# Policy Evaluation

Chen-Yu Wei

# Policy Evaluation

> Given: a policy $\pi$
>
> Evaluate $V^\pi(s)$ or $Q^\pi(s, a)$

**On-policy policy evaluation**: the learner can execute $\pi$ to evaluate $\pi$    $\pi_b(\cdot \mid s)$

**Off-policy/offline policy evaluation**: the learner can only execute some $\pi_b \neq \pi$, or can only access some existing dataset to evaluate $\pi$

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

$(s_1, a_1, r_1, s_2, a_2, \ldots)$

$\uparrow$ behavior policy

**Use cases:**

- Approximate policy iteration: $\pi^{(k)}(s) = \underset{a}{\mathrm{argmax}}\, Q^{\pi^{(k-1)}}(s, a)$

- Estimate the value of a policy before deploying it in the real world, e.g., COVID-related border measures, economic recovery policies, or policy changes in recommendation systems.

# Value Iteration for $V^\pi$ / $Q^\pi$

$Q^*(s,a)$   $V^*(s)$

$\pi^*(s) = \arg\max_a Q^*(s,a)$

**Input:** $\pi$

For $k = 1, 2, \ldots$

$V^{(k)}(s) \xrightarrow[k \to \infty]{} V^\pi(s)$

$$\forall s, \qquad V^{(k)}(s) \leftarrow \sum_a \pi(a|s)\left( R(s,a) + \gamma \sum_{s'} P(s'|s,a)\, V^{(k-1)}(s') \right)$$

**Input:** $\pi$

For $k = 1, 2, \ldots$

$Q^{(k)} \to Q^\pi$

$$\forall s, a, \qquad Q^{(k)}(s,a) \leftarrow R(s,a) + \gamma \sum_{s',a'} P(s'|s,a)\, \pi(a'|s')\, Q^{(k-1)}(s',a')$$

# On-Policy Policy Evaluation

# LSPE and TD

Collecting samples $\{(s_i, r_i, s_i')\}_{i=1}^n$ using $\pi$

For $k = 1, 2, \ldots$

$$\theta_k \leftarrow \operatorname*{argmin}_\theta \sum_{i=1}^n \left( V_\theta(s_i) - r_i - \gamma V_{\theta_{k-1}}(s_i') \right)^2$$

Least-Square Policy Evaluation (LSPE)

$$V(s) \underset{\theta_k}{\simeq} \left[ \sum_a \pi(a|s) \left( R(s,a) + \gamma \underset{s' \sim P(\cdot|s,a)}{\mathbb{E}} [V_{\theta_{k-1}}(s')] \right) \right]$$

For $i = 1, 2, \ldots$

Draw $a_i \sim \pi(\cdot \,|s_i)$

$(s, r, s')$

Observe reward $r_i$ and next state $s_{i+1}$

$$\theta_i \leftarrow \theta_{i-1} - \alpha \nabla_\theta \left( V_\theta(s_i) - r_i - \gamma V_{\theta_{i-1}}(s_{i+1}) \right)^2 \Big|_{\theta = \theta_{i-1}}$$

$s$  $r$  $s'$

Temporal difference learning

TD learning

TD(0)   TD($\lambda$)

linear: $V_\theta(s) = \phi(s)^T \theta \Rightarrow \nabla_\theta V_\theta(s) = \phi(s)$

$$2 \left( V_\theta(s_i) - r_i - \gamma V_{\theta_{i-1}}(s_{i+1}) \right) \nabla_\theta V_\theta(s_i) \Big|_{\theta = \theta_{i-1}}$$

# LSPEQ and ⟨TDQ⟩ *TD*

Collecting samples $\{(s_i, a_i, r_i, s_i')\}_{i=1}^n$ using $\pi$

For $k = 1, 2, \dots$

$$\theta_k \leftarrow \underset{\theta}{\text{argmin}} \sum_{i=1}^n \left( Q_\theta(s_i, a_i) - r_i - \gamma \sum_{a'} \pi(a'|s_i') \, Q_{\theta_{k-1}}(s_i', a') \right)^2$$

For $i = 1, 2, \dots$

Draw $a_i \sim \pi(\cdot | s_i)$, observe reward $r_i$ and next state $s_{i+1}$

$$\theta_i \leftarrow \theta_{i-1} - \alpha \nabla_\theta \left( Q_\theta(s_i, a_i) - r_i - \gamma \sum_{a'} \pi(a'|s_i') \, Q_{\theta_{i-1}}(s_i', a') \right)^2$$

