Exploration: UCB-VI

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

- Recap
- UCB-VI for tabular MDPs
- UCB-VI for linear MDPs



Why we don't want to treat MDPs as big contextual bandits

Without context:

Without context:

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

Without context:

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

With tabular context ($|\mathcal{X}|$ distinct contexts):

Without context:

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

With tabular context ($|\mathcal{X}|$ distinct contexts):

$$\pi_t(x_t) = \arg\max_k \hat{\mu}_t^{(k)}(x_t) + \sqrt{\ln(2TK|\mathcal{X}|/\delta)/2N_t^{(k)}(x_t)}$$

V/Q functions in Finite horizon MDP



$$\sum_{\tau=h}^{H-1} r(s_{\tau}, a_{\tau}) \left| s_{h} = s \right]$$
$$r(s_{\tau}, a_{\tau}) \left| (s_{h}, a_{h}) = (s, a_{\tau}) \right|$$

V/Q functions in Finite horizon MDP



Recall: $V_h^{\pi}(s) \le H$, $Q_h^{\pi}(s, a) \le H$

$$\sum_{\tau=h}^{H-1} r(s_{\tau}, a_{\tau}) \left| s_{h} = s \right]$$

$$r(s_{\tau}, a_{\tau}) \left| (s_h, a_h) = (s, a) \right|$$

V/Q functions in Finite horizon MDP



Recall: $V_h^{\pi}(s) \leq$

Bellman Consistency Equation:

 $Q_h^{\pi}(s,a) = r(s,a)$

$$\sum_{\tau=h}^{d-1} r(s_{\tau}, a_{\tau}) \left| s_{h} = s \right]$$

$$r(s_{\tau}, a_{\tau}) \left| (s_h, a_h) = (s, a) \right|$$

$$\leq H, \qquad Q_h^{\pi}(s,a) \leq H$$

) +
$$\mathbb{E}_{s' \sim P(s,a)} \left[V_{h+1}^{\pi}(s') \right]$$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,a)$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^* = \{\pi_0^*, \pi_1^*, \dots, \pi_{H-1}^*\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,$

$$V_{H-1}^{\star} = \max_{a} Q_{H-1}^{\star}$$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

 $_{-1}(s,a) = Q_{H-1}^{\star}(s,\pi_{H-1}^{\star}(s))$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,$

$$V_{H-1}^{\star} = \max_{a} Q_{H-1}^{\star}(s, a) = Q_{H-1}^{\star}(s, \pi_{H-1}^{\star}(s))$$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

2. Assuming we have computed V_{h+1}^{\star} , $h \leq H - 2$, i.e., assuming we know how to perform optimally starting at h + 1, then:

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,a)$

$$V_{H-1}^{\star} = \max_{a} Q_{H-1}^{\star}(s, a) = Q_{H-1}^{\star}(s, \pi_{H-1}^{\star}(s))$$

 $Q_h^\star(s,a) = r(s,a)$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

2. Assuming we have computed V_{h+1}^{\star} , $h \leq H - 2$, i.e., assuming we know how to perform optimally starting at h + 1, then:

$$a) + \mathbb{E}_{s' \sim P(s,a)} V_{h+1}^{\star}(s')$$

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,a)$

$$V_{H-1}^{\star} = \max_{a} Q_{H-1}^{\star}(s, a) = Q_{H-1}^{\star}(s, \pi_{H-1}^{\star}(s))$$

 $Q_h^\star(s,a) = r(s,a)$

$$\pi_h^\star(s) = \arg\max_a Q_h^\star(s,$$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

2. Assuming we have computed V_{h+1}^{\star} , $h \leq H - 2$, i.e., assuming we know how to perform optimally starting at h + 1, then:

$$a) + \mathbb{E}_{s' \sim P(s,a)} V_{h+1}^{\star}(s')$$

a),

VI = DP is a backwards in time approach for computing the optimal policy: $\pi^{\star} = \{\pi_0^{\star}, \pi_1^{\star}, \dots, \pi_{H-1}^{\star}\}$

