

Model-based RL

Recap: Planning algorithm for computing π^\star

We assumed that $P(s' | s, a), r(s, a) \forall s, a, s'$ are **known**



Recap: Planning algorithm for computing π^\star

We assumed that $P(s' | s, a), r(s, a) \forall s, a, s'$ are **known**

Value Iteration:

$$Q^{t+1}(s, a) \leftarrow r(s, a) + \max_a \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a'), \forall s, a$$

$\rightarrow Q^\star$

Recap: Planning algorithm for computing π^\star

We assumed that $P(s' | s, a), r(s, a) \forall s, a, s'$ are **known**

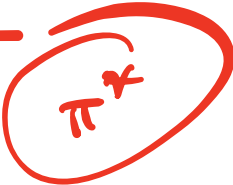
Value Iteration:

$$Q^{t+1}(s, a) \leftarrow r(s, a) + \max_a \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q^t(s', a'), \forall s, a$$

$$\pi \approx \pi^\star$$

Policy Iteration:

$$\pi^{t+1}(s) = \arg \max_a Q^{\pi^t}(s, a), \text{ for all } s$$

$$\pi^{t+1} \succeq \pi^t$$


Recap: Value-based Learning

When $P(s' | s, a)$ is unknown, Q-learning aims to learn Q^* directly

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \left(r(s, a) + \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)$$

Q-Target

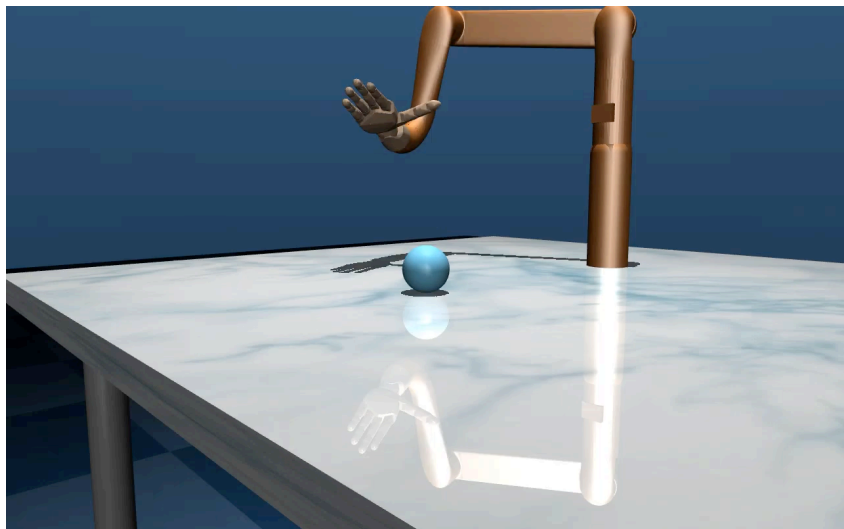
where $a \sim \epsilon$ -greedy(\hat{Q}), and $s' \sim P(\cdot | s, a)$, $r = r(s, a)$



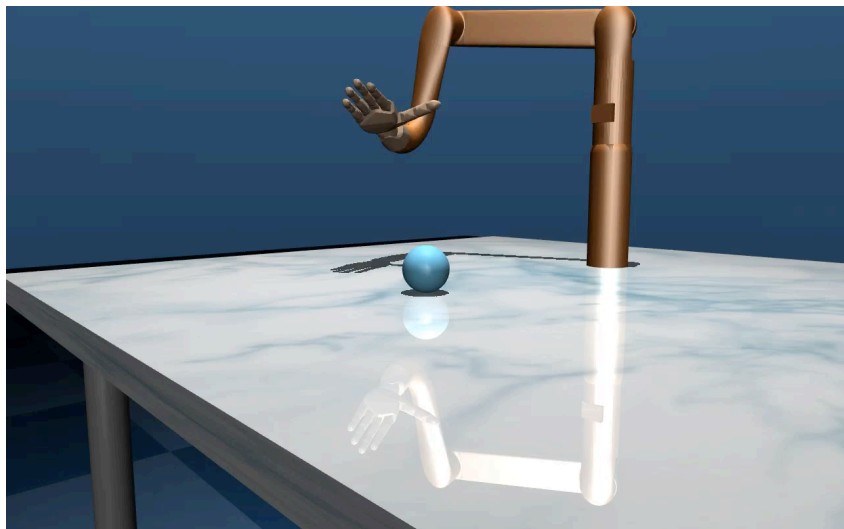
Questions for Today:

