

Learning the value functions

Sébastien Gros

Cybernetic, NTNU
Elec. Eng., Chalmers

TUM lectures on RL

Outline

- 1 Learning V_π
- 2 Learning Q_π
- 3 Learning Q_*



Bellman equations

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

$$Q_\pi(x, u) = L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u]$$

- Policy evaluation $V_{k+1}(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_k(x_+) | x, \pi]$ (sweep over x)
- Policy improvement $\pi'(x) = \arg \min_u Q_\pi(x, u)$
- Dynamic programming

$$Q_+(x, u) \leftarrow L(x, u) + \gamma \mathbb{E}[V(x_+) | x, u]$$

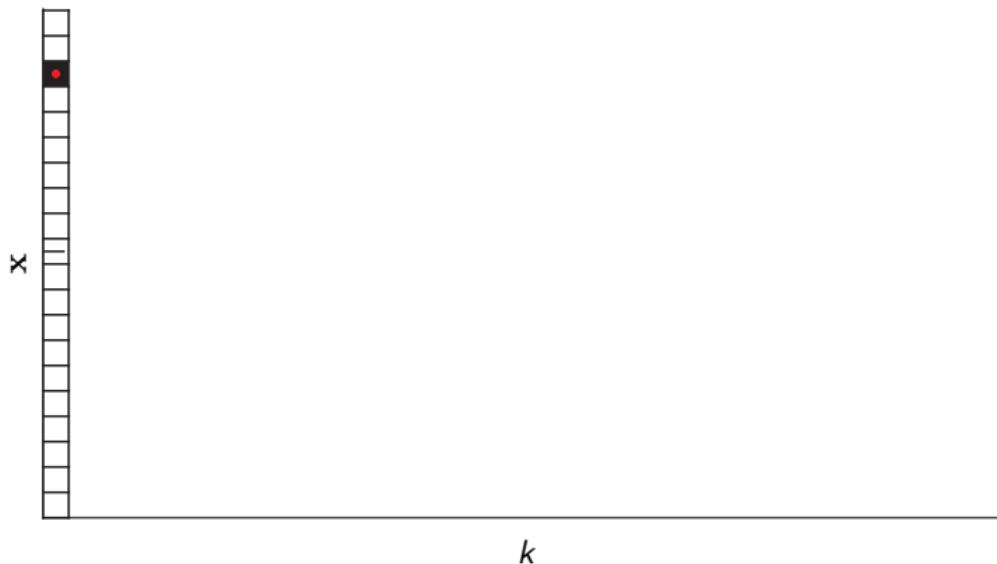
$$V(x) \leftarrow \min_u Q_+(x, u)$$

All these techniques require a model of the real system $\mathbb{P}[x_+ | x, u]$. What can we do from data?

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

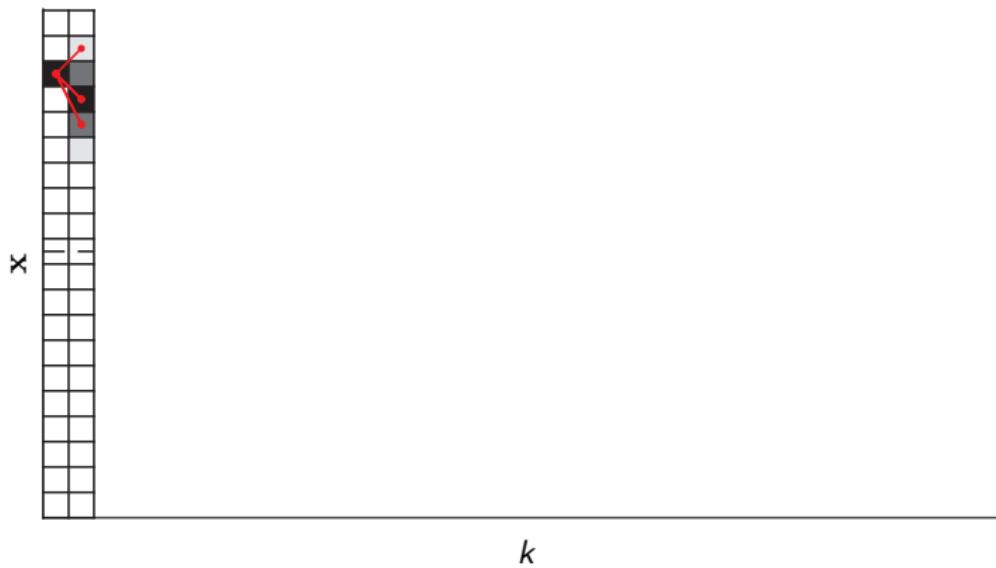
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

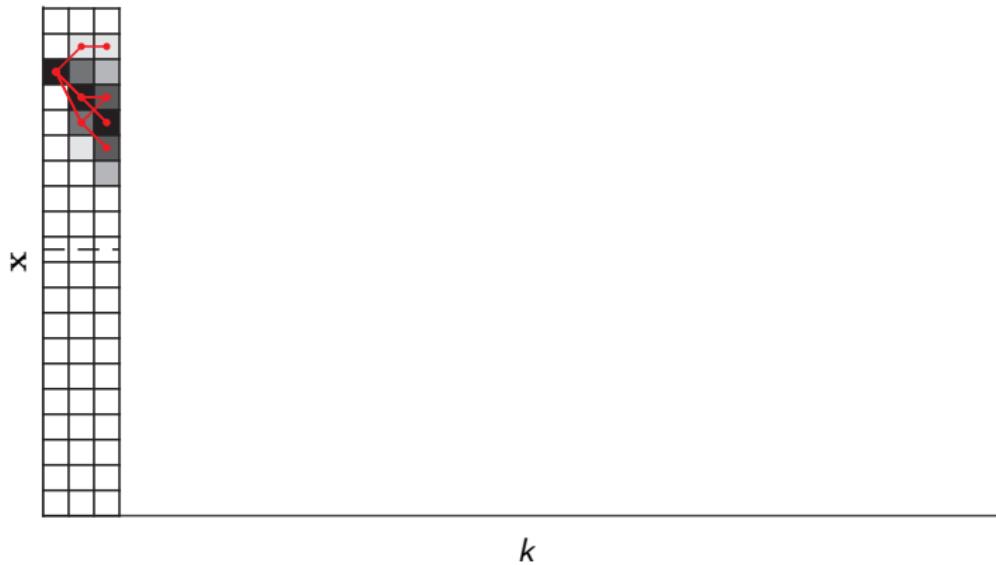


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

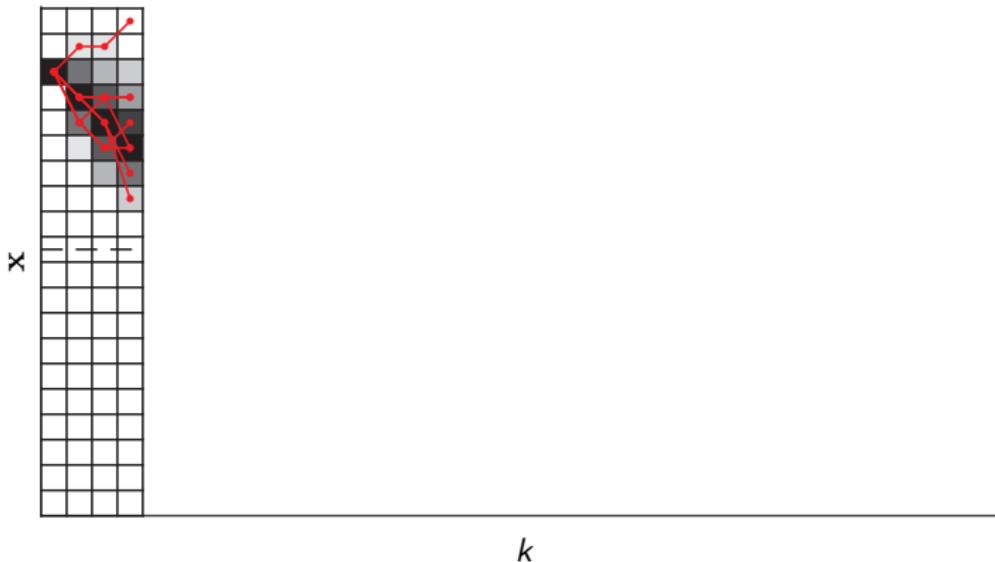
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

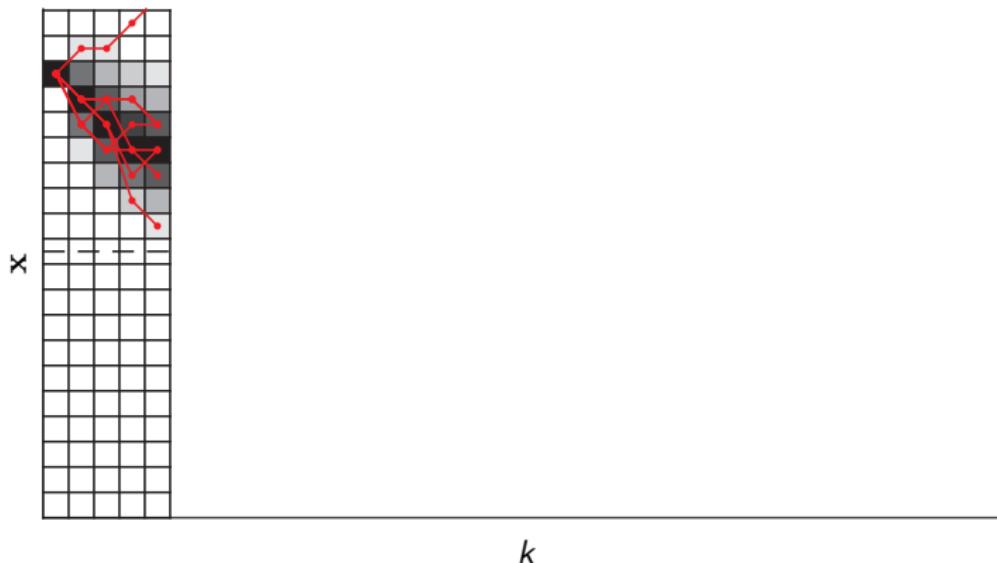


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

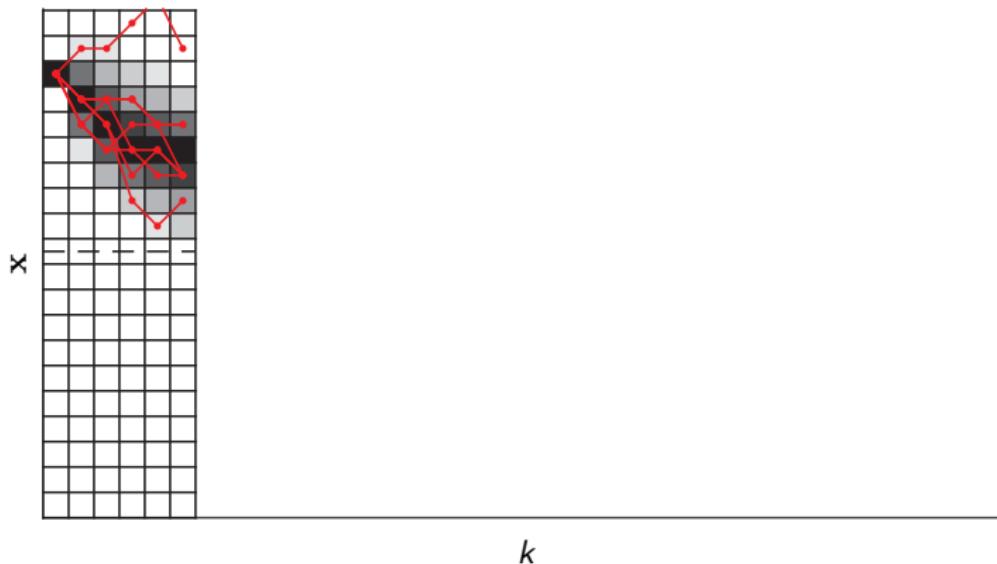


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

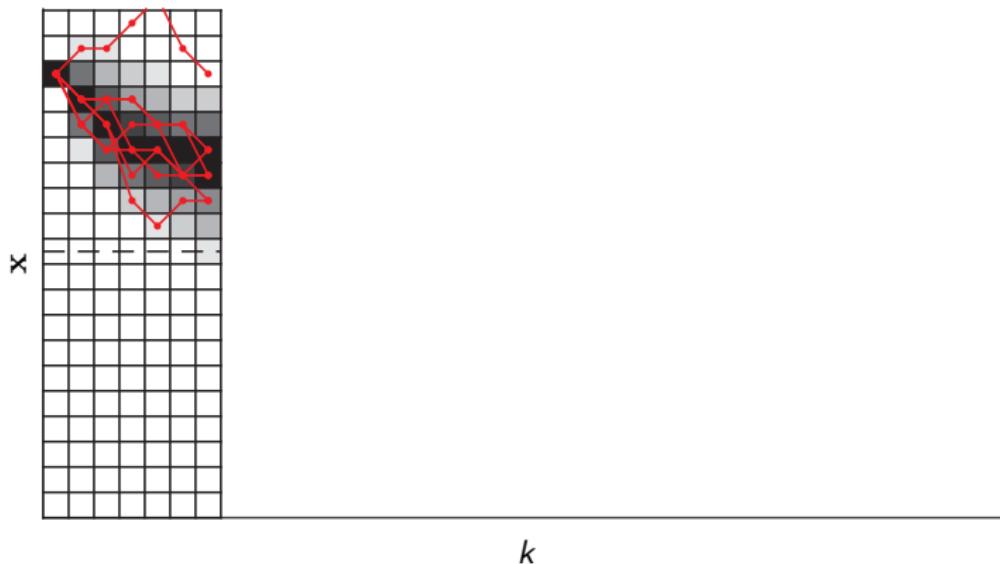


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

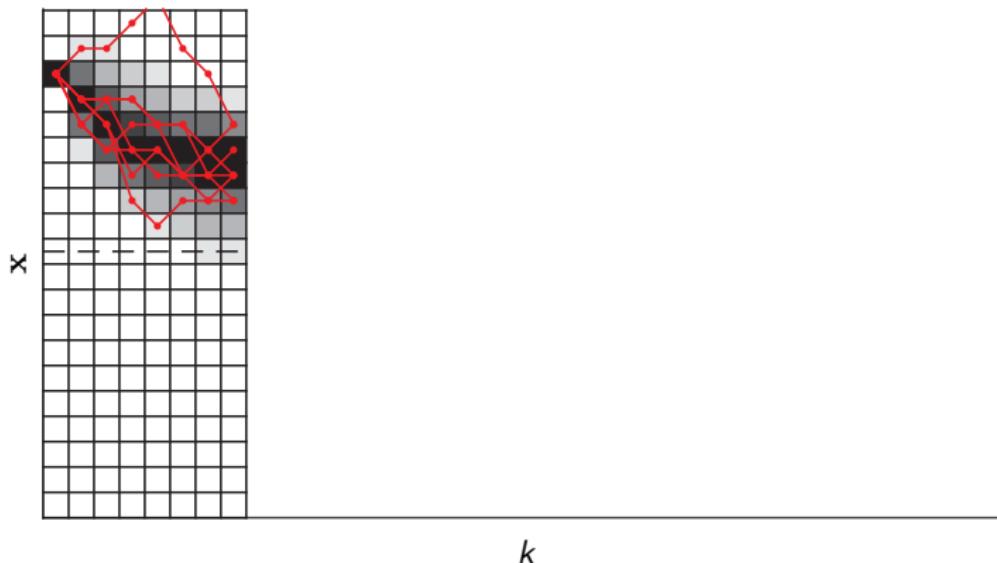


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
 - Input u can only move x on the grid
 - Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
 - Discount $\gamma = 0.9$

Dynamics: $\mathbf{x}_+ = \mathbf{x} + \mathbf{u} + \text{round}(\mathbf{e})$ where $\mathbf{e} \sim \mathcal{N}(0, 1)$

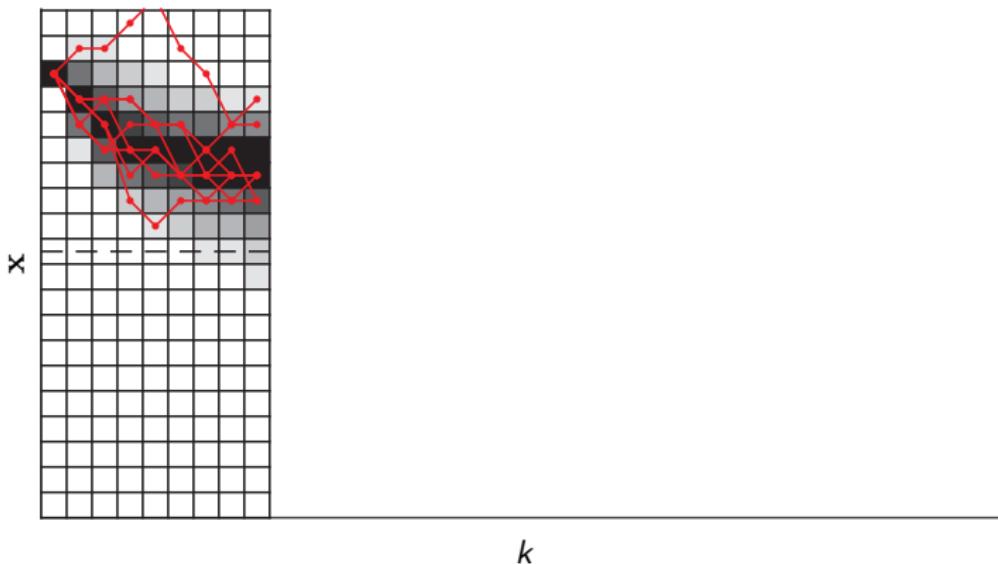


e.g. for
 $\mathbf{u} = \text{round}\left(-\frac{\mathbf{x}}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

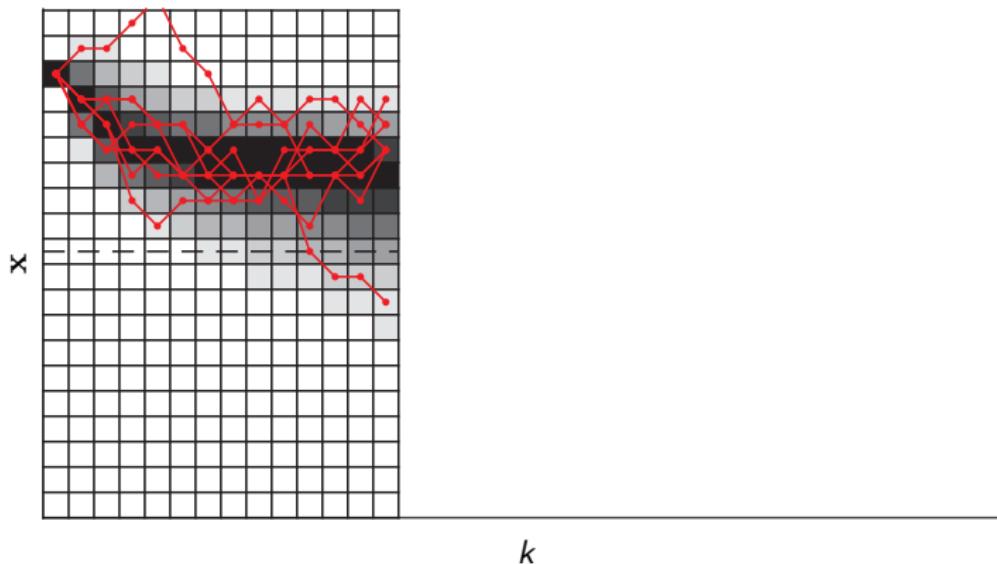


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

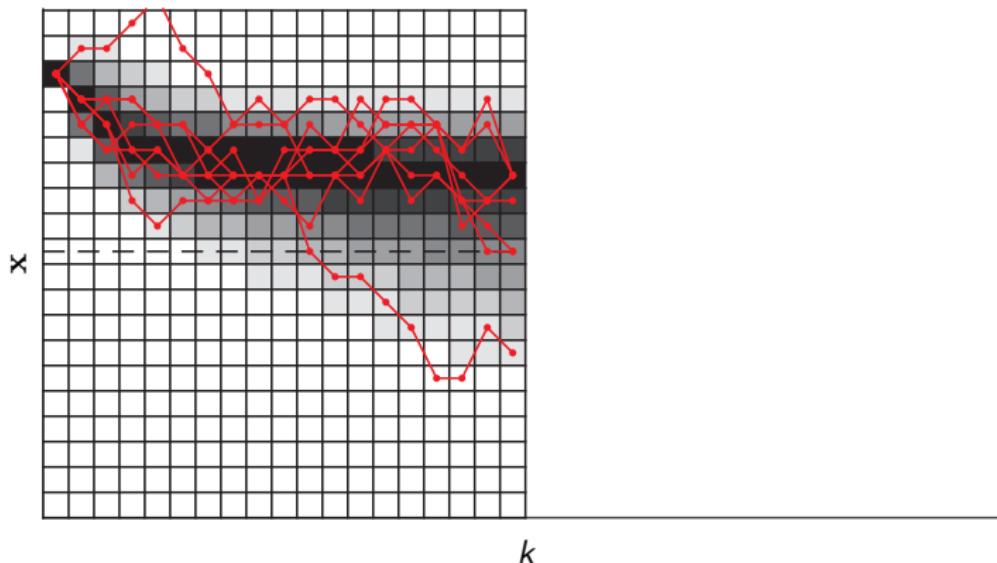


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

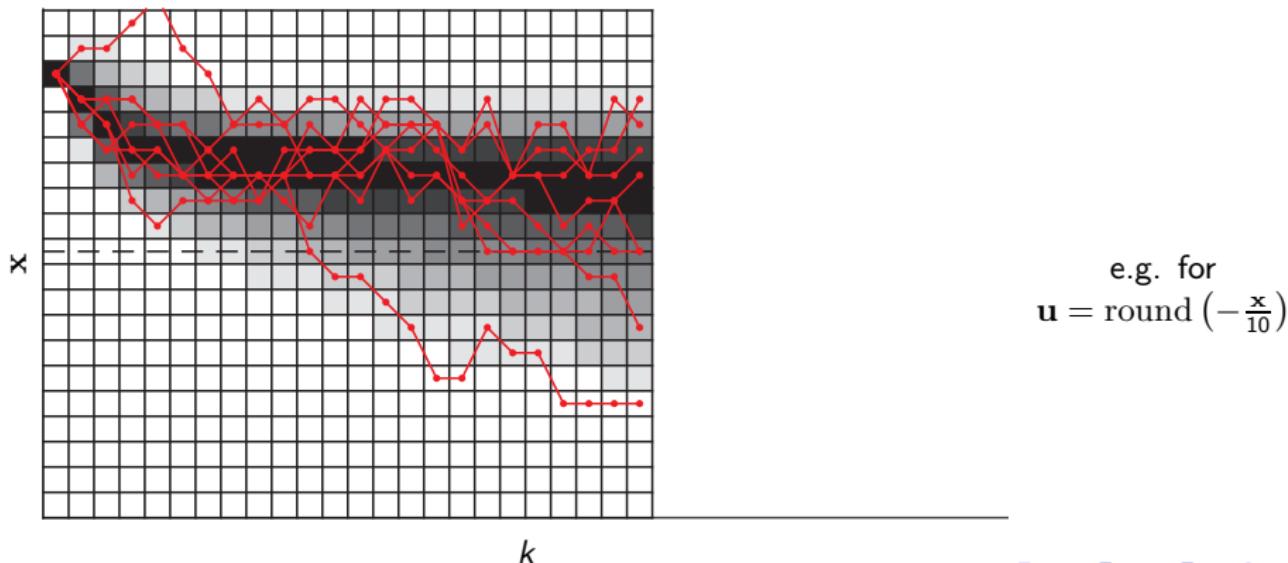


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

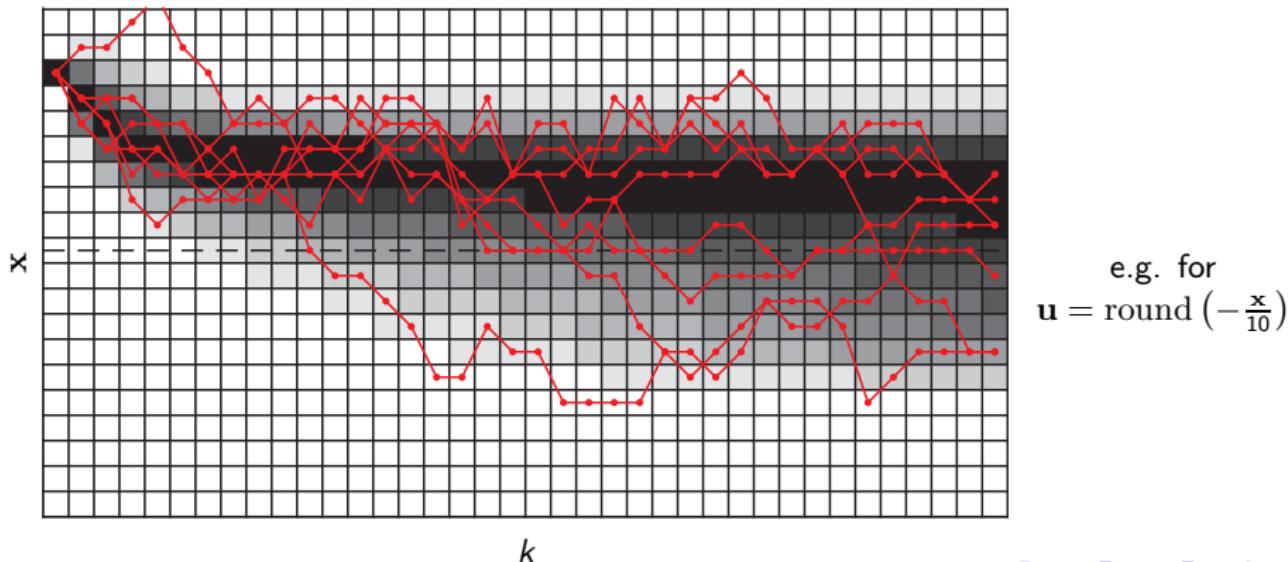
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Outline

