

CS292F StatRL Lecture 6

RL Algorithms + Bandits

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

Recap: Lecture 5

- Model-free RL that combines MC and DP
 - Temporal Difference methods and their variants
 - Function approximation ← MDP's large $|S|$ is large
 - TD methods (TD prediction / policy evaluation) as one that minimizes the square Bellman error a semi-gradient update.
 - LSTD and its fast computation
- Policy gradients methods

Recap: TD-method in a nutshell --- attempting to simulate Bellman equations with roll-out data

Bellman eqn. $V^\pi = r^\pi + \gamma P^\pi V^\pi \Leftrightarrow V^\pi(s) = \underbrace{r^\pi(s)}_{\text{anti-}\pi(s)} + \gamma \underbrace{E_{S' \sim P(\cdot|s)} [V^\pi(s')]}_{\text{sim } P(\cdot|s)}$

TD-prediction $\hat{V}_{t+1}^\pi(s) = \hat{V}_t^\pi(s) + \alpha [\underbrace{R_{t+1} + \gamma \hat{V}_t^\pi(S_{t+1})}_{C_T} - \hat{V}_t^\pi(S_t)]$

$E[C_T] = r^\pi(S_t) + \gamma E_{S_{t+1}} [\hat{V}_t^\pi(S_{t+1})]$

Setting $\alpha = 1$ $\alpha \ll 1$

"Bootstrap"

Use $\hat{V}_t^\pi(S_{t+1})$. not unbiased, lower variance

"MC"

$C_T = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$
unbiased but large variance

TD-methods: what I did not cover

- n-step TD
 - TD(λ) / Eligible traces
 - Off-policy methods and the “deadly triad”
 - Semi-gradient vs actual SGD?
 - Projected Bellman error and Gradient-TD
1. off policy
2. bootstrap
3. function approx.

(A large body of associated work on these.

Read Sutton and Barto Ch 7, CH 11-12 and the references therein.)

This lecture

- Continue with policy gradient methods
- Start exploration
 - Multi-armed bandits

Recap: Policy class and policy gradient methods

- Policy $\pi \in \Pi$ $\pi: S \rightarrow \Delta(A)$

- Parametric policy class:

$$\Pi = \{ \pi_\theta \mid \theta \in \mathbb{R}^d \}$$

- Goal: optimize the value

$$\theta^* = \underset{\theta}{\operatorname{argmax}} V^{\pi_\theta}(\mu)$$

$$V^{\pi_{\theta^*}} \neq V^*$$

- Policy gradient methods

- aim at learning the policy parameter by SGD.

$$\theta_{t+1} = \theta_t + \alpha \cdot \nabla_{\theta} V^{\pi_{\theta}}(\mu) \in G_D$$

$$\theta_{t+1} = \theta_t + \alpha \cdot g_t \quad (\text{SGD})$$

Soft max $d = |S| \cdot |A|$

$$p(a|s) = \frac{\exp(\theta_{s,a})}{\sum_a \exp(\theta_{s,a})}$$

$$\pi_\theta(\phi) = \frac{\exp(\theta^\top \phi(s,a))}{\sum_a \exp(\theta^\top \phi(s,a))}$$

$$\begin{aligned} & \textcircled{2} \mathbb{E}[\|g_t - \nabla_{\theta} V^{\pi_{\theta}}(\mu)\|^2 \mid \theta_t \in G_D] \\ & \mathbb{E}[g_t \mid \theta_t] = \nabla_{\theta} V^{\pi_{\theta}}(\mu) \end{aligned}$$

How to estimate the gradient?

$$R(\tau) = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots$$

- Policy gradient theorem:

Claim 1: $\nabla V^{\pi_0}(\mu) = \mathbb{E}_{\substack{\pi_0 \\ \text{trajectory}}} \left[R(\tau) \sum_{t=0}^{\infty} \nabla \log \pi_0(A_t | S_t) \right]$

① Reinforce π_0
 ① Roll out a trajectory τ with π_0
 ② $R(\tau) \rightarrow \sum_{t=0}^{\infty} \nabla \log \pi_0(A_t | S_t)$

