# Q-Learning

# Recap: Bellman Optimality

**Bellman Optimality**

$$Q^\star(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q^\star(s', a'), \forall s, a$$

**VI: An iterative approach for estimating $Q^\star$**

$$Q \Leftarrow \mathcal{T} Q$$

$$(\mathcal{T}Q)(s,a) = r(s,a) + \gamma \cdot \underset{s' \sim p(\cdot|s,a)}{\mathbb{E}} \max_{a'} Q(s', a')$$

# Recap: Bellman Optimality

**Bellman Optimality**

$$Q^\star(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q^\star(s', a'), \forall s, a$$

**VI: An iterative approach for estimating $Q^\star$**

$$Q \Leftarrow \mathcal{T} Q$$

1. Need to know the transition

2. Only works for discrete small MDPs

# Recap: Bellman Optimality

$Q(sa)$
$-(r(sa)$
$+\gamma \mathbb{E}_{s'} \max_{a'} Q(s'a')$
$= 0$

**Q: if there is some $Q(s, a)$, such that the following holds:**

$$Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} Q(s', a'), \forall s, a$$

$\triangle$

**is this $Q = Q^\star$?**

$Q^\star = T Q^\star$

$Q = T Q$

$\cdot Q^\star \xrightarrow{T} Q^\star \xrightarrow{r} Q^\star \rightarrow \cdots$

$\cdot Q \xrightarrow{T} Q \xrightarrow{r} Q \rightarrow \cdots$

# Today

$V^{\pi}$

**Given MDP** $\mathscr{M} = (S, A, r, P, \gamma)$,

**how to estimate** $Q^{\star}(s, a), \forall s$ **WITHOUT knowing** $P$
**i.e., how to learn** $Q^{\star}$ **(thus,** $\pi^{\star}$**) from experience**

$Q \rightarrow Q^{\star}$

$\pi \leftarrow \arg\max_a Q(s, a)$

# Motivation

Computing a near-optimal policy to achieve the long-term goals w/o knowing or explicitly modeling the world
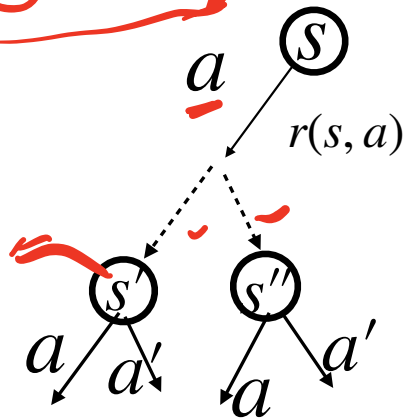
# Outline:

1. Q Learning

2. Revisit TD: Off-policy TD Learning

# Q Learning
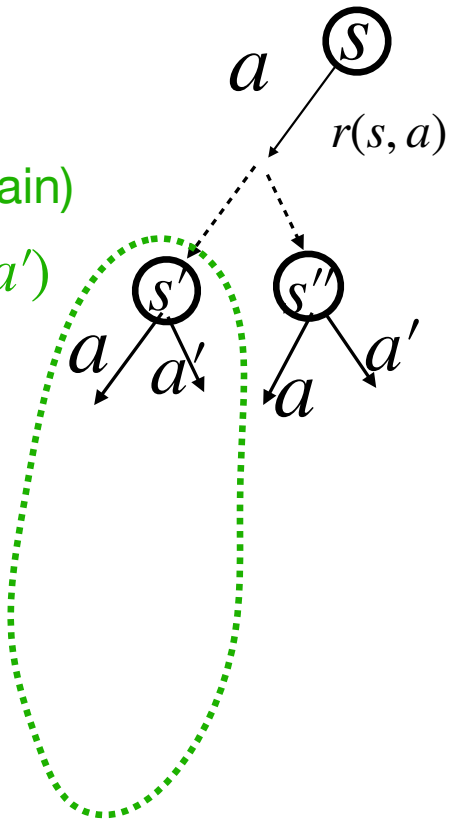


$\widehat{Q}$

$Q^*(s, a)$

$V^*(s')$

$a$   $S$

$r(s, a)$

$s'$   $s''$

$a$ $a'$   $a$ $a'$

# Q Learning



$a$   (S)

