

CMPT 419/983: Theoretical Foundations of Reinforcement Learning

Lecture 3

Sharan Vaswani

September 22, 2023

Recap

- **Stochastic Linear Bandits:** For arm $a \in [K]$, $\mu_a = \langle X_a, \theta^* \rangle$.
- On pulling arm a , we observe reward $R_t = \mu_{a_t} + \eta_t$, $\mathbb{E}[\eta_t] = 0$ and η_t is conditional 1 sub-Gaussian, i.e. for $\lambda \in \mathbb{R}$, $\mathbb{E}[\exp(\lambda\eta_t)|\mathcal{H}_{t-1}] \leq \exp(\lambda^2/2)$.

Algorithm LinUCB: **Blue** for finite number of arms, **Red** for infinitely many arms

1: **Input:** $\{\beta_t\}_{t=2}^{T+1}$, $V_0 = \lambda I_d \in \mathbb{R}^{d \times d}$, $b = 0 \in \mathbb{R}^d$

2: **For each arm** $a \in [K]$, initialize $U_a(1, \delta) := \infty$.

3: **for** $t = 1 \rightarrow T$ **do**

4: Select arm $a_t = \arg \max_{a \in [K]} U_a(t, \delta) = \arg \max_{a \in \mathcal{A}} U_a(t, \delta)$

5: Observe reward R_t and update:

$$V_t = V_{t-1} + X_t X_t^T \quad ; \quad b_t = b_{t-1} + R_t X_t \quad ; \quad \hat{\theta}_t = V_t^{-1} b_t$$

$$U_a(t+1) = \max_{\theta \in \mathcal{C}_{t+1}} \langle \theta, X_a \rangle = \langle X_a, \hat{\theta}_t \rangle + \sqrt{\beta_{t+1}} \|X_a\|_{V_t^{-1}}$$

6: **end for**

Recap

Claim: Assuming (i) $\|\theta^*\| \leq 1$, (ii) $\|X_a\| \leq 1$ for all a and (iii) $R_t \in [0, 1]$, UCB with $\sqrt{\beta_t} = \sqrt{d \log\left(\frac{\lambda d + t}{\lambda d}\right) + 2 \log(1/\delta) + \sqrt{\lambda}}$ achieves the following worst-case bound on the regret,

$$\text{Regret}(\text{LinUCB}, T) \leq O\left(d\sqrt{T} \log(T)\right)$$

Last time we showed the following results: if

$$G := \left\{ \forall t \in [T] \mid \theta^* \in \mathcal{C}_t := \left\{ \theta \mid \left\| \theta - \hat{\theta}_{t-1} \right\|_{V_{t-1}}^2 \leq \beta_t \right\} \right\},$$

$$(1): \text{Regret}(\text{LinUCB}, T) \leq 2 \sqrt{T \beta_T \mathbb{E} \left[\sum_{t=1}^T \|X_t\|_{V_{t-1}}^2 \mid G \right]} + T \Pr[G^c]$$

$$(2): \sum_{t=1}^T \|X_t\|_{V_{t-1}}^2 \leq 2d \log\left(\frac{\lambda d + T}{\lambda d}\right)$$

Today, we will prove: **(3):** For $\sqrt{\beta_t} = \sqrt{d \log\left(\frac{\lambda d + t}{\lambda d}\right) + 2 \log(T) + \sqrt{\lambda}}$, $\Pr[G^c] \leq \frac{1}{T}$, and thus finish the proof.

Digression – (Super)-Martingales

Martingale: Sequence of random variables for which, at a particular time, the conditional expectation of the next value in the sequence is equal to the present value, regardless of all prior values.

A sequence of random variables – M_1, M_2, \dots is a discrete-time martingale if for all t ,

$$\mathbb{E}[|M_t|] < \infty \quad ; \quad \mathbb{E}[M_t | M_1, M_2, \dots, M_{t-1}] = M_{t-1}$$

Example 1: An unbiased random walk

Example 2: Gambler's fortune: Suppose M_t is a gambler's fortune after t tosses of a fair coin, where the gambler wins \$1 if the coin comes up heads and loses \$1 if it comes up tails.