# TD with Linear Function Approximation

$A \succeq 0 :$ A is psd

$A \succeq B \iff A - B$ is psd

BC : $R(s,a) + \gamma \underset{s' \sim P(\cdot|s,a)}{\mathbb{E}} \max_{a'} \widetilde{Q}(s',a') = \phi(s,a)^\top \theta^* \quad \forall \widetilde{Q}$

imply $\left( \text{set } \widetilde{Q} = Q^* \right)$

Let $\mu$ be the stationary state distribution under policy $\pi$. Furthermore, assume

(1) $V^\pi(s) = \phi(s)^\top \theta^\star$      (realizability assumption)

(2) $\mathbb{E}_{s \sim \mu}[\phi(s)\phi(s)^\top] \succeq \rho I$ for some $\rho > 0$    (coverage assumption)

Then the following TD update:

Realizability in $Q^*$ :

$Q^*(s,a) = \phi(s,a)^\top \theta^*$

For $i = 1, 2, \ldots$

In fact, even if the samples are generated as $a_i \sim \pi(\cdot|s_i)$, $r_i = R(s_i, a_i)$, $s_{i+1} \sim P(\cdot|s_i, a_i)$

    Sample $s \sim \mu$,   $a \sim \pi(\cdot|s)$,   $r \sim R(s,a)$,   $s' \sim P(\cdot|s,a)$

    $\theta_i \leftarrow \theta_{i-1} - \alpha_i (\phi(s)^\top \theta_{i-1} - r - \gamma \phi(s')^\top \theta_{i-1}) \phi(s)$

converges to $\theta^\star$ with properly chosen $\alpha_i$.

$$V^{\pi}(s) = \phi(s)^T \theta^*$$

$$\|\theta_i - \alpha g - \theta^*\|^2 = \|\theta_i - \theta^*\|^2 - 2\alpha(\theta_i - \theta^*)^T g + \alpha^2 \|g\|^2$$

condition $\theta_i$

$$\|\theta_{i+1} - \theta^*\|^2 = \left\| \theta_i - \alpha\left(\phi(s)^T \theta_i - r - \gamma \phi(s')^T \theta_i\right)\phi(s) - \theta^* \right\|^2 \quad \text{where} \quad \left( s \sim \mu, \ a \sim \pi(\cdot|s), \ \mathbb{E}(r) = R(s,a) \atop s' \sim P(\cdot|s,a) \right)$$

$$= \|\theta_i - \theta^*\|^2 - 2\alpha(\theta_i - \theta^*)^T \left(\phi(s)^T \theta_i - r - \gamma \phi(s')^T \theta_i\right)\phi(s) + \alpha^2 \|g\|^2$$

$$\mathbb{E}\left[\|\theta_{i+1} - \theta^*\|^2\right] = \|\theta_i - \theta^*\|^2 \underbrace{- 2\alpha(\theta_i - \theta^*)^T \left[\mathbb{E}\left[ \overbrace{\phi(s)^T\theta_i - r - \gamma\phi(s')^T\theta_i}^{g}\right]\phi(s)\right]}_{} + \alpha^2 \mathbb{E}\left(\|g\|^2\right)$$

$$+ \mathbb{E}\left[\;\;\;\;\;\right]$$

$$\text{blue} = \mathbb{E}\left[-V^{\pi}(s) + r + \gamma V^{\pi}(s')\right]$$
$$= \mathbb{E}\left(-V^{\pi}(s) + \sum_a \pi(a|s)\left(R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)}(V^{\pi}(s'))\right)\right)$$
$$= 0$$

$$\Rightarrow \mathbb{E}\left[\|\theta_{i+1} - \theta^*\|^2\right] = \|\theta_i - \theta^*\|^2 - 2\alpha(\theta_i - \theta^*)^T \left[\mathbb{E}\left[\left(\phi(s)^T\theta_i - r - \gamma\phi(s')^T\theta_i\right)\phi(s)\right] - \mathbb{E}\left[\left(\phi(s)^T\theta^* - r - \gamma\phi(s')^T\theta^*\right)\phi(s)\right]\right] + \alpha^2 \mathbb{E}\left(\|g\|^2\right)$$

$$\left(V^{\pi}(s) - r - \gamma V^{\pi}(s')\right)\phi(s)$$
$$\mathbb{E}_{a \sim \pi(s)}\; \mathbb{E}_{s' \sim P(\cdot|s,a)}\left[V^{\pi}(s')\right]$$