1. Start at H - 1,

 $Q_{H-1}^{\star}(s,a) = r(s,$

$$V_{H-1}^{\star} = \max_{a} Q_{H-1}^{\star}(s, a) = Q_{H-1}^{\star}(s, \pi_{H-1}^{\star}(s))$$

$$Q_h^{\star}(s,a) = r(s,a) + \mathbb{E}_{s' \sim P(s,a)} V_{h+1}^{\star}(s')$$
$$\pi_h^{\star}(s) = \arg\max_a Q_h^{\star}(s,a), \quad V_h^{\star} = \max_a Q_h^{\star}(s,a)$$

a)
$$\pi_{H-1}^{\star}(s) = \arg\max_{a} Q_{H-1}^{\star}(s, a)$$

2. Assuming we have computed V_{h+1}^{\star} , $h \leq H - 2$, i.e., assuming we know how to perform optimally starting at h + 1, then:

Summary on Finite horizon MDP

 $\mathcal{M} = \{S, A, r, P, H\},\$ $r: S \times A \mapsto [0,1], H \in \mathbb{N}, P: S \times A \mapsto \Delta(S)$

Comparing to the infinite horizon, discounted MDP:

- 1. Policy will be time dependent
- 2. DP takes H steps to compute π^{\star}
 - total computation time is $O(H|S|^2|A|)$
 - no need to use contraction argument and no discount factor
- 3. Extension to non-stationary setting works immediately:

(i.e. with a non-stationary transition model: $P_0(s' \mid s, a), P_1(s' \mid s, a), \dots P_{H-1}(s' \mid s, a)$)





- UCB-VI for tabular MDPs
- UCB-VI for linear MDPs

Why we don't want to treat MDPs as big contextual bandits

Q: given a discrete MDP, how many unique policies we have?



Q: given a discrete MDP, how many unique policies we have?

 $\left(|A|^{|S|} \right)^{H}$



Q: given a discrete MDP, how many unique policies we have?

So treating each policy as an "arm" and running UCB gives us regret $\tilde{O}(\sqrt{|A|^{|S|H}N})$

 $\left(|A|^{|S|} \right)^{H}$



Q: given a discrete MDP, how many unique policies we have?

This seems bad, so are MDPs just super hard or can we do better?

 $\left(|A|^{|S|} \right)^{H}$

So treating each policy as an "arm" and running UCB gives us regret $\tilde{O}(\sqrt{|A|^{|S|H}N})$



$$S = \{a, b\},\$$

All state transitions happen with probability 1/2 for all actions

 $A = \{1, 2\}, H = 2$

Reward function: r(a,1) = r(b,1) = 0r(a,2) = r(b,2) = 1

$$S = \{a, b\},\$$

All state transitions happen with probability 1/2 for all actions

Reward function

Suppose we have a lot of data already on a policy $\pi^{(1)}$ that always takes action 1 and a policy $\pi^{(2)}$ that always takes action 2 (note $\pi^{(2)} = \pi^*$)

 $A = \{1, 2\}, H = 2$

n:
$$r(a,1) = r(b,1) = 0$$

 $r(a,2) = r(b,2) = 1$

$$S = \{a, b\},\$$

All state transitions happen with probability 1/2 for all actions

Reward function

Suppose we have a lot of data already on a policy $\pi^{(1)}$ that always takes action 1 and a policy $\pi^{(2)}$ that always takes action 2 (note $\pi^{(2)} = \pi^{\star}$)

What do we know about a policy $\pi^{(3)}$ which always takes action 1 in the first time step, and

 $A = \{1, 2\}, H = 2$

n:
$$r(a,1) = r(b,1) = 0$$

 $r(a,2) = r(b,2) = 1$

always takes action 2 at the second time step?

All state transitions happen with probability 1/2 for all actions

Reward functior

Suppose we have a lot of data already on a policy $\pi^{(1)}$ that always takes action 1 and a policy $\pi^{(2)}$ that always takes action 2 (note $\pi^{(2)} = \pi^{\star}$)

What do we know about a policy $\pi^{(3)}$ which always takes action 1 in the first time step, and always takes action 2 at the second time step?

Everything: we have a lot of data on every state-action reward and transition!

 $S = \{a, b\}, A = \{1, 2\}, H = 2$

n:
$$r(a,1) = r(b,1) = 0$$

 $r(a,2) = r(b,2) = 1$

All state transitions happen with probability 1/2 for all actions

Reward functior

Suppose we have a lot of data already on a policy $\pi^{(1)}$ that always takes action 1 and a policy $\pi^{(2)}$ that always takes action 2 (note $\pi^{(2)} = \pi^{\star}$)

What do we know about a policy $\pi^{(3)}$ which always takes action 1 in the first time step, and always takes action 2 at the second time step?

Everything: we have a lot of data on every state-action reward and transition!