Claim 2:
$$\begin{aligned} \nabla V^{\pi_0}(\mu) &= \mathbb{E}_{\pi_0} \left[\sum_{t=1}^{\infty} \gamma^t Q^{\pi_0}(S_t, A_t) \nabla \log \pi_0(A_t | S_t) \right] \\ &= \sum_S \underbrace{V_{\mu}^{\pi_0}(s)}_a \sum_a \pi_0(a|s) \cdot Q^{\pi_0}(s, a) \nabla \log \pi_0(a|s) \\ &= \frac{1}{1-\gamma} \mathbb{E}_{\text{und}_{\mu}^{\pi_0}} \mathbb{E}_{A \sim \pi_0(\cdot|S)} \left[Q^{\pi_0}(SA) \nabla \log \pi_0(A|S) \right] \end{aligned}$$

Claim 3:
$$\nabla V^{\pi_0}(\mu) = \sum_S \underbrace{V_{\mu}^{\pi_0}(s)}_a \sum_a \pi_0(a|s) \cdot \underbrace{(Q^{\pi_0}(s, a) - V^{\pi_0}(s))}_{A^{\pi_0}(SA)} \nabla \log \pi_0(a|s)$$

Proof of Policy Gradient Theorem

Claim 1: $\nabla V^{\pi_{\theta}}(\mu) = \nabla \sum_{\tau} R(\tau) \cdot P_{\mu}^{\pi_{\theta}}(\tau)$

$= \sum_{\tau} R(\tau) \cdot \nabla P_{\mu}^{\pi_{\theta}}(\tau)$

↳ "score function trick" !

$= \sum_{\tau} R(\tau) \cdot \underbrace{P_{\mu}^{\pi_{\theta}}(\tau)} \cdot \underbrace{\nabla \log P_{\mu}^{\pi_{\theta}}(\tau)}$

$= \sum_{\tau} R(\tau) \cdot P_{\mu}^{\pi_{\theta}}(\tau) \nabla_{\theta} \log \left(\underbrace{\mu(s_0)}_{\pi_{\theta}(A_0|s_0)} \cdot \underbrace{P(s_1|A_0,s_0)}_{\pi_{\theta}(A_1|s_1)} \cdot \underbrace{P(s_2|A_1,s_1)}_{\vdots} \right)$

$= \sum_{\tau} R(\tau) \cdot P_{\mu}^{\pi_{\theta}}(\tau) \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(A_t|s_t)$

* this works for non-stationary ~~π~~ π_{θ}

* works for POMDP

* works for cases when MDP model it self is wrong

□

Proof of Policy Gradient Theorem

Claim 2. $\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \nabla_{\theta} \sum_{a_0} \pi_{\theta}(a_0 | s_0) \cdot Q^{\pi_{\theta}}(s_0, a_0)$

$(x \cdot y)' = x'y + xy'$

$$= \sum_{a_0} \nabla \bar{u}_{\theta}(a_0 | s_0) Q^{\pi_{\theta}}(s_0, a_0) + \sum_{a_0} \pi_{\theta}(a_0 | s_0) \cdot \nabla Q^{\pi_{\theta}}(s_0, a_0)$$

$$= \sum_{a_0} \pi_{\theta}(a_0 | s_0) \nabla (\log \pi_{\theta}(a_0 | s_0) \cdot Q^{\pi_{\theta}}(s_0, a_0))$$

$$+ \sum_{a_0} \pi_{\theta}(a_0 | s_0) \nabla \left(r(s_0, a_0) + \gamma \sum_{s_1} P(s_1 | s_0, a_0) V^{\pi_{\theta}}(s_1) \right)$$

$$= \sum_{a_0} \pi_{\theta}(a_0 | s_0) \left[\nabla (\log \pi_{\theta}(a_0 | s_0)) Q^{\pi_{\theta}}(s_0, a_0) \right. \\ \left. + \gamma E \left[\nabla V^{\pi_{\theta}}(s_1) \right] \right]$$

$$\gamma \sum_{s_1} P(s_1 | s_0, a_0) \nabla V^{\pi_{\theta}}(s_1)$$

$$\gamma E \left[\nabla V^{\pi_{\theta}}(s_1) | s_0, a_0 \right]$$

$$\nabla V^{\pi_{\theta}}(s_0) = \sum_{a_0} \pi_{\theta}(a_0 | s_0) \nabla V^{\pi_{\theta}}(s_0)$$

$$= E \left[\nabla (\log \pi_{\theta}(A_0 | S_0)) \cdot Q^{\pi_{\theta}}(S_0, A_0) \right] + \gamma E_{S_0, \mu} E_{S_1, \pi_{\theta}^{a_0}} \left[\nabla V^{\pi_{\theta}}(s_1) \right] + \gamma^2 \sum_{s_2} \dots$$