Can we **learn the transition** from data and then compute its optimal policy;
and what performance guarantee we can get?

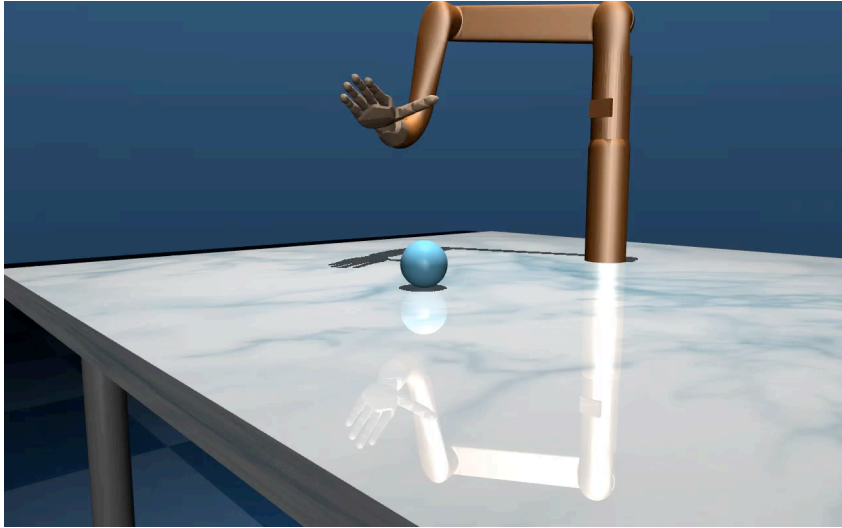
Motivation for Model-based Approach



Motivation for Model-based Approach

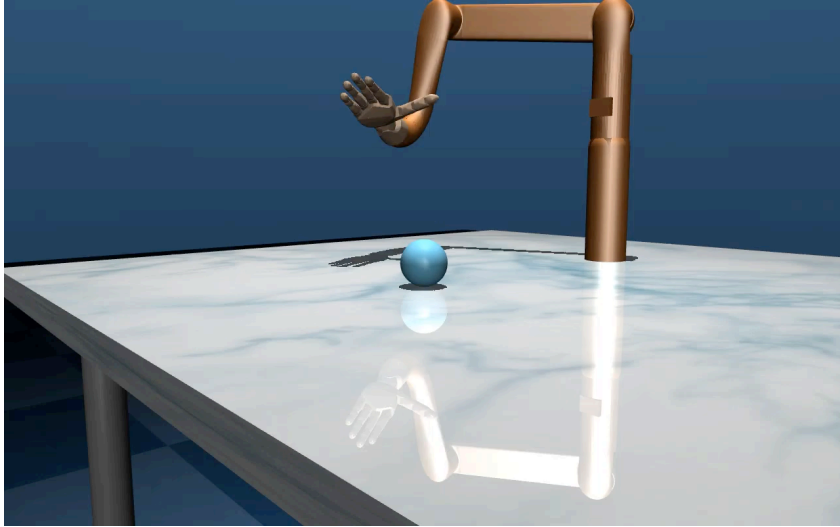


Motivation for Model-based Approach



While we cannot model the exact analytical dynamics, we can learn it from data $\{s, a, s'\}$

Motivation for Model-based Approach



While we cannot model the exact analytical dynamics, we can learn it from data $\{s, a, s'\}$

