Bandits: Upper Confidence Bound Algorithm

Lucas Janson and Sham Kakade CS/Stat 184: Introduction to Reinforcement Learning Fall 2022

- Feedback from last lecture
- Recap
- Confidence intervals for the arms
- Upper Confidence Bound (UCB) algorithm
- UCB regret analysis



1. Thank you to everyone who filled out the forms!

- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**

- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**

• I will slow down! (please let me know in today's feedback how the pace is)



- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**
 - I will slow down! (please let me know in today's feedback how the pace is)
 - I will review the last 15 minutes from last lecture (which was always) supposed to be the first part of today's lecture)



- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**
 - I will slow down! (please let me know in today's feedback how the pace is)
 - I will review the last 15 minutes from last lecture (which was always) supposed to be the first part of today's lecture)
- 3. Common points of confusion:



- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**
 - I will slow down! (please let me know in today's feedback how the pace is)
 - I will review the last 15 minutes from last lecture (which was always) supposed to be the first part of today's lecture)
- 3. Common points of confusion:
 - Union bound trick (I will review this today in more detail)



- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**
 - I will slow down! (please let me know in today's feedback how the pace is)
 - I will review the last 15 minutes from last lecture (which was always) supposed to be the first part of today's lecture)
- 3. Common points of confusion:
 - Union bound trick (I will review this today in more detail)
 - Why we bound the regret with high probability instead of in expectation **II**. (I will also talk about this a bit)





- 1. Thank you to everyone who filled out the forms!
- 2. Main feedback: **Too fast!**
 - I will slow down! (please let me know in today's feedback how the pace is)
 - I will review the last 15 minutes from last lecture (which was always) supposed to be the first part of today's lecture)
- 3. Common points of confusion:
 - Union bound trick (I will review this today in more detail)
 - Why we bound the regret with high probability instead of in expectation **II**. (I will also talk about this a bit)
- 4. Use <u>bit.ly</u> and QR code for feedback form: done!







- Recap
- Confidence intervals for the arms
- Upper Confidence Bound (UCB) algorithm
- UCB regret analysis



Pure greedy and pure exploration achieve linear regret

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:
 - balance exploration with exploitation

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:
 - balance exploration with exploitation
 - Achieve sublinear regret of $\tilde{O}(T^{2/3})$

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:
 - balance exploration with exploitation
 - Achieve sublinear regret of $\tilde{O}(T^{2/3})$
 - Exploration is non-adaptive

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:
 - balance exploration with exploitation
 - Achieve sublinear regret of $\tilde{O}(T^{2/3})$
 - Exploration is non-adaptive
- Today: can we do better than a rate of $T^{2/3}$?

- Pure greedy and pure exploration achieve linear regret
- Explore-then-commit (ETC) and ε -greedy:
 - balance exploration with exploitation
 - Achieve sublinear regret of $\tilde{O}(T^{2/3})$
 - Exploration is non-adaptive
- Today: can we do better than a rate of $T^{2/3}$?
- First, review a couple points of common confusion from last lecture

in a collection X_1, \ldots, X_n ,

 $X_i \leq B(\delta)$

to one that holds uniformly over a colle $X_i \leq B_n(\delta)$

Union bound trick: want to translate a bound that holds for each random variable X

w/p
$$1 - \delta$$
, $\forall i$
ection of random variables X_1, \dots, X_n , i.e.,
) $\forall i$, w/p $1 - \delta$



in a collection X_1, \ldots, X_n ,

 $X_i \leq B(\delta)$

to one that holds uniformly over a colle $X_i \leq B_n(\delta)$

Union bound trick: want to translate a bound that holds for each random variable X

w/p
$$1 - \delta$$
, $\forall i$
ection of random variables X_1, \dots, X_n , i.e.,
) $\forall i$, w/p $1 - \delta$ What is B_n ?





in a collection X_1, \ldots, X_n ,

 $X_i \leq B(\delta)$

to one that holds *uniformly* over a colle $X_i \leq B_n(\delta)$

i.e., $\mathbb{P}(X_i \leq B(\delta)) \geq 1 - \delta, \forall i$

Union bound trick: want to translate a bound that holds for each random variable X

$$\begin{array}{ll} & \text{w/p } 1 - \delta, \quad \forall i \\ \text{ection of random variables } X_1, \dots, X_n, \text{ i.e.,} \\ &) \quad \forall i, \quad \text{w/p } 1 - \delta & \text{What is } B_n? \\ & \longrightarrow & \mathbb{P}(\forall i, X_i \leq B_n(\delta)) \geq 1 - \delta \end{array}$$