1 Learning V_π

2 Learning Q_π

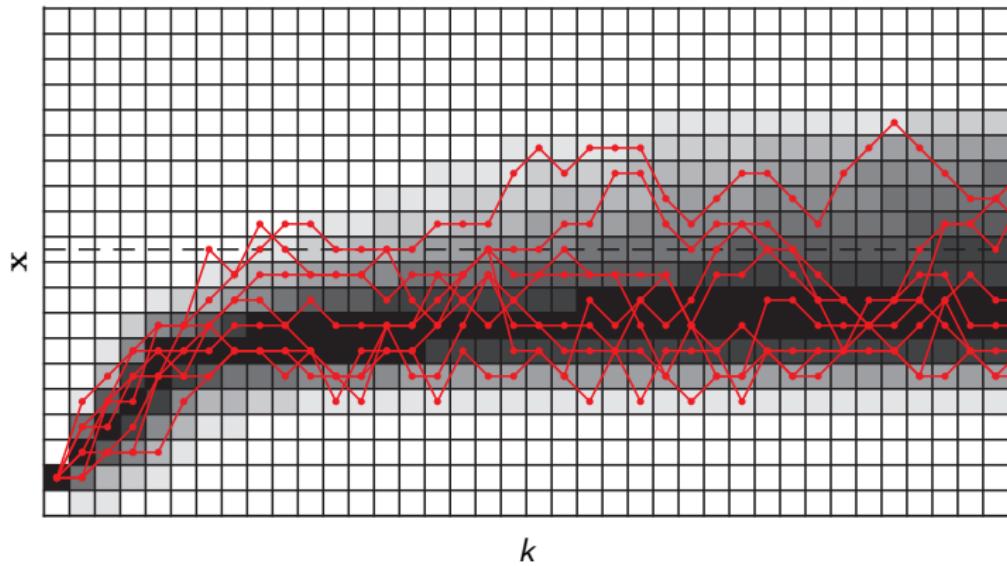
3 Learning Q_*

Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 155.18$$

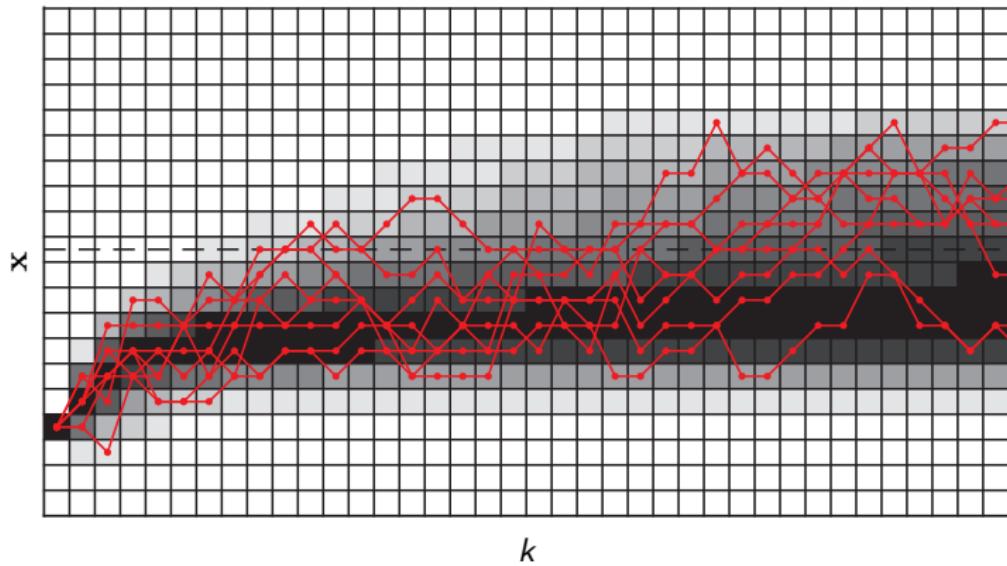


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 99.94$$

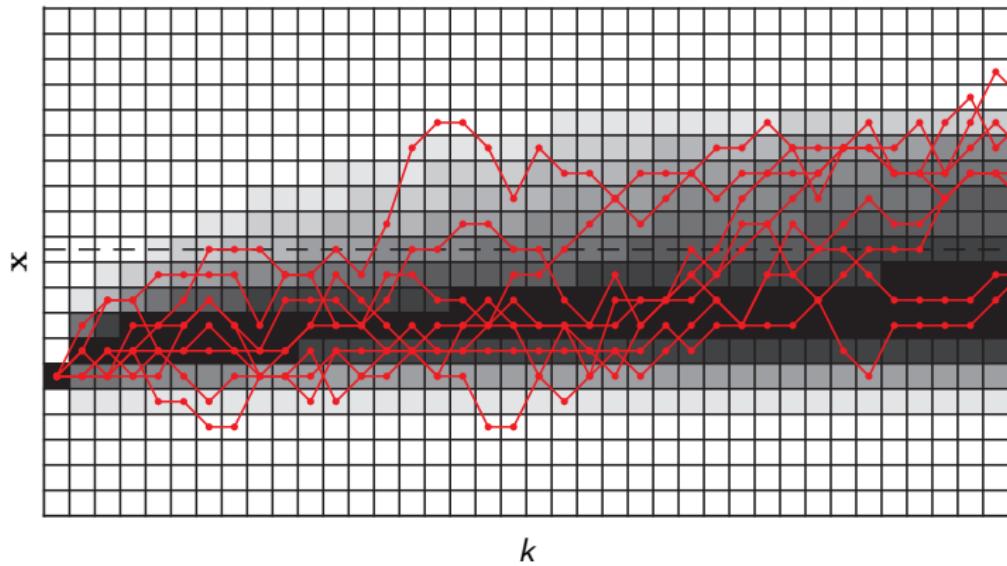


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 65.41$$

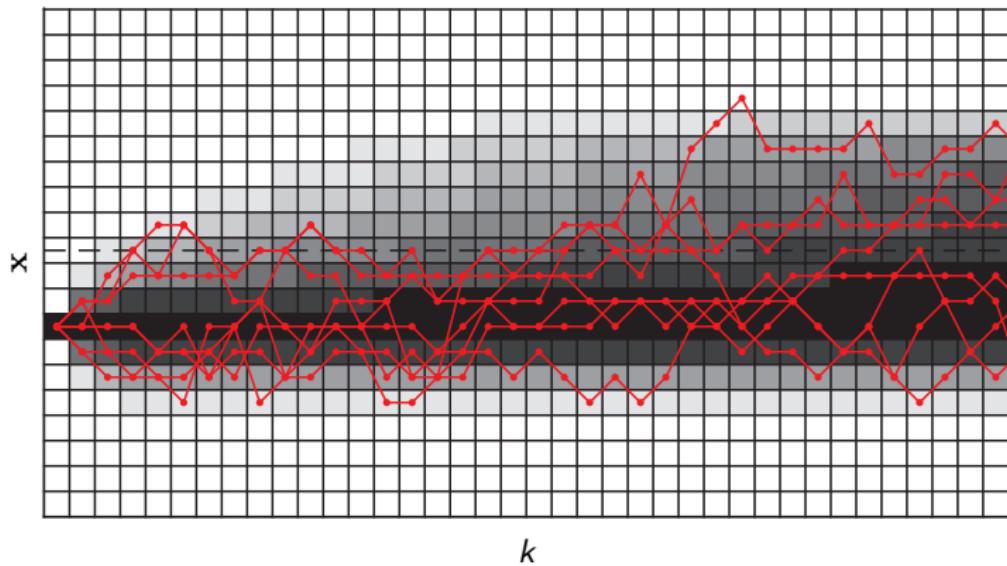


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 47.5$$

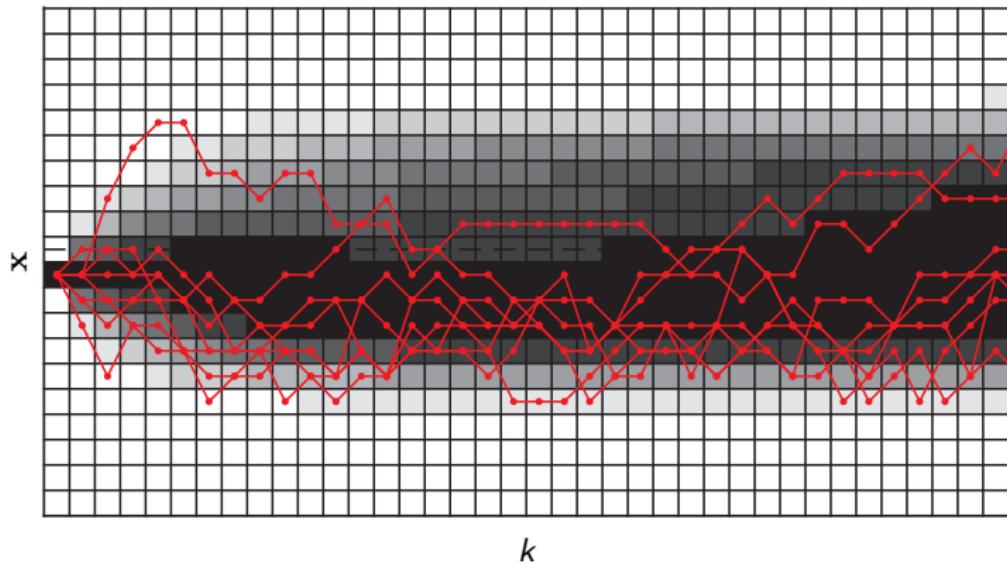


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 28.23$$

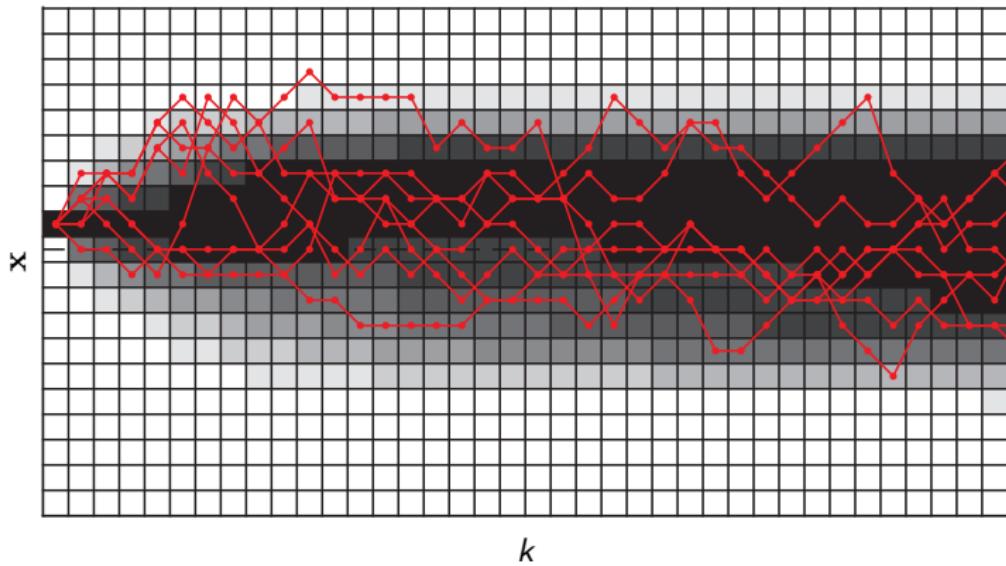


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 28.23$$

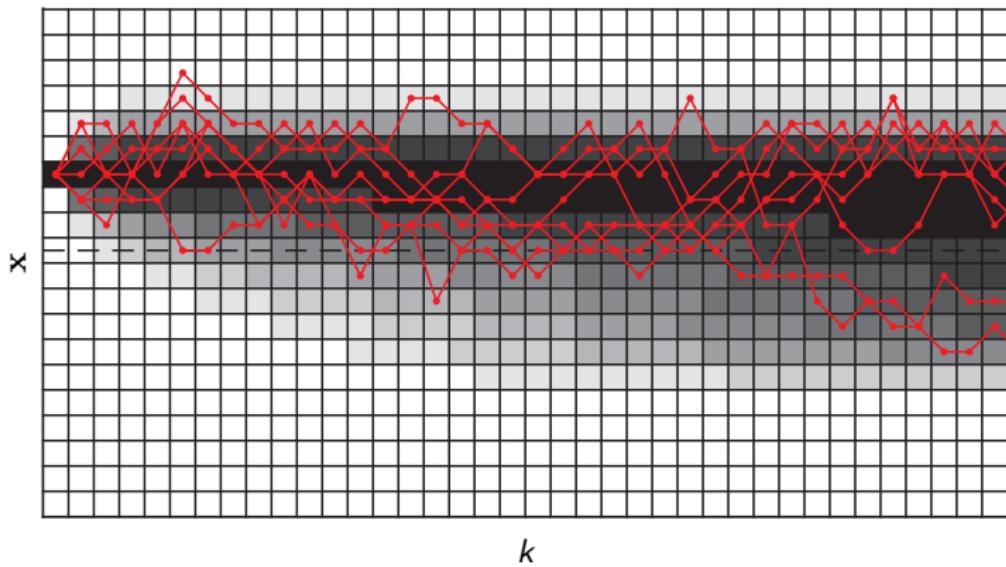


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 47.5$$

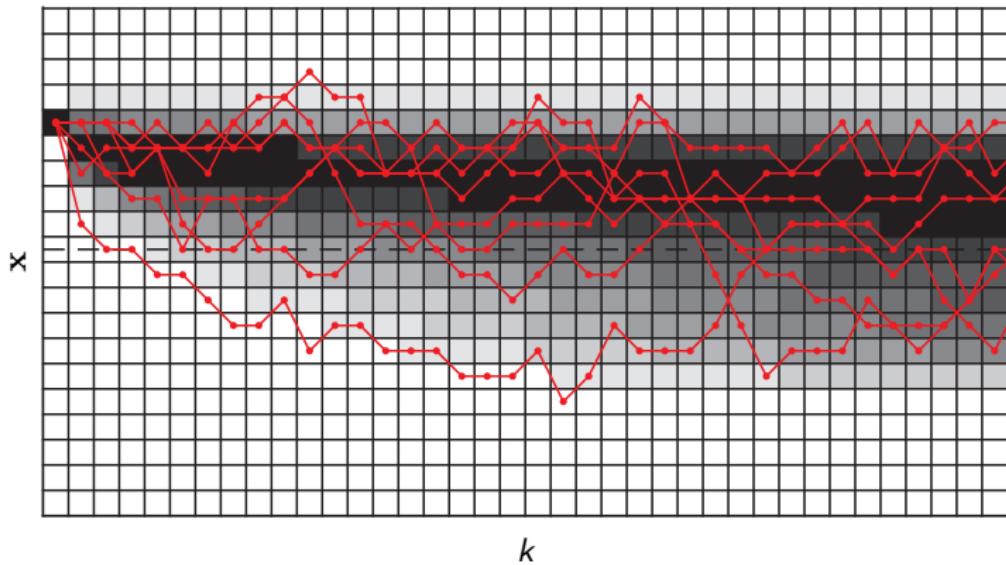


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 65.41$$

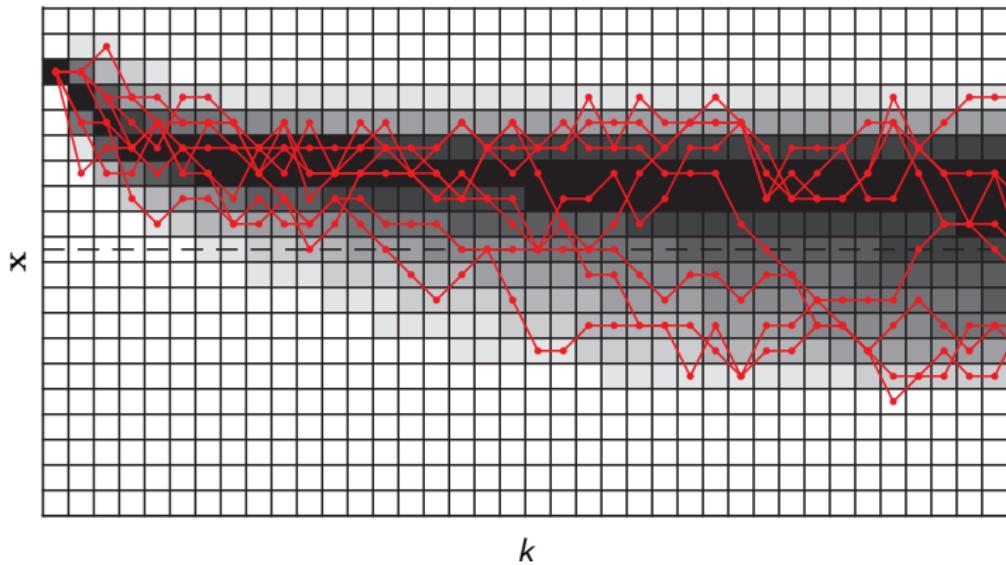


Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 99.94$$



Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

- Reminder **policy evaluation**: for a given policy π , iterate:

$$V_{\pi}(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_{\pi}(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

starting from any $V(\mathbf{x})$.

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

- Reminder **policy evaluation**: for a given policy π , iterate:

$$V_\pi(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_\pi(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

starting from any $V(\mathbf{x})$.

- Evaluation of $\mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$ requires a model:

$$\mathbb{E}[V_\pi(\mathbf{x}_+) \mid \mathbf{x}, \pi] = \int V_\pi(\mathbf{x}_+) \underbrace{\mathbb{P}[\mathbf{x}_+ \mid \mathbf{x}, \pi]}_{\text{Model}} d\mathbf{x}_+$$

Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

- Reminder **policy evaluation**: for a given policy π , iterate:

$$V_{\pi}(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_{\pi}(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

starting from any $V(\mathbf{x})$.

- Evaluation of $\mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$ requires a model:

$$\mathbb{E}[V_{\pi}(\mathbf{x}_+) \mid \mathbf{x}, \pi] = \int V_{\pi}(\mathbf{x}_+) \underbrace{\mathbb{P}[\mathbf{x}_+ \mid \mathbf{x}, \pi]}_{\text{Model}} d\mathbf{x}_+$$

- How can we evaluate $V_{\pi}(\mathbf{x})$ from data rather than from the model?

Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Monte Carlo (MC) evaluation:

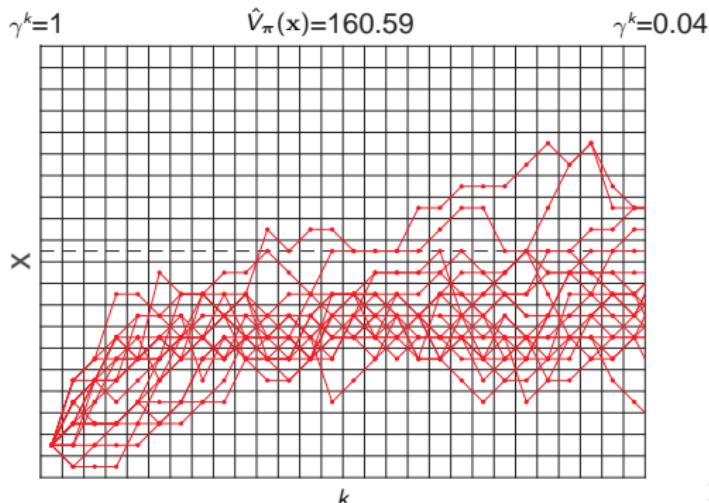
- Start the system in state x
- Unroll N trajectories using $\pi(x)$
- Collect trajectories $x^{(1)}, \dots, x^{(N)}$
- Return:

$$R_i = \sum_{k=0}^{\infty} \gamma^k L(x_k^{(i)}, \pi(x_k^{(i)}))$$

- Evaluate empirical mean:

$$\hat{V}_\pi(x) \approx \frac{1}{N} \sum_{k=1}^N R_i$$

E.g. for $N = 20$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Monte Carlo (MC) evaluation:

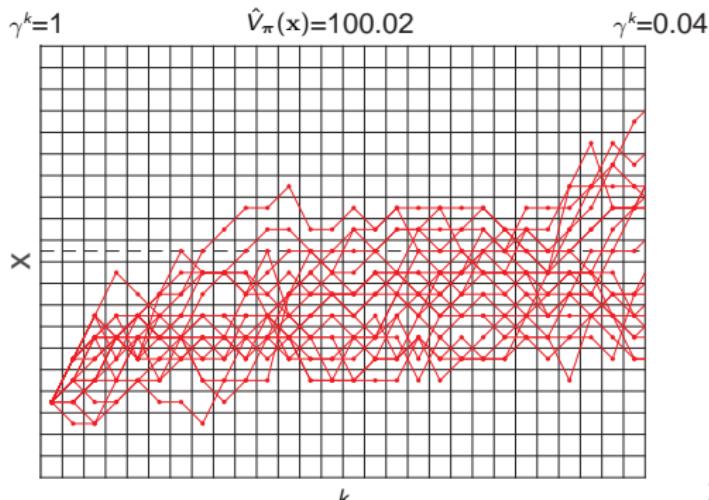
- Start the system in state x
- Unroll N trajectories using $\pi(x)$
- Collect trajectories $x^{(1)}, \dots, x^{(N)}$
- Return:

$$R_i = \sum_{k=0}^{\infty} \gamma^k L(x_k^{(i)}, \pi(x_k^{(i)}))$$

- Evaluate empirical mean:

$$\hat{V}_\pi(x) \approx \frac{1}{N} \sum_{k=1}^N R_i$$

E.g. for $N = 20$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Monte Carlo (MC) evaluation:

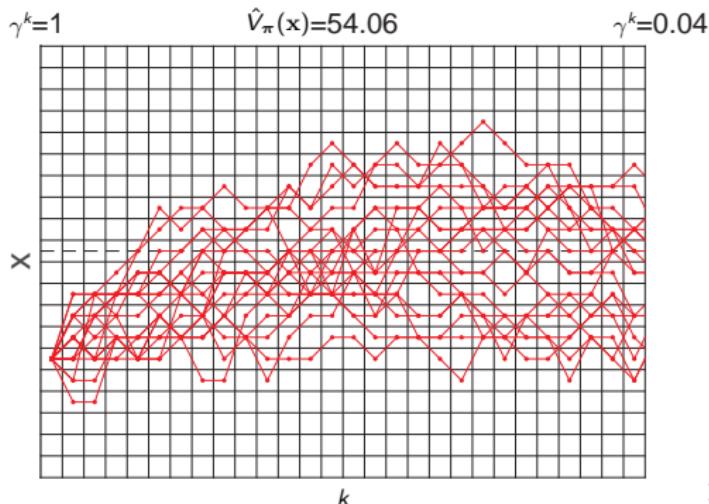
- Start the system in state x
- Unroll N trajectories using $\pi(x)$
- Collect trajectories $x^{(1)}, \dots, x^{(N)}$
- Return:

$$R_i = \sum_{k=0}^{\infty} \gamma^k L(x_k^{(i)}, \pi(x_k^{(i)}))$$

- Evaluate empirical mean:

$$\hat{V}_\pi(x) \approx \frac{1}{N} \sum_{k=1}^N R_i$$

E.g. for $N = 20$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Monte Carlo (MC) evaluation:

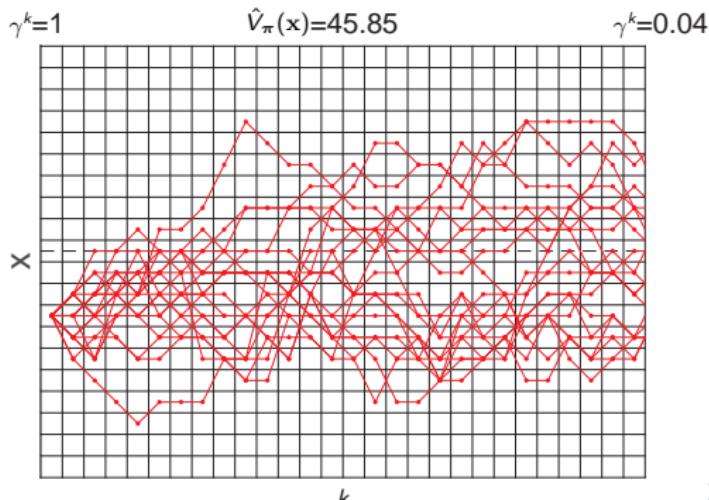
- Start the system in state x
- Unroll N trajectories using $\pi(x)$
- Collect trajectories $x^{(1)}, \dots, x^{(N)}$
- Return:

$$R_i = \sum_{k=0}^{\infty} \gamma^k L(x_k^{(i)}, \pi(x_k^{(i)}))$$

- Evaluate empirical mean:

$$\hat{V}_\pi(x) \approx \frac{1}{N} \sum_{k=1}^N R_i$$

E.g. for $N = 20$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Monte Carlo (MC) evaluation:

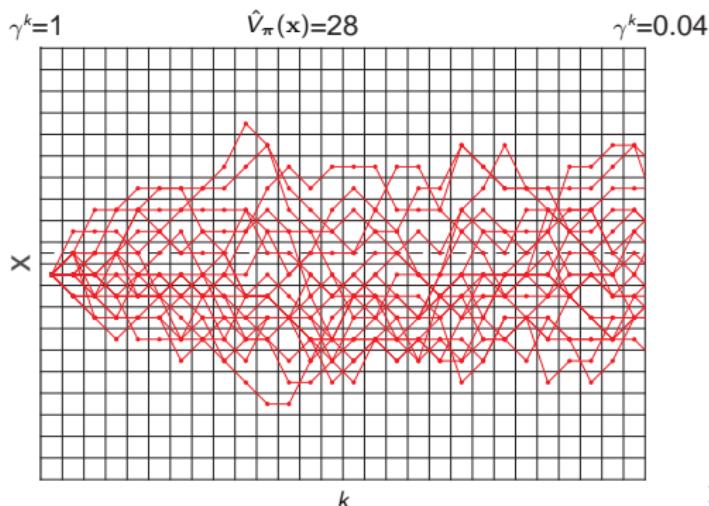
- Start the system in state x
- Unroll N trajectories using $\pi(x)$
- Collect trajectories $x^{(1)}, \dots, x^{(N)}$
- Return:

$$R_i = \sum_{k=0}^{\infty} \gamma^k L(x_k^{(i)}, \pi(x_k^{(i)}))$$

- Evaluate empirical mean:

$$\hat{V}_\pi(x) \approx \frac{1}{N} \sum_{k=1}^N R_i$$

E.g. for $N = 20$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

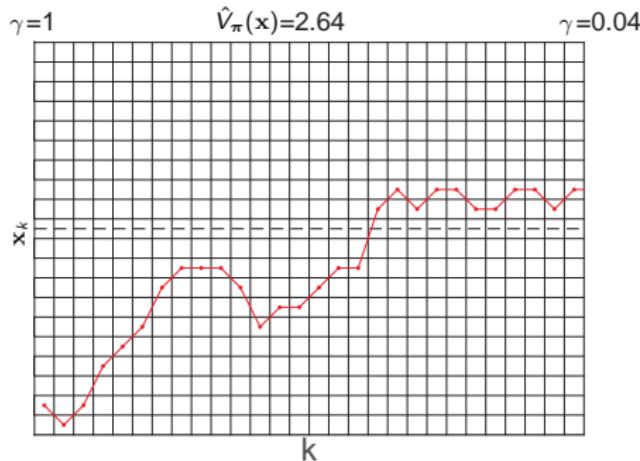
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a **noisy estimation** of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

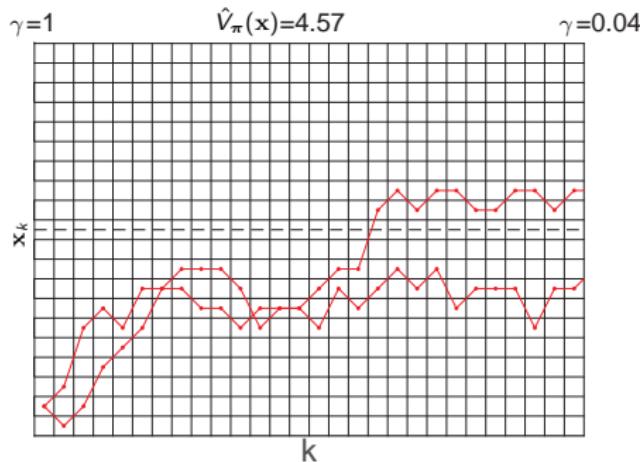
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a **noisy estimation** of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

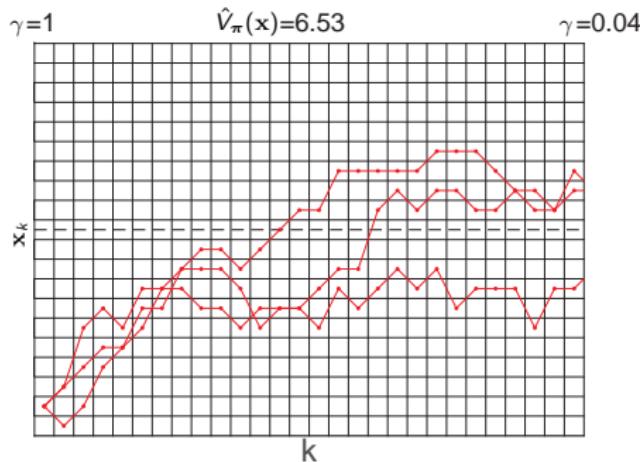
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a noisy estimation of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

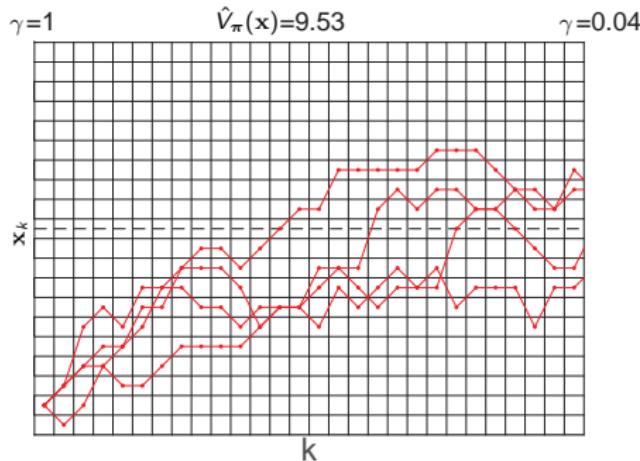
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a noisy estimation of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

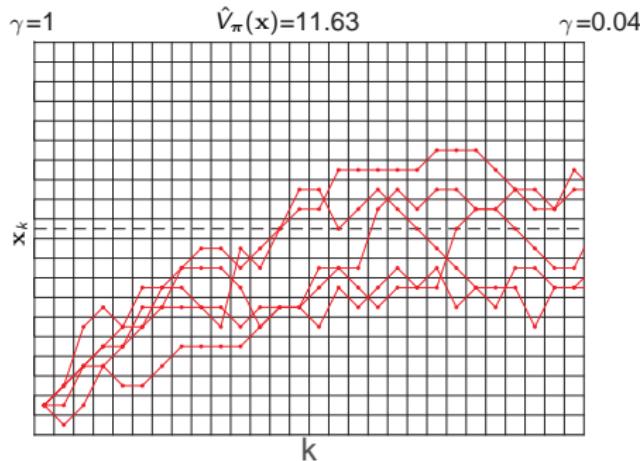
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a noisy estimation of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

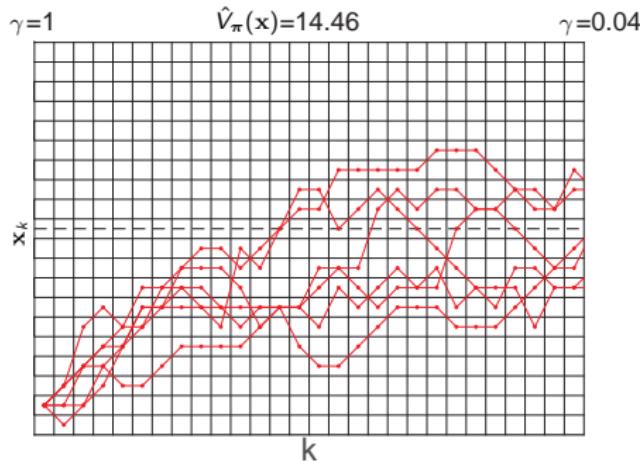
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a noisy estimation of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

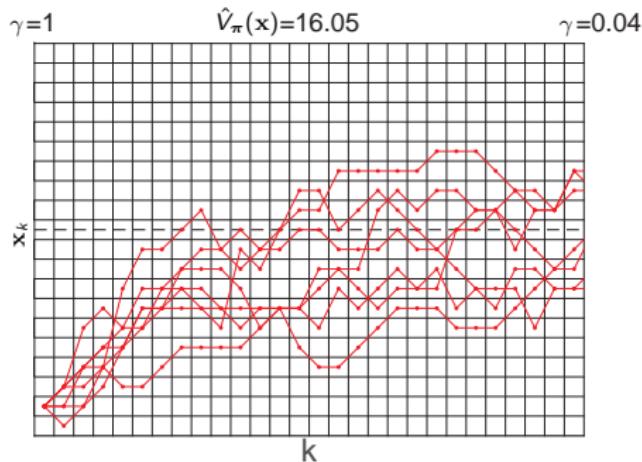
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Incremental Monte Carlo (MC)

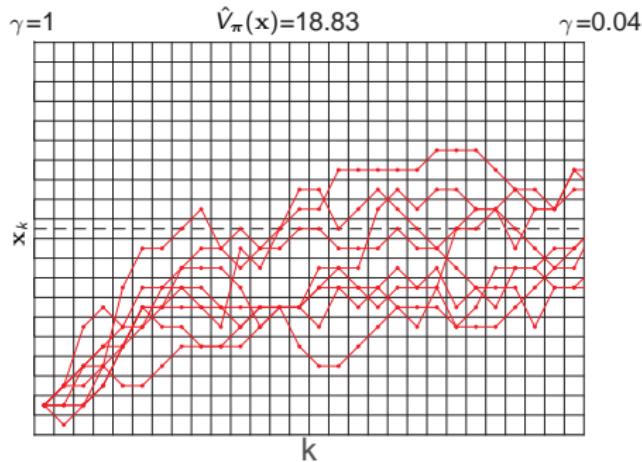
- Start the system in state $x \equiv x_0$
- Use $u_k = \pi(x_k)$
- Update $\hat{V}_\pi(x_0)$ as:

$$\hat{V}_\pi(x_0) \leftarrow \hat{V}_\pi(x_0) + \alpha (R - \hat{V}_\pi(x_0))$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k))$$

Note: R is a noisy estimation of $V_\pi(x_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

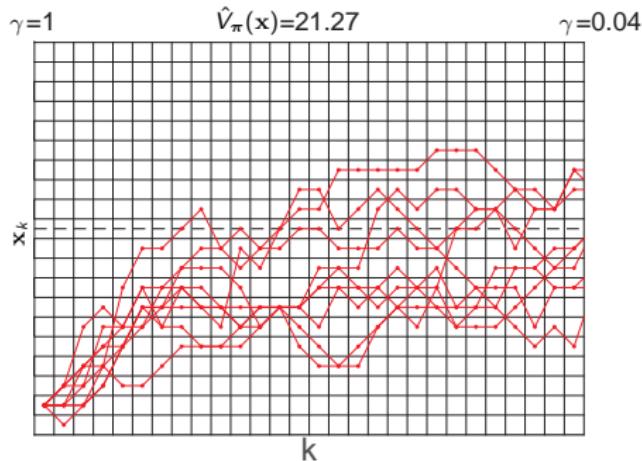
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a noisy estimation of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

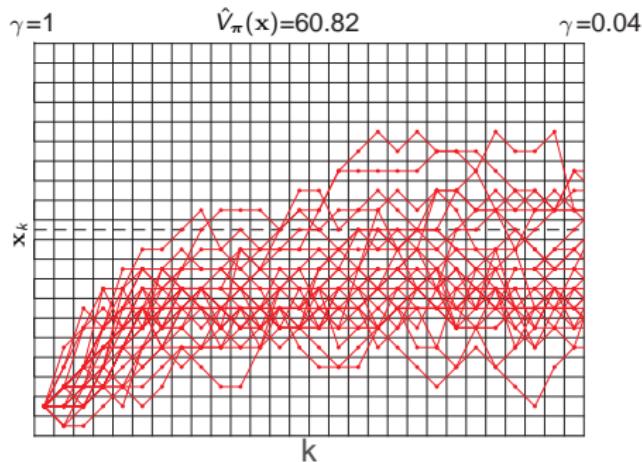
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

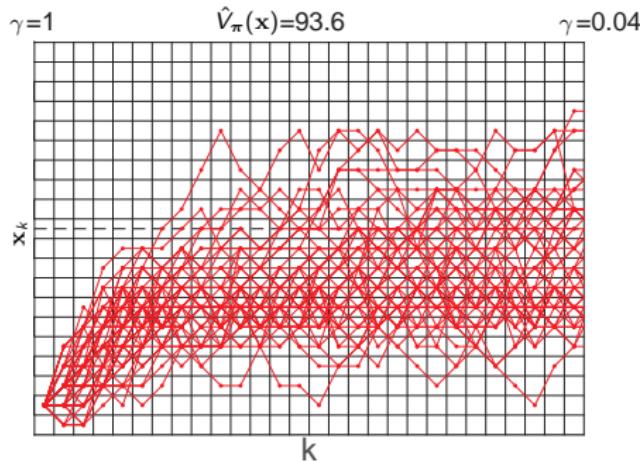
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

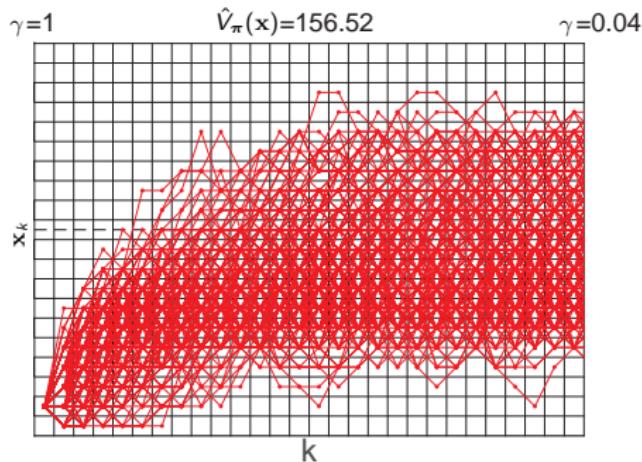
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

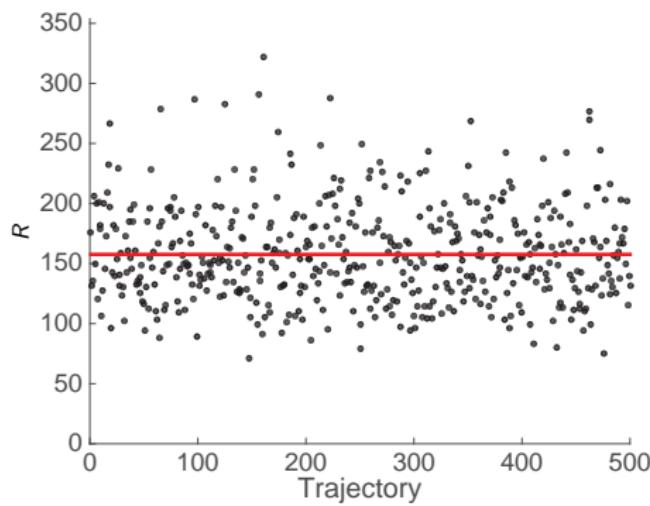
- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental Monte Carlo (MC)

- Start the system in state $\mathbf{x} \equiv \mathbf{x}_0$
- Use $\mathbf{u}_k = \pi(\mathbf{x}_k)$
- Update $\hat{V}_\pi(\mathbf{x}_0)$ as:

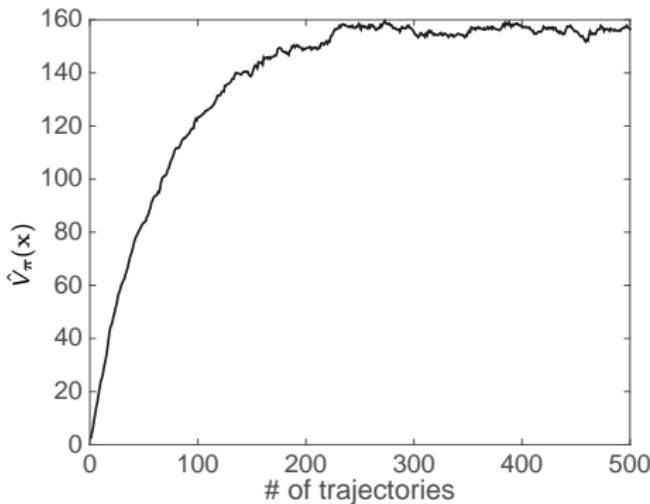
$$\hat{V}_\pi(\mathbf{x}_0) \leftarrow \hat{V}_\pi(\mathbf{x}_0) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}_0) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Note: R is a **noisy estimation** of $V_\pi(\mathbf{x}_0)$

Can build “on-the-fly” version

E.g. $\alpha = 5 \cdot 10^{-3}$



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental MC

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}) \right)$$

$$R = \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental MC

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}) \right)$$

$$R = L(\mathbf{x}, \pi(\mathbf{x})) + \sum_{k=1}^{\infty} \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k))$$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental MC

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}) \right)$$

$$R = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$$

Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental MC

$$\hat{V}_{\pi}(\mathbf{x}) \leftarrow \hat{V}_{\pi}(\mathbf{x}) + \alpha \left(R - \hat{V}_{\pi}(\mathbf{x}) \right)$$

$$R = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$$

Note that:

$$\mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) \mid \mathbf{x}_1 \right] = V_{\pi}(\mathbf{x}_1)$$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Incremental MC

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha \left(R - \hat{V}_\pi(\mathbf{x}) \right)$$

$$R = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))$$

Note that:

$$\mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) \mid \mathbf{x}_1 \right] = V_\pi(\mathbf{x}_1)$$

Temporal difference uses

$$R \approx L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \hat{V}_\pi(\mathbf{x}_1)$$

Learning the Value function

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_{\pi}(\mathbf{x}_k) \leftarrow \hat{V}_{\pi}(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_{\pi}(\mathbf{x}_{k+1}) - \hat{V}_{\pi}(\mathbf{x}_k)$$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Note: **Policy Evaluation** does

$$\begin{aligned} V_\pi(\mathbf{x}) &\leftarrow V_\pi(\mathbf{x}) + \bar{\delta} \\ \bar{\delta}(\mathbf{x}) &= \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x} \end{aligned}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

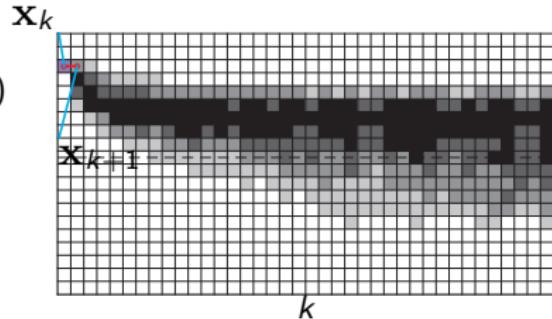
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

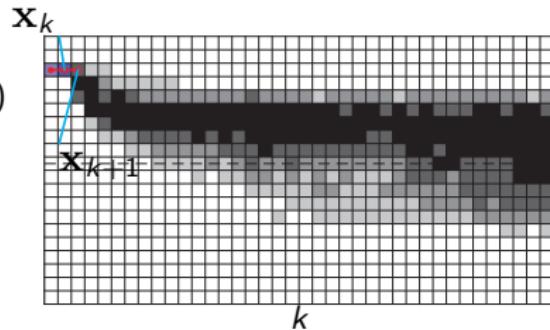
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

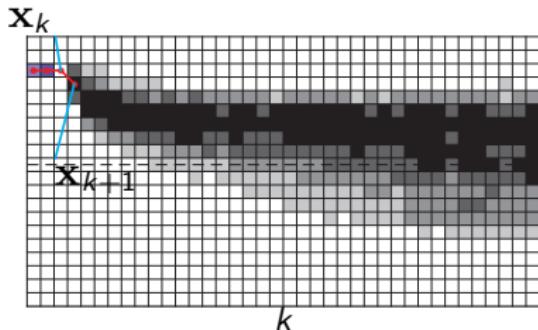
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

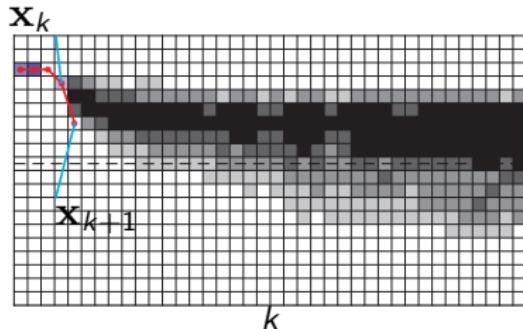
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

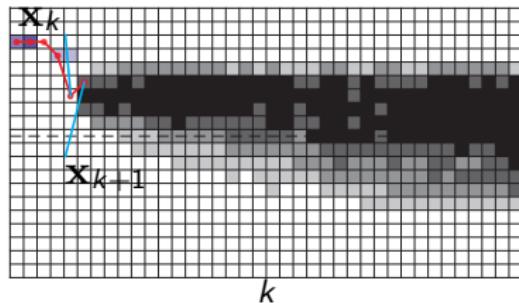
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

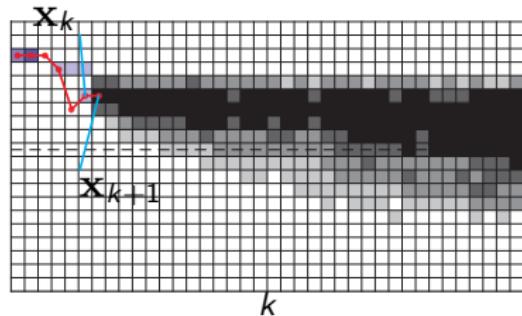
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

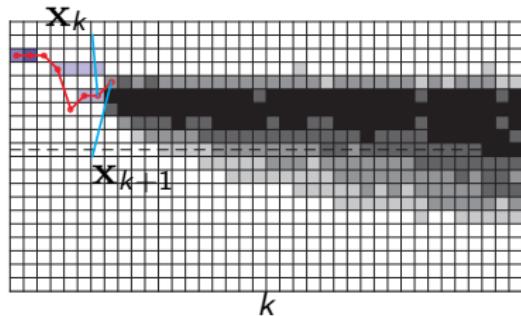
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

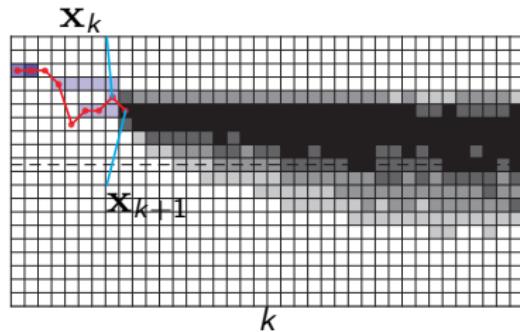
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

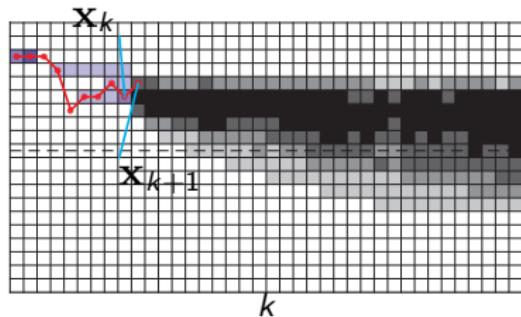
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