$r(s, a)$

(Bootstrapping again)

$$\max_{a'} \hat{Q}(s', a')$$

$\approx V^*(s')$

$a$ (s') $a$ (s'')

$a$ $d'$ $a$ $a'$

# Q Learning



$$Q^{\star}(s, a) \approx \underbrace{r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')}_{\text{target}}$$

$$\approx V^{*}(s')$$
$$= \max_{a'} Q^{*}(s', a')$$

(Bootstrapping again)

$$\max_{a'} \hat{Q}(s', a')$$

# Q Learning

$$Q^{\star}(s, a) \approx r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')$$

target

$a$   $s$

$r(s, a)$

(Bootstrapping again)

$\max_{a'} \hat{Q}(s', a')$

$s'$    $s''$

$a$   $a'$    $a'$

$a$

Target

$\hat{Q}(s, a)$

# Q Learning

$$Q^{\star}(s, a) \approx r(s, a) + \gamma \max_{a'} \hat{Q}(s', a') \neq Q^{\star}(s, a)$$

$$\underbrace{\phantom{r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')}}_{\text{target}}$$

$a$  $\textcircled{s}$

$r(s, a)$

(Bootstrapping again)

$\max_{a'} \hat{Q}(s', a')$

$a \quad \textcircled{s'} \quad a$

$a' \qquad \textcircled{s''} \quad a'$

$a$

Target

$\hat{Q}(s, a)$  $\widehat{Q}(sa)$

# Q Learning

$$Q^\star(s, a) \approx \underbrace{r(s, a) + \gamma \max_{a'} \hat{Q}(s', a')}_{\text{target}}$$

$a$

$r(s, a)$

(Bootstrapping again)
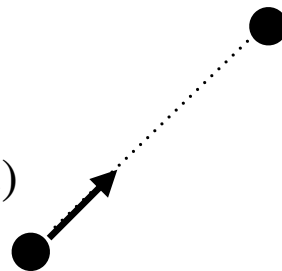
$$\max_{a'} \hat{Q}(s', a')$$

$a$ $a'$

$a$ $a'$

Target

$\hat{Q}(s, a)$

Q-learning update: move to the target with a small step

$V_\theta^\uparrow$  $Q^\uparrow$

# Q Learning

$Q^*$  $\underset{c}{\arg\max} Q^*(s,a)$

Given a **one-step transition** $(s, a, r, s')$ where $r = r(s, a)$, $s' \sim P(\cdot \mid s, a)$:

$a \sim \bar{a}(\cdot \mid s)$

Q-learning updates the guess at $(s, a)$ as follows:

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

Q-Target
(Constructed via Bootstrapping!)

# Q Learning

**Choice one**: trust current estimator $\hat{Q}$, always use $\arg\max_a \hat{Q}(s, a)$

# Q Learning

**How to collect data?**

**Choice one**: trust current estimator $\hat{Q}$, always use $\arg\max_a \hat{Q}(s, a)$

Issue: cannot explore (i.e., need to try something that hasn't been tried)

# Q Learning

**Choice one**: trust current estimator $\hat{Q}$, always use $\arg\max\limits_{a} \hat{Q}(s, a)$

Issue: cannot explore (i.e., need to try something that hasn't been tried)

**Choice two (quite effective in practice)**: $\epsilon$-greedy

$\epsilon \in (0, 1)$

$$Q^*(s_a)$$
$$= w^T \begin{pmatrix} s \\ a \end{pmatrix}$$

$$\hat{Q}(s, a_1) = 100$$
$$\hat{Q}(s, a_0) = -0.1$$

# Q Learning

## How to collect data?

**Choice one**: trust current estimator $\hat{Q}$, always use $\arg\max_a \hat{Q}(s, a)$

Issue: cannot explore (i.e., need to try something that hasn't been tried)

**Choice two (quite effective in practice)**: $\epsilon$-**greedy**

W/ prob $\epsilon$, select action uniform randomly ← Explore

W/ prob $1 - \epsilon$, select greedy action $\arg\max_a \hat{Q}(s, a)$ ← Exploitation