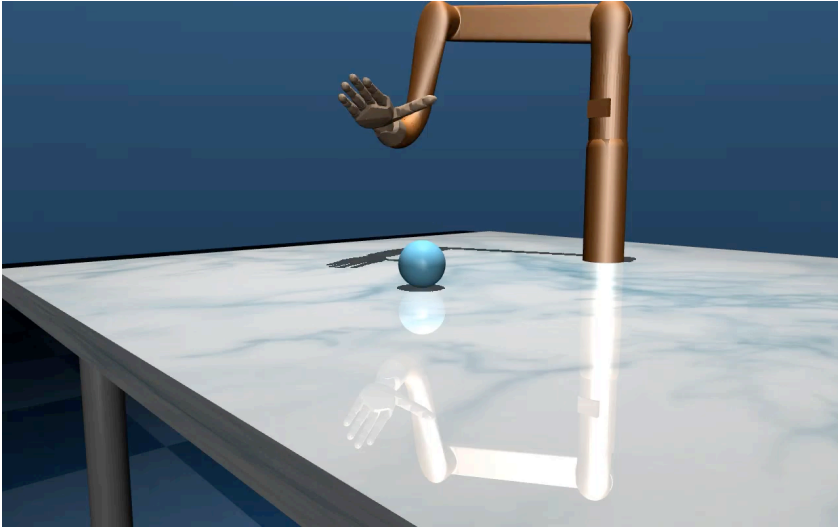
Then we do planning: e.g.,

$$\hat{\pi}^* = \text{VI}(\hat{P}, r)$$

Δ

$\hat{P} \approx P$

Motivation for Model-based Approach



While we cannot model the exact analytical dynamics, we can learn it from data $\{s, a, s'\}$

Then we do planning: e.g.,

$$\hat{\pi}^* = \text{VI}(\hat{P}, r)$$

(Often in practice we iterate the above process)

Motivation for Model-based Approach

Potential benefits of learning model over Q-learning

Motivation for Model-based Approach

Potential benefits of learning model over Q-learning

There are cases where model is much easier to learn than value function

Motivation for Model-based Approach

Potential benefits of learning model over Q-learning

There are cases where model is much easier to learn than value function

Once model is learned, we can optimize different rewards (i.e., multi-task)



Outline:

1. Simulation lemma:

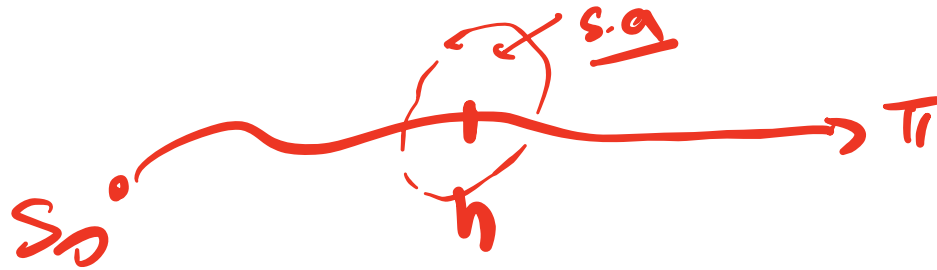
What is the performance of π under (\hat{P}, r)

2. Algorithm: estimate \hat{P} from data
and compute $\hat{\pi}^*$ — the optimal policy of \hat{P}

3. Analyzing the performance $\hat{\pi}^*$ under (P, r)

State-action distribution

Given a policy π and s_0 , we denote $\mathbb{P}_h^\pi(s, a | s_0)$ as the **prob of reaching (s, a) at time h , given we start at s_0 at $h = 0$**



State-action distribution

Given a policy π and s_0 , we denote $\mathbb{P}_h^\pi(s, a | s_0)$ as the **prob of reaching (s, a) at time h , given we start at s_0 at $h = 0$**

$$V^\pi(s_0)$$

$$= \mathbb{E}_{s \sim d_{s_0}^\pi} r(s, a)$$

Denote $d_{s_0}^\pi$ as the **average state-action distribution**:

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}_h^\pi(s, a | s_0)$$

$$\sum_{s, a} d_{s_0}^\pi(s, a) = 1$$



A key fundamental question in Model-based RL:

Given two transitions \hat{P} and P , how would π behave differently in \hat{P} and P ?



$$\hat{P} \approx P$$

A key fundamental question in Model-based RL:

Given two transitions \hat{P} and P , how would π behave differently in \hat{P} and P ?