If we treat the MDP as a contextual bandit, we treat $\pi^{(3)}$ as a new "arm" about which we know nothing...

 $S = \{a, b\}, A = \{1, 2\}, H = 2$

n:
$$r(a,1) = r(b,1) = 0$$

 $r(a,2) = r(b,2) = 1$



All state transitions happen with probability 1/2 for all actions

Reward function

Suppose we have a lot of data already on a policy $\pi^{(1)}$ that always takes action 1 and a policy $\pi^{(2)}$ that always takes action 2 (note $\pi^{(2)} = \pi^{\star}$)

What do we know about a policy $\pi^{(3)}$ which always takes action 1 in the first time step, and always takes action 2 at the second time step?

Everything: we have a lot of data on every state-action reward and transition!

If we treat the MDP as a contextual bandit, we treat $\pi^{(3)}$ as a new "arm" about which we know nothing...

 $|A|^{|S|H} = 2^4 = 16$ $S = \{a, b\}, A = \{1, 2\}, H = 2$

n:
$$r(a,1) = r(b,1) = 0$$

 $r(a,2) = r(b,2) = 1$







- Why we don't want to treat MDPs as big contextual bandits
 - UCB-VI for tabular MDPs
 - UCB-VI for linear MDPs





Use all previous data to estimate transitions $\widehat{P}_{1}^{n}, \ldots, \widehat{P}_{H-1}^{n}$



- Use all previous data to estimate transitions $\widehat{P}_{1}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$



- Use all previous data to estimate transitions $\widehat{P}_{1}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$
- Optimistic planning with learned model: $\pi^n = \text{Value-Iter}\left(\{\widehat{P}_h^n, r_h + b_h^n\}_{h=1}^{H-1}\right)$



- Use all previous data to estimate transitions $\widehat{P}_{1}^{n}, \ldots, \widehat{P}_{H-1}^{n}$
 - Design reward bonus $b_h^n(s, a), \forall s, a, h$
- Optimistic planning with learned model: $\pi^n = \text{Value-Iter}\left(\{\widehat{P}_h^n, r_h + b_h^n\}_{h=1}^{H-1}\right)$
- Collect a new trajectory by executing π^n in the real world $\{P_h\}_{h=0}^{H-1}$ starting from s_0


$$\mathcal{D}_h^n = \{s_h^i\}$$

Model Estimation

 ${}_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

Model Estimation

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h, \quad N_h^n(s,a,s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s,a,s')\}, \forall s, h \in \mathbb{N}\}$$

Model Estimation



$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h$$

Let's also maintain some statistics using these datasets:

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h, \quad N_h^n(s,a,s') = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i) = (s,a,s')\}, \forall s, h \in \mathbb{N}\}$$

Estimate model \widehat{P}

$$\widehat{P}_{h}^{n}(s'|s,a) = \frac{N_{h}^{n}(s,a,s')}{N_{h}^{n}(s,a)}$$

Model Estimation

$$\widehat{P}_{h}^{n}(s'|s,a), \forall s,a,s',h$$
:



Let us consider the very beginning of episode *n*:

 $\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \Lambda$

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h,$$

Let us consider the very beginning of episode *n*:

 $\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \Lambda$

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h,$$



Let us consider the very beginning of episode *n*:

 $\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, I$

$$N_h^n(s,a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_h^i, a_h^i) = (s,a)\}, \forall s, a, h,$$



Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$



Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

 $\widehat{V}_{H}^{n}(s) = 0, \forall s$



Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$



Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$

$$\widehat{V}_{h}^{n}(s) = \max_{a} \ \widehat{Q}_{h}^{n}(s,a),$$



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

Let us consider the very beginning of episode *n*:

$$\mathcal{D}_{h}^{n} = \{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=1}^{n-1}, \forall h, \ N_{h}^{n}(s, a) = \sum_{i=1}^{n-1} \mathbf{1}\{(s_{h}^{i}, a_{h}^{i}) = (s, a)\}, \forall s, a, h,$$

$$\widehat{V}_{H}^{n}(s) = 0, \forall s \qquad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}, \forall s, a$$

$$\widehat{V}_{h}^{n}(s) = \max_{a} \ \widehat{Q}_{h}^{n}(s,a),$$



$$\pi_h^n(s) = \arg\max_a \widehat{Q}_h^n(s, a), \forall s \qquad \left\| \widehat{V}_h^n \right\|_{\infty} \leq H, \forall$$