REINFORCE Algorithm

REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$

Repeat forever:

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \theta)$

For each step of the episode $t = 0, \dots, T - 1$:

$G \leftarrow$ return from step t

$\theta \leftarrow \theta + \alpha \gamma^t G \nabla_{\theta} \ln \pi(A_t | S_t, \theta)$

"Claim 1 of the Thm"

Variance of the gradient estimate and convergence

$\min_{\theta} f(\theta)$

- This is a non-convex optimization problem

$\max_{\theta} V^{Tr}(\theta)$

$t \leq \frac{\beta(V^*(\theta) - V(\theta))}{G^2}$

f is β -smooth

∇f is β Lipschitz

$\|\nabla f(\theta) - \nabla f(\theta')\| \leq \beta \|\theta - \theta'\|$

- Convergence to a stationary point.

- Choice of learning rate:
- Then (Lemma 9.8 AJKS)

$\alpha = \frac{1}{\beta} \sqrt{\frac{2}{\beta T}}$
 $\alpha = \sqrt{\frac{2}{\beta T}}$

$E[g(\theta) | \theta] = \nabla f(\theta)$

$E[\|g(\theta) - E[g(\theta) | \theta]\| | \theta] \leq G^2$ for all θ

$$\frac{1}{T} \sum E[\|\nabla V^{Tr}(\theta_t)\|^2] \leq \frac{2\beta(V^* - V^{Tr}(\theta_0))}{T} + \sqrt{\frac{2G}{T}}$$

- Does it converge to the global optimal solution?
 - Yes, sometimes (under additional assumptions, see Ch 10 AJKS).

REINFORCE with a baseline

- REINFORCE with a given (arbitrary) baseline

$$\theta^t = \theta + \alpha \gamma^t \cdot G \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | S_t)$$

REINFORCE
with baseline f

$$\theta^t = \theta + \alpha \gamma^t [G - f(S_t)] \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | S_t)$$

- Choose it to approximate the value function.

- Why does it work?

\sum_t suffices

$$E [f(S_t) \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | S_t)] = 0$$

$$E [E_{\pi} [f(S_{t+1}) \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) | S_t]]$$

$$= E \left[\sum_a \pi_{\theta}(a | S_t) \cdot f(S_{t+1}) \cdot \nabla_{\theta} \log \pi_{\theta}(a | S_t) \mid S_t \right] = E \left[\sum_a \pi_{\theta}(a | S_t) \frac{\nabla_{\theta} \pi_{\theta}(a | S_t)}{\pi_{\theta}(a | S_t)} \cdot f(S_{t+1}) \right]$$

Exercise: Calculate the variance when the baseline is chosen to be the correct value function.

$$f(S): S \rightarrow \mathbb{R}$$

$$f = \sum \pi_{\theta}$$

$$f = \sum \pi_{\theta}$$

$$E \left[\nabla_{\theta} \left(\sum_a \pi_{\theta}(a | S_t) \right) \mid f(S_t) \right]$$

$$\parallel \nabla_{\theta} \pi_{\theta}(a | S_t) \parallel$$

Actor-Critic: Learn the baseline and use the baseline for "bootstrapping"

- **Actor:** The policy that performs actions.
- **Critic:** A value function approximation that evaluates the "actor"

for $t = 0, 1, 2, \dots$

$$\text{TD error } \delta = \underbrace{R_{t+1}}_G + \underbrace{\gamma V_w(S_{t+1})}_{\text{baseline}} - \underbrace{V_w(S_t)}_{\text{baseline}}$$

$$\text{update "Critic" by TD: } w \leftarrow w + \alpha \cdot \delta \cdot \nabla_w V_w(S_t)$$

$$\text{update "Actor" by PG: } \theta \leftarrow \theta + \alpha \cdot \delta \cdot \nabla_{\theta} \log \pi_{\theta}(A_t | S_t)$$