Should we bound $\operatorname{Regret}_{T}$ with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

Should we bound $\operatorname{Regret}_{T}$ with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

• High probability bound on Regret_T

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T

Should we bound $\operatorname{Regret}_{T}$ with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values

pret $_T$ but really bad Regret $_T$ values

Should we bound $\operatorname{Regret}_{T}$ with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$

pret $_T$ but really bad Regret $_T$ values

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - Bounds *typical* behavior of Regret_T Ι.
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T

Should we bound $\operatorname{Regret}_{T}$ with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T
 - iv. Leaves possibility that the behavior of $\mathbb{E}[\text{Regret}_{T}]$ is atypical

ret $_T$ but really bad Regret $_T$ values

 $egret_T$ avior of $\mathbb{E}[\mathsf{Regret}_T]$ is *atypica*

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T
 - iv. Leaves possibility that the behavior of $\mathbb{E}[\text{Regret}_{T}]$ is atypical
- E.g., $\operatorname{Regret}_{T} \sim T^{100} \times \operatorname{Bernoulli}(10^{-10})$:

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T
 - iv. Leaves possibility that the behavior of $\mathbb{E}[\text{Regret}_{T}]$ is atypical
- E.g., $\text{Regret}_{T} \sim T^{100} \times \text{Bernoulli}(10^{-10})$:

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T
 - iv. Leaves possibility that the behavior of $\mathbb{E}[\text{Regret}_{T}]$ is atypical
- E.g., $\text{Regret}_{T} \sim T^{100} \times \text{Bernoulli}(10^{-10})$:
 - Regret_T = 0 w/p 0.99999999999
 - $\mathbb{E}[\text{Regret}_T] = 10^{-10} T^{100}$ ($\approx 10^{20}$ when T = 2)

Should we bound Regret_T with high probability (i.e., probability $\geq 1 - \delta$), or should we bound $\mathbb{E}[\operatorname{Regret}_{T}]$?

- High probability bound on Regret_T
 - i. Bounds *typical* behavior of Regret_T
 - ii. Leaves possibility of some rare but really bad Regret_T values
- Bound on $\mathbb{E}[\text{Regret}_T]$
 - iii. Bounds average behavior of Regret_T
 - iv. Leaves possibility that the behavior of $\mathbb{E}[\text{Regret}_{T}]$ is atypical
- E.g., $\text{Regret}_{T} \sim T^{100} \times \text{Bernoulli}(10^{-10})$:
 - Regret_T = 0 w/p 0.99999999999
 - $\mathbb{E}[\text{Regret}_T] = 10^{-10} T^{100}$ ($\approx 10^{20}$ when T = 2)
Feedback from last lecture



- Confidence intervals for the arms
- Upper Confidence Bound (UCB) algorithm
- UCB regret analysis



Intuition: maintain confidence intervals for mean of each arm and use them to focus exploration on most promising arms

Intuition: maintain confidence intervals for mean of each arm and use them to focus exploration on most promising arms

First: how to construct confidence intervals?

First: how to construct confidence intervals? Recall Hoeffding inequality:

$$|\hat{\mu} - \mu| \leq \sqrt{|\hat{\mu}|^2}$$

- Intuition: maintain confidence intervals for mean of each arm and use them to focus exploration on most promising arms

Sample mean of N i.i.d. samples on [0,1] satisfies

 $\sqrt{\frac{\ln(2/\delta)}{2N}} \text{ w/p } 1 - \delta$

First: how to construct confidence intervals? Recall Hoeffding inequality:

$$|\hat{\mu} - \mu| \le \sqrt{\frac{\ln(2/\delta)}{2N}} \text{ w/p } 1 - \delta$$

- Intuition: maintain confidence intervals for mean of each arm and use them to focus exploration on most promising arms

Sample mean of N i.i.d. samples on [0,1] satisfies

Worked for ETC b/c exploration phase was i.i.d., but in general the rewards from a given arm are *not* i.i.d. due to adaptivity of action selections