UCBVI: Put All Together

For $n = 1 \rightarrow N$: 1. Set $N_h^n(s, a) = \sum_{k=1}^n \mathbf{1}\{(s_h^i, a_h^i) = (s, a)\}, \forall s, a, h$ $i = 1_{n-1}$ 2. Set $N_h^n(s, a, s') = \sum_{k=1}^n \mathbf{1}\{(s_h^i, a_h^i, s_{h+1}^i)\}$ i=13. Estimate \widehat{P}^n : $\widehat{P}^n_h(s'|s,a) = \frac{N_h^n(s,a)}{N_h^n(s,a)}$

4. Plan: $\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right)$, with

5. Execute π^n : { $s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n$ }

$$= (s, a, s')\}, \forall s, a, a', h$$

$$a, s'), \forall s, a, s', h$$

$$b_h^n(s, a) = cH \sqrt{\frac{\log(SAHN/\delta)}{N_h^n(s, a)}}$$

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



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

1. What if $\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0})$ is small?

Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$

$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $(s_0) - V_0^{\pi''}(s_0)$ is small?

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

Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$

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

2. What if
$$\widehat{V}_0^n$$

$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $V_0^n(s_0) - V_0^{\pi^n}(s_0)$ is small?

 $(s_0) - V_0^{\pi''}(s_0)$ is large?

Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$

2. What if
$$\widehat{V}_{0}^{n}(s_{0}) = V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + b_{h}^$$

$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $V_0(s_0) - V_0^{\pi''}(s_0)$ is small?

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

 $V_0^n(s_0) - V_0^{\pi^n}(s_0)$ is large? $\cdot (\widehat{P}_{h}^{n}(\cdot | s, a) - P_{h}(\cdot | s, a)) \cdot \widehat{V}_{h+1}^{n}$ must be large





Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$



$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $S(s_0) - V_0^{\pi^n}(s_0)$ is small?

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

Not obvious $\begin{array}{l}
2. \text{ What if } \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \text{ is large}? \\
\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{k=1}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right] \text{ must be large}$





Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$



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

$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $\tilde{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$ is small?

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

Not obvious $\begin{array}{l}
2. \text{ What if } \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \text{ is large}? \\
\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{k=1}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right] \text{ must be large}$





Upper bound per-episode regret:

1. What if
$$\widehat{V}_0^n$$

Not obvious

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot | s,a) - P_{h}(\cdot | s,a)) \cdot \widehat{V}_{h+1}^{n} \right] \text{ must be large}$$

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

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}\sqrt{SAN}\right)$$

$$V_0^{\star}(s_0) - V_0^{\pi^n}(s_0) \le \widehat{V}_0^n(s_0) - V_0^{\pi^n}(s_0)$$

 $\int_{0}^{n} (s_0) - V_0^{\pi^n}(s_0)$ is small?

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







- UCB-VI for tabular MDPs
 - UCB-VI for linear MDPs



Why we don't want to treat MDPs as big contextual bandits

S & A could be large or even continuous, hence poly(S,A) is not acceptable

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

S & A could be large or even continuous, hence poly(S,A) is not acceptable

 $P_h(s' \mid s, a) = \mu_h^{\star}(s') \cdot \phi(s, a)$

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

),
$$\mu_h^{\star} \in S \mapsto \mathbb{R}^d, \phi \in S \times A \mapsto \mathbb{R}^d$$

S & A could be large or even continuous, hence poly(S,A) is not acceptable

$$P_{h}(s'|s,a) = \mu_{h}^{\star}(s') \cdot \phi(s,a), \quad \mu_{h}^{\star} \in S \mapsto \mathbb{R}^{d}, \phi \in S \times A \mapsto \mathbb{R}^{d}$$
$$r(s,a) = \theta_{h}^{\star} \cdot \phi(s,a), \quad \theta_{h}^{\star} \in \mathbb{R}^{d}$$

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

S & A could be large or even continuous, hence poly(S,A) is not acceptable

$$P_{h}(s'|s,a) = \mu_{h}^{\star}(s') \cdot \phi(s,a), \quad \mu_{h}^{\star} \in S \mapsto \mathbb{R}^{d}, \phi \in S \times A \mapsto \mathbb{R}^{d}$$
$$r(s,a) = \theta_{h}^{\star} \cdot \phi(s,a), \quad \theta_{h}^{\star} \in \mathbb{R}^{d}$$

Finite horizon time-dependent episodic MDP $\mathcal{M} = \{S, A, H, \{r\}_h, \{P\}_h, s_0\}$