Choice Three: $\pi(a|s) = \exp\left(\hat{Q}(s, a)\right) \Big/ \sum_{a'} \exp(\hat{Q}(s, a'))$

# Q TD Learning

Initialize $\hat{Q}(s, a) = 0, \forall s, a.$  Set initial state $s \in \mathcal{S}$

While True:

# TD Learning

Initialize $\hat{Q}(s, a) = 0, \forall s, a.$  Set initial state $s \in \mathcal{S}$

While True:

$(s, a, r, s')$

    Take action $a$ based on $\epsilon$-greedy of $\hat{Q}$, get reward $r$ and next state $s' \sim P(\,\cdot\,|s, a)$

# TD Learning

Initialize $\hat{Q}(s, a) = 0, \forall s, a$. Set initial state $s \in \mathcal{S}$

While True:

$\quad$ Take action $a$ based on $\epsilon$-greedy of $\hat{Q}$, get reward $r$ and next state $s' \sim P(\,\cdot\,|s, a)$

$\quad$ Form Q-target $r + \gamma \max_{a'} \hat{Q}(s', a')$ $\leftarrow$ Bootstrapping

# TD Learning

Initialize $\hat{Q}(s, a) = 0, \forall s, a$.  Set initial state $s \in \mathcal{S}$

While True:

  Take action $a$ based on $\epsilon$-greedy of $\hat{Q}$, get reward $r$ and next state $s' \sim P(\cdot \mid s, a)$

  Form Q-target $r + \gamma \max_{a'} \hat{Q}(s', a')$

  Update for $s, a$: $\hat{Q}(s, a) \Leftarrow \hat{Q}(s, a) + \eta \left( r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)$
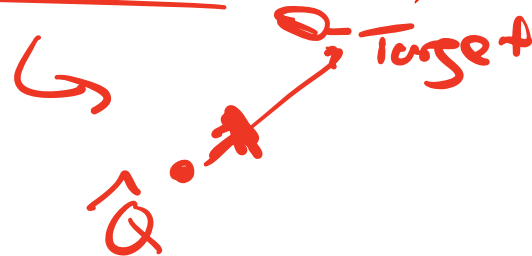
# TD Learning

Initialize $\hat{Q}(s, a) = 0, \forall s, a$. Set initial state $s \in \mathcal{S}$

While True:

> Take action $a$ based on $\epsilon$-greedy of $\hat{Q}$, get reward $r$ and next state $s' \sim P(\cdot \mid s, a)$
>
> Form Q-target $r + \gamma \max_{a'} \hat{Q}(s', a')$
>
> Update for $s, a$: $\hat{Q}(s, a) \Leftarrow \hat{Q}(s, a) + \eta \left( r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)$
>
> Set $s \Leftarrow s'$

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

$$\ell_{be}(\hat{Q}(s,a)) := \left( \hat{Q}(s,a) - y \right)^2, \text{ where } y = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} \hat{Q}(s',a')$$

$\hat{Q} = T\hat{Q}$

Bootstrapping

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

$$\ell_{be}(\hat{Q}(s,a)) := \left( \hat{Q}(s,a) - y \right)^2, \text{ where } y = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} \hat{Q}(s',a')$$

**This keeps changing as we learning**

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

$$\ell_{be}(\hat{Q}(s,a)) := \left( \hat{Q}(s,a) - y \right)^2, \text{ where } y = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} \hat{Q}(s',a')$$

$$\nabla \ell_{be}(x) \big|_{x=\hat{Q}(s,a)} := 2 \left( \hat{Q}(s,a) - y \right)$$

$s' \sim P(\cdot|sa)$

**This keeps changing as we learning**

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

$$\ell_{be}(\hat{Q}(s,a)) := \left(\hat{Q}(s,a) - y\right)^2, \text{ where } y = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} \hat{Q}(s',a')$$

$$\nabla \ell_{be}(x)|_{x=\hat{Q}(s,a)} := 2\left(\hat{Q}(s,a) - y\right)$$

**This keeps changing as we learning**

$$\widetilde{\nabla} \ell_{be}(x)|_{x=\hat{Q}(s,a)} := 2\left(\hat{Q}(s,a) - \left(r + \gamma \max_{a'} \hat{Q}(s',a')\right)\right)$$