Super-Martingale: A sequence of random variables – M_1, M_2, \dots is a discrete-time super-martingale if for all t ,

$$\mathbb{E}[|M_t|] < \infty \quad ; \quad \mathbb{E}[M_t | M_1, M_2, \dots, M_{t-1}] \leq M_{t-1}$$

Linear UCB – Regret Analysis

Claim: If (i) $\|\theta^*\| \leq 1$ and (ii) $\|X_a\| \leq 1$ for all a , for $\sqrt{\beta_t} = \sqrt{d \log\left(\frac{\lambda d + t}{\lambda d}\right) + 2 \log(T) + \sqrt{\lambda}}$ and $G := \{\forall t \in [T] \mid \theta^* \in \mathcal{C}_t := \left\{ \theta \mid \left\| \theta - \hat{\theta}_{t-1} \right\|_{V_{t-1}}^2 \leq \beta_t \right\}\}$, $\Pr[G^c] \leq \frac{1}{T}$.

Proof: Define $S_t := \sum_{s=1}^t \eta_s X_s$ and $K_t := \sum_{s=1}^t X_s X_s^\top$. We will prove the claim in 4 steps:

- (i) $\left\| \theta - \hat{\theta}_t \right\|_{V_t} \leq \|S_t\|_{V_t^{-1}} + \sqrt{\lambda}$.
- (ii) $M_t(z) = \exp\left(\langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2\right)$ is a non-negative super-martingale with $M_0(z) = 1$.
- (iii) Use the fact that a mixture of super-martingales given by $\bar{M}_t = \int_z M_t(z) h(z) dz$ is also a non-negative super-martingale for any probability density function $h(z)$.
- (iv) Use the maximal inequality for super-martingales to bound $\Pr\left[\sup_{t \in [T]} \log(\bar{M}_t(z)) \geq \log(1/\delta)\right]$ and hence bound $\left\| \theta - \hat{\theta}_t \right\|_{V_t}$.

Linear UCB – Regret Analysis

Part (i): If $S_t := \sum_{s=1}^t \eta_s X_s$ and $K_t := \sum_{s=1}^t X_s X_s^\top$, then $\|\theta^* - \hat{\theta}_t\|_{V_t} \leq \|S_t\|_{V_t^{-1}} + \sqrt{\lambda}$.

Proof:

$$\begin{aligned} b_t &= \sum_{s=1}^t X_s R_s = \sum_{s=1}^t X_s [\langle X_s, \theta^* \rangle + \eta_s] \\ &= \sum_{s=1}^t X_s^\top X_s \theta^* + \sum_{s=1}^t X_s \eta_s = S_t + \sum_{s=1}^t X_s^\top X_s \theta^*. \end{aligned}$$

$$\implies \hat{\theta}_t = V_t^{-1} b_t = V_t^{-1} S_t + V_t^{-1} \left[\sum_{s=1}^t X_s^\top X_s \right] \theta^* = V_t^{-1} S_t + V_t^{-1} K_t \theta^*$$

$$\begin{aligned} \|\theta^* - \hat{\theta}_t\|_{V_t} &= \|V_t^{-1} S_t + (V_t^{-1} K_t - I_d) \theta^*\|_{V_t} \leq \|S_t\|_{V_t^{-1}} + \sqrt{\theta^{*\top} (V_t^{-1} K_t - I_d) \underbrace{(K_t - V_t)}_{=-\lambda I_d} \theta^*} \\ &= \|S_t\|_{V_t^{-1}} + \sqrt{\lambda} \sqrt{\theta^{*\top} (I_d - V_t^{-1} K_t) \theta^*} \quad (\text{Since } \theta^{*\top} [V_t^{-1} K_t] \theta^* \geq 0) \end{aligned}$$

$$\implies \|\theta^* - \hat{\theta}_t\|_{V_t} \leq \|S_t\|_{V_t^{-1}} + \sqrt{\lambda} \|\theta^*\| \leq \|S_t\|_{V_t^{-1}} + \sqrt{\lambda} \quad \square$$

Linear UCB – Regret Analysis

Part (ii): If $S_t := \sum_{s=1}^t \eta_s X_s$ and $K_t := \sum_{s=1}^t X_s X_s^\top$, $M_t(z) = \exp\left(\langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2\right)$ is a non-negative super-martingale with $M_0(z) = 1$.