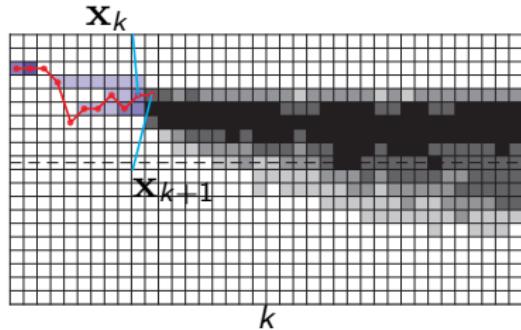
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

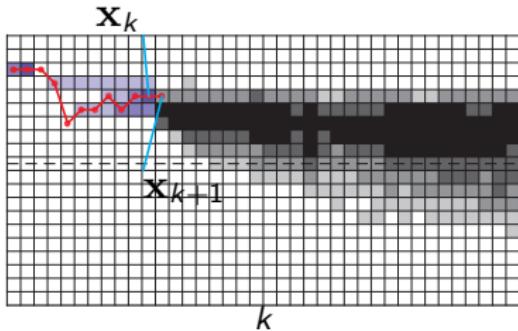
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

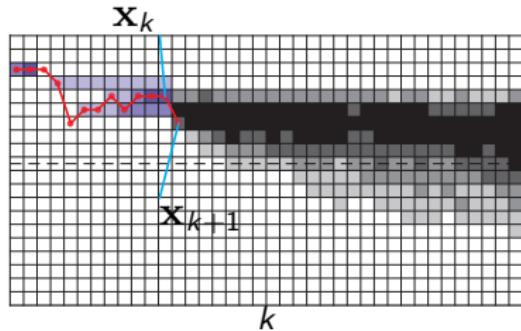
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$



hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size

Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

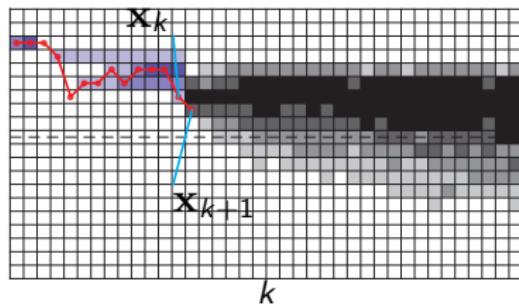
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

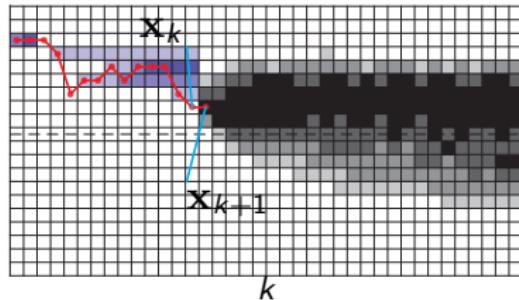
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

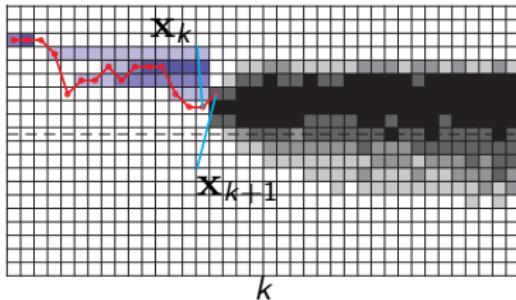
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

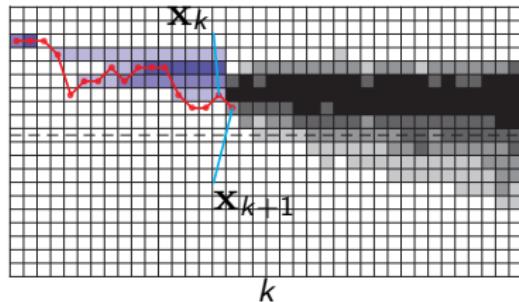
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

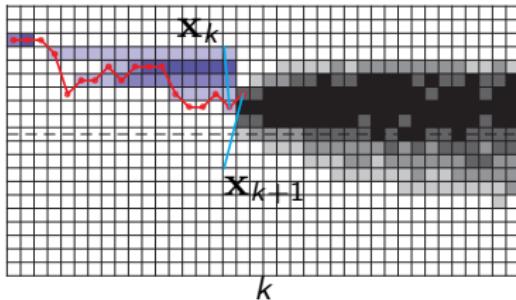
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

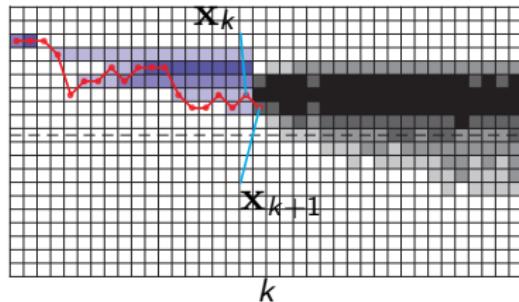
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

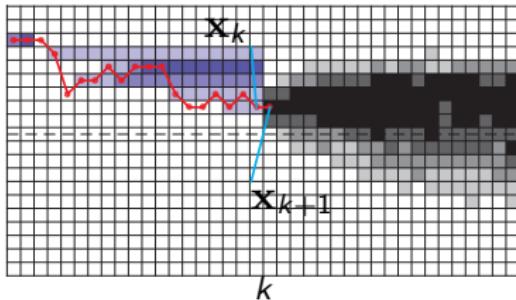
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta(x_k, x_{k+1})$$

$$\delta(x_k, x_{k+1}) = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

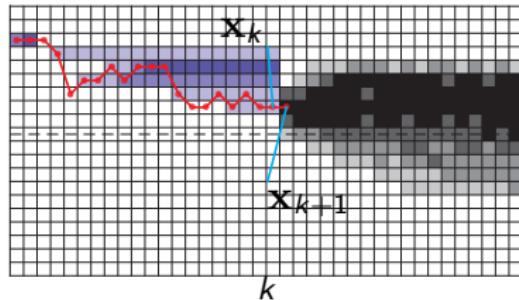
Note: **Policy Evaluation** does

$$V_\pi(x) \leftarrow V_\pi(x) + \bar{\delta}$$

$$\bar{\delta}(x) = \mathbb{E} [\delta(x, x_+) \mid x], \quad \forall x$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

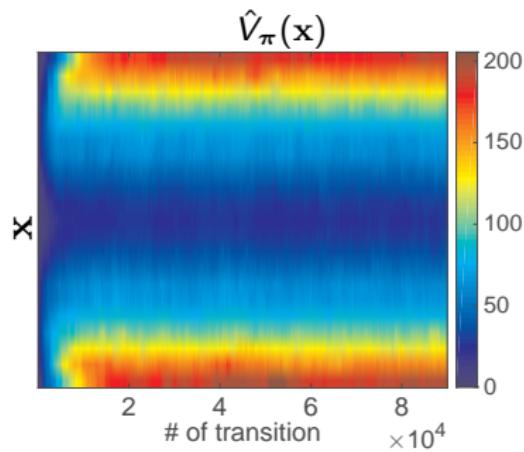
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



10^5 trajectories of length $N = 10$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

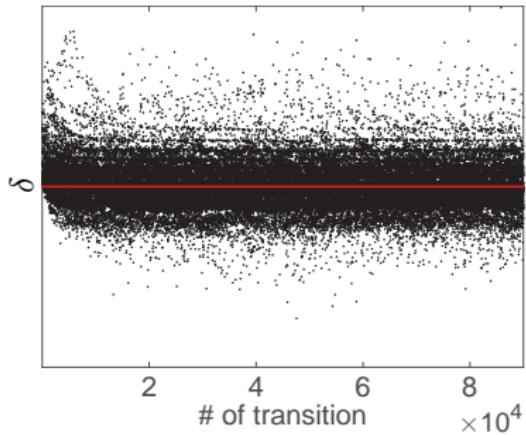
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



10^5 trajectories of length $N = 10$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

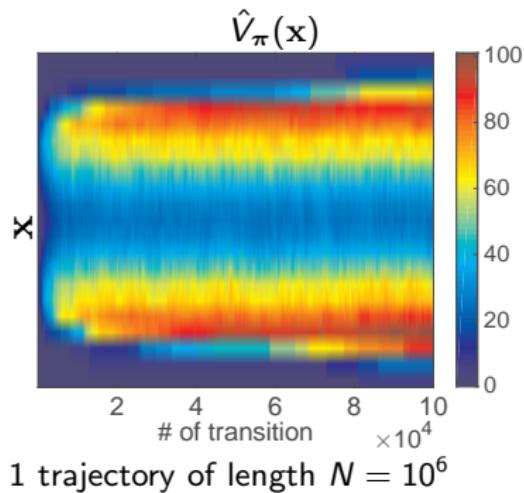
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

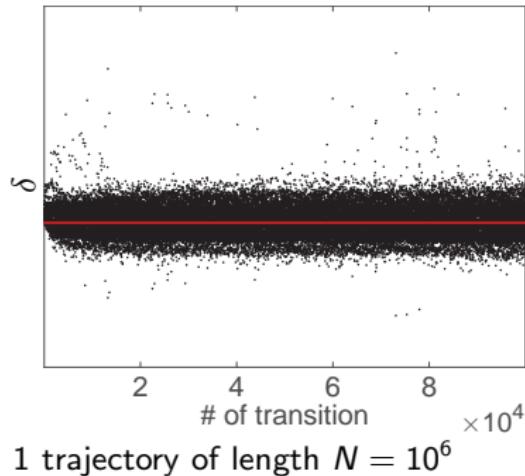
Note: **Policy Evaluation** does

$$V_\pi(\mathbf{x}) \leftarrow V_\pi(\mathbf{x}) + \bar{\delta}$$

$$\bar{\delta}(\mathbf{x}) = \mathbb{E} [\delta(\mathbf{x}, \mathbf{x}_+) \mid \mathbf{x}], \quad \forall \mathbf{x}$$

hence TD uses

- δ as a “noisy” version of $\bar{\delta}$
- $\alpha < 1$ as a “noise filter” & step size



1 trajectory of length $N = 10^6$

Learning the Value function

Value function for policy π :

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Temporal Difference (TD) uses

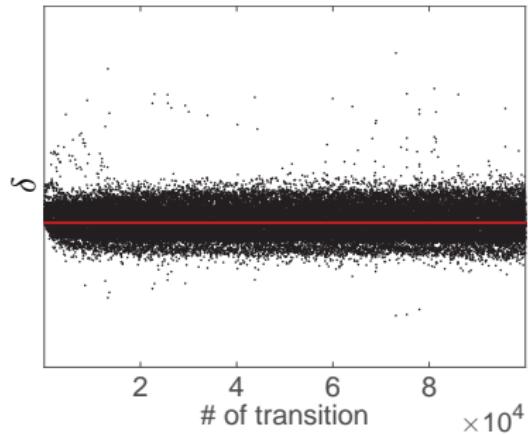
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta(\mathbf{x}_k, \mathbf{x}_{k+1})$$

$$\delta(\mathbf{x}_k, \mathbf{x}_{k+1}) = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

When $\hat{V}_\pi = V_\pi$

- $\mathbb{E}[\delta \mid \mathbf{x}] = 0$ holds
- TD-error δ has zero-mean over trajectories, i.e.

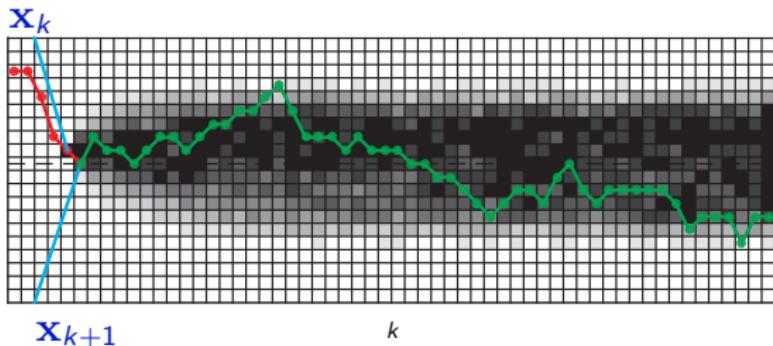
$$\mathbb{E}[\delta(\mathbf{x}_k, \mathbf{x}_{k+1})] = 0$$



Depth of Backup

Backup at time k

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta$$



Incremental Monte Carlo (IMC)

$$\delta = L(x_k, \pi(x_k)) + \gamma \sum_{i=1}^{\infty} \gamma^i L(x_{k+i}, \pi(x_{k+i})) - \hat{V}_\pi(x_k)$$

i.e. ∞ -deep backup

Temporal Difference (TD)

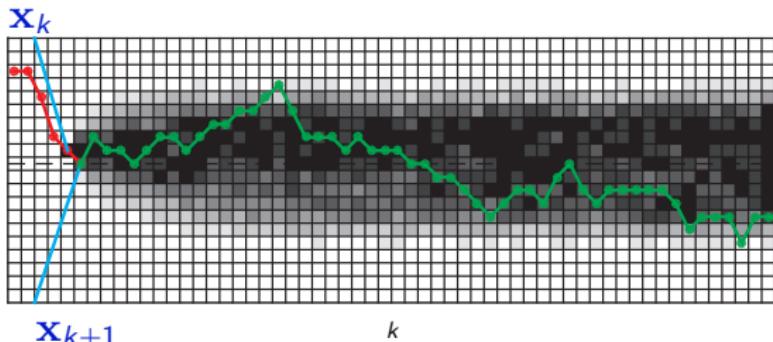
$$\delta = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

i.e. 1-step deep backup

Depth of Backup

Backup at time k

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta$$



Incremental Monte Carlo (IMC)

$$\delta = L(x_k, \pi(x_k)) + \gamma \sum_{i=1}^{\infty} \gamma^i L(x_{k+i}, \pi(x_{k+i})) - \hat{V}_\pi(x_k)$$

i.e. ∞ -deep backup

Is there something in
between 1 and ∞ ?

Temporal Difference (TD)

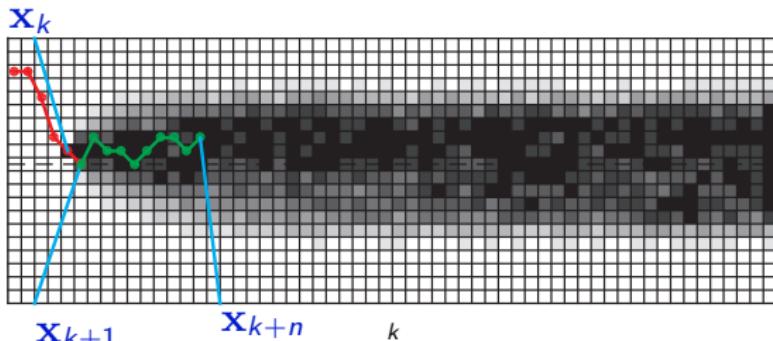
$$\delta = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

i.e. 1-step deep backup

Depth of Backup

Backup at time k

$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta$$



Incremental Monte Carlo (IMC)

$$\delta = L(x_k, \pi(x_k)) + \gamma \sum_{i=1}^{\infty} \gamma^i L(x_{k+i}, \pi(x_{k+i})) - \hat{V}_\pi(x_k)$$

i.e. ∞ -deep backup

Is there something in
between 1 and ∞ ?

Temporal Difference (TD)

$$\delta = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

i.e. 1-step deep backup

n -step Temporal Difference (TD) uses

$$\delta = L(x_k, \pi(x_k)) + \gamma^n \hat{V}_\pi(x_{k+n}) + \sum_{i=1}^{n-1} \gamma^i L(x_{k+i}, \pi(x_{k+i})) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive

TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

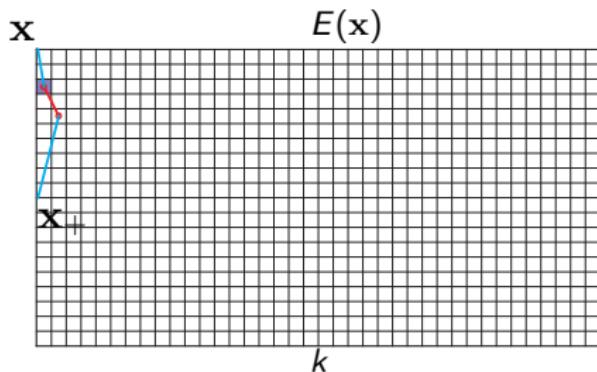
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

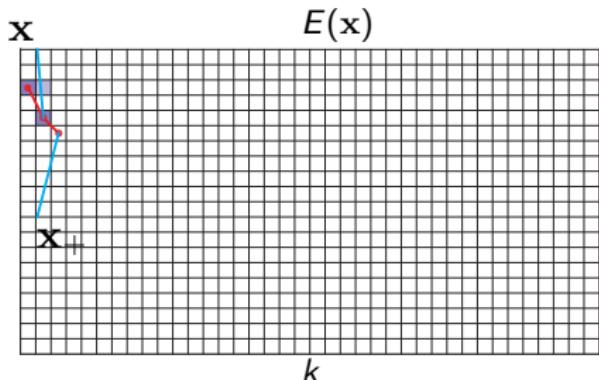
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

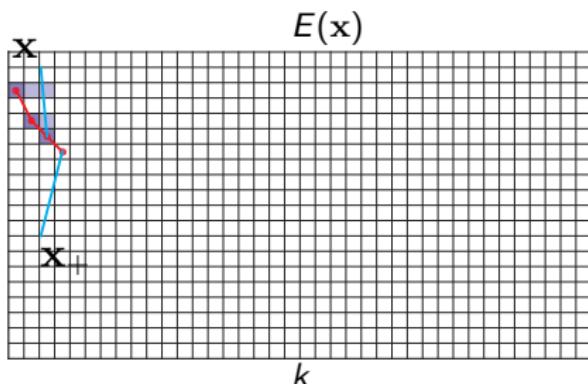
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

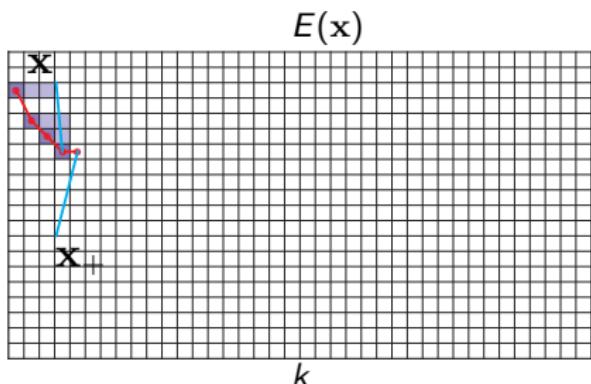
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

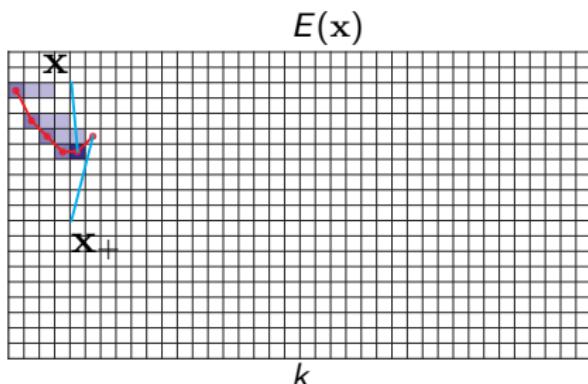
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

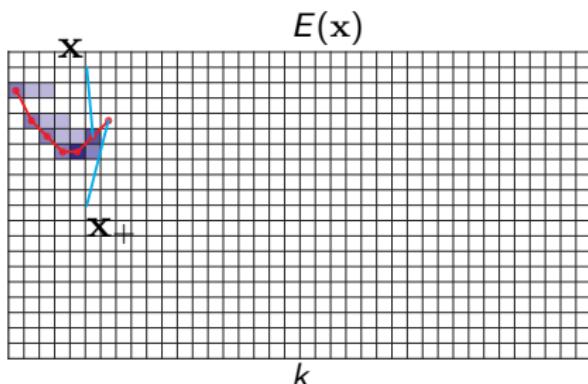
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

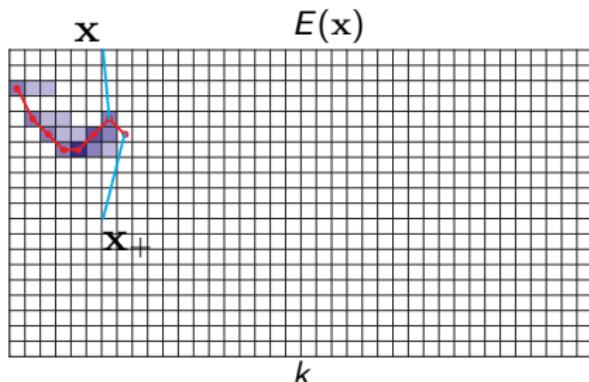
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

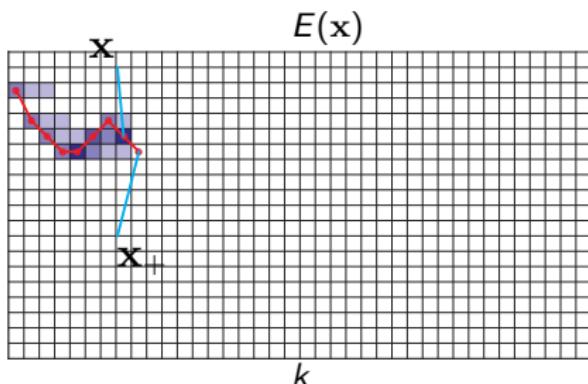
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

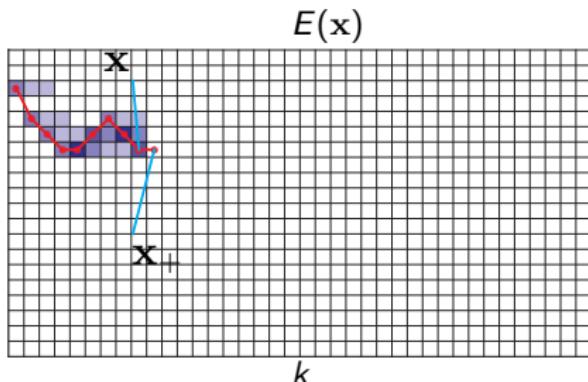
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

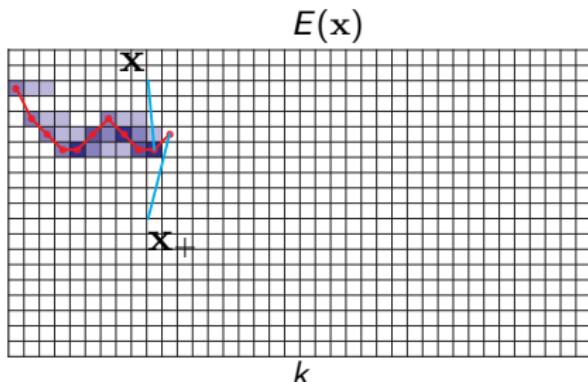
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

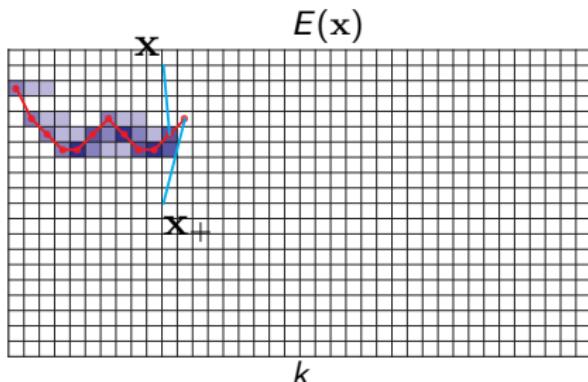
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

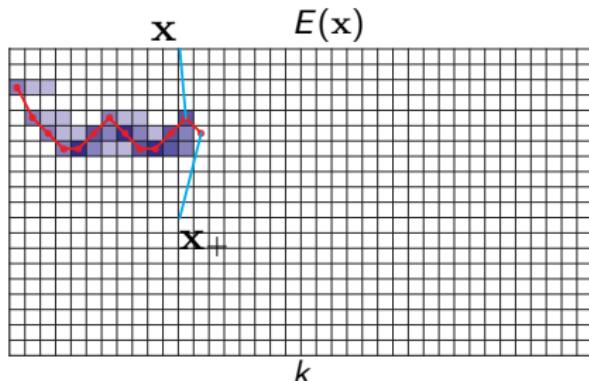
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