Constructing confidence intervals

Constructing confidence intervals

Notation:

Constructing confidence intervals

Notation:

Let $N_t^{(k)} = \sum_{\tau=k}^{t-1} 1_{\{a_{\tau}=k\}}$ be the number of times arm k is pulled before time t $\tau = 0$

Constructing confidence intervals Notation: Let $N_t^{(k)} = \sum_{k=1}^{t-1} 1_{\{a_r = k\}}$ be the number of times arm k is pulled before time t $\tau = 0$ Let $\hat{\mu}_t^{(k)} = \frac{1}{N_t^{(k)}} \sum_{\tau=0}^{t-1} 1_{\{a_\tau=k\}} r_\tau$ be the sample mean reward of arm k up to time t

Constructing confidence intervals Notation: Let $N_t^{(k)} = \sum_{k=1}^{t-1} 1_{\{a_r = k\}}$ be the number of times arm k is pulled before time t $\tau = 0$ Let $\hat{\mu}_t^{(k)} = \frac{1}{N_t^{(k)}} \sum_{\tau=0}^{t-1} 1_{\{a_\tau = k\}} r_\tau$ be the sample mean reward of arm k up to time t

So want Hoeffding to give us something like $\int \frac{\ln(2/\delta)}{2N^{(k)}} \text{ w/p } 1 - \delta$

$$\left| \hat{\mu}_t^{(k)} - \mu \right| \leq 1$$

Constructing confidence intervals Notation: Let $N_t^{(k)} = \sum_{a_r=k}^{t-1} 1_{\{a_r=k\}}$ be the number of times arm k is pulled before time t $\tau = 0$ Let $\hat{\mu}_t^{(k)} = \frac{1}{N_t^{(k)}} \sum_{\tau=0}^{t-1} 1_{\{a_\tau=k\}} r_\tau$ be the sample mean reward of arm k up to time t

So want Hoeffding to give us something like



But this is generally FALSE (unless a_t chosen very simply, like exploration phase of ETC)

The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$,

(all *arm* indexing (*k*) now in <u>superscripts;</u> <u>sub</u>scripts reserved for time index *t*) The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$,



The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t) $\hat{\mu}_{\star}^{(k)}$ is the sample mean of a random number $N_{\star}^{(k)}$ of returns



The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t)

 $\hat{\mu}_{t}^{(k)}$ is the sample mean of a random number $N_{t}^{(k)}$ of returns in general $N_{r}^{(k)}$ will depend on those returns themselves



 $\hat{\mu}_{\star}^{(k)}$ is the sample mean of a random number $N_{\star}^{(k)}$ of returns in general $N_{r}^{(k)}$ will depend on those returns themselves (i.e., how often we select arm k depends on the historical returns of arm k)

- The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t)



 $\hat{\mu}_{\star}^{(k)}$ is the sample mean of a random number $N_{\star}^{(k)}$ of returns in general $N_{\star}^{(k)}$ will depend on those returns themselves (i.e., how often we select arm k depends on the historical returns of arm k)

- The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t)
- Solution: First, imagine an infinite sequence of hypothetical i.i.d. draws from $\nu^{(k)}$: $\widetilde{r}_{0}^{(k)}, \widetilde{r}_{1}^{(k)}, \widetilde{r}_{2}^{(k)}, \widetilde{r}_{3}^{(k)}, \ldots$



 $\hat{\mu}_{t}^{(k)}$ is the sample mean of a random number $N_{t}^{(k)}$ of returns in general $N_{t}^{(k)}$ will depend on those returns themselves (i.e., how often we select arm k depends on the historical returns of arm k)

- The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t)
- Solution: First, imagine an infinite sequence of hypothetical i.i.d. draws from $\nu^{(k)}$: $\widetilde{r}_{0}^{(k)}, \widetilde{r}_{1}^{(k)}, \widetilde{r}_{2}^{(k)}, \widetilde{r}_{3}^{(k)}, \ldots$
- Then we can think of every time we pull arm k, just pulling the next $\tilde{r}_{i}^{(k)}$ off this list,