\hat{P} is known
↓

Denote:

$$\hat{V}^{\pi}(s_0) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi, \hat{P} \right]; \quad V^{\pi}(s_0) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi, P \right];$$

A key fundamental question in Model-based RL:

Given two transitions \hat{P} and P , how would π behave differently in \hat{P} and P ?



Denote:

$$|\hat{V}^{\pi}(s_0) - V^{\pi}(s_0)|$$

$$\hat{V}^{\pi}(s_0) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi, \hat{P} \right]; \quad V^{\pi}(s_0) = \mathbb{E} \left[\sum_{h=0}^{\infty} \gamma^h r(s_h, a_h) \mid \pi, P \right];$$

What is the difference between $\hat{V}^{\pi}(s_0)$ & $V^{\pi}(s_0)$?

In other words, how does the model error propagate to values

Simulation Lemma

$$\underline{\gamma \in (0, 1]}$$

Simulation Lemma:

$$\left| \widehat{V}^\pi(s_0) - V^\pi(s_0) \right| = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left\| \widehat{P}(\cdot | s, a) - P(\cdot | s, a) \right\|_1$$

$$\sum_{s'} \left| \widehat{P}(s' | sa) - P(s' | sa) \right| \stackrel{\|x\|_1}{=} \sum |x_i|$$



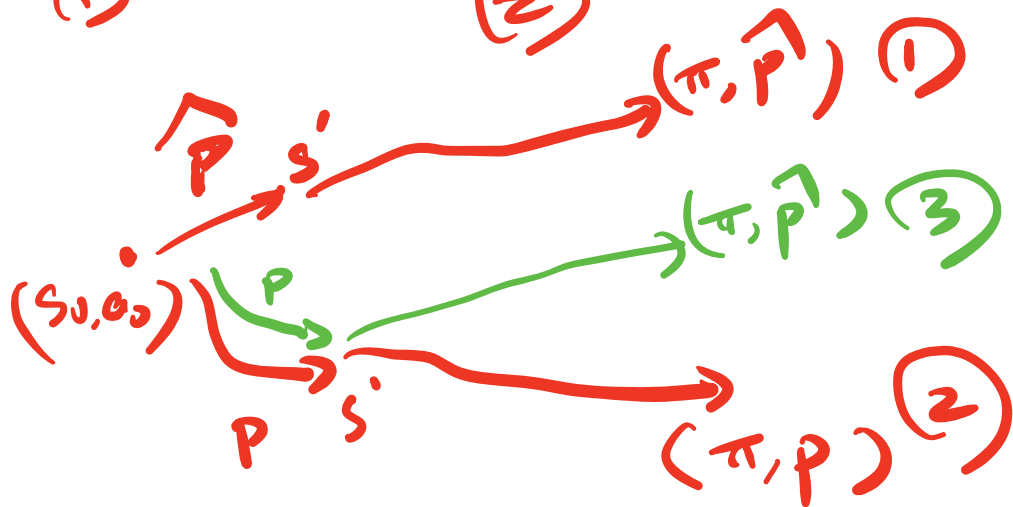
Simulation Lemma:

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

Simulation Lemma:

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\underbrace{\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1)}_{(1)} - \underbrace{\mathbb{E}_{s_1 \sim P(s_0,a_0)} V^\pi(s_1)}_{(2)} \right]$$



Simulation Lemma:

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1) - \mathbb{E}_{s_1 \sim P(s_0,a_0)} V^\pi(s_1) \right]$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\underbrace{\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1)}_{\textcircled{1}} - \underbrace{\mathbb{E}_{s_1 \sim P(s_0,a_0)} \widehat{V}^\pi(s_1)}_{\textcircled{2}} + \underbrace{\mathbb{E}_{s_1 \sim P(s_0,a_0)} \widehat{V}^\pi(s_1)}_{\textcircled{3}} - \underbrace{\mathbb{E}_{s_1 \sim P(s_0,a_0)} V^\pi(s_1)}_{\textcircled{4}} \right]$$

Simulation Lemma:

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1) - \mathbb{E}_{s_1 \sim P(s_0,a_0)} V^\pi(s_1) \right]$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1) - \mathbb{E}_{s_1 \sim P(s_0,a_0)} \widehat{V}^\pi(s_1) + \mathbb{E}_{s_1 \sim P(s_0,a_0)} \widehat{V}^\pi(s_1) - \mathbb{E}_{s_1 \sim P(s_0,a_0)} V^\pi(s_1) \right]$$

$$= \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0)} \left[\mathbb{E}_{s_1 \sim \widehat{P}(s_0,a_0)} \widehat{V}^\pi(s_1) - \mathbb{E}_{s_1 \sim P(s_0,a_0)} \widehat{V}^\pi(s_1) \right] \leftarrow \textcircled{1} - \textcircled{3}$$

$$+ \gamma \mathbb{E}_{a_0 \sim \pi(\cdot|s_0), s_1 \sim P(s_0,a_0)} \left[\widehat{V}^\pi(s_1) - V^\pi(s_1) \right] \leftarrow \textcircled{3} - \textcircled{2}$$



Summary so far:

$$\gamma \in [0, 1)$$

$$\widehat{V}^\pi \in [0, \frac{1}{1-\gamma}]$$

Simulation Lemma:

$$\widehat{V}^\pi(s_0) - V^\pi(s_0) = \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right]$$

$$\leq \frac{\gamma}{(1-\gamma)^2} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left\| \widehat{P}(\cdot | s, a) - P(\cdot | s, a) \right\|_1$$

$$\left| \mathbb{E}_{x \sim p} f(x) - \mathbb{E}_{x \sim q} f(x) \right|$$
$$\leq \left(\max_x |f(x)| \right) \|P - Q\|_1$$

Summary so far:

Simulation Lemma:

$$\begin{aligned}\widehat{V}^\pi(s_0) - V^\pi(s_0) &= \frac{\gamma}{1-\gamma} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left[\mathbb{E}_{s' \sim \widehat{P}(s,a)} \widehat{V}^\pi(s') - \mathbb{E}_{s' \sim P(s,a)} \widehat{V}^\pi(s') \right] \\ &\leq \frac{\gamma}{(1-\gamma)^2} \mathbb{E}_{s,a \sim d_{s_0}^\pi} \left\| \widehat{P}(\cdot | s, a) - P(\cdot | s, a) \right\|_1\end{aligned}$$

Total model disagreement over the real traces

Outline:



1. Simulation lemma:

What is the performance of π under any estimator \widehat{P}

2. Algorithm: estimate $(\widehat{P}, \widehat{r})$ from data
and compute $\widehat{\pi}^*$ — the optimal policy of $(\widehat{P}, \widehat{r})$

3. Analyzing the performance $\widehat{\pi}^*$ under (P, r)

A Model-based Algorithm

Assume reward r is known (just for analysis simplicity):

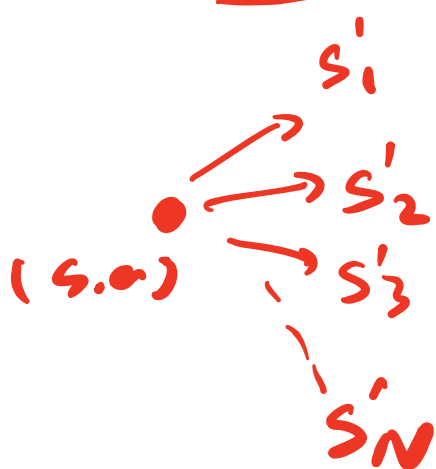
$$\underline{v(s, \omega)}$$

A Model-based Algorithm

Assume reward r is known (just for analysis simplicity):

