored parts of the state space

Performance of UCB₁

■ Theorem [Auer et al. 2002]

- Suppose optimal mean payoff is $\mu^* = \max_a \mu_a$
- And for each arm let $\Delta_a = \mu^* - \mu_a$
- Then it holds that

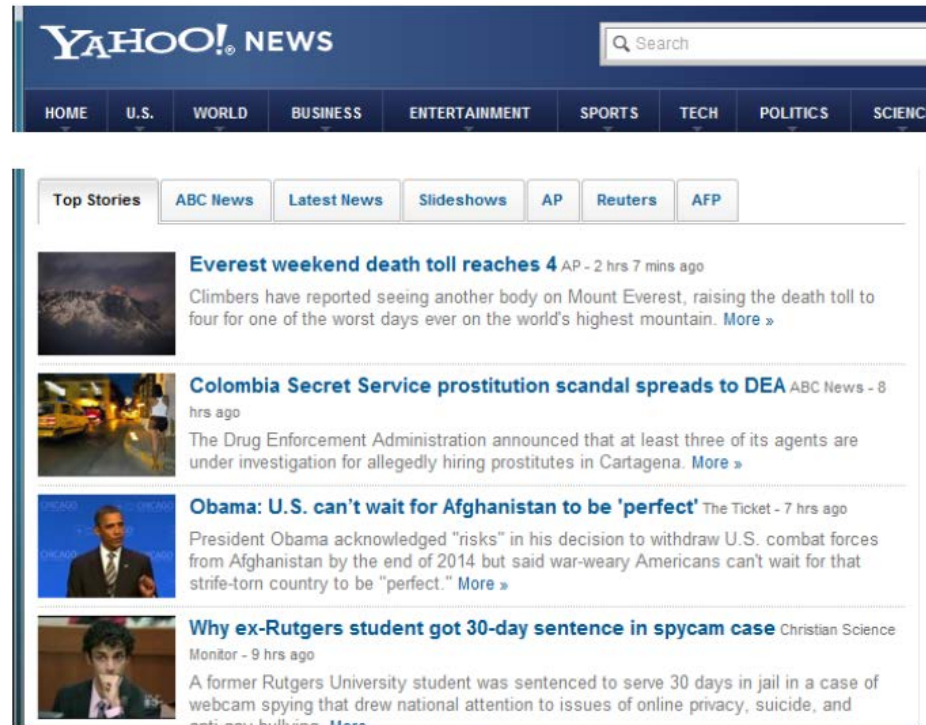
$$E[R_T] = \underbrace{\left[8 \sum_{a: \mu_a < \mu^*} \frac{\ln T}{\Delta_a} \right]}_{O(k \ln T)} + \underbrace{\left(1 + \frac{\pi^2}{3} \right) \left(\sum_{i=a}^k \Delta_a \right)}_{O(k)}$$

- So: $O\left(\frac{R_T}{T}\right) = k \frac{\ln T}{T}$

Summary so far

- k -armed bandit problem as a formalization of the exploration-exploitation tradeoff
- Analog of online optimization (e.g., SGD, BALANCE), but with **limited feedback**
- **Simple algorithms are able to achieve no regret (in the limit)**
 - Epsilon-greedy
 - UCB (Upper Confidence Sampling)

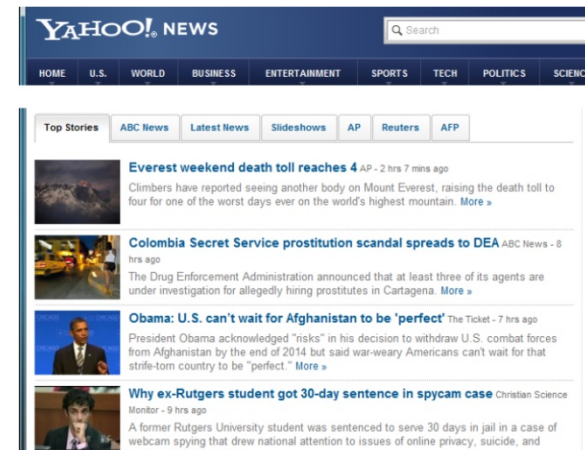
News Recommendation



- Every round receive **context** [Li et al., WWW '10]
 - **Context:** User features, articles view before
- **Model for each article's click-through rate**

News Recommendation

- **Feature-based exploration:**
 - **Select articles to serve users based on contextual information about the user and the articles**
 - **Simultaneously adapt article selection strategy based on user-click feedback to maximize total number of user clicks**



Contextual Bandits

- **Contextual bandit algorithm in round t**
 - **(1)** Algorithm observes user \mathbf{u}_t and a set \mathbf{A} of arms together with their features $\mathbf{x}_{t,a}$
 - Vector $\mathbf{x}_{t,a}$ summarizes both the user \mathbf{u}_t and arm \mathbf{a}
 - We call vector $\mathbf{x}_{t,a}$ the **context**
 - **(2)** Based on payoffs from previous trials, algorithm chooses arm $\mathbf{a} \in \mathbf{A}$ and receives payoff $r_{t,a}$
 - **Note only feedback for the chosen \mathbf{a} is observed**
 - **(3)** Algorithm improves arm selection strategy with each observation $(\mathbf{x}_{t,a}, \mathbf{a}, r_{t,a})$

