

Week 12

- ↳ Research talk tomorrow 4:15 in Kravis 62
- ↳ Office Hours Monday + Thursday 12:30-2:00
- ↳ Midterm on 11/25, practice exams next week :-)
- ↳ Project Proposal due 12/2

Reinforcement Learning

State, action, reward, ...

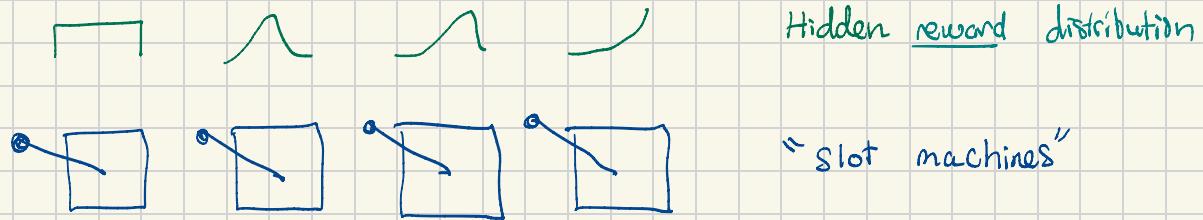
Key: Exploration vs. Exploitation

This Week: theoretical understanding

Simple model: multi-armed bandits

Question:

What can we say about draws from these distributions?



Goal: Pull arms to maximize reward

↳ clinical trials

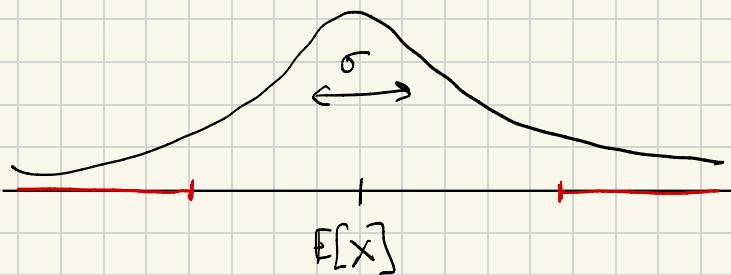
↳ recommendation systems

Random Variables

X r.v. $\Pr(X=x)$ = prob X takes value x

$$\mathbb{E}[X] = \sum_x x \Pr(X=x)$$

$$\text{Var}(x) = \mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$



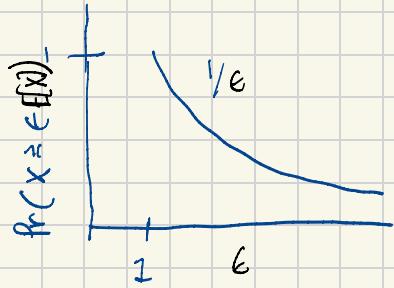
Goal: Bound probability of extreme values

Markov's Inequality

Consider non-negative X

$$\text{For } \epsilon > 0, \quad \Pr(X \geq \epsilon) \leq \frac{\mathbb{E}[X]}{\epsilon}$$

$$\mathbb{E}[X]$$



Proof:

$$\begin{aligned} \mathbb{E}[X] &= \sum_x \Pr(X=x) x \\ &= \sum_{x: x \geq \epsilon} \Pr(X=x) x + \sum_{x: x < \epsilon} \Pr(X=x) x \\ &\geq \sum_{x: x \geq \epsilon} \Pr(X=x) \epsilon + \sum_{x: x < \epsilon} \Pr(X=x) \cdot 0 \\ &= \epsilon \sum_{x: x \geq \epsilon} \Pr(X=x) \\ &= \epsilon \Pr(X \geq \epsilon) \end{aligned}$$

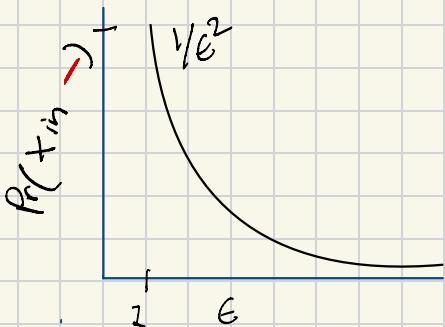
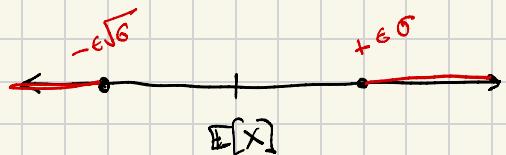
Q: Fix ϵ . Build a rv where $\Pr(X \geq \epsilon) = \frac{\mathbb{E}[X]}{\epsilon}$

