

Exploration in Tabular MDPs

Sham Kakade and Kianté Brantley

CS 2824: Foundations of Reinforcement Learning

Announcements

1. HW2: Out 02/24, Due 03/13

2. Course Project Website

<https://harvard-cs2824-s26.github.io/CS2824projects.html>

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$

(Handwritten notes: $\pi^n, \hat{p}, r+\gamma$ with arrows pointing to $\widehat{V}_0^n(s_0)$)

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

Then π^n is close to π^\star , i.e., we are doing exploitation

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

Then π^n is close to π^\star , i.e., we are doing exploitation

2. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$?

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

Then π^n is close to π^\star , i.e., we are doing exploitation

2. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$?

$$\epsilon \leq \underbrace{\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)} \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[\underbrace{b_h^n(s,a)} + \underbrace{(\widehat{P}_h^n(\cdot | s,a) - P_h(\cdot | s,a)) \cdot \widehat{V}_{h+1}^n} \right]$$

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

Then π^n is close to π^\star , i.e., we are doing exploitation

2. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$?

$$\epsilon \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[b_h^n(s,a) + (\widehat{P}_h^n(\cdot | s,a) - P_h(\cdot | s,a)) \cdot \widehat{V}_{h+1}^n \right]$$

We collect data at steps where bonus is large or model is wrong, i.e., exploration

1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

$f: V^*$

Given a fixed function $f: S \mapsto [0, H]$, w/ prob $1 - \delta$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O\left(H \sqrt{\ln(SAHN/\delta) / N_h^n(s, a)}\right), \forall s, a, h, N$$

$O\left(H \sqrt{\frac{1}{N_h^n(s, a)}}\right)$

1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function $f : S \mapsto [0, H]$, w/ prob $1 - \delta$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O\left(H \sqrt{\ln(SAHN/\delta) / N_h^n(s, a)}\right), \forall s, a, h, N$$

Bonus $b_h^n(s, a)$

1. Model Error using Hoeffding's inequality & Union Bound

$$\widehat{P}_h^n(s' | s, a) = \frac{N_h^n(s, a, s')}{N_h^n(s, a)}, \forall h, s, a, s'$$

Given a fixed function $f: S \mapsto [0, H]$, w/ prob $1 - \delta$:

$$\left| \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right)^\top f \right| \leq O\left(H \sqrt{\ln(SAHN/\delta) / N_h^n(s, a)}\right), \forall s, a, h, N$$


Bonus $b_h^n(s, a)$

From now on, assume this event being true

2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n(H) \right\}$$
$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$


2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$
$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \quad \forall s$

2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

$Q(s,a) \in \mathbb{H} \quad \forall s, a$

$V(s) \in \mathbb{H}$

$\widehat{Q} \geq Q^* \Rightarrow \widehat{V}_h \geq V_h^*$

Two cases

① $\widehat{Q} = H$

② $\widehat{Q} = r + b + \widehat{P}V$

Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \quad \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = \underbrace{r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n}_{\text{update}} - \underbrace{r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^*}_{\text{Bell eq.}}$$

2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \quad \forall s$

$$\begin{aligned} \widehat{Q}_h^n(s, a) - Q_h^*(s, a) &= \cancel{r_h(s, a)} + \cancel{b_h^n(s, a)} + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - \cancel{r_h(s, a)} - P_h(\cdot | s, a) \cdot V_{h+1}^* \\ &\geq \underline{b_h^n(s, a)} + \widehat{P}_h^n(\cdot | s, a) \cdot \underline{V_{h+1}^*} - P_h(\cdot | s, a) \cdot \underline{V_{h+1}^*} \end{aligned}$$

2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \quad \forall s$

$$\begin{aligned} \widehat{Q}_h^n(s, a) - Q_h^*(s, a) &= r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^* \\ &\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot V_{h+1}^* - P_h(\cdot | s, a) \cdot V_{h+1}^* \\ &= b_h^n(s, a) + \left| \widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right| \cdot V_{h+1}^* \leq 4\sqrt{\frac{1}{N_h^n(s, a)}} \approx b_h^n \end{aligned}$$