 $\hat{\mu}_{t}^{(k)}$ is the sample mean of a random number $N_{t}^{(k)}$ of returns in general $N_{\star}^{(k)}$ will depend on those returns themselves (i.e., how often we select arm k depends on the historical returns of arm k)

- The problem: Although $r_{\tau} \mid a_{\tau} = k$ is an i.i.d. draw from $\nu^{(k)}$, (all arm indexing (k) now in superscripts; subscripts reserved for time index t)
- Solution: First, imagine an infinite sequence of hypothetical i.i.d. draws from $\nu^{(k)}$: $\widetilde{r}_{0}^{(k)}, \widetilde{r}_{1}^{(k)}, \widetilde{r}_{2}^{(k)}, \widetilde{r}_{3}^{(k)}, \ldots$
- Then we can think of every time we pull arm k, just pulling the next $\tilde{r}_{i}^{(k)}$ off this list, We can think of every time we pair each, $f_{r_{\tau}}^{(k)}$, $i_{\tau}^{(k)} = \frac{1}{N_t^{(k)}} \sum_{i=0}^{N_t^{(k)}-1} \tilde{r}_i^{(k)}$ i.e., $r_{\tau} \mid a_{\tau} = k$ to simply equal to $\tilde{r}_{N_{\tau}^k}^{(k)}$, and hence $\hat{\mu}_t^{(k)} = \frac{1}{N_t^{(k)}} \sum_{i=0}^{N_t^{(k)}-1} \tilde{r}_i^{(k)}$

- Now Hoeffding applies to $\tilde{\mu}_n^{(k)}$ because *n* fixed/nonrandom

- Now Hoeffding applies to $\tilde{\mu}_n^{(k)}$ because *n* fixed/nonrandom
- and we know $\hat{\mu}_{t}^{(k)} = \tilde{\mu}_{n}^{(k)}$ for some $n \leq t$ (but which one is random)

and we know
$$\hat{\mu}_t^{(k)} = \tilde{\mu}_n^{(k)}$$
 for

- Now Hoeffding applies to $\tilde{\mu}_n^{(k)}$ because *n* fixed/nonrandom
 - some $n \leq t$ (but which one is *random*)
- Recall union bound in ETC analysis made Hoeffding hold simultaneously over $k \leq K$

and we know
$$\hat{\mu}_t^{(k)} = \tilde{\mu}_n^{(k)}$$
 for

$$\Rightarrow \mathbb{P}\left(\forall n \leq t, |\tilde{\mu}_n^{(k)} - \mu\right)$$

- Now Hoeffding applies to $\tilde{\mu}_n^{(k)}$ because *n* fixed/nonrandom
 - some $n \leq t$ (but which one is *random*)
- Recall union bound in ETC analysis made Hoeffding hold simultaneously over $k \leq K$
 - Hoeffding + union bound over $n \le t$: $n \le t$, $|\tilde{\mu}_n^{(k)} \mu^{(k)}| \le \sqrt{\ln(2t/\delta)/2n} \ge 1 \delta$

Constructing confidence intervals (cont'd) Hoeffding + union bound over $n \le t$: $\Rightarrow \mathbb{P}\left(\forall n \le t, |\tilde{\mu}_n^{(k)} - \mu^{(k)}| \le \sqrt{\ln(2t/\delta)/2n} \right) \ge 1 - \delta$

Hoeffding + union bound over $n \leq t$:

$$\Rightarrow \mathbb{P}\left(\forall n \leq t, |\tilde{\mu}_n^{(k)} - \mu\right)$$

$$\mathbb{P}\left(\left|\tilde{\mu}_{N_{t}^{(k)}}^{(k)} - \mu^{(k)}\right| \leq 1\right)$$

- $\mu^{(k)} | \leq \sqrt{\ln(2t/\delta)/2n} \Big) \geq 1 \delta$
- But since in particular $N_t^{(k)} \leq t$, this immediately implies