$$= \|\theta_i - \theta^*\|^2 - 2\alpha(\theta_i - \theta^*)^T \mathbb{E}\left[\left(\phi(s)^T(\theta_i - \theta^*) - \gamma\phi(s')^T(\theta_i - \theta^*)\right)\phi(s)\right] + \alpha^2 \mathbb{E}\left(\|g\|^2\right)$$

$$= \|\theta_i - \theta^*\|^2 - 2\alpha(\theta_i - \theta^*)^T \mathbb{E}\left[\phi(s)\left(\phi(s)^T - \gamma\phi(s')^T\right)\right](\theta_i - \theta^*) + \alpha^2 \mathbb{E}(\|g\|^2)$$

# Comparison

$\mathbb{E}(a)\mathbb{E}(a) \le \mathbb{E}(a^2)$

$s \sim \mu, \; a \sim \pi(\cdot|s)$   $s' \sim P(\cdot|s,a)$

On-Policy

Why does **Linear TD and Linear TDQ** converge (and converges to the correct solution) but **Linear Q-Learning** diverges?

Off-Policy

$$\mathbb{E}\left[\phi(s)\phi(s)^T - \gamma\phi(s)\phi(s')^T\right]$$

$$= \underset{s\sim\mu}{\mathbb{E}}\left[\phi(s)\phi(s)^T\right] - \gamma\underset{s\sim\mu}{\mathbb{E}}\left[\phi(s)\right]\underset{s\sim\mu}{\mathbb{E}}\left(\phi(s)\right)^T \succeq (1-\gamma)\underset{s\sim\mu}{\mathbb{E}}\left[\phi(s)\phi(s)^T\right]$$

$$\succeq (1-\gamma)\rho I$$

$$\mathbb{E}\left[\|\theta_{i+1}-\theta^*\|^2\right] \le \|\theta_i-\theta^*\|^2 - 2\alpha(\theta_i-\theta^*)^T\underbrace{\mathbb{E}\left[\phi(s)\left(\phi(s)-\gamma\phi(s')\right)^T\right]}_{\succeq cI}(\theta_i-\theta^*) + \alpha^2\|g\|^2$$

$$\le \|\theta_i-\theta^*\|^2 - 2\alpha c\|\theta_i-\theta^*\|^2 + \alpha^2\|g\|^2$$

$$\le (1-2\alpha c)\|\theta_i-\theta^*\|^2 + \alpha^2\|g\|^2$$

Q-learning  $(s,a,r,s')$ from arbitrary policy

$$Q^k(s,a) \leftarrow (1-\alpha)Q^{k-1}(s,a) + \alpha\left(r + \gamma\max_{a'}Q^{k-1}(s',a')\right)$$

TD  $(s,a,r,s') \; s\sim\mu, \; a\sim\pi(\cdot|s)$

$$Q^k(s,a) \leftarrow (1-\alpha)Q^{k-1}(s,a) + \alpha\left(r + \gamma\sum_{a'}\pi(a'|s')Q^{k-1}(s',a')\right)$$

# Comparison

Under coverage assumption
(i.e., the data $\{(s_i, a_i, r_i, s_i{}')\}$ sufficiently cover every state-action pair / feature space)

|  | **LSVI** | **Watkins's Q-Learning** | **On-Policy LSPE(Q) / TD(Q)** |
|---|---|---|---|
| Tabular | $Q^{(k)} \to Q^\star$ | $Q^{(k)} \to Q^\star$ | $V^{(k)} \to V^\pi$ / $Q^{(k)} \to Q^\pi$ under realizability |
| Linear Approx. | $Q^{(k)} \to Q^\star$ under Bellman completeness | Diverges even with Bellman completeness | |

# Monte Carlo Estimation

Start from $s_1 = s^\star$

Execute policy $\pi$ until the episode ends and obtain trajectory

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_\tau, a_\tau, r_\tau$$

Let $G = \sum_{h=1}^{\tau} \gamma^{h-1} r_h$

$\pi(a)$

$G$ is an unbiased estimator for $V^\pi(s^\star)$

**MC estimator**: unbiased, higher variance

**TD estimator**: biased, lower variance

a sample of arm $a$'s reward

$\widehat{R}(a)$

# A Family of Estimators

Suppose we have a **value function estimation** $V_\theta(s) \approx V^\pi(s)$

Suppose we also have a **trajectory** $s_1, a_1, r_1, \dots, s_\tau, a_\tau, r_\tau$ generated by $\pi$

Then the following are all valid estimators for $V^\pi(s_1)$ besides $V_\theta(s_1)$:

$G_1 = r_1 + \gamma V_\theta(s_2)$

$G_2 = r_1 + \gamma r_2 + \gamma^2 V_\theta(s_3)$

...