Proof: It is clear that $M_t(z) = \exp\left(\langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2\right)$ is non-negative and $M_0(z) = 1$. By our assumption on the noise, $\mathbb{E}[\exp(\omega \eta_t) | \mathcal{H}_{t-1}] \leq \exp\left(\frac{\omega^2}{2}\right)$. Setting $\omega = \langle z, X_t \rangle$, implies that

$$\mathbb{E}[\exp(\langle z, X_t \rangle \eta_t) | \mathcal{H}_{t-1}] \leq \exp\left(\frac{\|z\|_{X_t X_t^\top}^2}{2}\right) \implies \mathbb{E}\left[\exp\left(\langle z, X_t \rangle \eta_t - \frac{\|z\|_{X_t X_t^\top}^2}{2}\right) | \mathcal{H}_{t-1}\right] \leq 1 (*).$$

$$\begin{aligned}\mathbb{E}[M_t(z) | \mathcal{H}_{t-1}] &= \mathbb{E}\left[\exp\left(\langle z, S_{t-1} + \eta_t X_t \rangle - \frac{1}{2} \|z\|_{K_{t-1} + X_t X_t^\top}^2\right) | \mathcal{H}_{t-1}\right] \\ &= \mathbb{E}\left[\exp\left(\langle z, \eta_t X_t \rangle - \frac{1}{2} \|z\|_{X_t X_t^\top}^2\right) | \mathcal{H}_{t-1}\right] \exp\left(\langle z, S_{t-1} \rangle - \frac{1}{2} \|z\|_{K_{t-1}}^2\right) \\ &= M_{t-1}(z) \mathbb{E}\left[\exp\left(\langle z, \eta_t X_t \rangle - \frac{1}{2} \|z\|_{X_t X_t^\top}^2\right) | \mathcal{H}_{t-1}\right]\end{aligned}$$

$$\implies \mathbb{E}[M_t(z) | \mathcal{H}_{t-1}] \leq M_{t-1}(z) \quad (\text{Using } (*))$$

Linear UCB – Regret Analysis

Fact 1: For a probability density h , if $M_t(z)$ is a non-negative super-martingale with $M_0(z) = 1$, the “mixture” $\bar{M}_t := \int_z M_t(z) h(z) dz$ is also a non-negative super-martingale with $\bar{M}_0 = 1$.

Fact 2: For a non-negative super-martingale \bar{M}_t s.t. $\bar{M}_0 = 1$, for any $\epsilon > 0$, $\Pr[\sup_{t \in [T]} \bar{M}_t \geq \epsilon] \leq \frac{1}{\epsilon}$.

In order to construct \bar{M}_t , we will choose $h = \mathcal{N}(0, H^{-1})$ and $H = \lambda I_d$.

$$\bar{M}_t = \int_z M_t(z) h(z) dz = \frac{1}{\sqrt{(2\pi)^d \det[H^{-1}]}} \int_z \exp\left(\langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2 - \frac{1}{2} \|z\|_H^2\right) dz$$

From **Fact 1**, \bar{M}_t is a non-negative super-martingale, and hence using **Fact 2** with $\epsilon = 1/\delta$

$$\Pr\left[\sup_{t \in [T]} \bar{M}_t \geq \epsilon\right] = \Pr\left[\sup_{t \in [T]} \log(\bar{M}_t) \geq \log(\epsilon)\right] = \Pr\left[\sup_{t \in [T]} \log(\bar{M}_t) \geq \log(1/\delta)\right] \leq \delta$$

In the last part of the proof, we will relate \bar{M}_t to $\|S_t\|_{V_t^{-1}}$.

Linear UCB – Regret Analysis

Recall that $\bar{M}_t = \int_z M_t(z) h(z) dz = \frac{1}{\sqrt{(2\pi)^d \det[H^{-1}]}} \int_z \exp \left(\langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2 - \frac{1}{2} \|z\|_H^2 \right) dz$.