 $\sqrt{\ln(2t/\delta)/2N_t^{(k)}} \ge 1 - \delta$

Constructing confidence intervals (cont'd) Hoeffding + union bound over $n \leq t$: $\Rightarrow \mathbb{P}\left(\forall n \leq t, |\tilde{\mu}_n^{(k)} - \mu^{(k)}| \leq \sqrt{\ln(2t/\delta)/2n}\right) \geq 1 - \delta$ But since in particular $N_t^{(k)} \leq t$, this immediately implies $\mathbb{P}\left(|\tilde{\mu}_{N_{t}^{(k)}}^{(k)} - \mu^{(k)}| \leq \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right) \geq 1 - \delta$ And then since $\tilde{\mu}_{N_{t}^{(k)}}^{(k)} = \hat{\mu}_{t}^{(k)}$, we immediately get the kind of result we want: $\mathbb{P}\left(\left|\hat{\mu}_{t}^{(k)}-\mu^{(k)}\right|\leq\sqrt{1-\mu^{(k)}}\right)$

$$\left(\frac{\ln(2t/\delta)}{2N_t^{(k)}} \right) \ge 1 - \delta$$

$$\left(\frac{\ln(2t/\delta)}{2N_t^{(k)}}\right) \ge 1 - \delta$$

Constructing confidence intervals (cont'd) Hoeffding + union bound over $n \leq t$: $\Rightarrow \mathbb{P}\left(\forall n \leq t, |\tilde{\mu}_n^{(k)} - \mu'\right)$ But since in particular $N_t^{(k)} \leq t$, this immediately implies $\mathbb{P}\left(\left|\tilde{\mu}_{N_{t}^{(k)}}^{(k)} - \mu^{(k)}\right| \leq \sqrt{2}\right)$ And then since $\tilde{\mu}_{N_t^{(k)}}^{(k)} = \hat{\mu}_t^{(k)}$, we immediately get the kind of result we want: $\mathbb{P}\left(\left|\hat{\mu}_{t}^{(k)}-\mu^{(k)}\right|\leq\sqrt{\right.}$

$$^{(k)} | \leq \sqrt{\ln(2t/\delta)/2n} \ge 1 - \delta$$

$$\left(\frac{\ln(2t/\delta)/2N_t^{(k)}}{t} \right) \ge 1 - \delta$$

$$\left(\frac{\ln(2t/\delta)}{2N_t^{(k)}}\right) \ge 1 - \delta$$

<u>Summary</u>: to deal with problem of non-i.i.d. rewards that enter into $\hat{\mu}_{t}^{(k)}$, we used rewards' conditional i.i.d. property along with a union bound to get Hoeffding bound that is wider by just a factor of t in the log term

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time t from last equation: $\mathbb{P}\left(\left|\hat{\mu}_{t}^{(k)}-\mu^{(k)}\right| \leq \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right) \geq 1-\delta,$ i.e., $\hat{\mu}_t^{(k)} - \sqrt{\ln(2t/\delta)/2N_t^{(k)}}, \ \hat{\mu}_t^{(k)} + \sqrt{\ln(2t/\delta)/2N_t^{(k)}}$

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time *t* from last equation: $\mathbb{P}\left(|\hat{\mu}_t^{(k)} - \mu^{(k)}| \leq \sqrt{\ln(2t/\delta)/2N_t^{(k)}} \right) \geq 1 - \delta,$ i.e., $\left[\hat{\mu}_{t}^{(k)} - \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}, \hat{\mu}_{t}^{(k)} + \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right]^{\prime}$ Valid for any bandit algorithm! Of independent statistical interest for interpreting results

$$\geq 1 - \delta_{s}$$

for interpreting results

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time *t* from last equation: $\mathbb{P}\left(\left|\hat{\mu}_{t}^{(k)} - \mu^{(k)}\right| \leq \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right) \geq 1 - \delta,$ i.e., $\left[\hat{\mu}_{t}^{(k)} - \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}, \ \hat{\mu}_{t}^{(k)} + \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right]$ Valid for any bandit algorithm!
Of independent statistical interest for interpreting results

i.e.,
$$\hat{\mu}_t^{(k)} - \sqrt{\ln(2t/\delta)/2N_t^{(k)}}, \ \hat{\mu}_t^{(k)} + \sqrt{1}$$