LinUCB Algorithm (1)

- Payoff of arm \mathbf{a} : $E[r_{t,a}|x_{t,a}] = x_{t,a}^T \cdot \theta_a^*$
 - $\mathbf{x}_{t,a}$... d -dimensional feature vector
 - θ_a^* ... unknown coefficient vector we aim to learn
 - Note that θ_a^* are not shared between different arms!
- What's the difference between LinUCB, UCB1?
 - UCB1 directly estimates μ_a through experimentation (without any knowledge about arm \mathbf{a})
 - LinUCB estimates μ_a by regression $\mu_a = x_{t,a}^T \cdot \theta_a^*$
 - The hope is that we will be able to learn faster as we consider the context \mathbf{x}_a (user, ad) of arm \mathbf{a}

LinUCB Algorithm (2)

- Payoff of arm \mathbf{a} : $E[r_{t,a} | \mathbf{x}_{t,a}] = \mathbf{x}_{t,a}^T \cdot \boldsymbol{\theta}_a^*$
 - $\mathbf{x}_{t,a}$... d -dimensional feature vector
 - $\boldsymbol{\theta}_a^*$... unknown coefficient vector we aim to learn
- How to estimate $\boldsymbol{\theta}_a$?
 - \mathbf{D}_a ... $m \times d$ matrix of \mathbf{m} training inputs $[\mathbf{x}_{a,t}]$
 - \mathbf{b}_a ... \mathbf{m} -dim. vector of responses to \mathbf{a} (click/no-click)
 - Linear regression solution to $\boldsymbol{\theta}_a$ is then
 - $\hat{\boldsymbol{\theta}}_a = \arg \min_{\boldsymbol{\theta}} \sum_{\mathbf{m} \in \mathbf{D}_a} \left(\mathbf{x}_{t,a}^T \cdot \boldsymbol{\theta}_a - \mathbf{b}_a^{(m)} \right)^2$
 - Which is solved by: $\hat{\boldsymbol{\theta}}_a = \left(\mathbf{D}_a^T \mathbf{D}_a + \mathbf{I}_d \right)^{-1} \mathbf{D}_a^T \mathbf{b}_a$

\mathbf{I}_d is $d \times d$
identity matrix

LinUCB Algorithm (3)

- One can then show (using similar techniques as we used for UCB) that

$$\left| \mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a - \mathbf{E}[r_{t,a} | \mathbf{x}_{t,a}] \right| \leq \alpha \sqrt{\mathbf{x}_{t,a}^\top (\mathbf{D}_a^\top \mathbf{D}_a + \mathbf{I}_d)^{-1} \mathbf{x}_{t,a}}$$
$$\alpha = 1 + \sqrt{\ln(2/\delta)/2}$$

- So LinUCB arm selection rule is:

$$a_t \stackrel{\text{def}}{=} \arg \max_{a \in \mathcal{A}_t} \left(\underbrace{\mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a}_{\text{Estimated } \mu_a} + \alpha \sqrt{\underbrace{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}_{\text{Confidence interval: Standard deviation}}} \right)$$

$$\mathbf{A}_a \stackrel{\text{def}}{=} \mathbf{D}_a^\top \mathbf{D}_a + \mathbf{I}_d$$

LinUCB Algorithm (3)

Initialization:

For each arm a :

$$A_a = I_d$$

$$b_a = [0]_d$$

$$A_a \stackrel{\text{def}}{=} D_a^\top D_a + I_d$$

identity matrix $m \times m$
vector of zeros

Online algorithm:

For $t = 1, 2, 3, \dots T$:

Observe features of all arms $a : x_{t,a} \in R^d$

For each arm a :

$$\theta_a = A_a^{-1} b_a$$

regression coefficients

$$p_{t,a} = \theta_a^T x_{t,a} + \alpha \sqrt{x_{t,a}^T A_a^{-1} x_{t,a}}$$

confidence bound

Choose arm $a_t = \arg \max_a p_{t,a}$ choose arm

$$A_{a_t} = A_{a_t} + x_{t,a_t} x_{t,a_t}^T$$

update A for the chosen arm a_t

$$b_{a_t} = b_{a_t} + r_t x_{t,a_t}$$

updated b for the chosen arm a_t

LinUCB: Discussion

- **LinUCB** computational complexity is
 - **Linear** in the **number of arms** and
 - At most **cubic** in the **number of features**
- **LinUCB** works well for a **dynamic arm set** (arms come and go):
 - For example, in news article recommendation, for instance, editors add/remove articles to/from a pool

Yahoo! News Experiment

Featured | Entertainment | Sports | Life



McNair's final hours revealed

Police release 50 text messages that depict the late NFL player's alleged killer as losing control. » **Details**