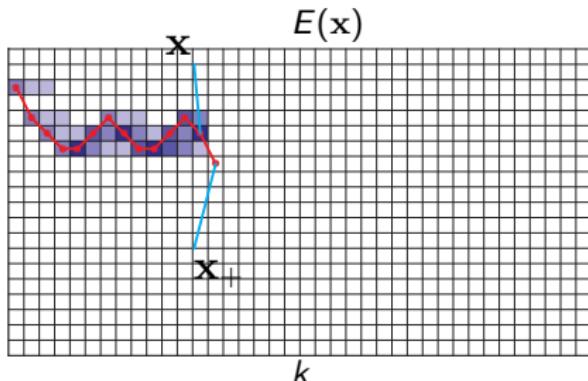
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

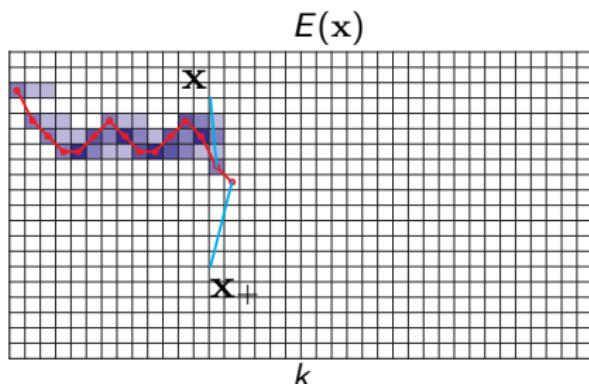
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

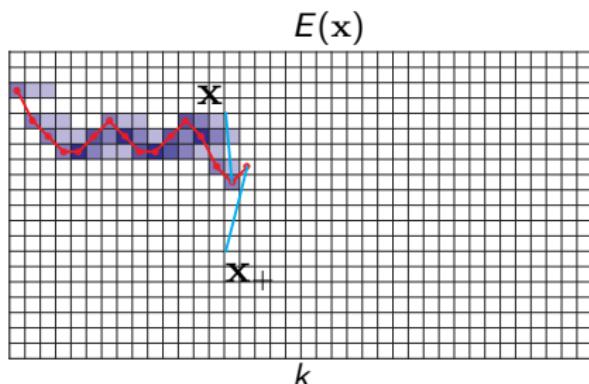
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

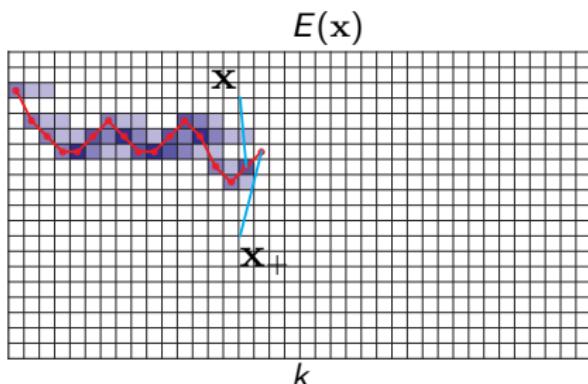
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

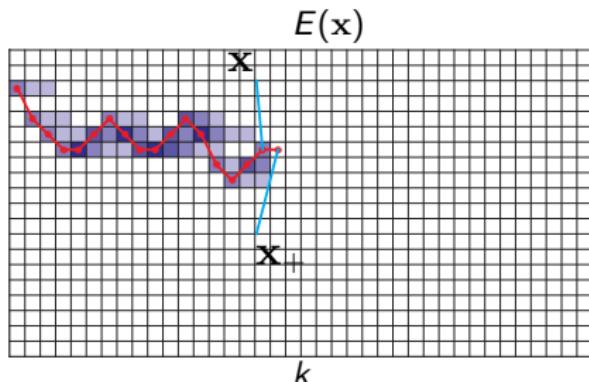
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

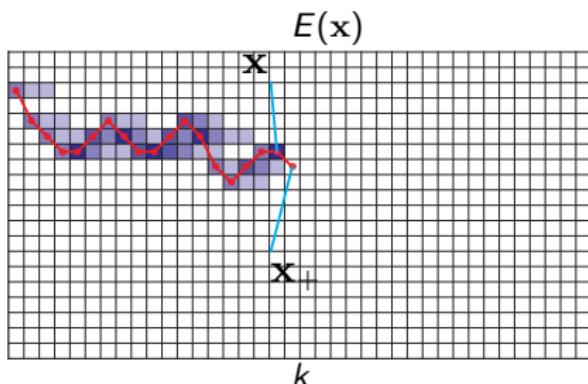
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

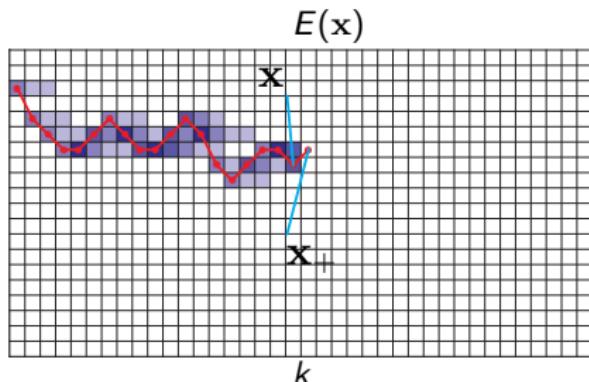
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

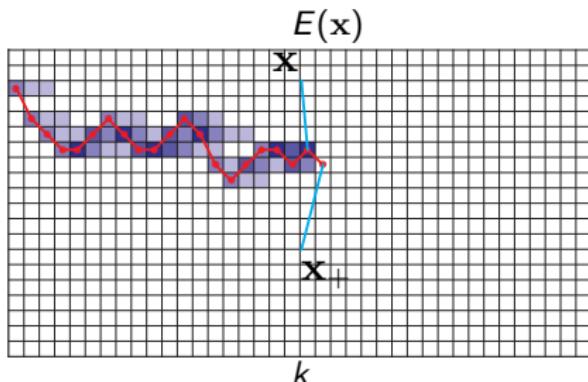
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

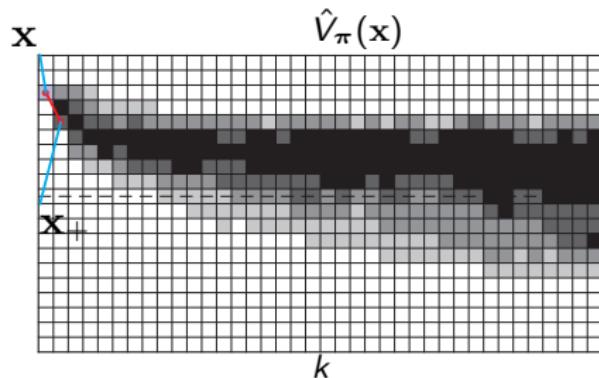
$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD



- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive

TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

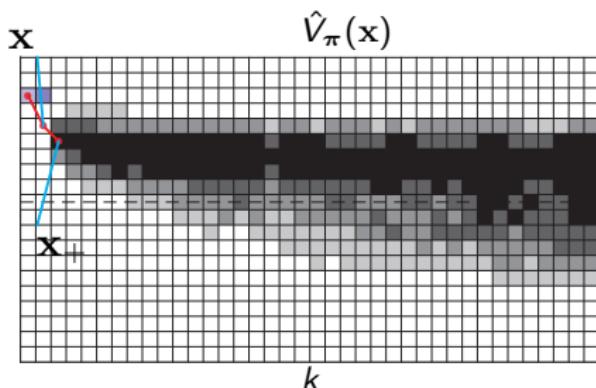
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

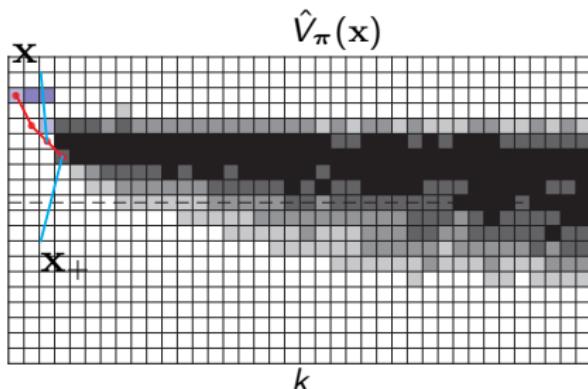
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

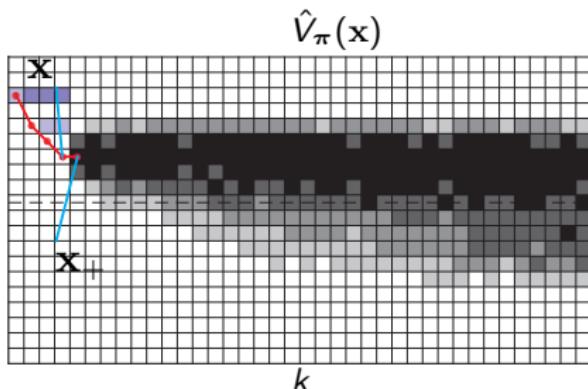
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

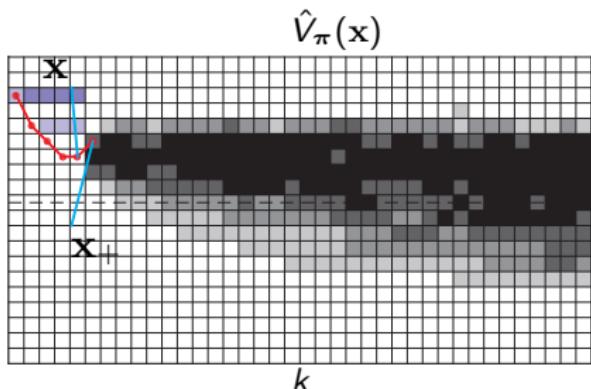
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

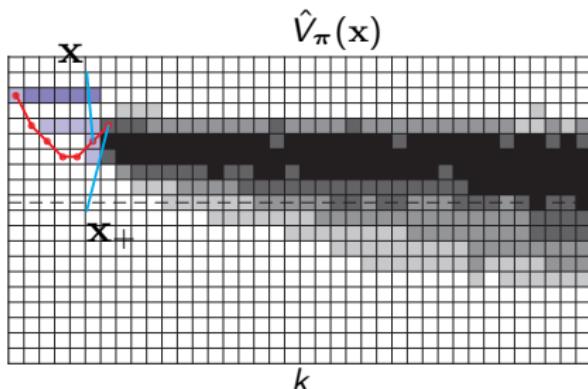
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

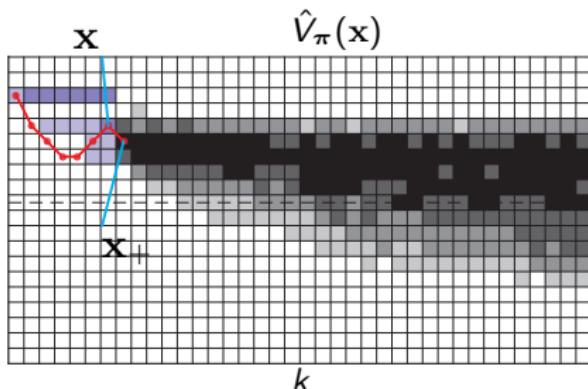
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

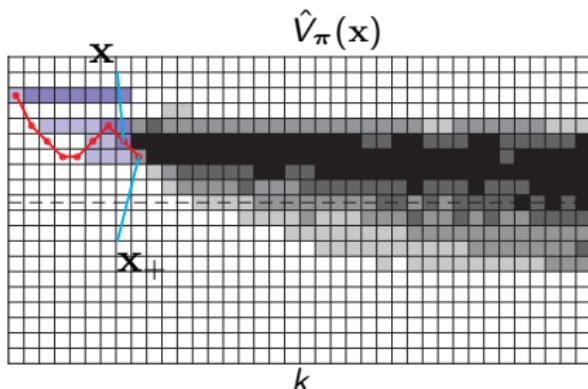
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

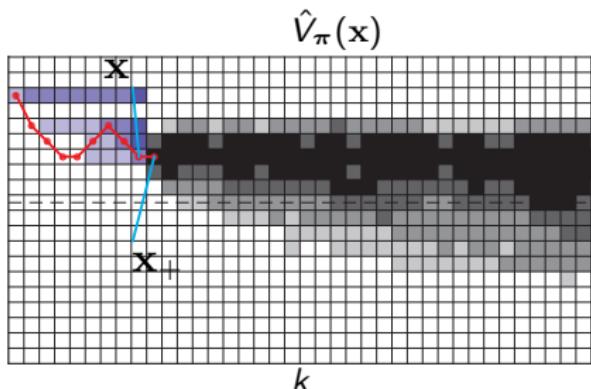
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

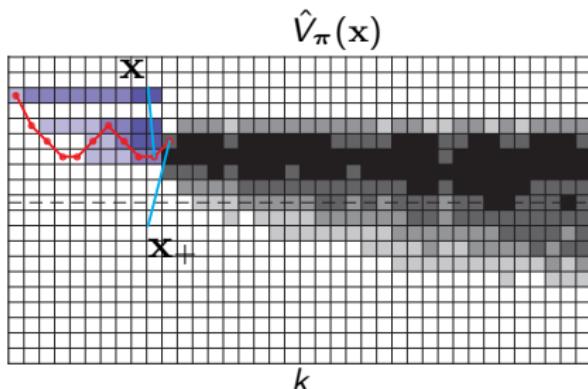
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

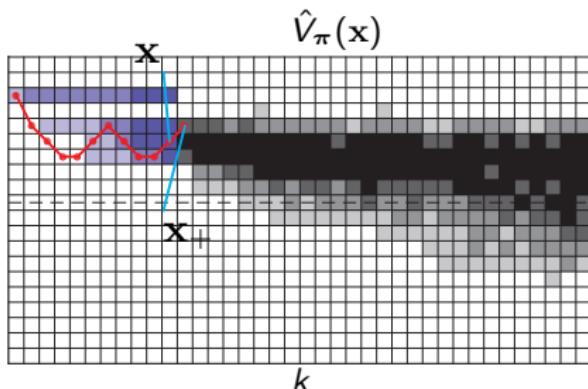
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

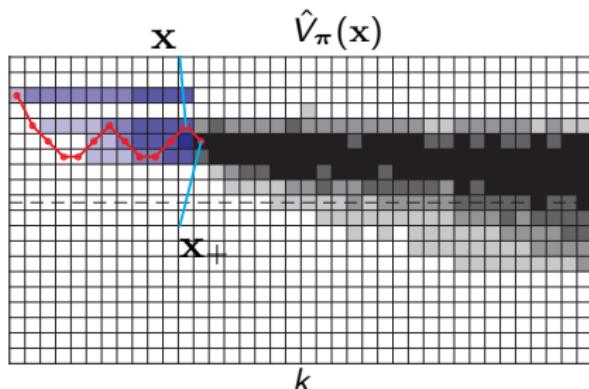
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

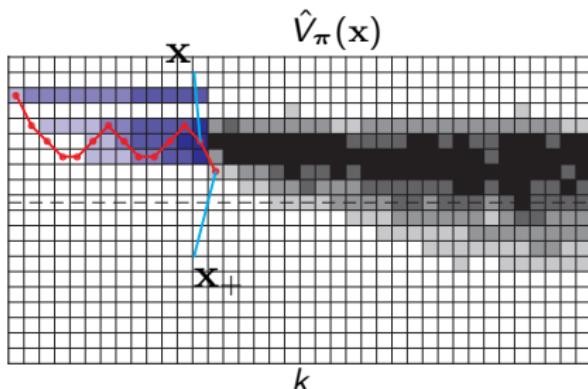
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

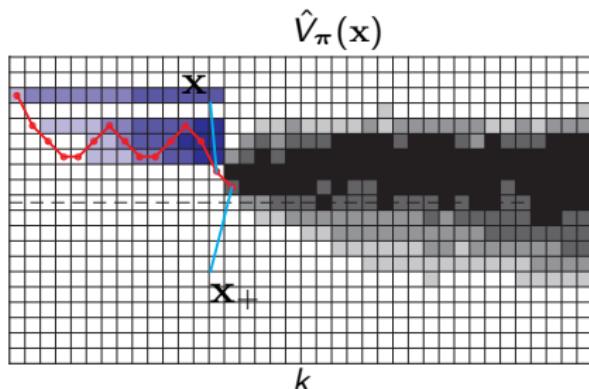
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

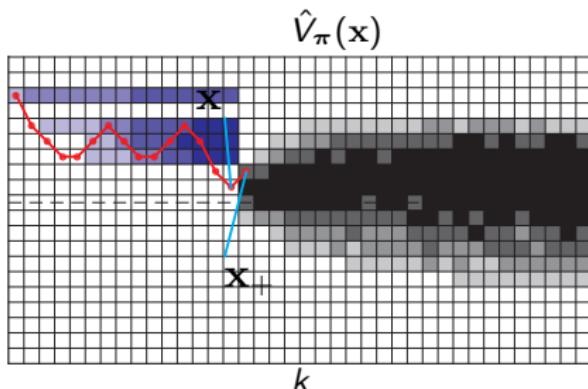
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

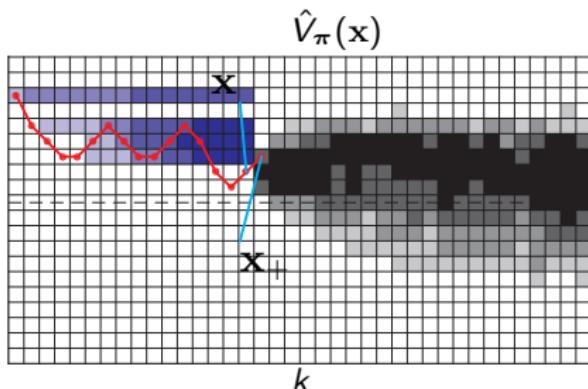
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

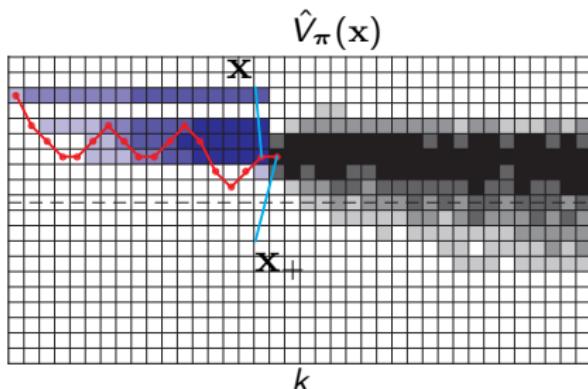
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = \mathbf{x}_k \\ 0 & \text{if } \textcolor{blue}{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

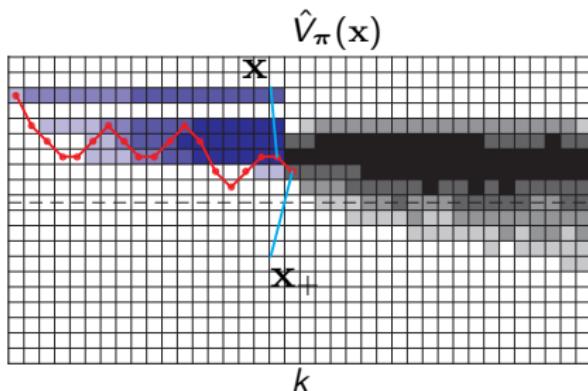
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(\mathbf{x}_k, \boldsymbol{\pi}(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



$\text{TD}(\lambda)$: trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

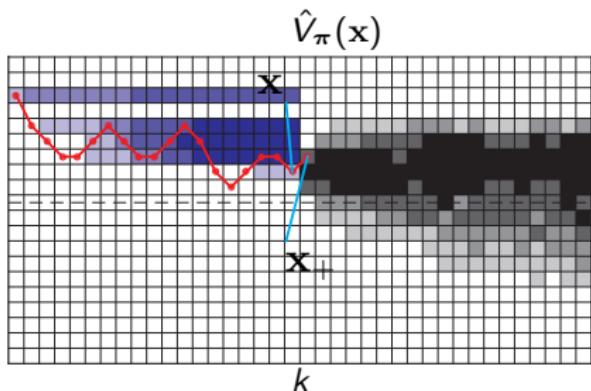
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

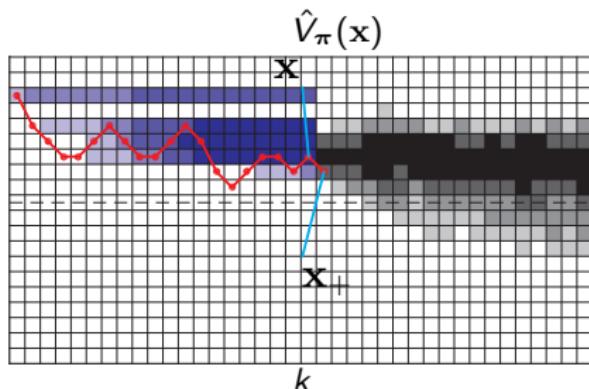
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow eligibility trace:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

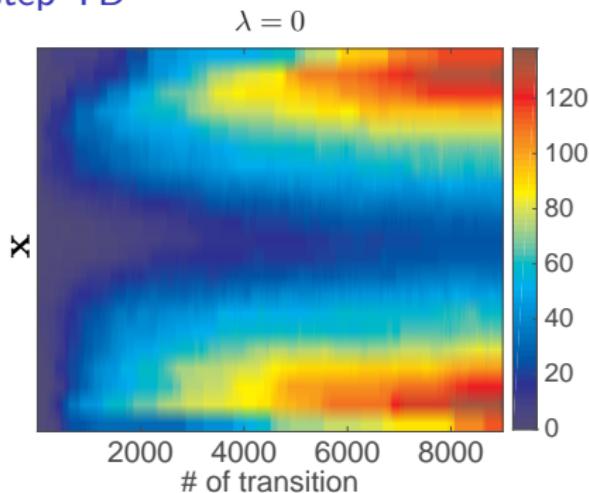
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

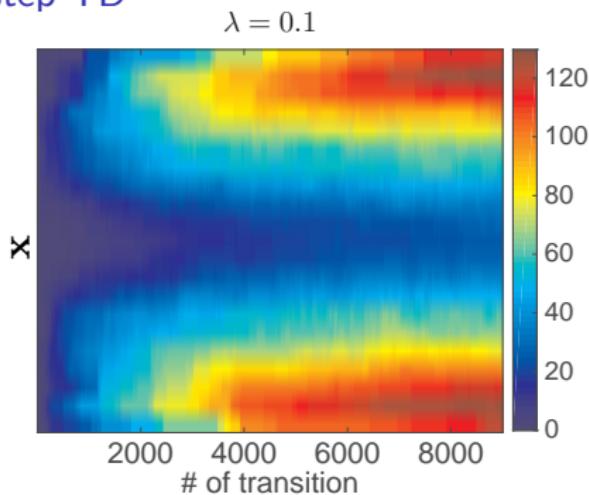
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

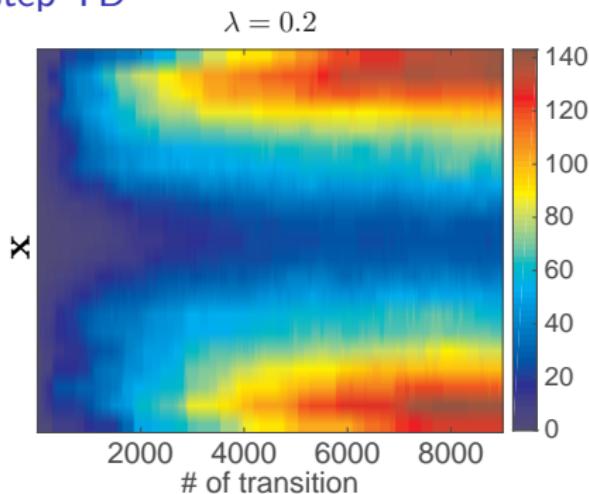
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\textcolor{blue}{x}) \leftarrow \gamma \lambda E(\textcolor{blue}{x}) + \begin{cases} 1 & \text{if } \textcolor{blue}{x} = x_k \\ 0 & \text{if } \textcolor{blue}{x} \neq x_k \end{cases}, \quad \forall \textcolor{blue}{x}$$

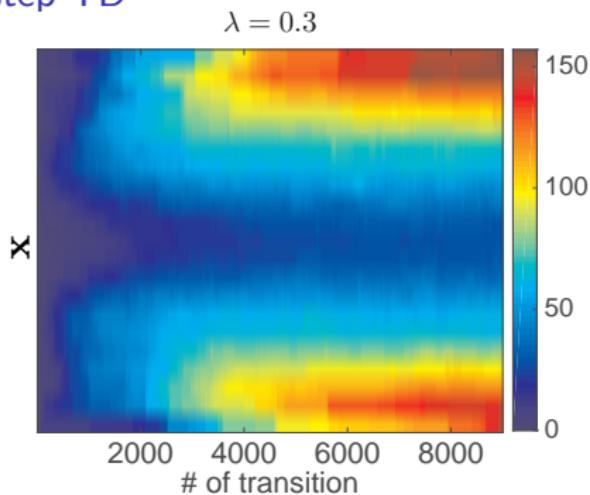
Value function is updated according to

$$\hat{V}_\pi(\textcolor{blue}{x}) \leftarrow \hat{V}_\pi(\textcolor{blue}{x}) + \alpha E(\textcolor{blue}{x}) \delta_k, \quad \forall \textcolor{blue}{x}$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