Simplifying the term inside exp,

$$\begin{aligned} \langle z, S_t \rangle - \frac{1}{2} \|z\|_{K_t}^2 - \frac{1}{2} \|z\|_H^2 &= \frac{1}{2} \|S_t\|_{(K_t+H)^{-1}}^2 - \frac{1}{2} \|z - (K_t + H)^{-1} S_t\|_{(K_t+H)}^2 \\ \implies \int_z M_t(z) h(z) dz &= \frac{\exp \left(\frac{1}{2} \|S_t\|_{V_t^{-1}}^2 \right)}{\sqrt{(2\pi)^d \det[H^{-1}]}} \int_z \exp \left(-\frac{1}{2} \|z - V_t^{-1} S_t\|_{V_t}^2 \right) dz \end{aligned}$$

The integral corresponds to the integral of the PDF for a multivariate Gaussian with mean $V_t^{-1} S_t$ and covariance V_t^{-1} . For a Gaussian with mean μ and covariance Σ^{-1} ,

$\frac{1}{\sqrt{(2\pi)^d \det[\Sigma^{-1}]}} \int_z \exp \left(-\frac{1}{2} \|z - \mu\|_{\Sigma}^2 \right) dz = 1$. Hence,

$$\bar{M}_t = \frac{\exp \left(\frac{1}{2} \|S_t\|_{V_t^{-1}}^2 \right)}{\sqrt{(2\pi)^d \det[H^{-1}]}} \sqrt{(2\pi)^d \det[V_t^{-1}]} = \sqrt{\frac{\det[H]}{\det[V_t]}} \exp \left(\frac{1}{2} \|S_t\|_{V_t^{-1}}^2 \right)$$

Linear UCB – Regret Analysis

Putting everything together, we know that for all $t \in [T]$, w.p $1 - \delta$, $\log(\bar{M}_t) \leq \log(1/\delta)$. Using the result from the previous slide, w.p $1 - \delta$, for all $t \in [T]$

$$\begin{aligned} \frac{1}{2} \|S_t\|_{V_t^{-1}}^2 + \frac{1}{2} \log\left(\frac{\det[H]}{\det[V_t]}\right) &\leq \log(1/\delta) \implies \|S_t\|_{V_t^{-1}} \leq \sqrt{\log\left(\frac{\det[V_t]}{\lambda^d}\right) + 2\log(1/\delta)} \\ &\implies \|S_t\|_{V_t^{-1}} \leq \sqrt{d \log\left(\frac{\lambda d + t}{\lambda d}\right) + 2\log(1/\delta)} \end{aligned}$$

From **Part (i)**, we know that,

$$\|\theta^* - \hat{\theta}_t\|_{V_t} \leq \|S_t\|_{V_t^{-1}} + \sqrt{\lambda} \leq \underbrace{\sqrt{d \log\left(\frac{\lambda d + t}{\lambda d}\right) + 2\log(1/\delta)}}_{:=\sqrt{\beta_t}} + \sqrt{\lambda}$$

We have shown that w.p. $1 - \frac{1}{T}$, for all $t \in [T]$ $\|\theta^* - \hat{\theta}_t\|_{V_t}^2 \leq \beta_t$, and hence $\Pr[G^c] \leq \frac{1}{T}$ \square

Improvements to LinUCB

- LinUCB results in $O(d\sqrt{T}\log(T))$ regret. Importantly, the same regret analysis works for infinitely many arms and even for a potentially changing set of actions \mathcal{A}_t .
- When the number of arms is finite, fixed and equal to K , a phase-based elimination algorithm can achieve $O(\sqrt{dT\log(KT)})$ regret (see [LS20, Chapter 22]).
- **Lower Bound:** For any bandit algorithm, there exists a linear bandit instance (with the set of actions \mathcal{A} equal to a unit hyper-cube or a unit sphere) such that $\text{Regret}(T) = \Omega(d\sqrt{T})$ (see [LS20, Chapter 24]).
- LinUCB maintain confidence intervals, and ensures optimism. An alternative set of strategies that work better in practice is *Posterior Sampling* of which Thompson Sampling is the most common (see [LS20, Chapter 36]).

Markov Decision Processes

Markov Decision Processes (MDPs)

- In bandit problems, the “state” of the environment does not change *as a result of an action*.
- Applications in robotics, operations research or conversational agents require explicitly modelling the current information available in a round.
- *Example 1*: A robot needs to model what is its position, velocity in order to take an action at the next round. This information is summarized as the “state” of the environment. The robot’s action changes its velocity, position and hence the “state”.
- *Example 2*: A conversational agent requires context (the past conversation, who it is speaking to) in order to decide what to respond to a particular user. The agent’s action can change the context of the conversation, and hence the “state”.
- Markov Decision Processes (MDPs) is the standard approach to sequential decision-making in such applications.