But analysis easier if CIs are *uniformly valid* over time t and arm k

$$\geq 1 - \delta$$

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time *t* from last equation: $\mathbb{P}\left(|\hat{\mu}_{t}^{(k)} - \mu^{(k)}| \leq \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right) \geq 1 - \delta,$ i.e., $\left[\hat{\mu}_{t}^{(k)} - \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}, \ \hat{\mu}_{t}^{(k)} + \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right]$ Valid for any bandit algorithm!
Of independent statistical interest for interpreting results

i.e.,
$$\hat{\mu}_t^{(k)} - \sqrt{\ln(2t/\delta)/2N_t^{(k)}}, \ \hat{\mu}_t^{(k)} + \sqrt{1}$$

But analysis easier if CIs are *uniformly valid* over time t and arm k

By same argument as last two slides using a union bound over Hoeffding applied to all $\tilde{\mu}_n^{(k)}$ for $n \leq T$, and noting that $N_t^{(k)} \leq T$ for all t < T, we get:

$$\geq 1 - \delta$$
,

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time *t* from last equation: $\mathbb{P}\left(\left|\hat{\mu}_t^{(k)} - \mu^{(k)}\right| \le \sqrt{\ln(2t/\delta)/2N_t^{(k)}}\right) \ge 1 - \delta,$ Г

i.e.,
$$\hat{\mu}_t^{(k)} - \sqrt{\ln(2t/\delta)/2N_t^{(k)}}, \ \hat{\mu}_t^{(k)} + \sqrt{1}$$

By same argument as last two slides using a union bound over Hoeffding applied to all $\tilde{\mu}_n^{(k)}$ for $n \leq T$, and noting that $N_t^{(k)} \leq T$ for all t < T, we get:

$$\mathbb{P}\left(\left|\forall t < T, \left|\hat{\mu}_{t}^{(k)} - \mu^{(k)}\right| \le \sqrt{\ln(2T/\delta)/2N_{t}^{(k)}}\right|\right) \ge 1 - \delta$$

$$\geq 1 - \delta$$

 $\cdot \sqrt{\ln(2t/\delta)/2N_t^{(k)}}$ Valid for any bandit algorithm! Of independent statistical interest for interpreting results

But analysis easier if CIs are *uniformly valid* over time t and arm k
Uniform confidence intervals

 $\mathbb{P}\left(\left|\hat{\mu}_{t}^{(k)} - \mu^{(k)}\right| \leq \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}\right)$ Г

i.e.,
$$\hat{\mu}_t^{(k)} - \sqrt{\ln(2t/\delta)/2N_t^{(k)}}, \ \hat{\mu}_t^{(k)} + \sqrt{2}$$

By same argument as last two slides using a union bound over Hoeffding applied to all $\tilde{\mu}_n^{(k)}$ for $n \leq T$, and noting that $N_t^{(k)} \leq T$ for all t < T, we get:

$$\mathbb{P}\left(\forall t < T, |\hat{\mu}_t^{(k)} - \mu^{(k)}| \le \sqrt{\ln(2T/\delta)/2N_t^{(k)}}\right) \ge 1 - \delta$$

By same argument made in ETC analysis, union bound over K makes coverage uniform over k: $\mathbb{P}\left(\forall k \leq K, t < T, | \hat{\mu}_t^{(k)} - \mu^{(k)} \right)$

So we have a valid $(1 - \delta)$ confidence interval (CI) for $\mu^{(k)}$ at time t from last equation:

$$\geq 1 - \delta$$
,

, $\hat{\mu}_{t}^{(k)} + \sqrt{\ln(2t/\delta)/2N_{t}^{(k)}}$ Valid for any bandit algorithm! Of independent statistical interest for interpreting results

But analysis easier if CIs are *uniformly valid* over time t and arm k

$$|k| \le \sqrt{\ln(2TK/\delta)/2N_t^{(k)}} \ge 1 - \delta$$







Feedback from last lecture



- Confidence intervals for the arms
 - Upper Confidence Bound (UCB) algorithm
 - UCB regret analysis



Upper Confidence Bound (UCB) algorithm For t = 0, ..., T - 1: Choose the arm with the highest upper confidence bound, i.e., $a_t = \arg \max_{k \in \{1, \dots, K\}} \hat{\mu}_t^{(k)} + \sqrt{\ln(2TK/\delta)/2N_t^{(k)}}$