2. Proving Optimism via Induction

Lemma [Optimism]: $\widehat{V}_h^n(s) \geq V_h^*(s), \forall n, h, s$

Recall Bonus-enhanced Value Iteration at episode n:

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Inductive hypothesis: $\widehat{V}_{h+1}^n(s) \geq V_{h+1}^*(s), \quad \forall s$

$$\widehat{Q}_h^n(s, a) - Q_h^*(s, a) = r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n - r_h(s, a) - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot V_{h+1}^* - P_h(\cdot | s, a) \cdot V_{h+1}^*$$

$$= b_h^n(s, a) + \left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot V_{h+1}^*$$

$$\geq b_h^n(s, a) - b_h^n(s, a) = 0, \quad \forall s, a$$

$$\widehat{Q}_h^n - Q_h^* \geq 0$$

$$\widehat{Q}_h^n \geq Q_h^*$$

$$\Rightarrow \widehat{V}_h^n \geq V_h^*$$

3. Upper Bounding Regret using Optimism

$$\text{per-episode regret} := V_0^*(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

The diagram includes several annotations: a blue circle around $V_0^*(s_0)$, a green circle around $\widehat{V}_0^n(s_0)$, and a blue circle around $V_0^{\pi^n}(s_0)$. A blue line points from the text $\pi^n, \hat{P}, r + b$ to the green circle. Another blue line points from the text π^n, ρ, Γ to the blue circle around $V_0^{\pi^n}(s_0)$.

This is something
we can control!
And this is related
to our policy π^n

Recall simulation lemma — the lemma measures the difference of a policy under two MDPs

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \underline{\pi}_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\underline{\widehat{V}}_0^n(s_0) - V_0^{\pi^n}(s_0) = \underline{\widehat{Q}}_0^n(s_0, \underline{\pi}^n(s_0)) - \underline{Q}_0^{\pi^n}(s_0, \underline{\pi}^n(s_0))$$

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\} \leq a$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0))$$

$$\leq \underbrace{r_0(s_0, \pi^n(s_0)) + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n}_{\text{update}} - \underbrace{r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}}_{\text{hel equ.}}$$

$\min(a, b) \leq a$

update

hel equ.

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0))$$

$$\leq \underbrace{r_0(s_0, \pi^n(s_0))}_{\text{blue}} + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - \underbrace{r_0(s_0, \pi^n(s_0))}_{\text{blue}} - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

$$= \underbrace{b_h^n(s_0, \pi^n(s_0))}_{\text{blue}} + \boxed{\widehat{P}_0^n(\cdot | s_0, \pi^n(s_0))} \cdot \widehat{V}_1^n - \boxed{P_0(\cdot | s_0, \pi^n(s_0))} \cdot V_1^{\pi^n}$$

+/-

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

step 1

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0))$$

$$\leq r_0(s_0, \pi^n(s_0)) + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

$$= b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

$$= b_h^n(s_0, \pi^n(s_0)) + \left(\widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \right) \cdot \widehat{V}_1^n + P_0(\cdot | s_0, \pi^n(s_0)) \cdot \left(\widehat{V}_1^n - V_1^{\pi^n} \right)$$

step 2

4. Upper bounding Regret via Simulation Lemma

$$\widehat{V}_H^n(s) = 0, \quad \widehat{Q}_h^n(s, a) = \min \left\{ r_h(s, a) + b_h^n(s, a) + \widehat{P}_h^n(\cdot | s, a) \cdot \widehat{V}_{h+1}^n, H \right\}$$

$$\widehat{V}_h^n(s) = \max_a \widehat{Q}_h^n(s, a), \quad \pi_h^n(s) = \arg \max_a \widehat{Q}_h^n(s, a), \forall s$$