Advantage representation of the policy gradient A2C

Async. Adv. Actor-Critic A3C

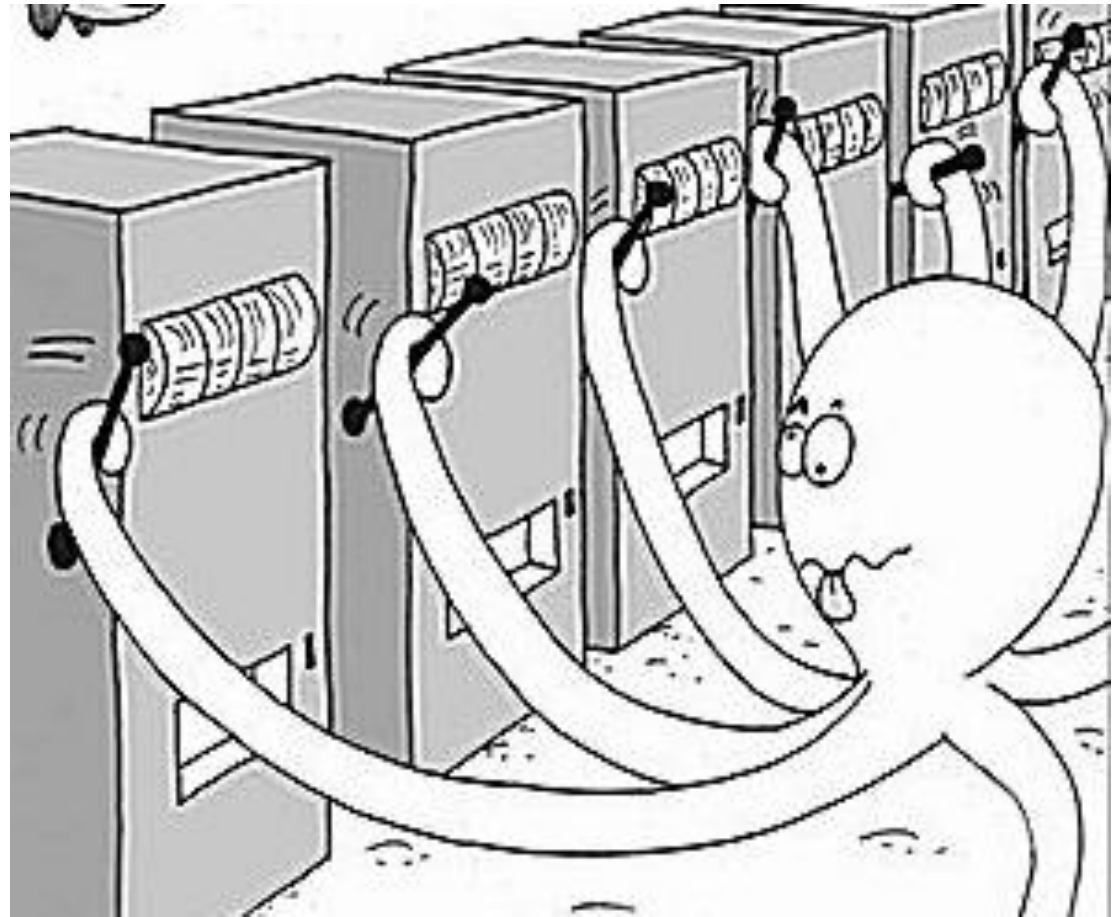
What I did not talk about

- Natural Policy Gradient
 - Extension of Natural Gradient (i.e., Fisher-Scoring)
 - Multiplicative form in updating the policy
 - Function approximation
- Global convergence of PG / NPG
- KL-divergence regularization / imitation learning / conservative policy iteration / trust-region policy iteration

Recap: Let us tackle different aspects of the RL problem one at a time

- Markov Decision Processes: (Lecture 1-3)
 - Dynamics are given no need to learn. planning only.
- RL algorithms: (Lecture 4-5)
 - Ideas on how to use “sampled” transitions / trials-and-errors to learn and plan at the same time (without touching upon exploration).
- **Strategic exploration: (Lecture 6- 9)**
 - Bandits: Explore-Exploit in simple settings
 - RL: Explore-Exploit in Learning MDPs