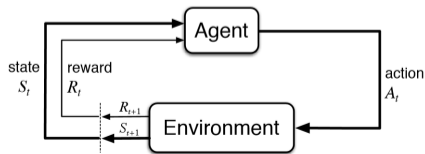
Markov Decision Processes (MDPs)

An MDP can be described by 5 elements: the state space (\mathcal{S}), action space (\mathcal{A}), starting state distribution (ρ), transition probabilities (\mathcal{P}) and rewards (r).

- State space \mathcal{S}
 - A state summarizes all the relevant information available to the agent. We will assume that the states are fully observable.
 - *Example:* Position of the rover on Mars, Inventory level of products.
 - States are mutually exclusive and exhaustive.
 - We will assume that the state space is discrete and finite, and $|\mathcal{S}| = S$.
- Starting state distribution $\rho \in \Delta_{\mathcal{S}}$:
 - $\rho(s)$ corresponds to the probability that the agent starts in state s . $\sum_{s \in \mathcal{S}} \rho(s) = 1$.
- Action space \mathcal{A} :
 - Consists of the actions an agent can take. The action space can be different in each state.
 - *Example:* Move north for the Mars rover, buy more stock of a particular product.
 - We will assume that \mathcal{A} is fixed, discrete and finite, and $|\mathcal{A}| = A$.

Markov Decision Processes (MDPs)

- Transition probabilities \mathcal{P} :
 - Model the inherent stochasticity in the system.
 - $\mathcal{P}(s'|s, a)$ is the probability of moving to a state s' when taking action a in state s .
 - $\sum_{s' \in \mathcal{S}} \mathcal{P}(s'|s, a) = 1$; $\mathcal{P}(s'|s, a) \geq 0$.
 - **Markov property:** $\mathcal{P}(s'|s, a)$ only depends on the current state s and action a .
 - In some examples, such as robotics, transitions can be deterministic.
 - If \mathcal{P} does not change, the transition probabilities are referred to as *stationary*.
- Rewards r : Model how much the agent has moved towards achieving its goal.
 - $r(s, a)$ is the reward obtained on taking action a in state s .
 - The reward can depend on s' , the state to which the agent transitioned to and is denoted as $r(s', a, s)$. In this case, $r(s, a) = \sum_{s' \in \mathcal{S}} r(s', a, s) \mathcal{P}(s'|s, a)$.



Protocol: At round (epoch) t , the agent observes state s_t and takes action a_t , transitions to state s_{t+1} and receives reward r_t .

Markov Decision Processes (MDPs)

Decision Rule: Describes the information and mechanism an agent uses to select an action in a given state and round. Can be classified as follows:

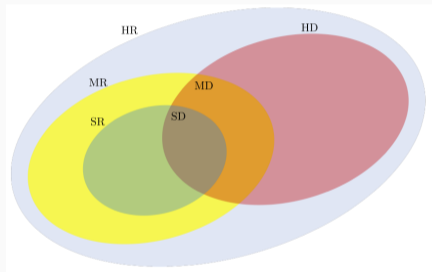
- *Information:* History dependent vs Markovian
 - A *history-dependent* decision rule uses some or all of the previous states and actions up to and including the current state when choosing an action.
 - A *Markovian* decision rule uses only the current state to select actions.
- *Mechanism:* Randomized vs Deterministic
 - A *randomized* decision rule maintains a probability distribution over the actions that can be taken in each state.
 - A *deterministic* decision rule corresponds to a degenerate distribution and consists of a deterministic mapping from states to actions.
- We define π_t to be the decision-rule at round t .
- A **policy** π is a sequence of decision rules, one for each round t , i.e. $\pi = (\pi_0, \pi_1, \pi_2, \dots)$.

Q: Why are history dependent policies computationally expensive to implement in general? **Ans:**

Need to look at the whole sequence of states and actions to decide which action to take.