For t = 0, ..., T - 1: $\hat{\mu}_{t}^{(2)} + \sqrt{\ln(2TK/\delta)/2N_{t}^{(2)}}$ $\hat{\mu}_{t}^{(1)} + \sqrt{\ln(2TK/\delta)/2N_{t}^{(1)}}$ $\hat{\mu}_t^{(2)} \bullet^{\mu^{(2)}}$ $\hat{\mu}_t^{(1)}$ $\hat{\mu}_t^{(1)} - \sqrt{\ln(2TK/\delta)/2N_t^{(1)}}$

Choose the arm with the highest upper confidence bound, i.e., $a_t = \arg \max_{k \in \{1, \dots, K\}} \hat{\mu}_t^{(k)} + \sqrt{\ln(2TK/\delta)/2N_t^{(k)}}$





For t = 0, ..., T - 1: $\hat{\mu}_{t}^{(1)} + \sqrt{\ln(2TK/\delta)/2N_{t}^{(1)}}$ $\hat{\mu}_t^{(2)} \bullet^{\mu^{(2)}}$ $\hat{\mu}_t^{(1)}$ $\hat{\mu}_t^{(1)} - \sqrt{\ln(2TK/\delta)/2N_t^{(1)}}$

Choose the arm with the highest upper confidence bound, i.e., $a_t = \arg \max_{k \in \{1, \dots, K\}} \hat{\mu}_t^{(k)} + \sqrt{\ln(2TK/\delta)/2N_t^{(k)}}$



 $\hat{\mu}_{t}^{(3)} + \sqrt{\ln(2TK/\delta)/2N_{t}^{(3)}}$ $\hat{\mu}_{t}^{(3)} \stackrel{\mu^{(3)}}{\bullet}$ $\hat{\mu}_t^{(3)} - \sqrt{\ln(2TK/\delta)/2N_t^{(3)}}$



(we can't see the $\mu^{(k)}$)





For t = 0, ..., T - 1: $\hat{\mu}_{t}^{(1)} + \sqrt{\ln(2TK/\delta)/2N_{t}^{(1)}}$ $\hat{\mu}_t^{(1)}$ $\hat{\mu}_t^{(1)} - \sqrt{\ln(2TK/\delta)/2N_t^{(1)}}$ $\hat{\mu}_{t}^{(2)} - \sqrt{\ln(2TK/\delta)/2N_{t}^{(2)}}$

Choose the arm with the highest upper confidence bound, i.e., $a_t = \arg \max_{k \in \{1, \dots, K\}} \hat{\mu}_t^{(k)} + \sqrt{\ln(2TK/\delta)/2N_t^{(k)}}$ $\hat{\mu}_t^{(3)} - \sqrt{\ln(2TK/\delta)/2N_t^{(3)}}$

(we can't see the $\mu^{(k)}$)





Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action



Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action

In UCB, this means constructing a CI (i.e., set of plausible values) for each $\mu^{(k)}$, and being greedy with respect to the <u>upper bound</u> of the CIs



- Since each upper bound is $\hat{\mu}_t^{(k)} + \sqrt{\ln(2KT/\delta)/2N_t^{(k)})}$, this means when we select
- $a_t = k$, at least one of the two terms is large, i.e., either

Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action

In UCB, this means constructing a CI (i.e., set of plausible values) for each $\mu^{(k)}$, and being greedy with respect to the <u>upper bound</u> of the CIs



Since each upper bound is
$$\hat{\mu}_t^{(k)} + \sqrt{\ln t}$$

 $a_t = k$, at least one of the two terms is large, i.e., either 1. $\sqrt{\ln(2KT/\delta)/2N_t^{(k)}}$ large, i.e., we haven't explored arm k much (exploration)

Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action

In UCB, this means constructing a CI (i.e., set of plausible values) for each $\mu^{(k)}$, and being greedy with respect to the <u>upper bound</u> of the CIs

 $n(2KT/\delta)/2N_{\star}^{(k)}$, this means when we select





Since each upper bound is
$$\hat{\mu}_t^{(k)} + \sqrt{\ln t}$$

 $a_t = k$, at least one of the two terms is large, i.e., either 1. $\sqrt{\ln(2KT/\delta)/2N_t^{(k)}}$ large, i.e., we haven't explored arm k much (exploration)

Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action

In UCB, this means constructing a CI (i.e., set of plausible values) for each $\mu^{(k)}$, and being greedy with respect to the <u>upper bound</u> of the CIs