Slot machines and Multi-arm bandits



Multi-arm bandits: Problem setup

$$|S|=1$$

- No state. k-actions $a \in \mathcal{A} = \{1, 2, \dots, k\}$
- You decide which arm to pull in every iteration

$$A_1, A_2, \dots, A_T$$

- You collect a cumulative payoff of $\sum_{t=1}^T R_t$
- The goal of the agent is to maximize the expected payoff.
 - For future payoffs?
 - For the expected cumulative payoff?

$$\max_{A_1, \dots, A_T} \mathbb{E} \left[\sum_{t=1}^T R_t \right]$$

$$A_t = \text{function}(A_1, \dots, A_{t-1}, R_1, \dots, R_{t-1})$$
$$A_i = a^* = \arg \max_a \mathbb{E}(R_t | A_i = a)$$

How do we measure the performance of an online learning agent?

- The notion of “Regret”:
 - I wish I have done things differently.
 - Comparing to the best actions in the hindsight, how much worse did I do.

- For MAB, the regret is defined as follow

$$\underbrace{T \max_{a \in [k]} \mathbb{E}[R_t | a]}_{\text{grade}} - \underbrace{\sum_{t=1}^T \mathbb{E}_{a \sim \pi} [\mathbb{E}[R_t | a]]}_{\text{agent}} = o(T) \text{ (No-regret learning alg)}$$

$$T \max_a \mathbb{E}[R_t(a)] - \sum_{t=1}^T \mathbb{E}[R_t | A_t]$$

Greedy strategy

Handwritten notes:

$A_1=1$ $R_1=0 \sim \text{Ber}(0.7)$ $a=1$ $Q_1(1)=0$

$A_2=0$ $R_2=1 \sim \text{Ber}(0.5)$ $Q_2(0)=1$

Keep choosing $A_1=0$

- Expected reward

$$r(s, a) = \underline{Q(s, a)} = \mathbb{E}[R_t \mid A_t = a].$$

- Estimate the expected reward

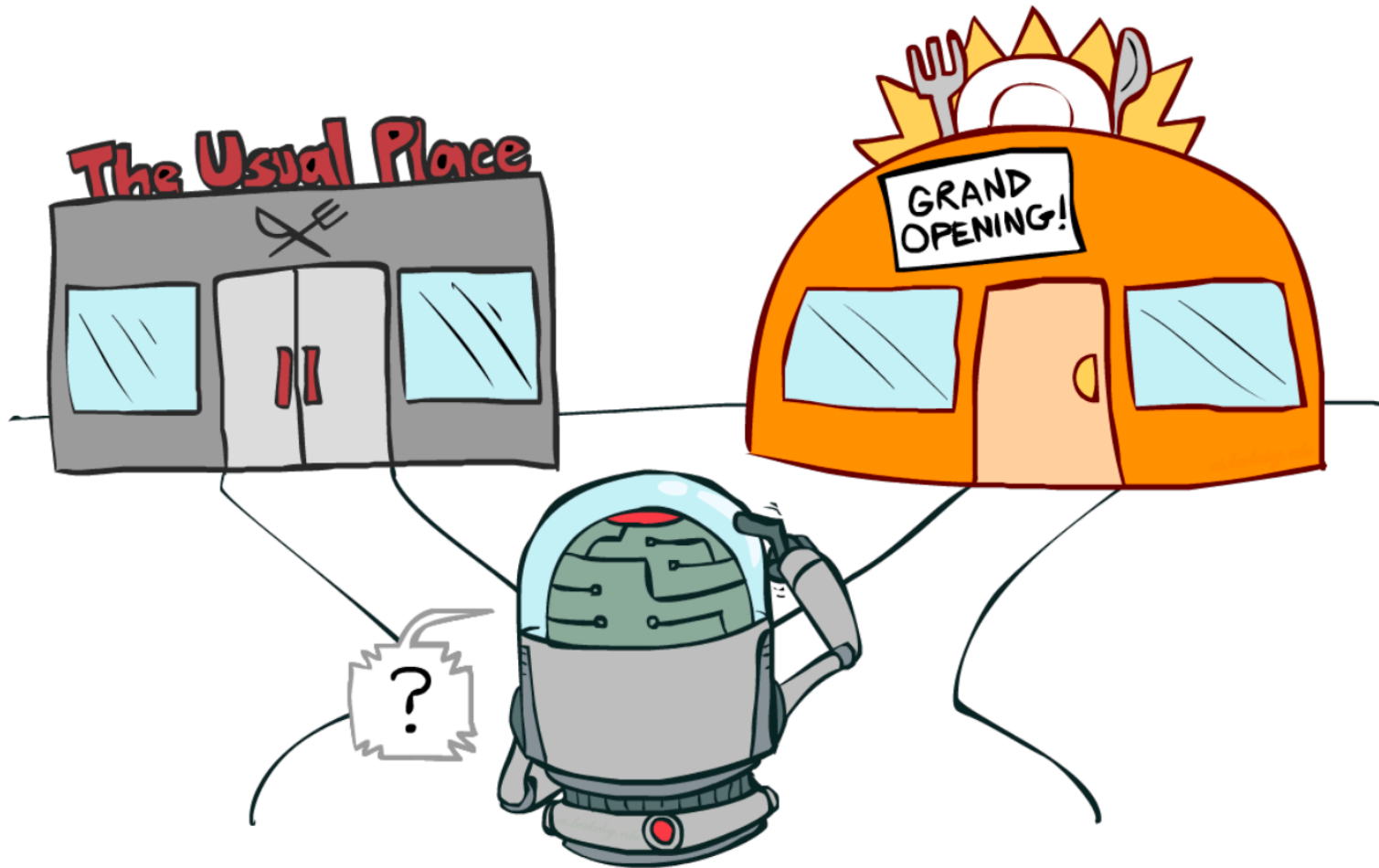
$$\hat{Q}_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t}$$