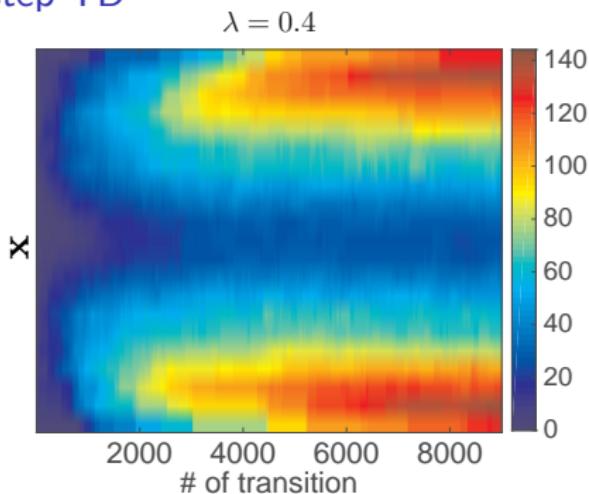
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

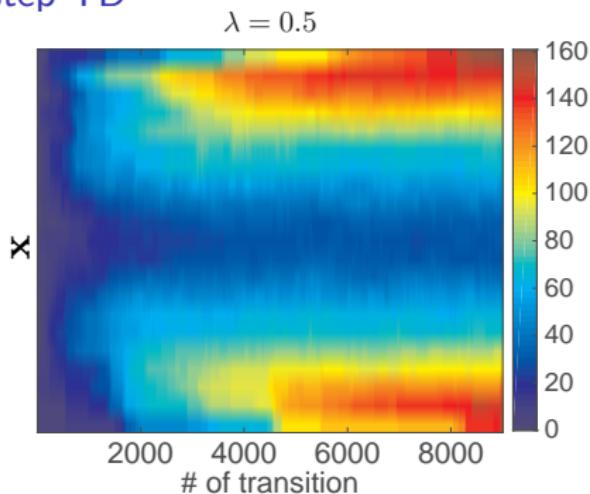
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

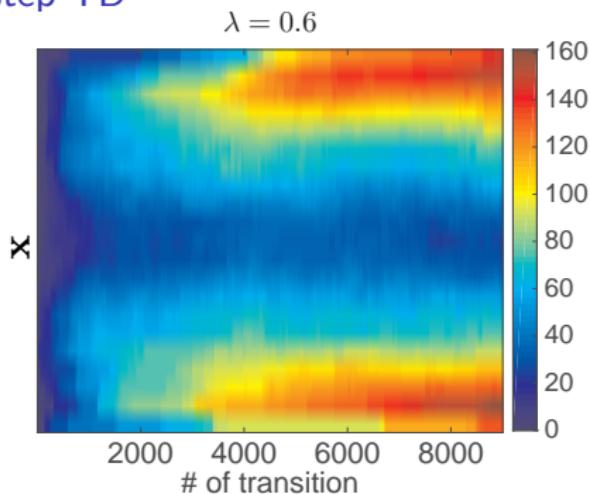
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

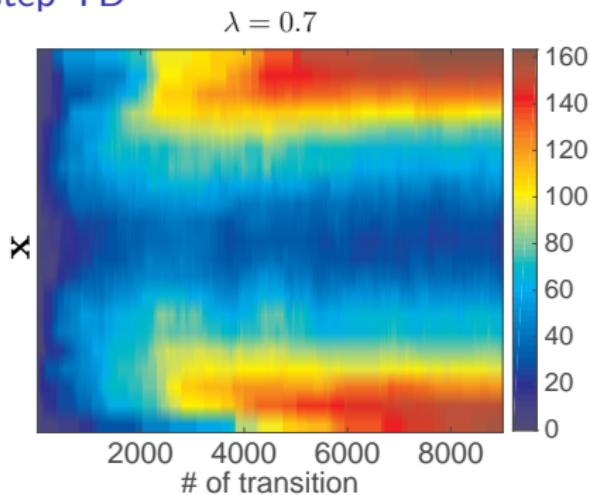
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

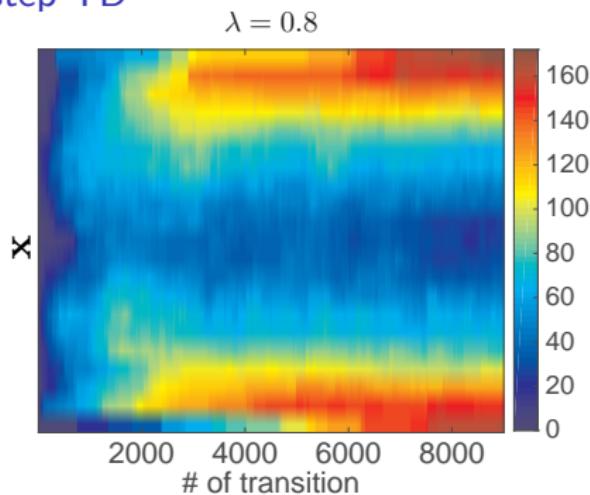
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

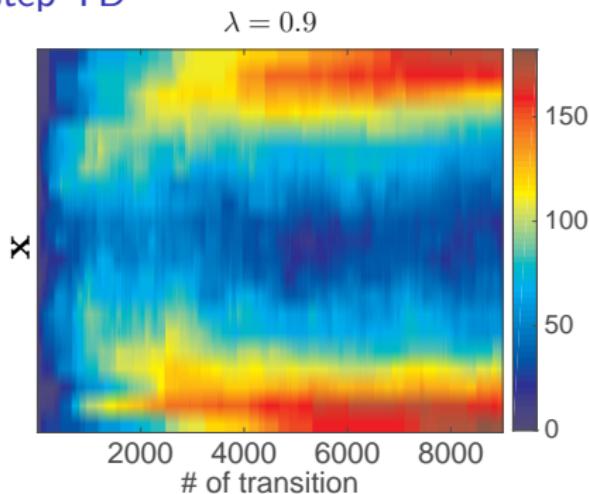
Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

TD(λ) method - Implementing n -step TD

- MC or n -steps TD is hard to implement
- Requires keeping traces of the state trajectories, memory intensive



TD(λ): trace of state visits \rightarrow **eligibility trace:**

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}, \quad \forall \mathbf{x}$$

Value function is updated according to

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Outline

1 Learning V_π

2 Learning Q_π

3 Learning Q_*

Learning the action-value function - Off policy

Action-value function for policy π :

$$Q_{\pi}(\mathbf{x}, \mathbf{u}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Learning the action-value function - Off policy

Action-value function for policy π :

$$Q_\pi(\mathbf{x}, \mathbf{u}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Policy Evaluation does

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\hat{Q}_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) \mid \mathbf{x}, \mathbf{u} \right]$$

Learning the action-value function - Off policy

Action-value function for policy π :

$$Q_\pi(\mathbf{x}, \mathbf{u}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Policy Evaluation does

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\hat{Q}_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) \mid \mathbf{x}, \mathbf{u} \right]$$

Temporal Difference (TD)

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) \leftarrow (1 - \alpha) \hat{Q}_\pi(\mathbf{x}, \mathbf{u}) + \alpha \left(L(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) \right)$$

where \mathbf{u} can be arbitrary.

Learning the action-value function - Off policy

Action-value function for policy π :

$$Q_\pi(\mathbf{x}, \mathbf{u}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Policy Evaluation does

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\hat{Q}_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) \mid \mathbf{x}, \mathbf{u} \right]$$

Temporal Difference (TD)

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) \leftarrow (1 - \alpha) \hat{Q}_\pi(\mathbf{x}, \mathbf{u}) + \alpha \left(L(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) \right)$$

where \mathbf{u} can be arbitrary. Hence TD for Q-learning does:

$$\begin{aligned} \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \underbrace{\hat{Q}_\pi(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1}))}_{:= \hat{V}_\pi(\mathbf{x}_{k+1})} - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \end{aligned}$$

This is labelled an off-policy TD evaluation

Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

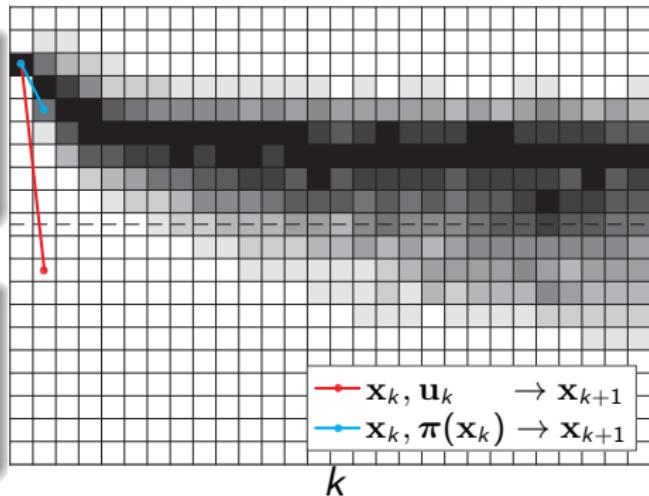
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

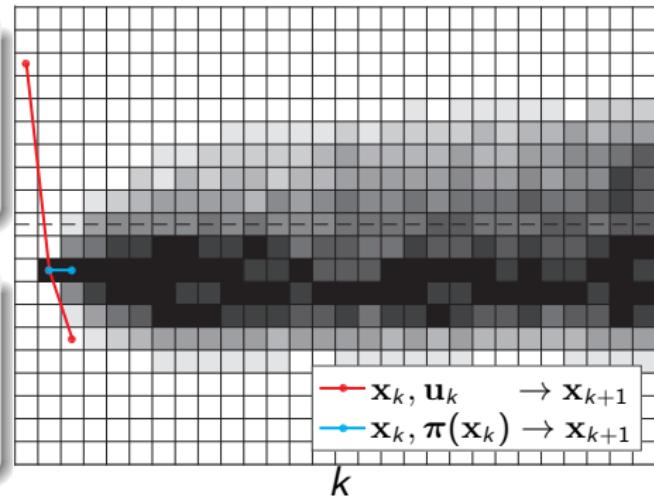
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

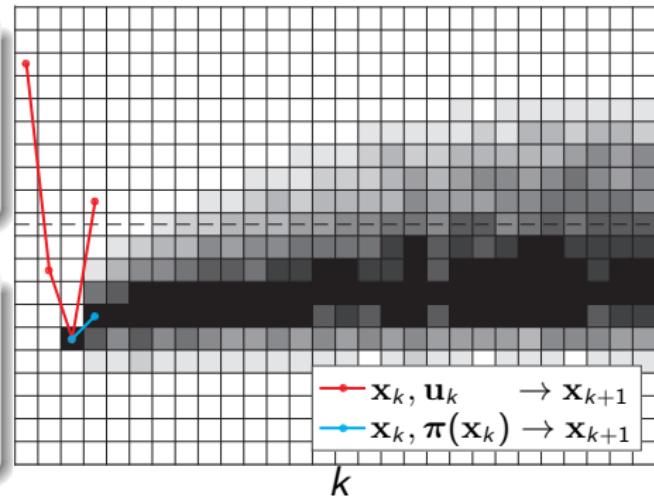
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta_k$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(x_k, u_k) \leftarrow \hat{Q}_\pi(x_k, u_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(x, u_k) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{Q}_\pi(x_k, u_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

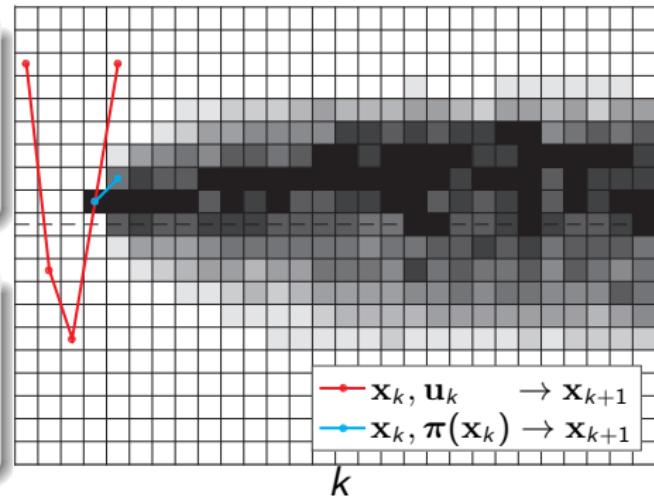
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

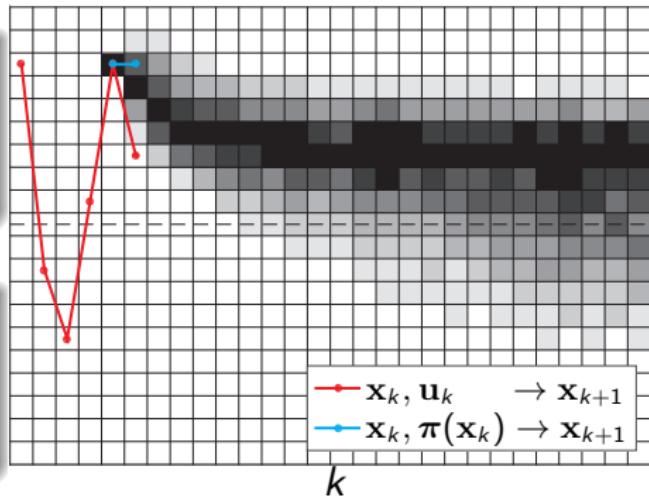
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

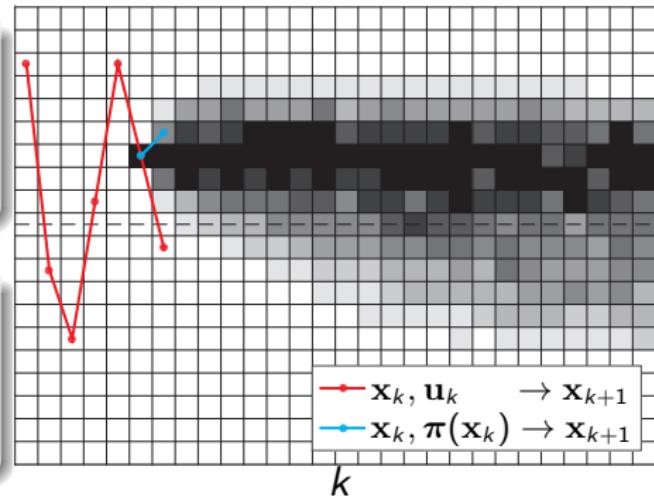
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

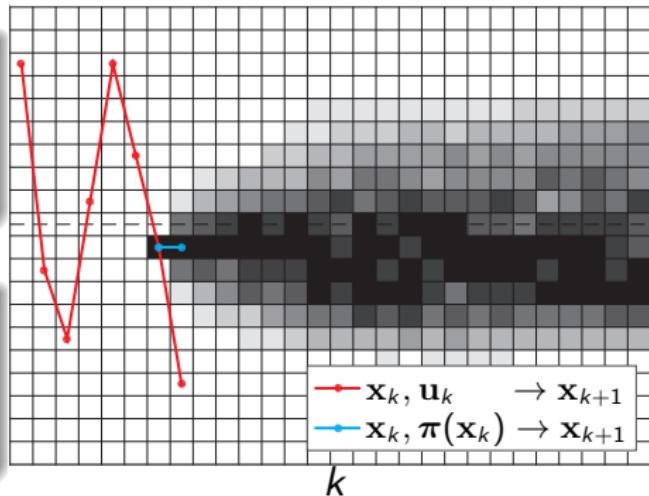
$$\hat{V}_\pi(x_k) \leftarrow \hat{V}_\pi(x_k) + \alpha \delta_k$$

$$\delta_k = L(x_k, \pi(x_k)) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{V}_\pi(x_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(x_k, u_k) \leftarrow \hat{Q}_\pi(x_k, u_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(x, u_k) + \gamma \hat{V}_\pi(x_{k+1}) - \hat{Q}_\pi(x_k, u_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

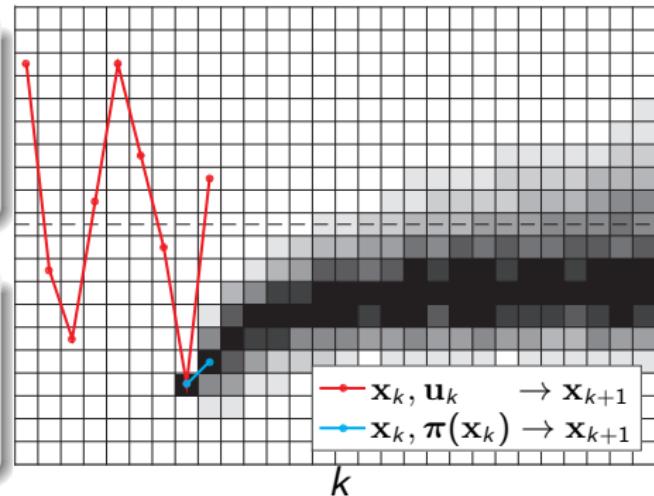
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

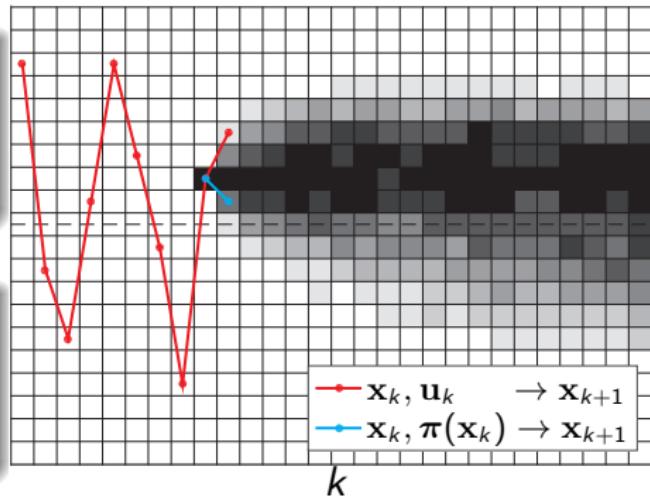
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

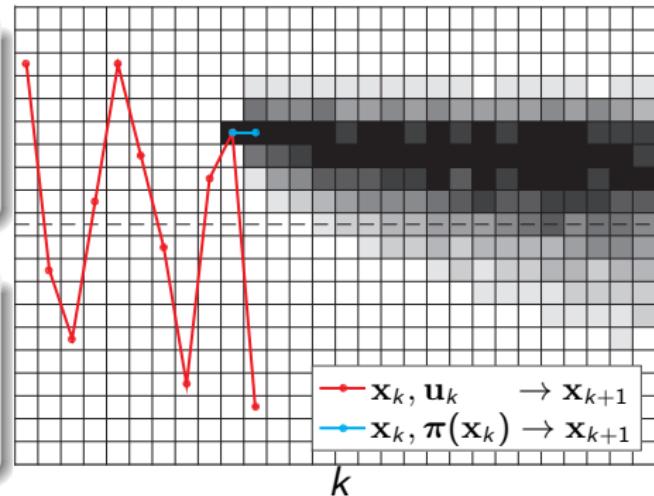
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

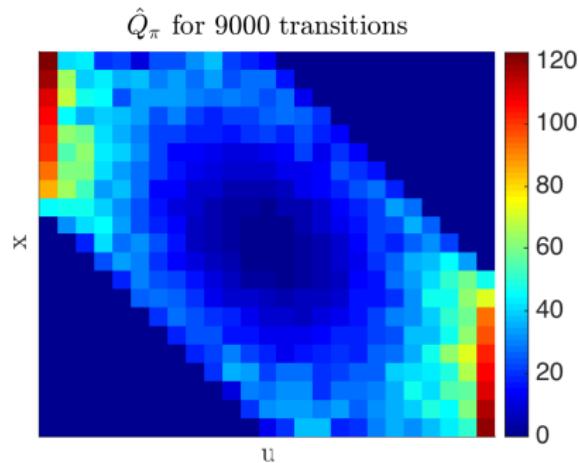
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

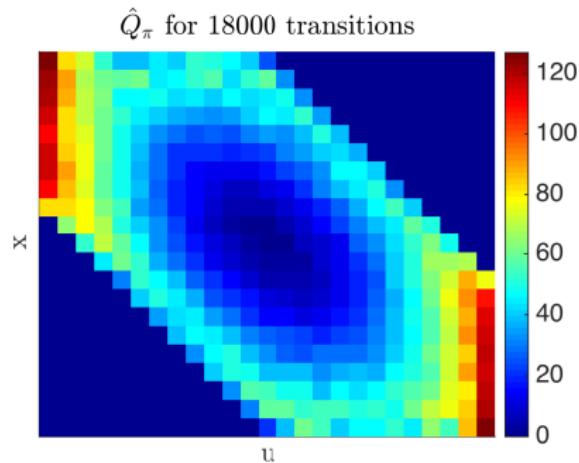
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

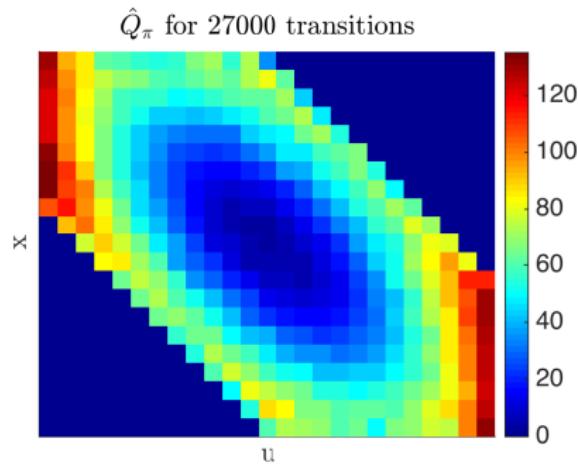
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

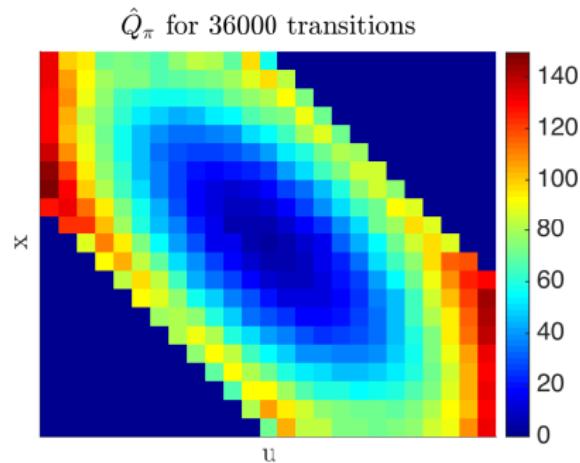
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

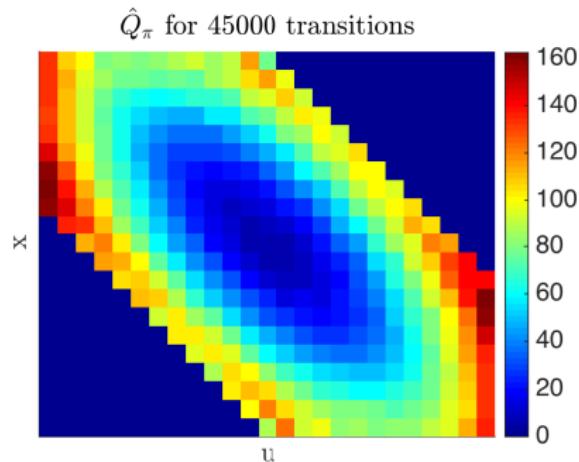
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

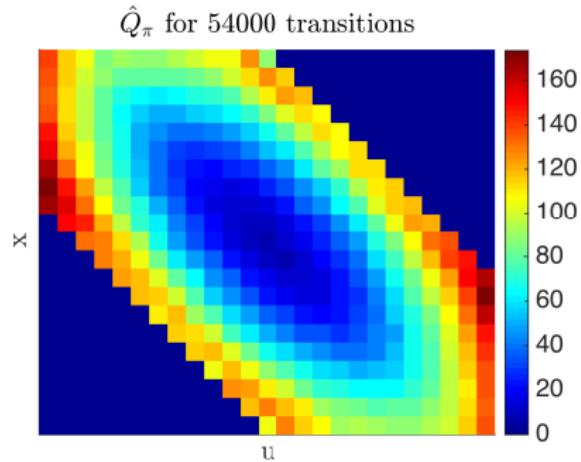
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