Lemma [Simulation lemma]:

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) = \widehat{Q}_0^n(s_0, \pi^n(s_0)) - Q_0^{\pi^n}(s_0, \pi^n(s_0))$$

$$\leq r_0(s_0, \pi^n(s_0)) + b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - r_0(s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

$$= b_h^n(s_0, \pi^n(s_0)) + \widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) \cdot \widehat{V}_1^n - P_0(\cdot | s_0, \pi^n(s_0)) \cdot V_1^{\pi^n}$$

$$= b_h^n(s_0, \pi^n(s_0)) + \left(\widehat{P}_0^n(\cdot | s_0, \pi^n(s_0)) - P_0(\cdot | s_0, \pi^n(s_0)) \right) \cdot \widehat{V}_1^n + P_0(\cdot | s_0, \pi^n(s_0)) \cdot \left(\widehat{V}_1^n - V_1^{\pi^n} \right)$$

$$= \sum_{h=0}^{H-1} \mathbb{E}_{s, a \sim d_h^{\pi^n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

4. Upper bounding Regret via Simulation Lemma

$$\text{per-episode regret} := V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$$

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

$$|(\widehat{P} - P^*)V^+| \leq H \sqrt{\frac{1}{N_h^{\text{tr}}(s,a)}}$$

4. Upper bounding Regret via Simulation Lemma

per-episode regret $:= V_0^\star(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent
(this is different from V_h^\star) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

4. Upper bounding Regret via Simulation Lemma

per-episode regret $:= V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent (this is different from V_h^*) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\|a \circ b\| \leq \|a\|_f \cdot \|b\|_g \quad v \in [0, 1]$$

$$\left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\|V\|_\infty \leq 1$$

4. Upper bounding Regret via Simulation Lemma

per-episode regret $:= V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent
(this is different from V_h^*) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \| \widehat{V}_{h+1}^n \|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{ with prob } 1 - \delta$$

4. Upper bounding Regret via Simulation Lemma

per-episode regret := $V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent (this is different from V_h^*) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + \underbrace{(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n}_{\text{purple bracket}} \right]$$

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$b_h^n = H \sqrt{\frac{f_h^{h-1}}{N_h^n(s, a)}} \leq H \sqrt{\frac{S \ln(\dots)}{N_h^n(s, a)}}$$

$$\left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{ with prob } 1 - \delta$$

4. Upper bounding Regret via Simulation Lemma

per-episode regret := $V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent (this is different from V_h^*) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$\leq 2 \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$\left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{ with prob } 1 - \delta$$

4. Upper bounding Regret via Simulation Lemma

per-episode regret := $V_0^*(s_0) - V_0^{\pi_n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi_n}(s_0)$ But \widehat{V}_h^n is data-dependent (this is different from V_h^*) !!!

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + (\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a)) \cdot \widehat{V}_{h+1}^n \right]$$

Let's do Holder's inequality

$$\leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[b_h^n(s, a) + H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right]$$

$$\leq 2 \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}} \right] = 2H \sqrt{S \ln(SAHN/\delta)} \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi_n}} \left[\sqrt{\frac{1}{N_h^n(s, a)}} \right]$$

$$\left(\widehat{P}_h^n(\cdot | s, a) - P_h(\cdot | s, a) \right) \cdot \widehat{V}_{h+1}^n \leq \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \|\widehat{V}_{h+1}^n\|_\infty$$

$$\leq H \|P_h(\cdot | s, a) - \widehat{P}_h^n(\cdot | s, a)\|_1 \leq H \sqrt{\frac{S \ln(SAHN/\delta)}{N_h^n(s, a)}}, \forall s, a, h, n, \text{ with prob } 1 - \delta$$

5. Final Step

Remember we had two failure events for bounding transitions errors.