Feature map ϕ is known to the learner! (We assume reward is known, i.e., θ^{\star} is known)

 $V_H^{\star}(s) = 0, \forall s,$

$$V_H^{\star}(s) = 0, \forall s,$$

$$Q_{h}^{\star}(s,a) = r_{h}(s,a) + \mathbb{E}_{s' \sim P_{h}(\cdot|s,a)} V_{h+1}^{\star}(s')$$

$$V_{H}^{\star}(s) = 0, \forall s,$$

$$Q_{h}^{\star}(s, a) = r_{h}(s, a) + \mathbb{E}_{s' \sim P_{h}(\cdot | s, a)} V_{h+1}^{\star}(s')$$

$$= \theta_{h}^{\star} \cdot \phi(s, a) + \left(\mu_{h}^{\star} \phi(s, a)\right)^{\top} V_{h+1}^{\star}$$

$$\begin{aligned} V_{H}^{\star}(s) &= 0, \forall s, \\ Q_{h}^{\star}(s,a) &= r_{h}(s,a) + \mathbb{E}_{s' \sim P_{h}(\cdot \mid s,a)} V_{h+1}^{\star}(s') \\ &= \theta_{h}^{\star} \cdot \phi(s,a) + \left(\mu_{h}^{\star} \phi(s,a)\right)^{\top} V_{h+1}^{\star} \\ &= \phi(s,a)^{\top} \left(\theta_{h}^{\star} + (\mu_{h}^{\star})^{\top} V_{h+1}^{\star}\right) \end{aligned}$$

$$\begin{aligned} V_{H}^{\star}(s) &= 0, \forall s, \\ Q_{h}^{\star}(s,a) &= r_{h}(s,a) + \mathbb{E}_{s' \sim P_{h}(\cdot \mid s,a)} V_{h+1}^{\star}(s') \\ &= \theta_{h}^{\star} \cdot \phi(s,a) + \left(\mu_{h}^{\star}\phi(s,a)\right)^{\top} V_{h+1}^{\star} \\ &= \phi(s,a)^{\top} \left(\theta_{h}^{\star} + (\mu_{h}^{\star})^{\top} V_{h+1}^{\star}\right) \\ &= \phi(s,a)^{\top} w_{h} \end{aligned}$$

$$V_{H}^{\star}(s) = 0, \forall s,$$

$$Q_{h}^{\star}(s, a) = r_{h}(s, a) + \mathbb{E}_{s' \sim P_{h}(\cdot | s, a)} V_{h+1}^{\star}(s')$$

$$= \theta_{h}^{\star} \cdot \phi(s, a) + (\mu_{h}^{\star} \phi(s, a))^{\top} V_{h+1}^{\star}$$

$$= \phi(s, a)^{\top} (\theta_{h}^{\star} + (\mu_{h}^{\star})^{\top} V_{h+1}^{\star})$$

$$= \phi(s, a)^{\top} w_{h}$$

$$V_{h}^{\star}(s) = \max \phi(s, a)^{\top} w_{h}, \quad \pi_{h}^{\star}(s) = \arg \max \phi(s, a)^{\top} w_{h}$$

a

a

Planning in Linear MDP: Value Iteration $P_h(\cdot | s, a) = \mu_h^* \phi(s, a), \quad \mu_h^* \in \mathbb{R}^{S \times d}, \phi(s, a) \in \mathbb{R}^d$ $r_h(s, a) = (\theta_h^*)^\top \phi(s, a), \quad \theta_h^* \in \mathbb{R}^d$

$$V_{H}^{\star}(s) = 0, \forall s,$$

$$Q_{h}^{\star}(s, a) = r_{h}(s, a) + \mathbb{E}_{s' \sim P_{h}(\cdot | s, a)} V_{h+1}^{\star}(s')$$

$$= \theta_{h}^{\star} \cdot \phi(s, a) + (\mu_{h}^{\star} \phi(s, a))^{\top} V_{h+1}^{\star}$$

$$= \phi(s, a)^{\top} (\theta_{h}^{\star} + (\mu_{h}^{\star})^{\top} V_{h+1}^{\star})$$

$$= \phi(s, a)^{\top} w_{h}$$

$$f_{h}^{\star}(s) = \max \phi(s, a)^{\top} w_{h}, \quad \pi_{h}^{\star}(s) = \arg \max \phi(s, a)^{\top} w_{h}$$

a

Indeed we can show that $Q_h^{\pi}(\cdot, \cdot)$ Is linear with respect to ϕ as well, for any π, h

a

At the beginning of iteration n:

UCBVI in Linear MDPs

At the beginning of iteration n:

1. Learn transition model $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$ from all previous data $\{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=0}^{n-1}$

UCBVI in Linear MDPs
At the beginning of iteration n:

1. Learn transition model $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$ from all previous data $\{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=0}^{n-1}$

2. Design reward bonus $b_h^n(s, a), \forall s, a$

UCBVI in Linear MDPs

At the beginning of iteration n:

1. Learn transition model $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$ from all previous data $\{s_{h}^{i}, a_{h}^{i}, s_{h+1}^{i}\}_{i=0}^{n-1}$

3. Plan: $\pi^{n+1} = Value$

UCBVI in Linear MDPs

2. Design reward bonus $b_h^n(s, a), \forall s, a$

e-Iter
$$\left(\left\{ \widehat{P}^n \right\}_h, \left\{ r_h + b_h^n \right\} \right)$$



Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s



Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

Given s, a, note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)} \left[\delta(s') \right] = P_h(\cdot | s, a) = \mu_h^* \phi(s, a)$

How to estimate $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$?

How to estimate $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$?

Given s, a, note that $\mathbb{E}_{s' \sim P_{h}(\cdot | s, a)}$

Penalized Linear Regression:

$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

$$P_{a}\left[\delta(s')\right] = P_{h}(\cdot \mid s, a) = \mu_{h}^{\star}\phi(s, a)$$

How to estimate $\{\widehat{P}_{h}^{n}\}_{h=0}^{H-1}$?

Given *s*, *a*, note that $\mathbb{E}_{s' \sim P_h(\cdot | s, c)}$

Penalized Linear Regression:

$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

$$A_{h}^{n} = \sum_{i=1}^{n-1} \phi(s_{h}^{i}, a_{h}^{i}) \phi(s_{h}^{i}, a_{h}^{i})^{\mathsf{T}} + \lambda I$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

$$_{a)}\left[\delta(s')\right] = P_{h}(\cdot \mid s, a) = \mu_{h}^{\star}\phi(s, a)$$

$$\widehat{\mu}_{h}^{n} = (A_{h}^{n})^{-1} \sum_{i=1}^{n-1} \delta(s_{h+1}^{i}) \phi(s_{h}^{i}, a_{h}^{i})^{\mathsf{T}}$$

How to estim

Given *s*, *a*, note that $\mathbb{E}_{s' \sim P_h(\cdot | s, a)}$

Penalized Linear Regression:

$$\min_{\mu} \sum_{i=1}^{n-1} \|\mu \phi(s_h^i, a_h^i) - \delta(s_{h+1}^i)\|_2^2 + \lambda \|\mu\|_F^2$$

$$A_{h}^{n} = \sum_{i=1}^{n-1} \phi(s_{h}^{i}, a_{h}^{i}) \phi(s_{h}^{i}, a_{h}^{i})^{\top} + \lambda I$$

hate
$$\{ \widehat{P}_{h}^{n} \}_{h=0}^{H-1} ?$$

Denote $\delta(s) \in \mathbb{R}^{|S|}$ with zero everywhere except the entry corresponding to s

$$_{a)}\left[\delta(s')\right] = P_{h}(\cdot \mid s, a) = \mu_{h}^{\star}\phi(s, a)$$

$$\widehat{\mu}_{h}^{n} = (A_{h}^{n})^{-1} \sum_{i=1}^{n-1} \delta(s_{h+1}^{i}) \phi(s_{h}^{i}, a_{h}^{i})^{\mathsf{T}}$$

 $\widehat{P}_{h}^{n}(\cdot | s, a) = \widehat{\mu}_{h}^{n} \phi(s, a)$

Chebyshev-like approach, similar to in linUCB:

How to choose $b_h^n(s, a)$?