$$= \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

- Choose $A_t \doteq \arg \max_a \hat{Q}_t(a)$,

What is the issue with this strategy?

Exploration vs. Exploitation



(Illustration from Dan Klein and Pieter Abbeel's course in UC Berkeley)

Exploration first strategy

- Let's spend the first N step exploring.
 - Play each action for N / k times.

$$\hat{Q}_t(a) = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

- For $t = N + 1, N+2, \dots, T$:

$$\underline{A_t \doteq \arg \max_a Q_t(a),}$$

Recap: Concentration inequalities --- finite-sample bounds of LLN and CLT

- **Hoeffding's inequality:** Assume X_1, \dots, X_n are independent and their support bounded:

$$S_n = X_1 + \dots + X_n$$
$$P(S_n - \mathbb{E}[S_n] \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right),$$

- Easy version, if $0 < X_i < B$, **with probability $1 - \delta$:**

$$|\bar{X} - \mathbb{E}[\bar{X}]| \leq \sqrt{\frac{B^2}{2n} \log(2/\delta)}$$

Regret analysis of Exploration First

Regret analysis of Exploration First

ϵ -Greedy strategy: one way to balance exploration and exploitation

- You choose with probability $1 - \epsilon$

$$A_t \doteq \operatorname{argmax}_a Q_t(a),$$

- With probability ϵ , choose an action **uniformly at random!**
 - Including the argmax.
- Carefully choose ϵ parameter.

A sketch of the analysis for ϵ -greedy

- In expectation, each arm is chosen for at least ϵt times.
- Condition on the number of times, apply Hoeffding's inequality / union bound for all t and a
- Regret bound is

$$\epsilon T + \sum_{t=1}^T C \sqrt{\frac{k}{\epsilon t}}$$

Optimism-in-the-face of uncertainty: Upper Confidence Bound algorithm

Martingale

- We say that a sequence of r.v. X_1, \dots, X_n, \dots is a Martingale if for any n

$$\mathbf{E}(|X_n|) < \infty$$

$$\mathbf{E}(X_{n+1} \mid X_1, \dots, X_n) = X_n.$$

- Example:
 - Random-walk: Total number of heads minus tails in n coin tosses

Azuma-Hoeffding's inequality

- **Azuma-Hoeffding's inequality:** Assume X_1, \dots, X_n are **Martingale differences**

$$S_n = X_1 + \dots + X_n$$

$$\mathbb{P} [S_n \geq \epsilon] \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

- Apply Azuma-Hoeffding's inequality to our problem

Regret analysis of UCB

Regret analysis of UCB

Summary of Exploration in Multi-Armed Bandits

- Explore-First
- ϵ -greedy
- UCB