$G_\tau = r_1 + \gamma r_2 + \gamma^2 r_3 + \cdots + \gamma^{\tau-1} r_\tau + \square$

*higher bias*

$\dfrac{w_1 G_1 + w_2 G_2 + \cdots + w_\tau G_\tau}{}$

$w_1 + w_2 + \cdots w_\tau = 1$

*higher variance*

$+ \odot$

$G_{\tau+1}$

$G_{\tau+2}$

Below, we will show

$TD(0) = TD \text{ learning}$

$\lambda \in (0, 1]$

$TD(1) = MC \text{ estimation}$
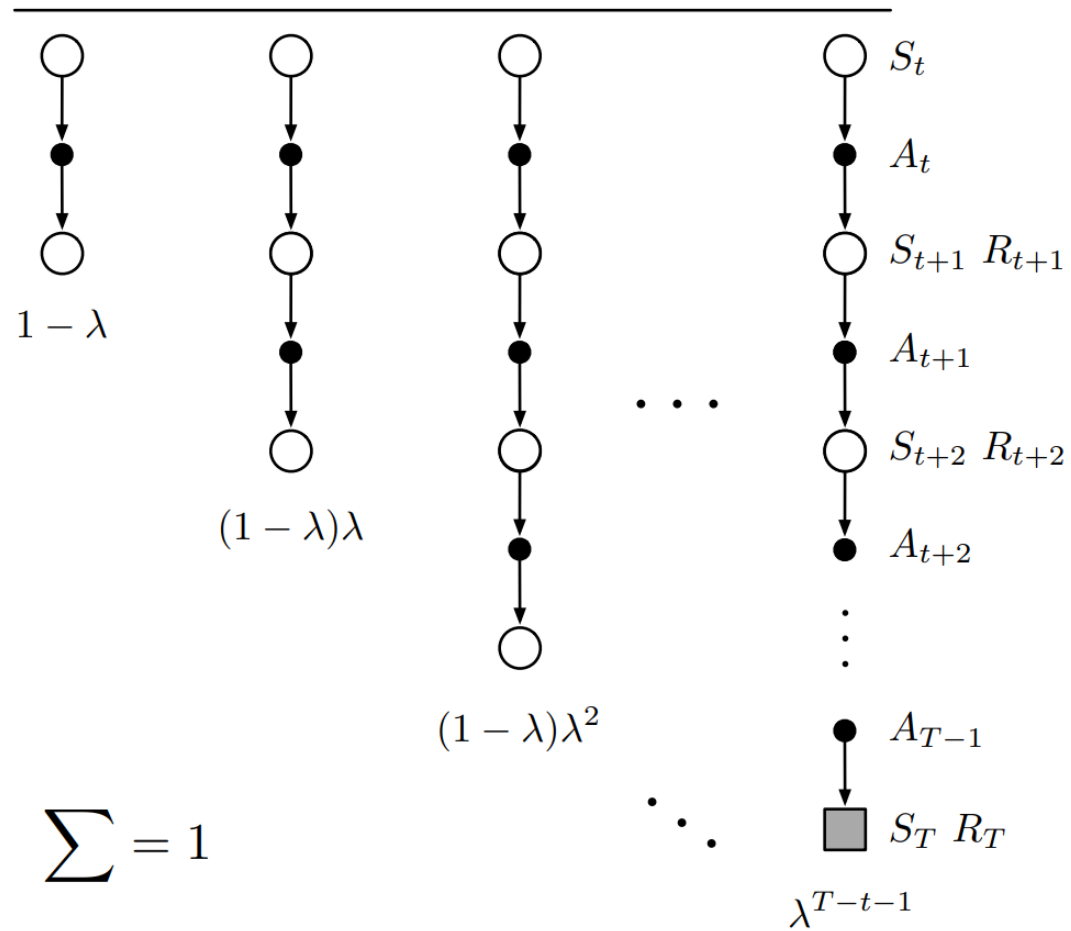
1. A way to combine these estimators