 $b_h^n(s,a) = \beta \sqrt{\phi(s,a)^{\mathsf{T}}(A_h^n)^{-1}\phi(s,a)}, \quad \beta = \widetilde{O}(dH)$

linUCB-VI: Put All Together

For $n = 1 \rightarrow N$: 1. Set $A_h^n = \sum_{k=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^{\top} + \lambda I$ $i=1 \qquad n-1 \\ \text{2. Set } \widehat{\mu}_h^n = (A_h^n)^{-1} \sum_{k=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top$ i=1

3. Estimate \widehat{P}^n : $\widehat{P}^n_h(\cdot | s, a) = \widehat{\mu}^n_h \phi(s, a)$

4. Plan: $\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right)$, with $b_h^n(s, a) = cdH_{\sqrt{\phi(s, a)^T(A_h^n)^{-1}\phi(s, a)}}$

5. Execute π^n : { $s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n$ }

linUCB-VI: Put All Together

For $n = 1 \rightarrow N$: 1. Set $A_h^n = \sum_{k=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^{\top} + \lambda I$ $i=1 \qquad n-1 \\ \text{2. Set } \widehat{\mu}_h^n = (A_h^n)^{-1} \sum_{k=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top$ i = 1

3. Estimate
$$\widehat{P}^n$$
 : $\widehat{P}^n_h(\cdot | s, a) = \widehat{\mu}^n_h \phi$

4. Plan:
$$\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right), v$$

5. Execute
$$\pi^{n}$$
: $\{s_{0}^{n}, a_{0}^{n}, r_{0}^{n}, \dots, s_{H-1}^{n}, a_{H-1}^{n}, r_{H-1}^{n}, s_{H}^{n}\}$

$$\mathbb{E}\left[\operatorname{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}d^{1.5}\sqrt{N}\right)$$

$$\mathbb{E}\left[\mathsf{Regret}_{N}^{n}, a_{0}^{n}, r_{0}^{n}, \dots, s_{H-1}^{n}, a_{H-1}^{n}, r_{H-1}^{n}, s_{H}^{n}\right]$$

$$\mathbb{E}\left[\mathsf{Regret}_{N}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}d^{1.5}\sqrt{N}\right)$$

b(s,a)

with $b_h^n(s, a) = c dH_{\sqrt{\phi(s, a)^T (A_h^n)^{-1} \phi(s, a)}}$

linUCB-VI: Put All Together

For $n = 1 \rightarrow N$: 1. Set $A_h^n = \sum_{k=1}^{n-1} \phi(s_h^i, a_h^i) \phi(s_h^i, a_h^i)^{\mathsf{T}} + \lambda I$ $i=1 \qquad n-1 \\ \text{2. Set } \widehat{\mu}_h^n = (A_h^n)^{-1} \sum_{k=1}^{n-1} \delta(s_{h+1}^i) \phi(s_h^i, a_h^i)^\top$ i = 1

3. Estimate
$$\widehat{P}^n$$
 : $\widehat{P}^n_h(\cdot | s, a) = \widehat{\mu}^n_h \phi$

4. Plan:
$$\pi^n = VI\left(\{\widehat{P}_h^n, r_h + b_h^n\}_h\right), v$$

5. Execute
$$\pi^n : \{s_0^n, a_0^n, r_0^n, \dots, s_{H-1}^n, a_{H-1}^n, r_{H-1}^n, s_H^n\}$$

$$\mathbb{E}\left[\operatorname{\mathsf{Regret}}_N\right] := \mathbb{E}\left[\sum_{n=1}^N \left(V^{\star} - V^{\pi^n}\right)\right] \le \widetilde{O}\left(H^2 d^{1.5} \sqrt{N}\right)$$

$$\mathbb{E}\left[\mathsf{Regret}_{N}^{n}, a_{0}^{n}, r_{0}^{n}, \dots, s_{H-1}^{n}, a_{H-1}^{n}, r_{H-1}^{n}, s_{H}^{n}\right] \\ \mathbb{E}\left[\mathsf{Regret}_{N}^{n}\right] := \mathbb{E}\left[\sum_{n=1}^{N} \left(V^{\star} - V^{\pi^{n}}\right)\right] \leq \widetilde{O}\left(H^{2}d^{1.5}\sqrt{N}\right)$$

b(s,a)

with $b_h^n(s, a) = c dH_{\sqrt{\phi(s, a)^T (A_h^n)^{-1} \phi(s, a)}}$

No *S*, *A* dependence!