# Interpret Q-learning as "SGD" on Bellman error

Q-learning is not the usual SGD, i.e., it is **not** running SGD to minimize a fixed loss function

Q-learning may be interpreted as running SGD on an **evolving** loss function (Bellman error)

$$\ell_{be}(\hat{Q}(s,a)) := \left( \hat{Q}(s,a) - y \right)^2, \text{ where } y = r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} \max_{a'} \hat{Q}(s',a')$$

$$\nabla \ell_{be}(x) \big|_{x=\hat{Q}(s,a)} := 2 \left( \hat{Q}(s,a) - y \right)$$

Unbiased
estimate of $y$

**This keeps changing as we learning**

$$\widetilde{\nabla} \ell_{be}(x) \big|_{x=\hat{Q}(s,a)} := 2 \left( \hat{Q}(s,a) - \left( r + \gamma \max_{a'} \hat{Q}(s',a') \right) \right)$$

$$\hat{Q}(sa) - \eta \cdot \widetilde{\nabla} \ell_{BE}(x) \big|_{x = \hat{\partial}(sa)}$$

Q-Target

# Q Learning Theory

[Informal] Assume the $\epsilon$-greedy strategy has non-trivial probability of visiting every state-action pair. Setting learning rate $\eta$ properly, we will have:

$$\frac{1}{\sqrt{n}} \quad \frac{1}{n} \qquad \hat{Q}(s,a) \to Q^{\star}(s,a), \forall s, a$$

when # of interactions approaches to $\infty$

(concrete convergence rates are known as well)

# Demo: Q-learning on CartPole

Note: Cartpole's state is continuous, so we will need Q-learning w/ function approximation, e.g., neural network (we will get there very soon)

1. Does Q learning eventually learn a good policy

2. How does the $\epsilon$ affect the learning

# Outline:

1. Q Learning

2. Revisit TD: Off-policy TD Learning

# TD Learning

Given $(s, a, r, s')$, where $a \sim \pi(\cdot \mid s), s' \sim P(\cdot \mid s, a)$, TD updates:

TD-Error

$$\hat{V}^\pi(s) \Leftarrow \hat{V}^\pi(s) + \eta \left( r + \gamma \hat{V}^\pi(s') - \hat{V}^\pi(s') \right)$$

TD-Target

**On-policy**: data is generated from the policy $\pi$ itself

**Off-policy: data is generated from policy $\pi_b$ where $\pi_b \neq \pi$**

# TD Learning

Given $(s, a, r, s')$, where $a \sim \pi( \cdot \,|\, s), s' \sim P( \cdot \,|\, s, a)$, TD updates:

$$\hat{V}^\pi(s) \Leftarrow \hat{V}^\pi(s) + \eta \left( r + \gamma \hat{V}^\pi(s') - \hat{V}^\pi(s') \right)$$

**On-policy**: data is generated from the policy $\pi$ itself

**Off-policy: data is generated from policy $\pi_b$ where $\pi_b \neq \pi$**

Q: is Q-learning off-policy or on-policy?

$Q\text{-leang}$

$\text{arg max} \, Q^*$

# Motivation for off-policy evaluation

**Counterfactual**: what would happen if I did something different?

$$\pi_b \qquad \pi_{new}$$

# Off-policy TD Learning

$a \sim \pi_b(\cdot|s)$

Setting: data is generated by $\pi_b$, but we want to estimate $V^\pi$ for some $\pi \neq \pi_b$

Key trick: importance weighting

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

$BE - \pi$

$a \sim \pi_b$

$= \sum_a \pi(a|s)(----) = \sum_a \pi_b(a|s) \frac{\pi(a|s)}{\pi_b(a|s)}(---)$

$\pi = \mathbb{E}_{a \sim \pi_b(\cdot|s)} \frac{\pi(a|s)}{\pi_b(a|s)}(\sim-)$

# Off-policy TD Learning

Setting: data is generated by $\pi_b$, but we want to estimate $V^\pi$ for some $\pi \neq \pi_b$