2. A more general policy evaluation method TD($\lambda$) based on these estimators

$\text{TD}(\lambda)$

$S_t$

$A_t$

$S_{t+1}\ R_{t+1}$

$1 - \lambda$

$A_{t+1}$

$S_{t+2}\ R_{t+2}$

$(1 - \lambda)\lambda$

$A_{t+2}$

$(1 - \lambda)\lambda^2$

$A_{T-1}$

$\sum = 1$

$S_T\ R_T$

$\lambda^{T-t-1}$

# Striking a Balance Between Bias and Variance

$$(1-\lambda) + (1-\lambda)\lambda + (1-\lambda)\lambda^2 + \cdots$$

$$\underline{G_\theta(\lambda)} = (1-\lambda)\big(G_1 + \lambda G_2 + \lambda^2 G_3 + \cdots\big)$$

$$= (1-\lambda)\big(r_1 + \gamma V_\theta(s_2)\big) + (1-\lambda)\lambda\big(r_1 + \gamma r_2 + \gamma^2 V_\theta(s_3)\big) + (1-\lambda)\lambda^2(\cdots) + \cdots$$

# TD($\lambda$)

$$\boxed{S_1, a_1, r_1}\ S_2, a_2, r_2, \cdots\ , S_t, a_t, r_t$$

$$G_{\theta_k}(0)$$

$$\text{TD(0):}\quad \theta_{k+1} \leftarrow \theta_k - \alpha \nabla_\theta \left( V_\theta(s_1) - r_1 - \gamma V_{\theta_k}(s_2) \right)^2$$

$$\text{TD}(\lambda):\quad \theta_{k+1} \leftarrow \theta_k - \alpha \nabla_\theta \left( V_\theta(s_1) - G_{\theta_k}(\lambda) \right)^2$$

$$S_1, a_1, r_1, S_2, a_2, r_2, \cdots\ , S_t, a_t, r_t$$

**Implementation details:**

How to make update before reaching the end of the episode?
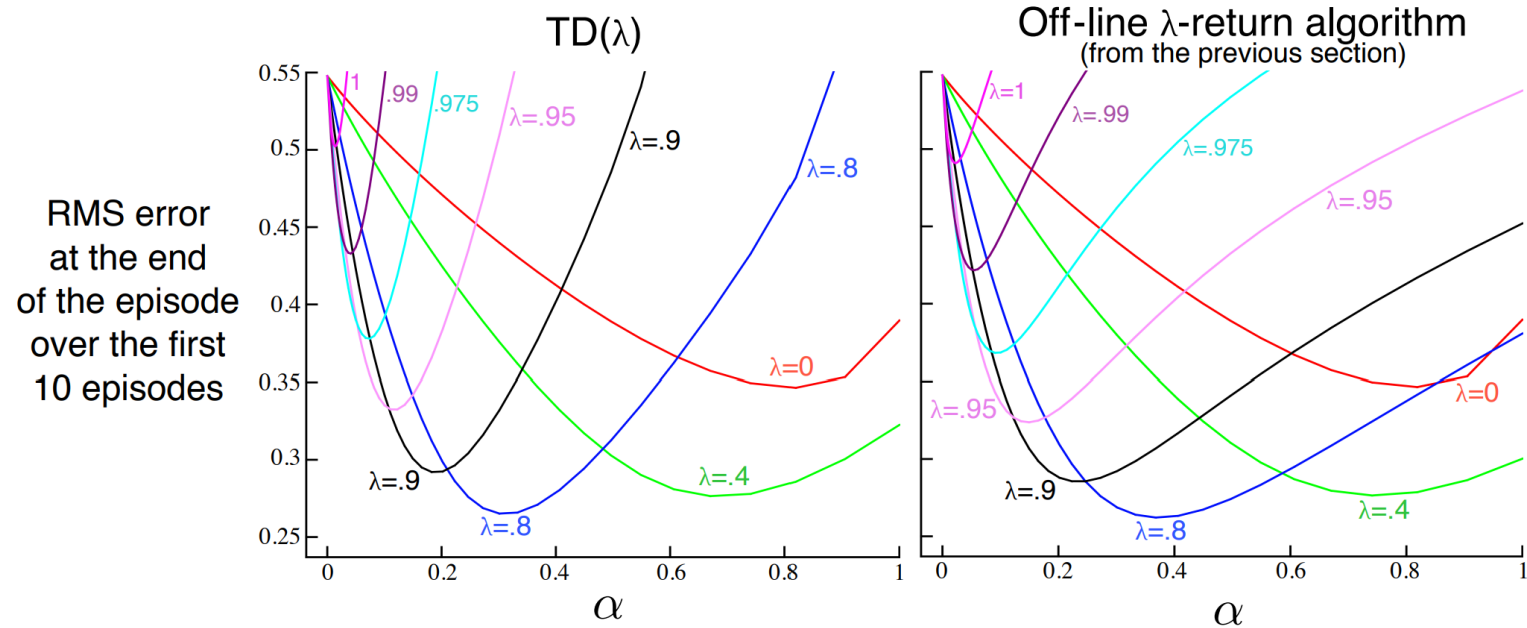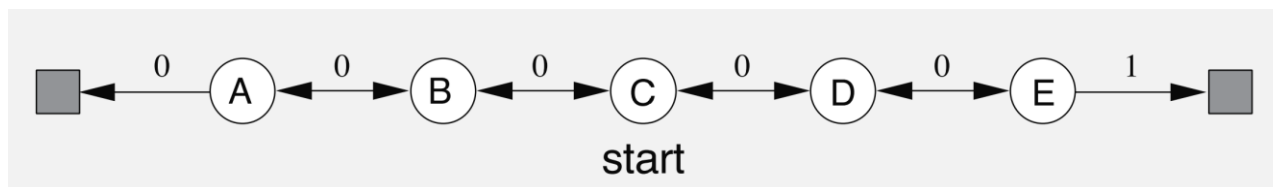
([Sutton and Barto](#) Chapter 12)