Chebyshov's Inequality

Idea: Use extra information for tighter bound

$$\sigma^2 = \text{Var}(x). \text{ For } \epsilon > 0,$$

$$\Pr(|x - \mathbb{E}[x]| \geq \epsilon \sigma) \leq \frac{1}{\epsilon^2}$$



$$\text{Proof: } z = (x - \mathbb{E}[x])^2$$

By Markov's,

$$\Pr(z \geq \epsilon^2) \leq \frac{1}{\epsilon^2}$$

$$\Leftrightarrow \Pr((x - \mathbb{E}[x])^2 \geq \epsilon^2 \mathbb{E}[(x - \mathbb{E}[x])^2]) \leq \frac{1}{\epsilon^2}$$

(\geq)

$$\Pr(|x - \mathbb{E}[x]| \geq \sqrt{\epsilon^2 \sigma}) \leq \frac{1}{\epsilon^2}$$

(\leq)

$$\Pr(|x - \mathbb{E}[x]| \geq \epsilon \sigma) \leq \frac{1}{\epsilon^2}$$

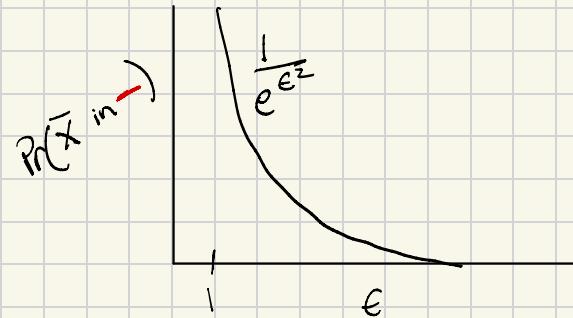
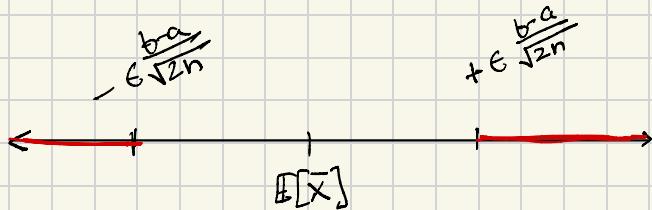
Hoeffding's Inequality

Idea: Formalize central limit theorem

Consider x_1, \dots, x_n where $a \leq x_i \leq b$

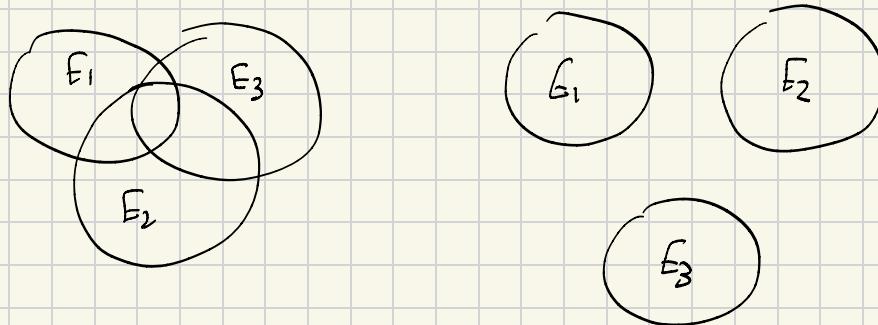
Let $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. For $\epsilon > 0$,

$$\Pr(|\bar{x} - \mathbb{E}[\bar{x}]| \geq \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$



Union Bound

$$\Pr(E_1 \cup E_2 \cup \dots \cup E_m) \leq \Pr(E_1) + \Pr(E_2) + \dots + \Pr(E_m)$$



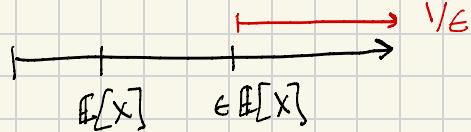
Logistics

↳ Midterm

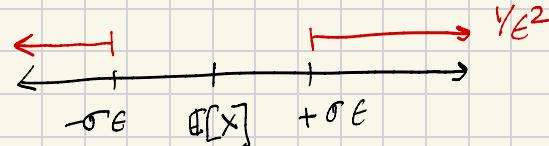
↳ Project

Concentration Inequalities

Markov's

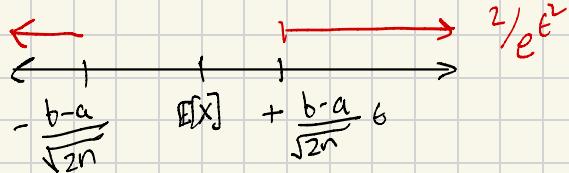


Chabyshev's



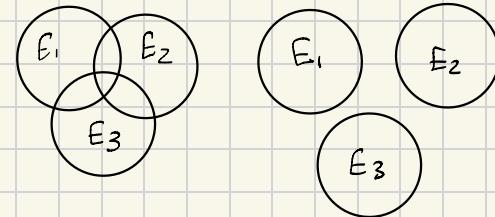
$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad a \leq X_i \leq b$$