Today's summary:

Today's summary:

UCB-VI algorithm for tabular MDPsUses UCB idea, but leverages MDP structure

Today's summary:

UCB-VI algorithm for tabular MDPs • Uses UCB idea, but leverages MDP structure

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





Lemma [Simulation lemma]:

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right]$$

$$(s,a) + b_h^n(s,a) + \widehat{P}_h^n(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^n, H \bigg\}$$

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

Lemma [Simulation lemma]:

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right]$$

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

$$(s,a) + b_h^n(s,a) + \widehat{P}_h^n(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^n, H \bigg\}$$

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

 $\pi^n(s_0)$

Lemma [Simulation lemma]:

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right]$$

$$\begin{aligned} \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) &= \widehat{Q}_{0}^{n}(s_{0}, \pi^{n}(s_{0})) - Q_{0}^{\pi^{n}}(s_{0}, \pi^{n}(s_{0})) \\ &\leq r_{0}(s_{0}, \pi^{n}(s_{0})) + b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - r_{0}(s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \end{aligned}$$

$$(s,a) + b_h^n(s,a) + \widehat{P}_h^n(\cdot | s,a) \cdot \widehat{V}_{h+1}^n, H \bigg\}$$

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



Lemma [Simulation lemma]: $\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s, t_{0}) \right]$ h=0

$$\begin{split} \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) &= \widehat{Q}_{0}^{n}(s_{0}, \pi^{n}(s_{0})) - Q_{0}^{\pi^{n}}(s_{0}, \pi^{n}(s_{0})) \\ &\leq r_{0}(s_{0}, \pi^{n}(s_{0})) + b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - r_{0}(s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \\ &= b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot V_{1}^{\pi^{n}} \end{split}$$

$$(s,a) + b_h^n(s,a) + \widehat{P}_h^n(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^n, H \bigg\}$$

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

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

$$\pi^n(s_0))$$



Lemma [Simulation lemma]:

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right]$$

$$\begin{split} \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) &= \widehat{Q}_{0}^{n}(s_{0}, \pi^{n}(s_{0})) - Q_{0}^{\pi^{n}}(s_{0}, \pi^{n}(s_{0})) \\ &\leq r_{0}(s_{0}, \pi^{n}(s_{0})) + b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - r_{0}(s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \\ &= b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot V_{1}^{\pi^{n}} \\ &= b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \left(\widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0}))\right) \cdot \widehat{V}_{1}^{n} + P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \left(\widehat{V}_{1}^{n} - V_{1}^{\pi^{n}}\right) \end{split}$$

$$(s,a) + b_h^n(s,a) + \widehat{P}_h^n(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^n, H \bigg\}$$

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



$$\begin{array}{l}
\begin{array}{l}
\begin{array}{l}
\begin{array}{l}
\begin{array}{l}
\begin{array}{l}
\begin{array}{l}
\end{array}\\
\end{array}\\
\end{array}\\
\hline V_{H}^{n}(s) = 0, \quad \widehat{Q}_{h}^{n}(s,a) = \min\left\{r_{h}(s,a) + b_{h}^{n}(s,a) + \widehat{P}_{h}^{n}(\cdot \mid s,a) \cdot \widehat{V}_{h+1}^{n}, H\right\}\\ \hline \widehat{V}_{h}^{n}(s) = \max_{a} \ \widehat{Q}_{h}^{n}(s,a), \quad \pi_{h}^{n}(s) = \arg\max_{a} \ \widehat{Q}_{h}^{n}(s,a), \forall s\end{array}$$

$$\widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) \leq \sum_{h=0}^{H-1} \mathbb{E}_{s,a \sim d_{h}^{\pi^{n}}} \left[b_{h}^{n}(s,a) + (\widehat{P}_{h}^{n}(\cdot \mid s,a) - P_{h}(\cdot \mid s,a)) \cdot \widehat{V}_{h+1}^{n} \right]$$

$$\begin{split} \widehat{V}_{0}^{n}(s_{0}) - V_{0}^{\pi^{n}}(s_{0}) &= \widehat{Q}_{0}^{n}(s_{0}, \pi^{n}(s_{0})) - Q_{0}^{\pi^{n}}(s_{0}, \pi^{n}(s_{0})) \\ &\leq r_{0}(s_{0}, \pi^{n}(s_{0})) + b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - r_{0}(s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \\ &= b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \widehat{V}_{1}^{n} - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot V_{1}^{\pi^{n}} \\ &= b_{h}^{n}(s_{0}, \pi^{n}(s_{0})) + \left(\widehat{P}_{0}^{n}(\cdot \mid s_{0}, \pi^{n}(s_{0})) - P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0}))\right) \cdot \widehat{V}_{1}^{n} + P_{0}(\cdot \mid s_{0}, \pi^{n}(s_{0})) \cdot \left(\widehat{V}_{1}^{n} - V_{1}^{\pi^{n}}\right) \end{split}$$

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

Simulation lemma]:

 $(,a))\cdot \widehat{V}_{h+1}^n$