 $n(2KT/\delta)/2N_{\star}^{(k)}$, this means when we select

2. $\hat{\mu}_{t}^{(k)}$ large, i.e., based on what we've seen so far, arm k is the best (exploitation)





Since each upper bound is
$$\hat{\mu}_t^{(k)} + \sqrt{\ln t}$$

 $a_t = k$, at least one of the two terms is large, i.e., either 1. $\sqrt{\ln(2KT/\delta)/2N_t^{(k)}}$ large, i.e., we haven't explored arm k much (exploration)

Optimism in the face of uncertainty is an important principle in RL It basically says to give each arm the benefit of the doubt, and basically act as if that arm is as good as it could plausibly be in choosing an action

In UCB, this means constructing a CI (i.e., set of plausible values) for each $\mu^{(k)}$, and being greedy with respect to the <u>upper bound</u> of the CIs

 $n(2KT/\delta)/2N_{\star}^{(k)}$, this means when we select

2. $\hat{\mu}_{t}^{(k)}$ large, i.e., based on what we've seen so far, arm k is the best (exploitation)

Note that the exploration here is *adaptive*, i.e., focused on most promising arms







- Feedback from last lecture
- Recap
- Confidence intervals for the arms
- Upper Confidence Bound (UCB) algorithm
 - UCB regret analysis



UCB Regret Analysis Strategy

UCB Regret Analysis Strategy

1. Bound regret at each time step

UCB Regret Analysis Strategy

- 1. Bound regret at each time step
- 2. Bound the sum of those bounds over time steps

UCB regret at each time step

Recall k^{\star} is optimal arm, so $\mu^{(k^{\star})}$ is true best arm mean. Thus time step t regret is:

 $\mu^{(k^{\star})} - \mu^{(a_t)}$



Sum of UCB per-time-step regrets

2.

1. per-time-step regret bound $\mu^{(k^{\star})} - \mu^{(a_t)} \leq \sqrt{2 \ln(2KT/\delta)/N_t^{(a_t)}}$ w/p $1 - \delta$

Finally, putting it all together, we get: $\operatorname{Regret}_T \leq 2K\sqrt{T}\sqrt{2\ln(KT/\delta)} \quad \text{w/p } 1 - \delta$

Finally, putting it all together, we get: $\operatorname{Regret}_{T} \leq 2K\sqrt{T}$ $= \tilde{O}(\sqrt{T})$

$$\sqrt{2 \ln(KT/\delta)}$$
 w/p 1 – δ
w/p 1 – δ

Finally, putting it all together, we get: $\operatorname{Regret}_{T} \leq 2K\sqrt{T}$ $= \tilde{O}(\sqrt{T})$

In fact, a more sophisticated analysis c

$$\sqrt{2 \ln(KT/\delta)}$$
 w/p 1 – δ
w/p 1 – δ

can get: Regret_T =
$$\tilde{O}(\sqrt{KT})$$
 w/p 1



 Feedback from last lecture Recap • Confidence intervals for the arms Upper Confidence Bound (UCB) algorithm • UCB regret analysis



Upper Confidence Bound (UCB) algorithm:

- Uses uncertainty quantification inside algorithm
- Applies principle of optimism in the face of uncertainty (OFU)
- Adaptive exploration
- Achieves regret of $\tilde{O}(\sqrt{TK})$

Upper Confidence Bound (UCB) algorithm:

- Uses uncertainty quantification inside algorithm
- Applies principle of optimism in the face of uncertainty (OFU)
- Adaptive exploration
- Achieves regret of $\tilde{O}(\sqrt{TK})$

Next time:

- Can't do better than $\tilde{\Omega}(\sqrt{TK})$ regret in worst case
- Instance-dependent regret is another way to analyze regret

Upper Confidence Bound (UCB) algorithm:

- Uses uncertainty quantification inside algorithm
- Applies principle of optimism in the face of uncertainty (OFU)
- Adaptive exploration
- Achieves regret of $\tilde{O}(\sqrt{TK})$

Next time:

- Can't do better than $\tilde{\Omega}(\sqrt{TK})$ regret in worst case Instance-dependent regret is another way to analyze
- regret

1-minute feedback form: <u>https://bit.ly/3RHtlxy</u>