1. Model fitting:

$\forall s, a$: collect N next states, $s'_i \sim P(\cdot | s, a)$, $i \in [N]$; set



$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

A Model-based Algorithm

Assume reward r is known (just for analysis simplicity):

1. Model fitting:

$\forall s, a$: collect N next states, $s'_i \sim P(\cdot | s, a)$, $i \in [N]$; set

$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

2. Planning w/ the learned model:

$$\widehat{\pi}^* = \text{PI} \left(\widehat{P}, r \right)$$

Detour: estimating mean of Bernoulli distribution

Given: we have a biased coin:

With probability p , it gives +1, and w/ prob $1-p$, it gives -1;

$$\hat{p} \approx p$$



Detour: estimating mean of Bernoulli distribution

Given: we have a biased coin:

With probability p , it gives +1, and w/ prob $1-p$, it gives -1;

To estimate p : We flip the coin N times independently, get N outcomes, $\{x_i\}_{i=1}^N$,

$$x_i \in \{-1, +1\}$$



Detour: estimating mean of Bernoulli distribution

Given: we have a biased coin:

With probability p , it gives +1, and w/ prob $1-p$, it gives -1;

To estimate p : We flip the coin N times independently, get N outcomes, $\{x_i\}_{i=1}^N$,

$$x_i \in \{-1, +1\}$$

p $\leftarrow \hat{p} = \frac{\sum_{i=1}^N \mathbf{1}\{x_i = +1\}}{N}$

Detour: estimating mean of Bernoulli distribution

Given: we have a biased coin:

With probability p , it gives +1, and w/ prob $1-p$, it gives -1;

To estimate p : We flip the coin N times independently, get N outcomes, $\{x_i\}_{i=1}^N$,

$$x_i \in \{-1, +1\}$$

$$\hat{p} = \frac{\sum_{i=1}^N \mathbf{1}\{x_i = +1\}}{N}$$

(Informal) we can show that $|\hat{p} - p| \lesssim \sqrt{\frac{1}{N}}$

(proof out of scope)

Model-based RL

1. Model fitting:

$\forall s, a$: collect N next states,

$s'_i \sim P(\cdot | s, a), i \in [N]$; set

$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

Estimating the prob of observing s' at (s, a)

2. Planning w/ the learned model:

$$\widehat{\pi}^* = \text{PI}(\widehat{P}, r)$$

Handwritten notes in red:

$$\widehat{P}(s' | s, a) - P(s' | s, a) \leq \sqrt{\frac{1}{N}}$$

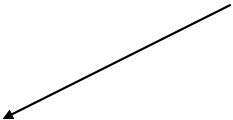
Model-based RL

1. Model fitting:

$\forall s, a$: collect N next states,
 $s'_i \sim P(\cdot | s, a), i \in [N]$; set

$$\widehat{P}(s' | s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N};$$

Estimating the prob of observing
 s' at (s, a)



We can show (informally):

2. Planning w/ the learned model:

$$\widehat{\pi}^* = \text{PI} \left(\widehat{P}, r \right)$$

$$\underbrace{\|\widehat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1}_{\triangle} \lesssim \sqrt{1/N}, \forall s, a$$

Summary so far:

By collecting data (next state) at every (s, a) , we build an estimator \hat{P} that is close to P
(e.g., possible to show error shrinks $\approx 1/\sqrt{N}$)

Outline:



1. Simulation lemma:

What is the performance of π under any estimator \hat{P}



2. Algorithm: estimate (\hat{P}, \hat{r}) from data
and compute $\hat{\pi}^*$ — the optimal policy of (\hat{P}, \hat{r})



3. Analyzing the performance $\hat{\pi}^*$ under (P, r)

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

△

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

△

$$|V^{\hat{\pi}^*} - V^{\pi^*}| \leq \epsilon \cdot \frac{1}{1 - \gamma}$$

$$V^{\hat{\pi}^*}$$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e. $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

$$V^*(s_0) \approx V^{\hat{\pi}^*}(s_0)$$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0)$$

$$V^*(s_0) - \hat{V}^{\hat{\pi}^*}(s_0) + \hat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0)$$

$$\leq V^*(s_0) - \hat{V}^{\pi^*}(s_0)$$

$$\leq \frac{1}{(1-\gamma)^2} \mathbb{E} \|\hat{P} - P\|_1$$

$$\leq \left(\frac{1}{1-\gamma}\right)^2 \cdot \epsilon$$

Simulation lemma

$$\leq \left(\frac{1}{1-\gamma}\right)^2 \epsilon$$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$$\begin{aligned} & V^*(s_0) - V^{\hat{\pi}^*}(s_0) \\ & \leq V^*(s_0) - \widehat{V}^{\pi^*}(s_0) + \widehat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0) \end{aligned}$$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$V^*(s_0) - V^{\hat{\pi}^*}(s_0)$ Q: why this is true?

$$\leq V^*(s_0) - \widehat{V}^{\pi^*}(s_0) + \widehat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0)$$

$\Rightarrow 0$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0)$$

Q: why this is true?

$$\leq \underbrace{V^*(s_0) - \widehat{V}^{\pi^*}(s_0)} + \underbrace{\widehat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0)} \quad (\text{Simulation lemma})$$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0)$$

Q: why this is true?

$$\leq V^*(s_0) - \widehat{V}^{\pi^*}(s_0) + \widehat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0) \quad (\text{Simulation lemma})$$

$$\leq \frac{1}{(1-\gamma)^2} \left[\underbrace{\mathbb{E}_{s,a \sim d_{s_0}^{\pi^*}} \|\widehat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1}_{\leq \epsilon} + \mathbb{E}_{s,a \sim d_{s_0}^{\hat{\pi}^*}} \|\widehat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \right]$$

Performance of the learned policy

Now assume that \hat{P} is nearly accurate, i.e., $\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \leq \epsilon, \forall s, a$

How good is $\hat{\pi}^* = \text{PI}(\hat{P}, r)$?

$$V^*(s_0) - V^{\hat{\pi}^*}(s_0)$$

Q: why this is true?

$$\leq V^*(s_0) - \widehat{V}^{\pi^*}(s_0) + \widehat{V}^{\hat{\pi}^*}(s_0) - V^{\hat{\pi}^*}(s_0) \quad (\text{Simulation lemma})$$

$$\leq \frac{1}{(1-\gamma)^2} \left[\mathbb{E}_{s,a \sim d_{s_0}^{\pi^*}} \|\widehat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 + \mathbb{E}_{s,a \sim d_{s_0}^{\hat{\pi}^*}} \|\widehat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 \right]$$

$$\leq \frac{2\epsilon}{(1-\gamma)^2}$$

Summary for Today:

1. A model-based RL **Algorithm for small-size MDP**

$$\widehat{P}(s'|s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N}, \forall s, a; \quad \widehat{\pi}^* = \mathbf{PI}(\widehat{P}, r)$$

Summary for Today:

1. A model-based RL **Algorithm for small-size MDP**

$$\widehat{P}(s'|s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N}, \forall s, a; \quad \widehat{\pi}^* = \mathbf{PI}(\widehat{P}, r)$$

2. **Simulation lemma** allows us to link model error to policy's performance

Summary for Today:

1. A model-based RL **Algorithm for small-size MDP**

$$\widehat{P}(s'|s, a) = \frac{\sum_{i=1}^N \mathbf{1}\{s'_i = s'\}}{N}, \forall s, a; \quad \widehat{\pi}^* = \mathbf{PI}(\widehat{P}, r)$$

2. **Simulation lemma** allows us to link model error to policy's performance

3. Good model leads to a good policy, up to some error amplification based on effective horizon

$$\frac{1}{(1-\gamma)^2} \cdot \epsilon$$