# TD(λ)



**Figure 12.6:** 19-state Random walk results (Example 7.1): Performance of TD($\lambda$) alongside that of the off-line $\lambda$-return algorithm. The two algorithms performed virtually identically at low (less than optimal) $\alpha$ values, but TD($\lambda$) was worse at high $\alpha$ values.

(Sutton and Barto Chapter 12)

# Summary: On-Policy Policy Evaluation

- Double time-scale: **LSPE, LSPEQ**, Single time-scale: **TD, TDQ**
- TD (TD(0)) update:

$$(s, a, r, s') \sim \pi$$

$$\theta_{i+1} \leftarrow \theta_i - \alpha \, \nabla_{\theta} \left( V_{\theta}(s) - r - \gamma V_{\theta_i}(s') \right)^2 \Big|_{\theta = \theta_i}$$

- In the linear case, when realizability and coverage hold, we can show $\theta_i \to \theta^\star$
- Monte Carlo Estimator
- An estimator with parameter $\lambda$ that balances variance and bias
- TD($\lambda$)

# Off-Policy Policy Evaluation

# Off-Policy LSPEQ / TDQ

Collecting samples $\{(s_i, a_i, r_i, s_i')\}_{i=1}^n$ using $\boldsymbol{\pi_b}$

For $k = 1, \ 2, \dots$

$$\theta_k \leftarrow \underset{\theta}{\text{argmin}} \sum_{i=1}^n \left( Q_\theta(s_i, a_i) - r_i - \gamma \sum_{a'} \pi(a'|s_i') \, Q_{\theta_{k-1}}(s_i', a') \right)^2$$

Bellman completeness + coverage will make it work

For $i = 1, \ 2, \dots$

Draw $a_i \sim \boldsymbol{\pi_b}(\cdot \,|s_i)$, observe reward $r_i$ and next state $s_{i+1}$

$$\theta_i \leftarrow \theta_{i-1} - \alpha \nabla_\theta \left( Q_\theta(s_i, a_i) - r_i - \gamma \sum_{a'} \pi(a'|s_i') \, Q_{\theta_{i-1}}(s_i', a') \right)^2$$

Like Q-learning, this is not stable

# Off-Policy LSPE

Collecting samples $\{(s_i, a_i, r_i, s_i')\}_{i=1}^n$ using $\boldsymbol{\pi_b}$

For $k = 1, \ 2, \dots$

$$\theta_k \leftarrow \underset{\theta}{\mathrm{argmin}} \sum_{i=1}^{n} \frac{\pi(a_i|s_i)}{\pi_b(a_i|s_i)} \left( V_\theta(s_i) - r_i + \gamma V_{\theta_{k-1}}(s_i') \right)^2$$

Bellman completeness + coverage will make it work

(Sutton and Barto Chapter 11.7 and 11.8 have more techniques to deal with the $V_\theta$ case)