Markov Decision Processes (MDPs)

- The **policy class** depends on the decision rule it uses.
- A policy can be in Π_{HR} , Π_{HD} , Π_{MR} , Π_{MD} depending on whether the decision rule is history-dependent (H) or Markovian (M); randomized (R) or deterministic (D).
- *Example:* If $\mathcal{H}_t = \{S_0, A_0, S_1, \dots, S_t\}$ is the history of interactions until round t , then, $\Pi_{HR} = \{\pi_0, \pi_1, \pi_2, \dots\}$ where $\pi_t : \mathcal{H}_t \rightarrow \Delta_A$,
- *Example:* $\Pi_{MD} = \{\pi_0, \pi_1, \pi_2, \dots\}$ where $\pi_t : S_t \rightarrow \mathcal{A}$.
- A policy is *stationary* if it uses the same decision rule in every round, i.e. $\pi = \{\pi_0, \pi_0, \dots\}$.
- We will only consider stationary policies that are Markovian, and define $\Pi_{SR} \subset \Pi_{MR} \subset \Pi_{HR}$ and $\Pi_{SD} \subset \Pi_{MD} \subset \Pi_{HD}$.



Markov Decision Processes (MDPs)

- Specifying ρ and choosing a policy π results in a stochastic process over the state and action space. We will denote this *trajectory* as (S_0, A_0, S_1, \dots) .
- When $\pi \in \Pi_{MR}$, the stochastic process is a discrete-time Markov chain.

Q: For a policy $\pi \in \Pi_{MR}$, calculate the probability of the trajectory $(s_0, a_0, s_1, a_1, \dots)$ **Ans:**
 $\Pr[(s_0, a_0, s_1, a_1, \dots)] = \rho(s_0) \pi_0(a_0|s_0) \mathcal{P}(s_1|s_0, a_0) \pi_1(a_1|s_1) \dots$

- The trajectory over states and actions generates a *reward process*:
 $(R_0, R_1, \dots) = (r(S_0, A_0), r(S_1, A_1), \dots)$.
- When the stochastic process over states-actions is a Markov chain, the corresponding reward process is a *Markov reward process*.

Q: How do we judge whether one reward process is “better” than the other?

We need some notion of *utility*. Common choice of utility functions is *additive*, i.e. the utility of a reward process (r_0, r_1, \dots) is given by: $\mathbb{E}[\sum_{i=0} U_i(R_i)]$ where $U_i : \mathbb{R} \rightarrow \mathbb{R}$ and the expectation is over the different trajectories produced by the policy.

Markov Decision Processes (MDPs)

Choosing a policy gives rise to a reward process, and we can design additive utility functions to compare different reward processes. This gives different optimality criterion w.r.t to policies:

- (a) For a finite *horizon* H , $\max_{\pi \in \Pi_{\text{HR}}} \mathbb{E} \left[\sum_{t=1}^H R_t \right]$. [**Finite Horizon Total Reward**]
- (b) For an infinite horizon, $\max_{\pi \in \Pi_{\text{HR}}} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right]$ where $\gamma \in (0, 1)$ is the discount factor.
[**Infinite Horizon Discounted Reward**]
- (c) For an infinite horizon, $\max_{\pi \in \Pi_{\text{HR}}} \lim_{T \rightarrow \infty} \frac{\mathbb{E} \left[\sum_{t=0}^T R_t \right]}{T}$. [**Infinite Horizon Average Reward**]

We will focus mainly on (b) infinite horizon discounted reward setting and towards the end, consider (a) finite horizon total reward setting.

Infinite Horizon Discounted Reward: The discount factor γ models the fact that near-term rewards are preferable to future rewards. For example, it models inflation meaning that a penny today is worth 10 in the future.

Infinite-horizon Discounted Setting

Objective: For a starting state s_0 , find policy $\pi \in \Pi_{\text{HR}}$ that maximizes the **value function** $v^\pi(s_0)$

$$v^\pi(s_0) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s_0 \right],$$

where the expectation is over the randomness in the reward process induced by policy π . For a starting state distribution ρ , the related objective is to maximize $v^\pi(\rho) := \mathbb{E}_{s \sim \rho} v^\pi(s)$.

Assumptions:

- The reward function does not change across rounds.
- The rewards are bounded in $[0, 1]$.

Q: What are the upper and lower-bounds on the value function? **Ans:** $\frac{1}{1-\gamma}$ and 0

Infinite-horizon Discounted Setting

Claim: For each $s \in \mathcal{S}$, for a given policy $\pi = (\pi_0, \pi_1, \dots) \in \Pi_{\text{HR}}$, there exists a policy $\pi' = (\pi'_0, \pi'_1, \dots) \in \Pi_{\text{MR}}$ with the same value, conditioned on $S_0 = s_0$.