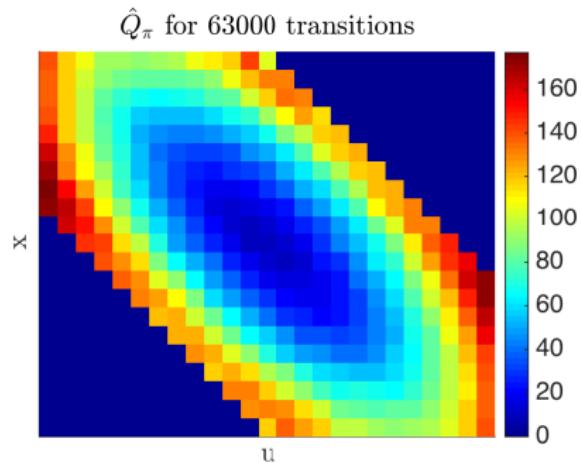
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

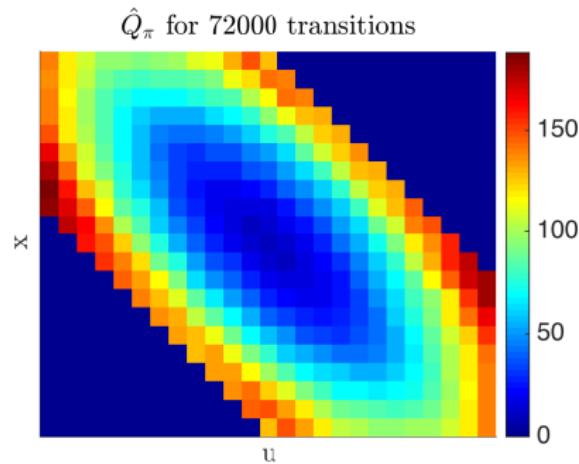
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

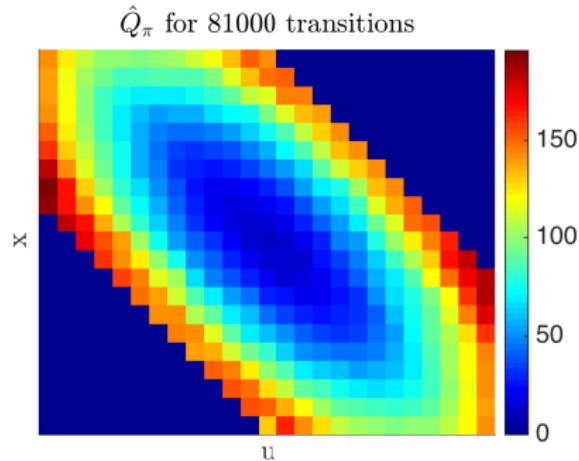
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

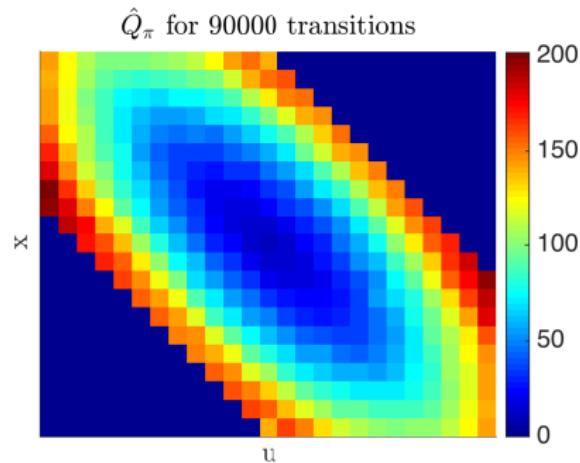
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

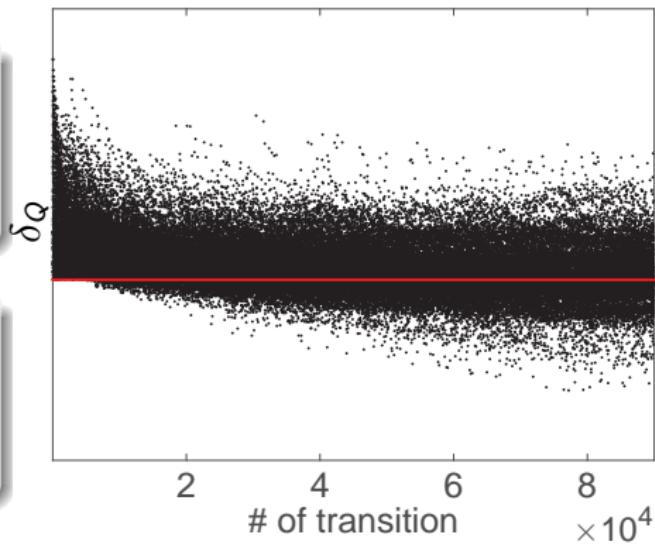
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Learning the action-value function - Off policy

Temporal Difference (TD) for \hat{V}_π

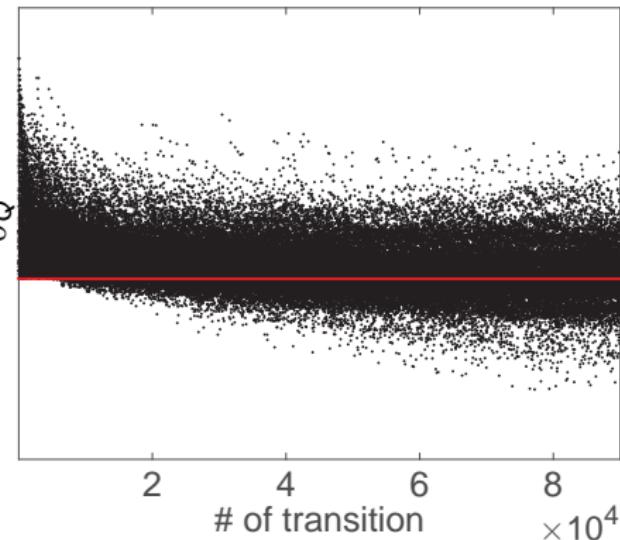
$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{V}_\pi(\mathbf{x}_k)$$

Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}, \mathbf{u}_k) + \gamma \hat{V}_\pi(\mathbf{x}_{k+1}) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)$$



Can we pick $\mathbf{u}_{0, \dots, \infty}$ arbitrarily? Not really...

- must visit all state-input \mathbf{x}, \mathbf{u} pairs infinitely many times to build

$$\hat{Q}_\pi(\mathbf{x}, \mathbf{u}) = Q_\pi(\mathbf{x}, \mathbf{u})$$

- “partial picture” is built from visiting only subparts of the state-input space

Outline

- 1 Learning V_π
- 2 Learning Q_π
- 3 Learning Q_*



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Dynamic Programming (DP)

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \leftarrow \mathbb{E} \left[L(\mathbf{x}, \mathbf{u}) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_+, \mathbf{u}') \mid \mathbf{x}, \mathbf{u} \right]$$

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Dynamic Programming (DP)

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \leftarrow \mathbb{E} \left[L(\mathbf{x}, \mathbf{u}) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_+, \mathbf{u}') \mid \mathbf{x}, \mathbf{u} \right]$$

Temporal Difference (TD)

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \leftarrow (1 - \alpha) \hat{Q}_*(\mathbf{x}, \mathbf{u}) + \alpha \left(L(\mathbf{x}, \mathbf{u}) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_+, \mathbf{u}') \right)$$

where \mathbf{u} can be arbitrary.

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Dynamic Programming (DP)

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \leftarrow \mathbb{E} \left[L(\mathbf{x}, \mathbf{u}) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_+, \mathbf{u}') \mid \mathbf{x}, \mathbf{u} \right]$$

Temporal Difference (TD)

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \leftarrow (1 - \alpha) \hat{Q}_*(\mathbf{x}, \mathbf{u}) + \alpha \left(L(\mathbf{x}, \mathbf{u}) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_+, \mathbf{u}') \right)$$

where \mathbf{u} can be arbitrary.

also reads as Off-policy TD control:

$$\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k$$

$$\delta_k = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_{k+1}, \mathbf{u}') - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)$$

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD evaluation

$$\begin{aligned}\hat{Q}_{\pi}(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_{\pi}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{Q}_{\pi}(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - \hat{Q}_{\pi}(\mathbf{x}_k, \mathbf{u}_k)\end{aligned}$$

TD control

$$\begin{aligned}\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \underbrace{\min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}_{k+1}, \mathbf{u}') - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)}_{:= \hat{V}_{\star}(\mathbf{x}_{k+1})}\end{aligned}$$

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD evaluation

$$\begin{aligned}\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{Q}_\pi(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)\end{aligned}$$

TD control

$$\begin{aligned}\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \underbrace{\min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_{k+1}, \mathbf{u}') - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)}_{:= \hat{V}_*(\mathbf{x}_{k+1})}\end{aligned}$$

- Selection of $\mathbf{u}_0, \dots, \infty$ must yield ∞ -many visits of all state-input pairs to get

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \rightarrow Q_*(\mathbf{x}, \mathbf{u})$$

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD evaluation

$$\begin{aligned}\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{Q}_\pi(\mathbf{x}_{k+1}, \pi(\mathbf{x}_{k+1})) - \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k)\end{aligned}$$

TD control

$$\begin{aligned}\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) &\leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q \\ \delta_k^Q &= L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \underbrace{\min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}_{k+1}, \mathbf{u}') - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)}_{:= \hat{V}_*(\mathbf{x}_{k+1})}\end{aligned}$$

- Selection of $\mathbf{u}_{0,\dots,\infty}$ must yield ∞ -many visits of all state-input pairs to get

$$\hat{Q}_*(\mathbf{x}, \mathbf{u}) \rightarrow Q_*(\mathbf{x}, \mathbf{u})$$

- Optimal policy is $\hat{\pi}_*(\mathbf{x}) = \arg \min_{\mathbf{u}} \hat{Q}_*(\mathbf{x}, \mathbf{u})$

Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

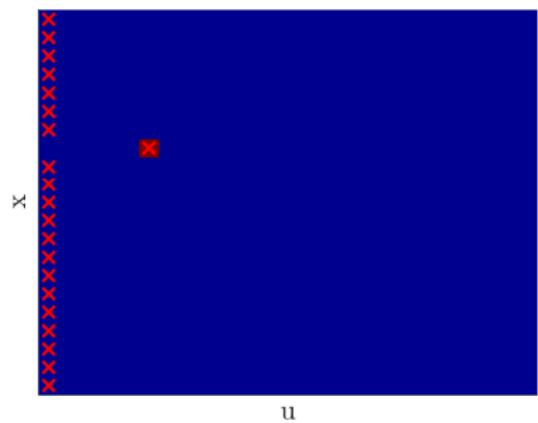
where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

\hat{Q}_* for 1 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

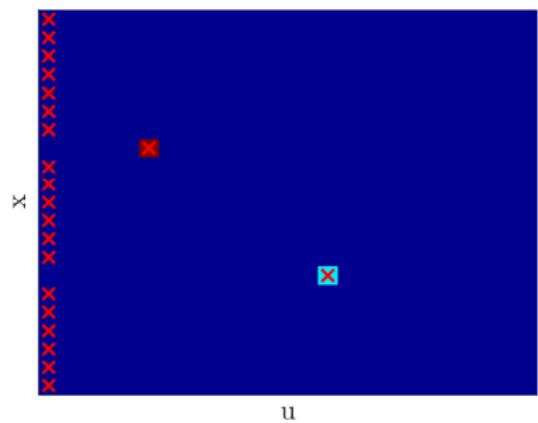
where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

\hat{Q}_{\star} for 2 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

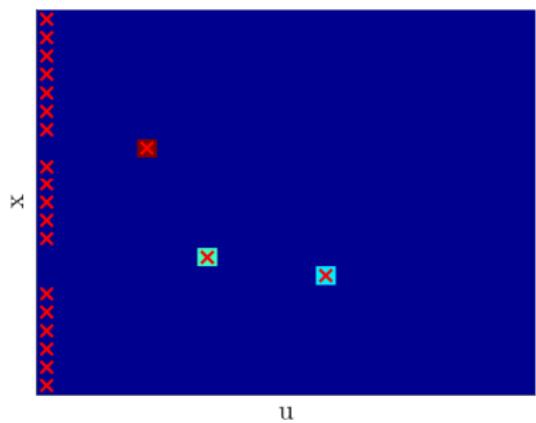
where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

\hat{Q}_* for 3 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

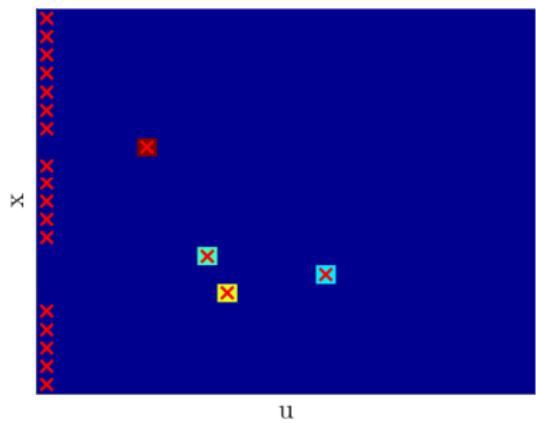
where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

\hat{Q}_* for 4 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

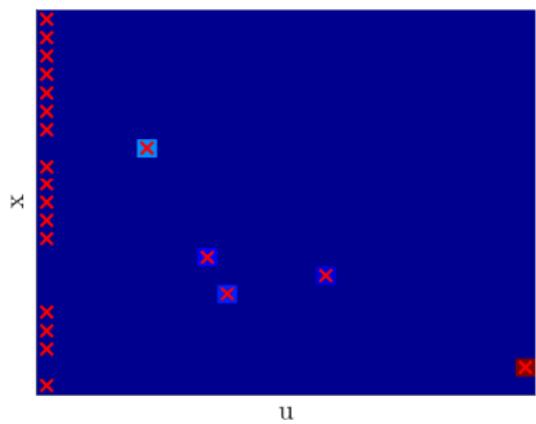
where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

\hat{Q}_* for 5 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

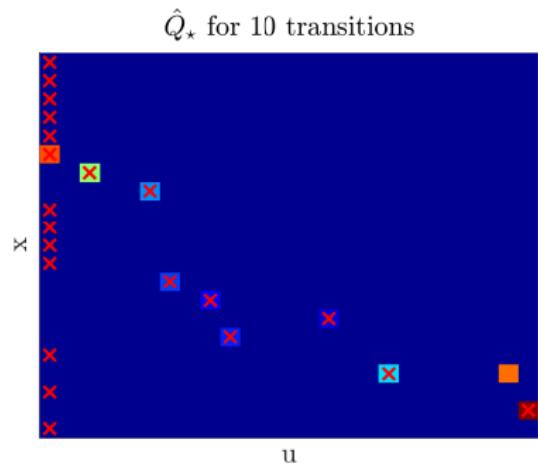
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_*(\mathbf{x}_{k+1}) - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_*(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

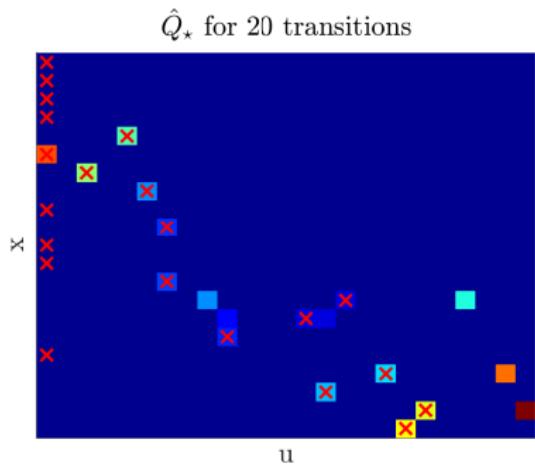
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

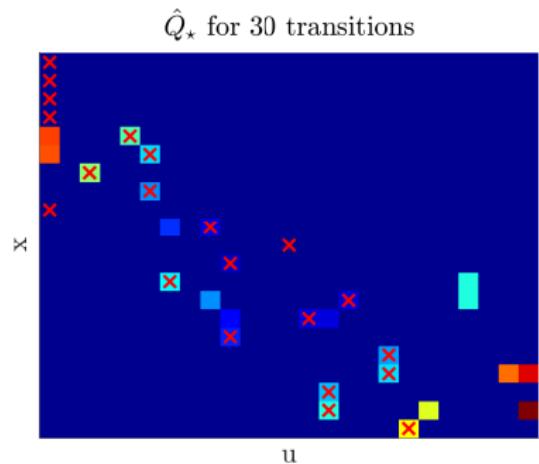
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

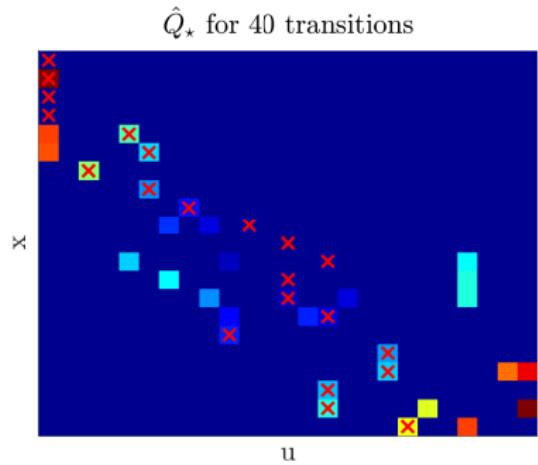
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

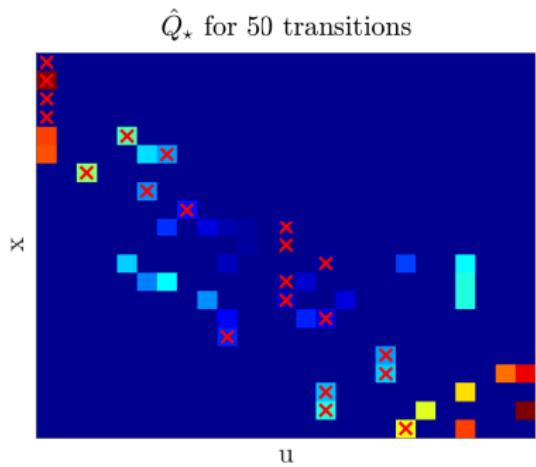
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

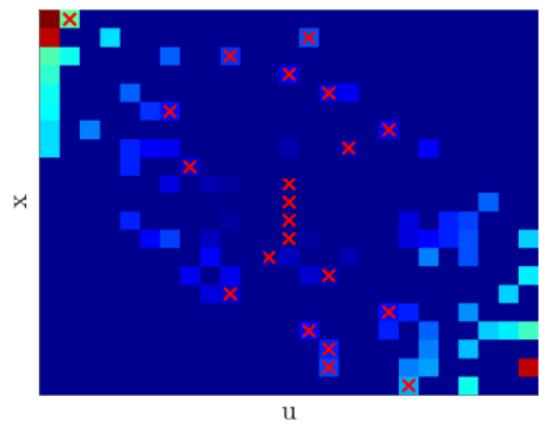
where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

\hat{Q}_* for 100 transitions



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

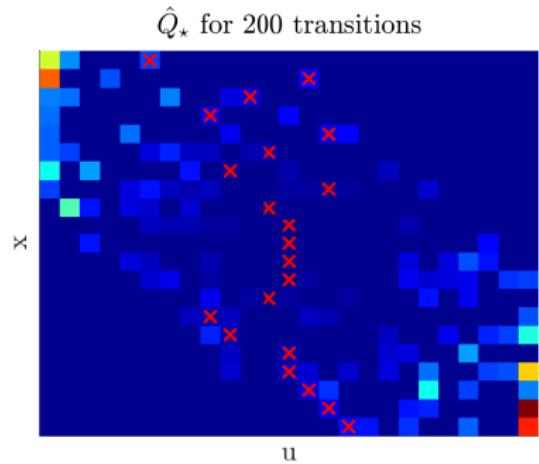
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

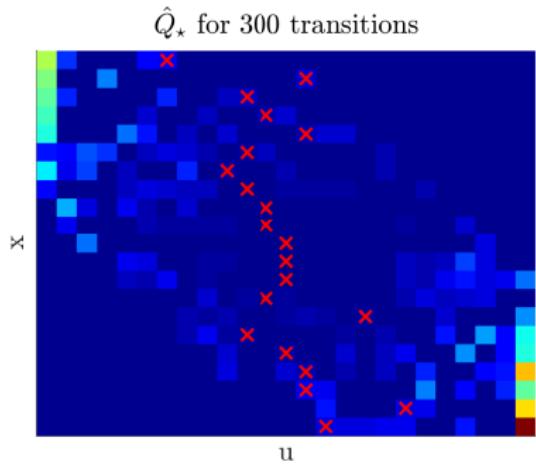
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

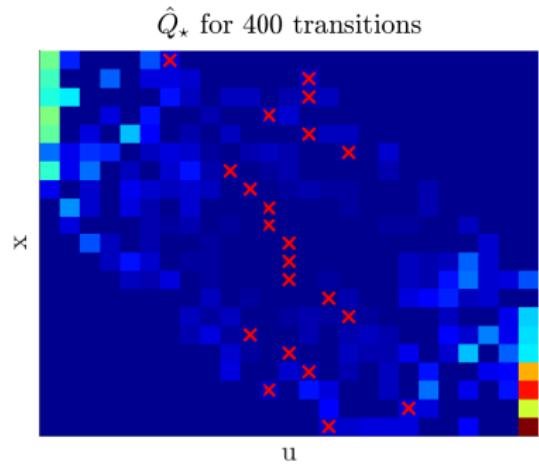
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_*(\mathbf{x}_{k+1}) - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_*(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

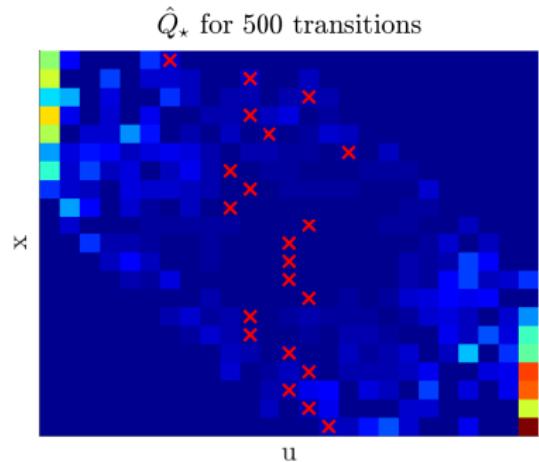
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

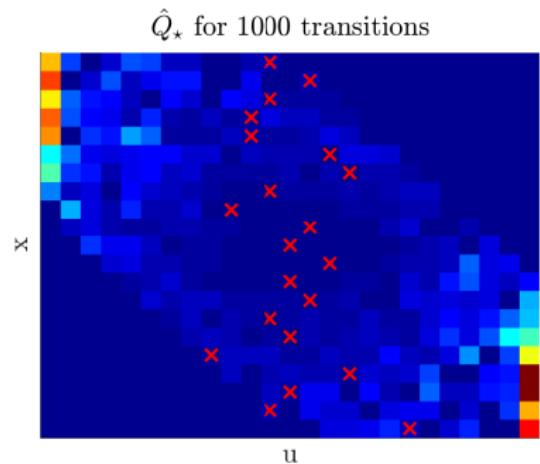
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

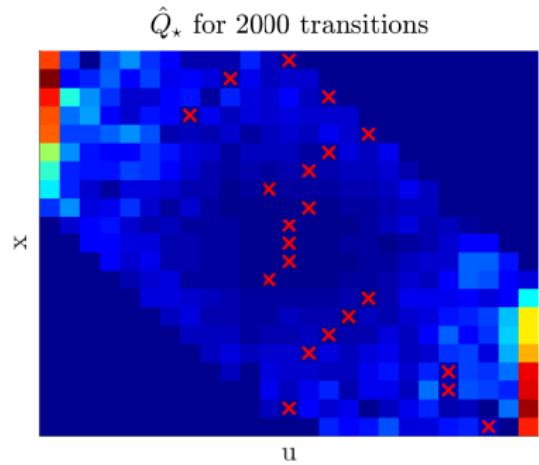
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

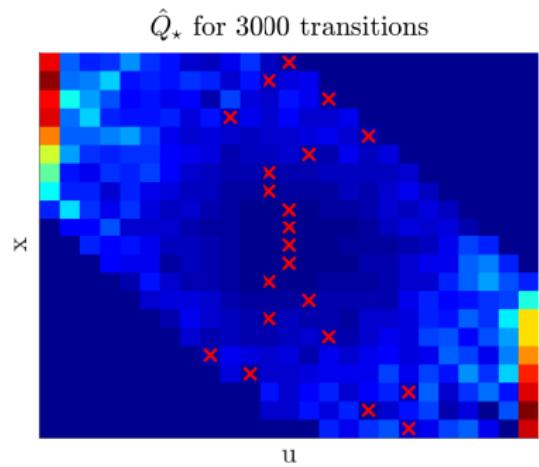
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_*(\mathbf{x}_{k+1}) - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_*(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