- UConn murder victim mourned

 Find Steve McNair murder case

**F1** Steve McNair's final hours revealed

**F2** Cindy Crawford stays fierce in a black mini

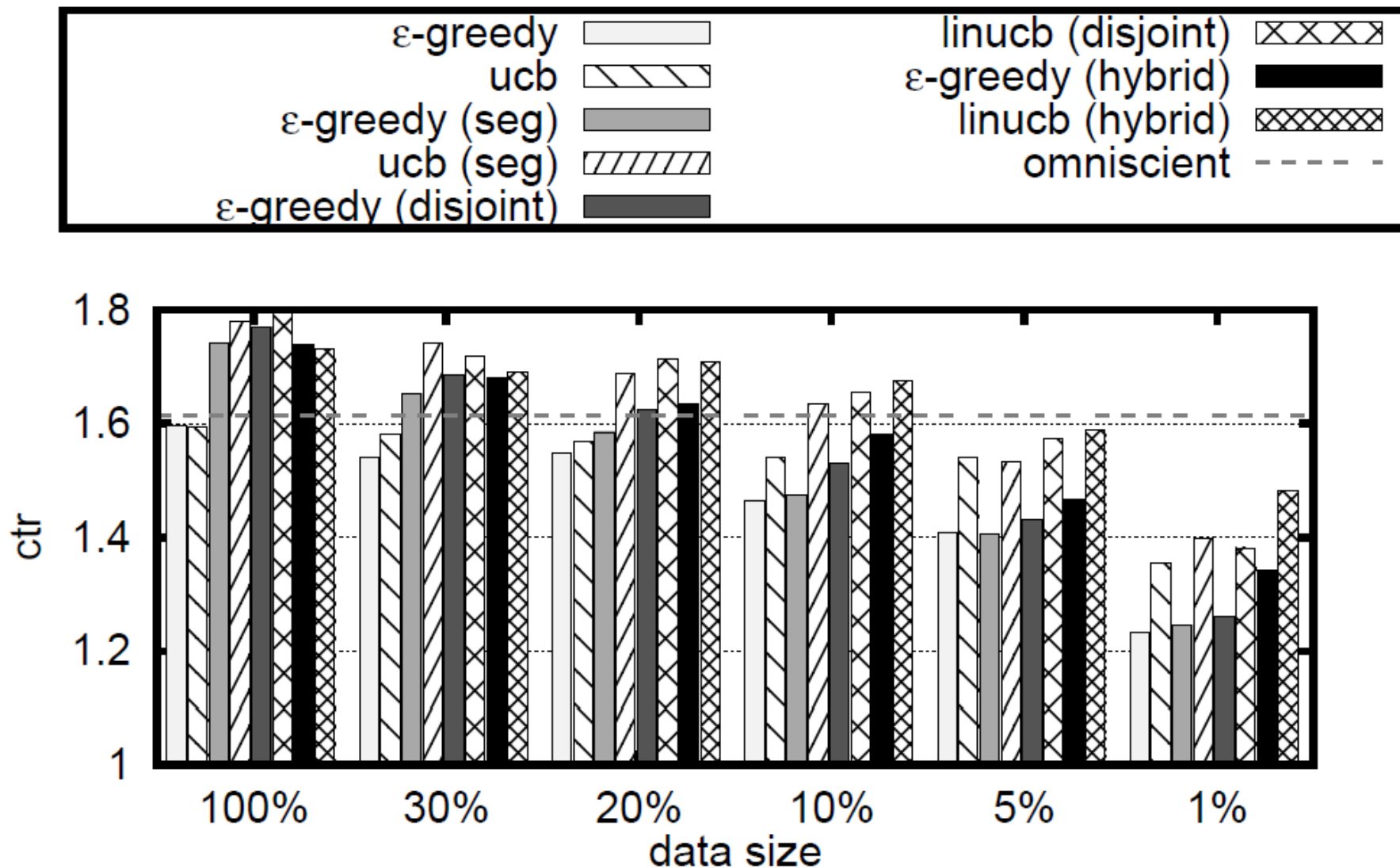
**F3** Watch for dozens of 'shooting stars' tonight

**F4** At team's big moment, star player isn't around

» More: **Featured** | **Buzz**

- What to put in slots F1, F2, F3, F4 to make the user click? Use **LinUCB**

Results



Example: A/B testing vs. Bandits

- Imagine you have two versions of the website and you'd like to test which one is better
 - Version A has engagement rate of 5%
 - Version B has engagement rate of 4%
- You want to establish with 95% confidence that version A is better
 - You'd need 22,330 observations (11,165 in each arm) to establish that
 - Use student's t-test to establish the sample size
 - Can bandits do better?

Example: Bandits vs. A/B testing

- **How long it does it take to discover $A > B$?**
 - **A/B test:** We need 22,330 observations. Assuming 100 observations/day, we need 223 days
 - **Bandits:** We use UCB1 and keep track of confidences for each version we stop as soon as A is better than B with 95% confidence.
How much do we save?
 - 175 days on the average!
 - **48 days vs. 223 days**
 - More at: <http://bit.ly/1pywka4>