- Since there exists a Markov policy that has the same value as every history-dependent policy, we only need to consider Π_{MR} when we optimize for the optimal policy.
- Markov policies only need to maintain the knowledge of the current state, and are hence computationally tractable.

Proof: Using the definition of the value function,

$$\begin{aligned} v^\pi(s_0) &:= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s_0 \right] = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) \mid S_0 = s_0 \right] \\ &= \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(s, a) \Pr^\pi [S_t = s, A_t = a \mid S_0 = s_0] \end{aligned}$$

Here, \Pr^π corresponds to the probability distribution induced by policy π .

Infinite-horizon Discounted Setting

Recall that $v^\pi(s_0) = \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(s, a) \Pr^\pi[S_t = s, A_t = a | S_0 = s_0]$

Construct $\pi' \in \Pi_{MR}$ as follows: $\pi'_t(A_t = a | S_t = s) = \Pr^\pi[A_t = a | S_t = s, S_0 = s_0]$.

We will prove (by induction) that $\Pr^\pi[S_t = s, A_t = a | S_0 = s_0] = \Pr^{\pi'}[S_t = s, A_t = a | S_0 = s_0]$, and hence, $v^\pi(s_0) = v^{\pi'}(s_0)$.

$$v^\pi(s_0) = \sum_{t=0}^{\infty} \gamma^t \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(s, a) \Pr^\pi[A_t = a | S_t = s, S_0 = s_0] \Pr^\pi[S_t = s | S_0 = s_0]$$

Base Case: For $t = 0$, $\sum_a \pi_0[A_0 = a | S_0 = s_0] = \sum_a \pi'_0(A_0 = a | S_0 = s_0)$ by def. of π' .

Inductive Hypothesis: For $t \geq 1$, assume that

$\Pr^\pi[S_t = s, A_t = a | S_0 = s_0] = \Pr^{\pi'}[S_t = s, A_t = a | S_0 = s_0]$. Let us now prove it for $t + 1$.

Infinite-horizon Discounted Setting

Recall that $\pi'_t(A_t = a|S_t = s) = \Pr^\pi[A_t = a|S_t = s_t, S_0 = s_0]$ by def. of π' , and by the inductive hypothesis, $\Pr^\pi[S_t = s, A_t = a|S_0 = s_0] = \Pr^{\pi'}[S_t = s, A_t = a|S_0 = s_0]$.

For a fixed (s, a) , using the definition of π' ,

$$\begin{aligned} & \Pr^\pi[A_{t+1} = a|S_{t+1} = s, S_0 = s_0] \Pr^\pi[S_{t+1} = s|S_0 = s_0] \\ &= \Pr^{\pi'}[A_{t+1} = a|S_{t+1} = s, S_0 = s_0] \Pr^\pi[S_{t+1} = s|S_0 = s_0] \end{aligned}$$

We need to show that $\Pr^\pi[S_{t+1} = s|S_0 = s_0] = \Pr^{\pi'}[S_{t+1} = s|S_0 = s_0]$. For an arbitrary $s' \in \mathcal{S}$,

$$\begin{aligned} \Pr^\pi[S_{t+1} = s'|S_0 = s_0] &= \sum_s \sum_a \mathcal{P}[s'|s, a] \Pr^\pi[S_t = s, A_t = a|S_0 = s_0] \\ &= \sum_s \sum_a \mathcal{P}[s'|s, a] \Pr^{\pi'}[S_t = s, A_t = a|S_0 = s_0] && \text{(Inductive Hypothesis)} \\ &= \Pr^{\pi'}[S_{t+1} = s'|S_0 = s_0] \end{aligned}$$

Hence, $\Pr^\pi[S_{t+1} = s, A_{t+1} = a|S_0 = s_0] = \Pr^{\pi'}[S_{t+1} = s, A_{t+1} = a|S_0 = s_0]$. Using the definition of v^π , $v^\pi(s_0) = v^{\pi'}(s_0)$. □

 Tor Lattimore and Csaba Szepesvári, *Bandit algorithms*, Cambridge University Press, 2020.