Hoeffding's

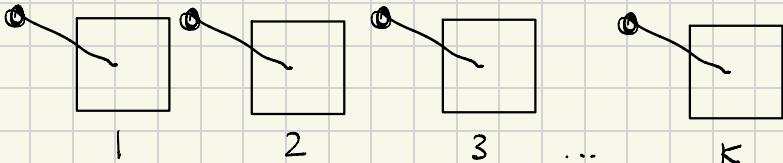


Union Bound

$$\Pr(E_1 + E_2 + \dots + E_m) \leq \Pr(E_1) + \Pr(E_2) + \dots + \Pr(E_m)$$



Multi-armed Bandits



T time steps

Choose action $a^{(t)} \in \{1, \dots, k\}$ at time t

Receive reward $r_{a(t)} \in [-1, 1]$

Expected reward of arm a : $\mu_a = E[r_a]$

Goal: Minimize average regret

$$R_T = \max_{a \in \{1, \dots, k\}} \mu_a - \frac{1}{T} \sum_{t=1}^T \mu_{a(t)}$$

Exploration vs. exploitation

Strategy: "optimism in the face of uncertainty"

Upper Confidence Bound (UCB) Algorithm

$$t = \underbrace{1, 2, \dots, K}_{\text{Try each arm}}, \underbrace{K+1, \dots, T}_{\text{Choose optimistically best arm}}$$

At time t ,

- $\tilde{\mu}_a^{(t)}$ = empirical average of arm a so far $= \frac{1}{n_a^{(t)}} \sum_{l=1}^t \mathbb{I}[a^{(l)} = a] r_a^{(l)}$
- $\epsilon_a^{(t)}$ = uncertainty of estimate so far

Informally, $|\mu_a - \tilde{\mu}_a^{(t)}| \leq \epsilon_a^{(t)}$ with high prob. e.g., w.p. $1-\delta$

Choose $a^{(t)} = \operatorname{argmax}_a \tilde{\mu}_a^{(t)} + \epsilon_a^{(t)}$

Lemma a: For all $t = k+1, \dots, T$ and $a = 1, \dots, k$ with probability $1-\delta$

$$|\mu_a - \hat{\mu}_a^{(t)}| \leq \epsilon_a^{(t)} = \sqrt{\frac{2 \log(\frac{2Tk}{\delta})}{n_a^{(t)}}}$$

Proof: Hoeffding's :

$$\Pr(|E[\bar{X}] - \bar{X}| \geq \epsilon \frac{b-a}{\sqrt{2n}}) \leq 2e^{-\epsilon^2}$$

$$b-a = 1-(-1) = 2, \quad n = n_a^{(t)}, \quad \hat{\mu}_a^{(t)} = \bar{X}$$

$$\begin{aligned} 2e^{-\epsilon^2} &= \frac{\delta}{Tk} \\ (\Leftrightarrow) \quad \frac{2Tk}{\delta} &= e^{\epsilon^2} \\ (\Leftrightarrow) \quad \sqrt{\log \frac{2Tk}{\delta}} &= \epsilon \end{aligned}$$

$$\epsilon \frac{b-a}{\sqrt{2n}} = \sqrt{\frac{2 \log \frac{2Tk}{\delta}}{n_a^{(t)}}}$$

wp $\frac{\delta}{Tk}$, $\Pr(|\mu_a - \hat{\mu}_a^{(t)}| \geq \sqrt{\frac{2 \log \frac{2Tk}{\delta}}{n_a^{(t)}}}) \leq \frac{\delta}{Tk}$

for one t and a

union bound

$$\Pr\left(\bigcup_{t=1}^T \bigcup_{a=1}^k |\mu_a - \hat{\mu}_a^{(t)}| \geq \epsilon_a^{(t)}\right) \leq \sum_{t=1}^T \sum_{a=1}^k \frac{\delta}{Tk} = \delta$$

Theorem: With prob $1-\delta$,

$$R_T = \max_a \mu_a - \frac{1}{T} \sum_{t=1}^T \mu_a^{(t)} \leq O\left(\frac{K}{T} + \sqrt{\log \frac{T K}{\delta}} \frac{\sqrt{K}}{\sqrt{T}}\right)$$
$$\approx O\left(\sqrt{\frac{K}{T}}\right) \quad \text{since 1) } K \ll T$$

2) $\log(\text{anything}) = \text{small}$

$$\log(\# \text{atoms}) \approx 82$$

Intuition:

- Increase # arms by 100x, only 10x increase in regret
- Increase # steps by 100x, only 10x decrease in regret