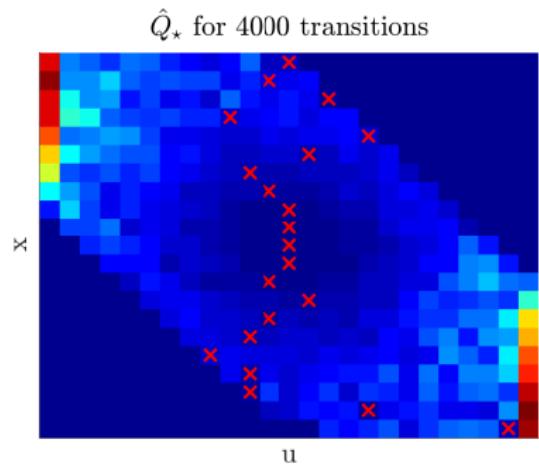
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

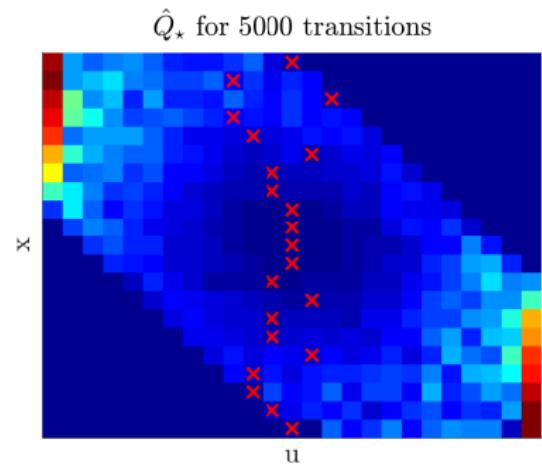
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_*(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

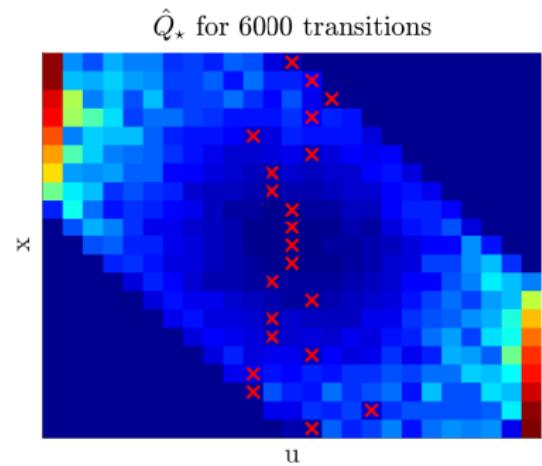
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_*(\mathbf{x}_{k+1}) - \hat{Q}_*(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_*(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_*(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

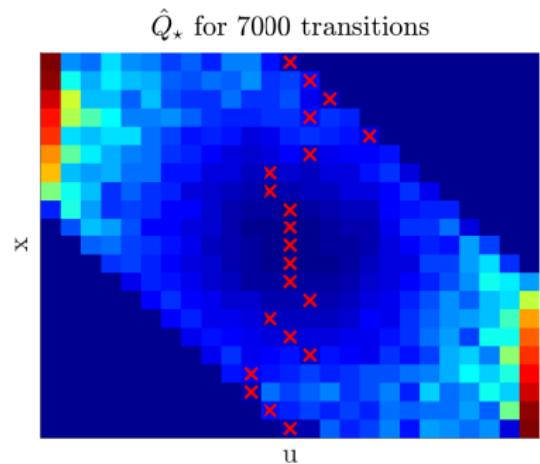
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

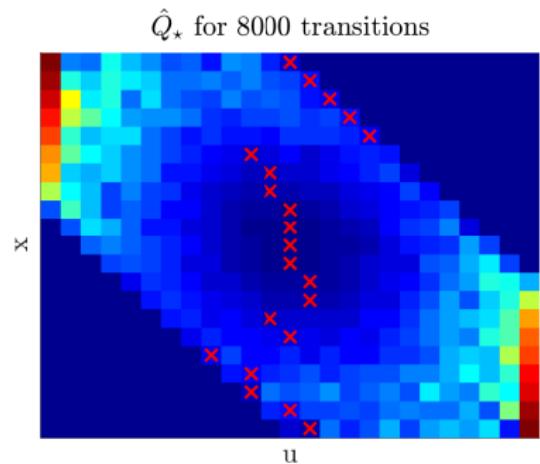
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

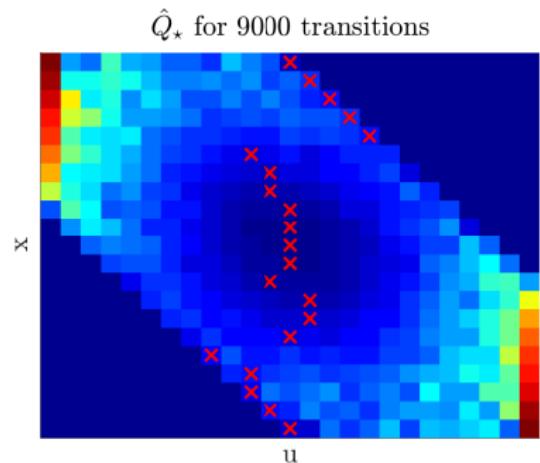
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

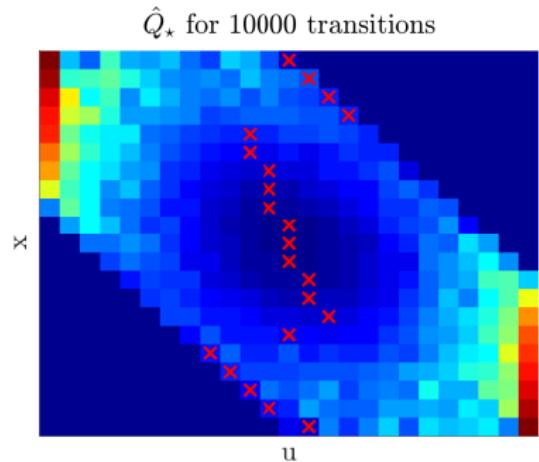
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

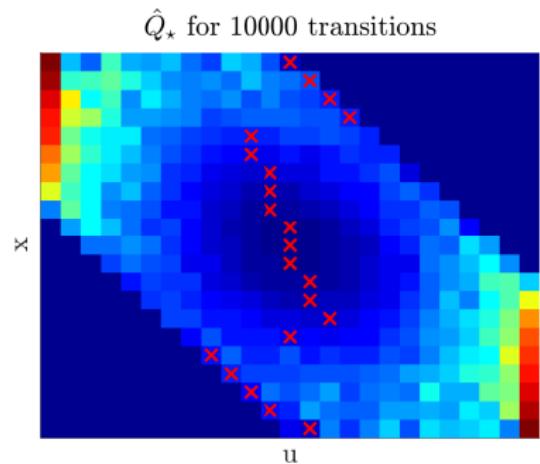
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

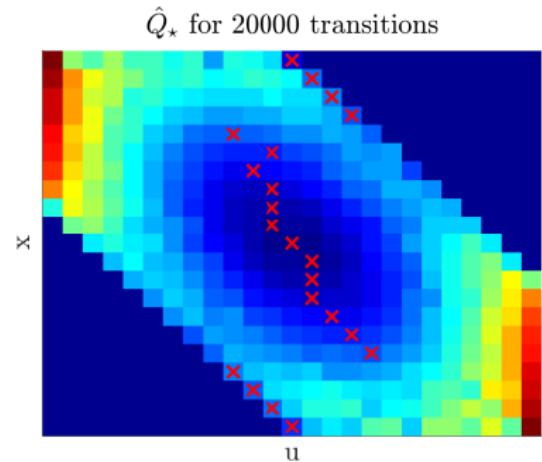
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

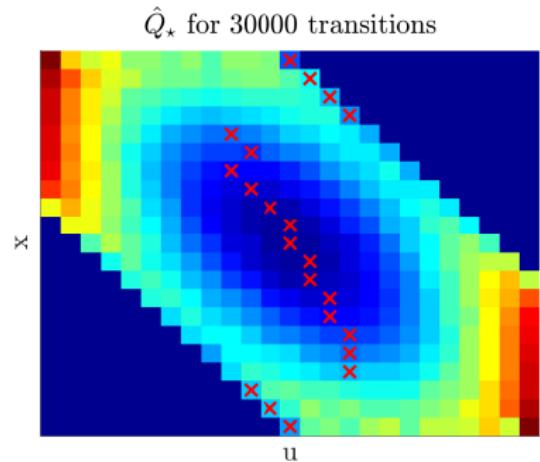
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

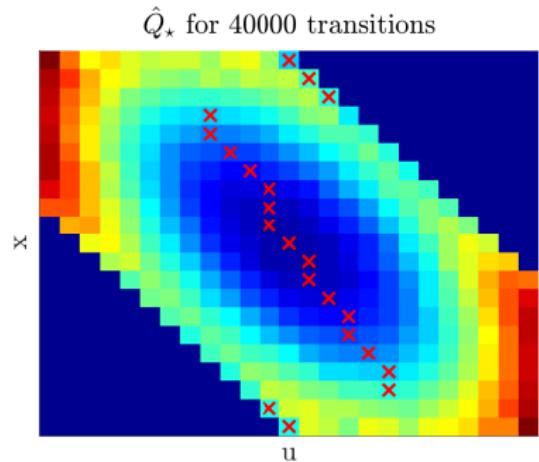
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

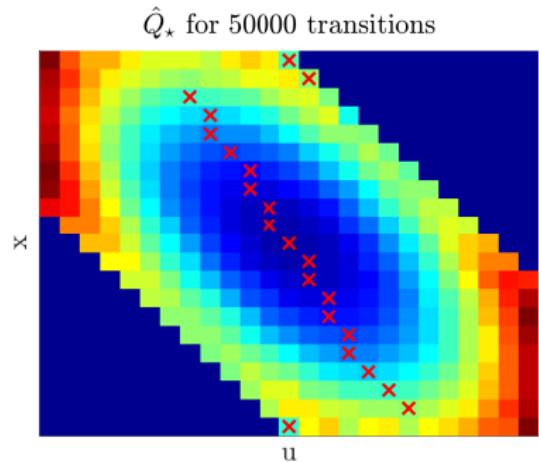
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

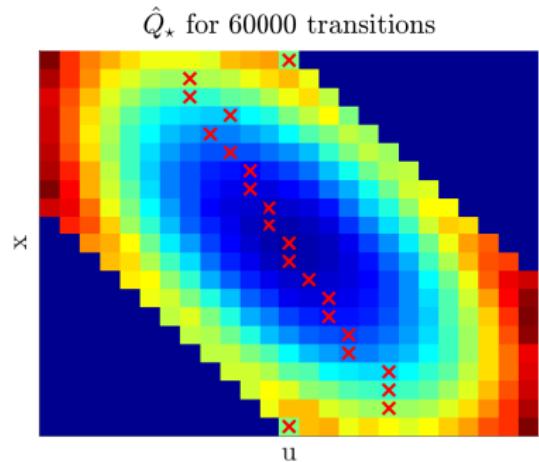
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

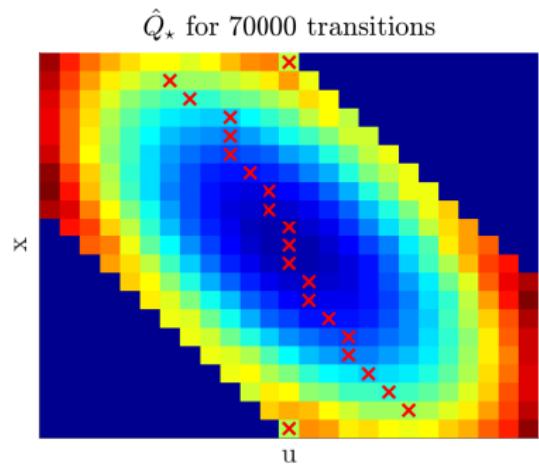
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

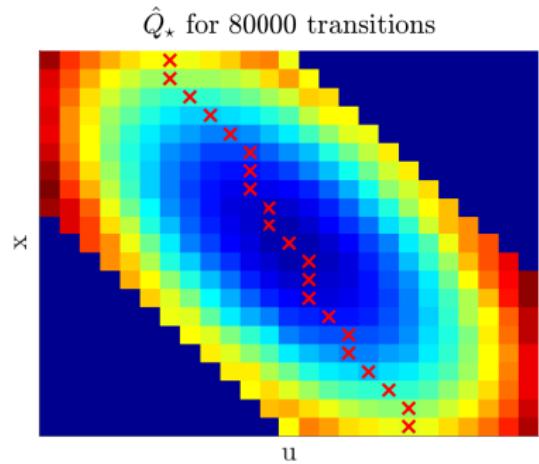
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_* (\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

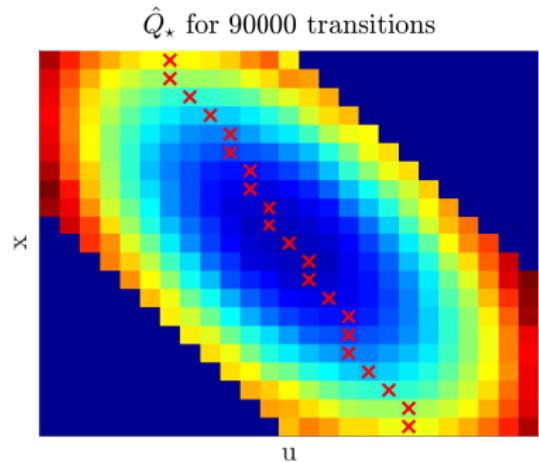
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_* (\mathbf{x}_{k+1}) - \hat{Q}_* (\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_* (\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_* (\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

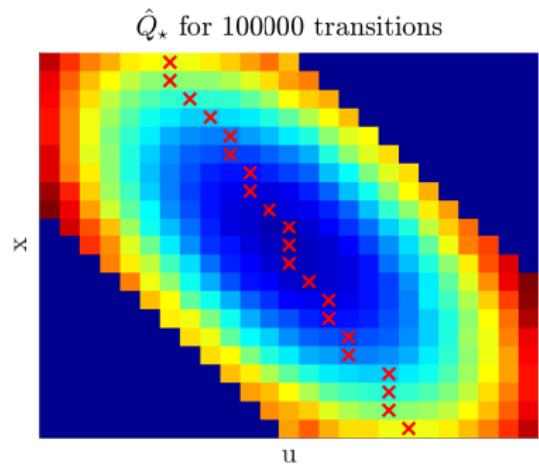
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Learning the optimal action-value function - Off-policy

Optimal action-value function for policy π :

$$Q_{\star}(\mathbf{x}, \mathbf{u}) = \min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_0 = \mathbf{u}, \mathbf{u}_{k>0} = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

TD control

$$\hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

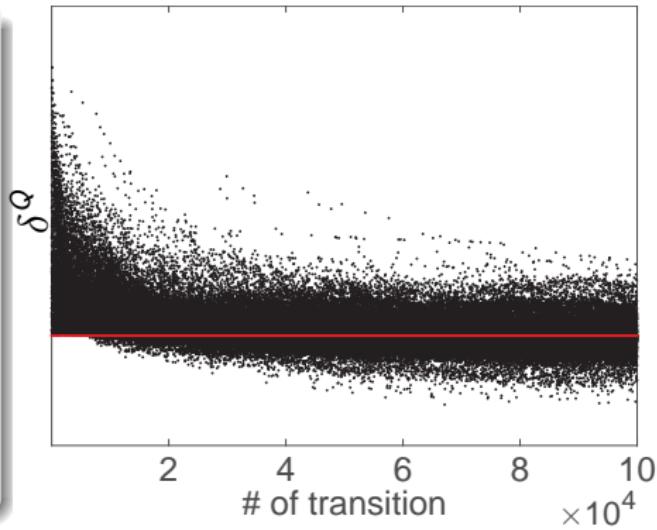
$$\delta_k^Q = L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \hat{V}_{\star}(\mathbf{x}_{k+1}) - \hat{Q}_{\star}(\mathbf{x}_k, \mathbf{u}_k)$$

where

$$\hat{V}_{\star}(\mathbf{x}) = \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$

Policy defined as:

$$\hat{\pi}(\mathbf{x}) = \arg \min_{\mathbf{u}'} \hat{Q}_{\star}(\mathbf{x}, \mathbf{u}')$$



Eligibility trace for learning the optimal action-value function

1-step backup

TD for V_π evaluation:

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

“ ∞ ” backup

TD(λ) for V_π evaluation:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}$$

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

Eligibility trace for learning the optimal action-value function

1-step backup

TD for V_π evaluation:

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

“ ∞ ” backup

TD(λ) for V_π evaluation:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}$$

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

TD for \hat{Q}_π - Off-policy

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k$$

where **u** can be chosen arbitrarily

Eligibility trace for learning the optimal action-value function

1-step backup

TD for V_π evaluation:

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

“ ∞ ” backup

TD(λ) for V_π evaluation:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}$$

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

TD for \hat{Q}_π - Off-policy

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k$$

where \mathbf{u} can be chosen arbitrarily

Can we use E for \hat{Q}_π ???

Eligibility trace for learning the optimal action-value function

1-step backup

TD for V_π evaluation:

$$\hat{V}_\pi(\mathbf{x}_k) \leftarrow \hat{V}_\pi(\mathbf{x}_k) + \alpha \delta_k$$

“ ∞ ” backup

TD(λ) for V_π evaluation:

$$E(\mathbf{x}) \leftarrow \gamma \lambda E(\mathbf{x}) + \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{x}_k \\ 0 & \text{if } \mathbf{x} \neq \mathbf{x}_k \end{cases}$$

$$\hat{V}_\pi(\mathbf{x}) \leftarrow \hat{V}_\pi(\mathbf{x}) + \alpha E(\mathbf{x}) \delta_k, \quad \forall \mathbf{x}$$

TD for \hat{Q}_π - Off-policy

$$\hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(\mathbf{x}_k, \mathbf{u}_k) + \alpha \delta_k$$

where \mathbf{u} can be chosen arbitrarily

Can we use E for \hat{Q}_π ???

Yes, but only on-policy!

Learning the action-value function - On policy (SARSA)

Off-policy Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(x_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(x_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(x_k, \mathbf{u}_k) + \gamma \hat{Q}_\pi(x_{k+1}, \pi(x_{k+1})) - \hat{Q}_\pi(x_k, \mathbf{u}_k)$$

where \mathbf{u}_k can be “arbitrary”

Remarks: if π is defined as:

$$\pi(x_k) = \min_u \hat{Q}_\pi(x_k, u)$$

throughout the iterations, then

$$\hat{Q}_\pi \rightarrow Q_\star \text{ and } \pi \rightarrow \pi_\star$$

under some conditions...

On-policy Temporal Difference (TD) for \hat{Q}_π

$$\hat{Q}_\pi(x_k, \mathbf{u}_k) \leftarrow \hat{Q}_\pi(x_k, \mathbf{u}_k) + \alpha \delta_k^Q$$

$$\delta_k^Q = L(x_k, \mathbf{u}_k) + \gamma \hat{Q}_\pi(x_{k+1}, \mathbf{u}_{k+1}) - \hat{Q}_\pi(x_k, \mathbf{u}_k)$$

where \mathbf{u}_k must be chosen from $\mathbf{u}_k = \pi(x_k) + d_k$

Disturbance d_k must decay over the iterations!

SARSA(λ)

SARSA with eligibility trace - “Emulates” n -step Monte Carlo:

$$E(x, u) \leftarrow \gamma \lambda E(x, u) + \begin{cases} 1 & \text{if } x = x_k, u = u \\ 0 & \text{otherwise} \end{cases}$$

$$\delta^Q = L(x_k, u_k) + \gamma \hat{Q}_\pi(x_{k+1}, u_{k+1}) - \hat{Q}_\pi(x_k, u_k)$$

$$\hat{Q}_\pi(x, u) \leftarrow \hat{Q}_\pi(x, u) + \alpha \delta^Q E(x, u)$$

Remarks

- All u_k must be chosen using $\pi(x_k) + d_k$ (c.f. ϵ -greedy)
- δ^Q computed based on latest sarsa transition: $x_k, u_k, L, x_{k+1}, u_{k+1}$
- Updates of E and \hat{Q}_π sweep through all state-action pairs x, u
- E keeps a trace of state-input pairs recently visited
- If π is updated as

$$\pi(x_k) = \arg \min_u \hat{Q}_\pi(x_k, u)$$

throughout the iterations, then

$$\hat{Q}_\pi \rightarrow Q_\star \text{ and } \pi \rightarrow \pi_\star$$

under some conditions...

Dynamic Programming vs. Temporal-Difference Learning

Full Backup

Iterative π Evaluation

$$V_{\pi}(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi(x)]$$

sweep on all x

Sampled Backup

TD learning

$$V_{\pi}(x_k) \leftarrow L(x_k, \pi(x_k)) + \gamma V_{\pi}(x_{k+1})$$

update on observed $x_k \rightarrow x_{k+1}$

Dynamic Programming vs. Temporal-Difference Learning

Full Backup

Iterative π Evaluation

$$V_\pi(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi(x)]$$

sweep on all x

Sampled Backup

TD learning

$$V_\pi(x_k) \leftarrow L(x_k, \pi(x_k)) + \gamma V_\pi(x_{k+1})$$

update on observed $x_k \rightarrow x_{k+1}$

Q-Evaluation

$$Q_\pi(x, u) \leftarrow L(x, u) + \gamma \mathbb{E}[Q_\pi(x_+, \pi(x_+)) | x, u]$$

sweep on all x, u

TD on Q

$$Q_\pi(x_k, u_k) \leftarrow L(x_k, u_k) + \gamma Q_\pi(x_{k+1}, \pi(x_{k+1}))$$

update on observed $x_k, u_k \rightarrow x_{k+1}, u_{k+1}$

Dynamic Programming vs. Temporal-Difference Learning

Full Backup

Iterative π Evaluation

$$V_\pi(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi(x)]$$

sweep on all x

Sampled Backup

TD learning

$$V_\pi(x_k) \leftarrow L(x_k, \pi(x_k)) + \gamma V_\pi(x_{k+1})$$

update on observed $x_k \rightarrow x_{k+1}$

Q-Evaluation

$$Q_\pi(x, u) \leftarrow L(x, u) + \gamma \mathbb{E}[Q_\pi(x_+, \pi(x_+)) | x, u]$$

sweep on all x, u

SARSA

$$Q_\pi(x_k, u_k) \leftarrow L(x_k, u_k) + \gamma Q_\pi(x_{k+1}, u_{k+1})$$

update on observed $x_k, u_k \rightarrow x_{k+1}, u_{k+1}$

Dynamic Programming vs. Temporal-Difference Learning

Full Backup

Iterative π Evaluation

$$V_\pi(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_\pi(\mathbf{x}_+) | \mathbf{x}, \pi(\mathbf{x})]$$

sweep on all \mathbf{x}

Sampled Backup

TD learning

$$V_\pi(\mathbf{x}_k) \leftarrow L(\mathbf{x}_k, \pi(\mathbf{x}_k)) + \gamma V_\pi(\mathbf{x}_{k+1})$$

update on observed $\mathbf{x}_k \rightarrow \mathbf{x}_{k+1}$

Q-Evaluation

$$Q_\pi(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[Q_\pi(\mathbf{x}_+, \pi(\mathbf{x}_+)) | \mathbf{x}, \mathbf{u}]$$

sweep on all \mathbf{x}, \mathbf{u}

SARSA

$$Q_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow L(\mathbf{x}_k, \mathbf{u}_k) + \gamma Q_\pi(\mathbf{x}_{k+1}, \mathbf{u}_{k+1})$$

update on observed $\mathbf{x}_k, \mathbf{u}_k \rightarrow \mathbf{x}_{k+1}, \mathbf{u}_{k+1}$

Q-Iteration

$$Q_\pi(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}_+} Q_\pi(\mathbf{x}_+, \mathbf{u}_+) | \mathbf{x}, \mathbf{u} \right]$$

sweep on all \mathbf{x}, \mathbf{u}

Q-learning

$$Q_\pi(\mathbf{x}_k, \mathbf{u}_k) \leftarrow L(\mathbf{x}_k, \mathbf{u}_k) + \gamma \min_{\mathbf{u}_+} Q_\pi(\mathbf{x}_{k+1}, \mathbf{u}_+)$$

update on observed $\mathbf{x}_k, \mathbf{u}_k \rightarrow \mathbf{x}_{k+1}$