Key trick: importance weighting

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

$$= \mathbb{E}_{a \sim \pi_b(\cdot|s)} \frac{\pi(a \mid s)}{\pi_b(a \mid s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

# Off-policy TD Learning

Setting: data is generated by $\pi_b$, but we want to estimate $V^\pi$ for some $\pi \neq \pi_b$

Key trick: importance weighting

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

$$= \mathbb{E}_{a \sim \pi_b(\cdot|s)} \frac{\pi(a \mid s)}{\pi_b(a \mid s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

Importance weight

# Off-policy TD Learning

Setting: data is generated by $\pi_b$, but we want to estimate $V^\pi$ for some $\pi \neq \pi_b$

Key trick: importance weighting

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

$$= \mathbb{E}_{a \sim \pi_b(\cdot|s)} \frac{\pi(a\,|\,s)}{\pi_b(a\,|\,s)} \left( r(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} V^\pi(s') \right)$$

Now action is
sampled from $\pi_b$

Importance weight

$a \sim \pi_b (\cdot | s)$

# Off-policy TD Learning

Given $(s, a, r, s')$, where $a \sim \pi_b( \cdot \mid s), s' \sim P( \cdot \mid s, a),$

Off-policy TD updates as follows:

$$\hat{V}^\pi(s) \Leftarrow \hat{V}^\pi(s) + \eta \frac{\pi(a \mid s)}{\pi_b(a \mid s)} \left( r + \gamma \hat{V}^\pi(s') - V^\pi(s) \right)$$

# Off-policy TD Learning

Given $(s, a, r, s')$, where $a \sim \pi_b( \cdot \mid s)$, $s' \sim P( \cdot \mid s, a)$,

Off-policy TD updates as follows:

$$\hat{V}^\pi(s) \Leftarrow \hat{V}^\pi(s) + \eta \frac{\pi(a \mid s)}{\pi_b(a \mid s)} \left( r + \gamma \hat{V}^\pi(s') - V^\pi(s) \right)$$

Case 1: $\pi(a \mid s)$ is large but $\pi_b(a \mid s)$ is small

$$\frac{\pi(a \mid s)}{\pi_b(a \mid s)} \nearrow$$

# Off-policy TD Learning

Given $(s, a, r, s')$, where $a \sim \pi_b( \cdot \,|\, s), s' \sim P( \cdot \,|\, s, a)$,

Off-policy TD updates as follows:

$$\hat{V}^{\pi}(s) \Leftarrow \hat{V}^{\pi}(s) + \eta \frac{\pi(a\,|\,s)}{\pi_b(a\,|\,s)} \left( r + \gamma \hat{V}^{\pi}(s') - V^{\pi}(s) \right)$$

Case 1: $\pi(a\,|\,s)$ is large but $\pi_b(a\,|\,s)$ is small

Case 2: $\pi(a\,|\,s)$ is small but $\pi_b(a\,|\,s)$ is large

$$\frac{\pi}{\pi_b} \approx 0$$

# Off-policy TD Learning is SGD on TD loss

Given $(s, a, r, s')$, where $a \sim \pi_b(\cdot \mid s), s' \sim P(\cdot \mid s, a)$. Off-policy TD updates:

$$\hat{V}^\pi(s) \Leftarrow \hat{V}^\pi(s) + \eta \frac{\pi(a \mid s)}{\pi_b(a \mid s)} \left( r + \gamma \hat{V}^\pi(s') - V^\pi(s) \right)$$

Check if it is doing one-step SGD on the TD loss:

$$\ell_{td}(\hat{V}^\pi(s)) = \left( \hat{V}^\pi(s) - y \right)^2 \text{ where } y = \mathbb{E}_{a \sim \pi(\cdot \mid s)} \left( r + \gamma \mathbb{E}_{s' \sim P(s,a)} \hat{V}^\pi(s') \right)$$

The off-policy TD update is one-step SGD on $\ell_{td}$ (more in HW2)

# Summary

Q-Learning: online algorithm that learns $Q^\star$ (bootstrapping)

Exploration & Exploitation tradeoff: $\epsilon$-greedy is an effective heuristic

Off-policy policy evaluation: importance weighting
(also known as inverse probability weighting in causal inference)