5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\mathbb{E} [\text{Regret}_N] = \mathbb{E} \left[\mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}\{\text{events don't hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right]$$

5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\begin{aligned}\mathbb{E} [\text{Regret}_N] &= \mathbb{E} \left[\mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}\{\text{events don't hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\leq \mathbb{E} \left[\mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \underbrace{\mathbb{P}(\text{events don't hold}) \cdot NH}\end{aligned}$$

5. Final Step

Remember we had two failure events for bounding transitions errors.

$$\begin{aligned}\mathbb{E} [\text{Regret}_N] &= \mathbb{E} \left[\mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{E} \left[\mathbf{1}\{\text{events don't hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] \\ &\leq \mathbb{E} \left[\mathbf{1}\{\text{events hold}\} \sum_{n=1}^N (V_0^*(s_0) - V_0^{\pi^n}(s_0)) \right] + \mathbb{P}(\text{events don't hold}) \cdot NH \\ &\leq \underbrace{H\sqrt{S \ln(SANH/\delta)}} \mathbb{E} \left[\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^h)}} \right] + \underbrace{2\delta NH}\end{aligned}$$

5. Final Step

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}}$$

5. Final Step

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \sum_{h=0}^{H-1} \underbrace{\sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}}}$$

5. Final Step

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \sum_{h=0}^{H-1} \sum_{s,a} \boxed{\sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \boxed{\sqrt{N_h^N(s,a)}}$$

$\sum_{i=1}^N \frac{1}{\sqrt{i}} \leq 2\sqrt{N}$

5. Final Step

$$\sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} = \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)}$$

Ca - inequality

$$\begin{aligned} & \sum_{s,a} \sqrt{N_h^N(s,a)} \\ & \leq \sqrt{\sum_{s,a} 1} \cdot \sqrt{\sum_{s,a} N_h^N(s,a)} \\ & \quad \parallel \\ & \sqrt{SA} \end{aligned}$$

5. Final Step

$$\begin{aligned} \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)} \\ &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN} \end{aligned}$$

5. Final Step

$$\begin{aligned} \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)} \\ &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN} \end{aligned}$$

$$\mathbb{E} [\text{Regret}_N] \leq \underbrace{2H^2 S \sqrt{AN \ln(SAHN/\delta)}} + 2\delta NH$$

5. Final Step

$$\begin{aligned}
 \sum_{n=1}^N \sum_{h=0}^{H-1} \frac{1}{\sqrt{N_h^n(s_h^n, a_h^n)}} &= \sum_{h=0}^{H-1} \sum_{s,a} \sum_{i=1}^{N_h^N(s,a)} \frac{1}{\sqrt{i}} \leq \sum_{h=0}^{H-1} \sum_{s,a} \sqrt{N_h^N(s,a)} \leq \sum_{h=0}^{H-1} \sqrt{SA \sum_{s,a} N_h^N(s,a)} \\
 &\leq \sum_{h=0}^{H-1} \sqrt{SAN} = H\sqrt{SAN}
 \end{aligned}$$

$$\mathbb{E} [\text{Regret}_N] \leq 2H^2S\sqrt{AN \ln(SAHN/\delta)} + 2\delta NH \quad \text{Set } \delta = 1/(HN)$$

$$\leq 2H^2S\sqrt{AN \cdot \ln(SAH^2N^2)} = \boxed{\tilde{O}\left(H^2S\sqrt{AN}\right)}$$

High-level Idea: Exploration or Exploitation Tradeoff

Upper bound per-episode regret: $V_0^\star(s_0) - V_0^{\pi^n}(s_0) \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$

1. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \epsilon$?

Then π^n is close to π^\star , i.e., we are doing exploitation

2. What if $\widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \geq \epsilon$?

$$\epsilon \leq \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_h^{\pi^n}} \left[b_h^n(s,a) + (\widehat{P}_h^n(\cdot | s,a) - P_h(\cdot | s,a)) \cdot \widehat{V}_{h+1}^n \right]$$

We collect data at steps where bonus is large or model is wrong, i.e., exploration