# TD Learning

Shusen Wang

# Outline

1. Mathematically derive the TD target.

2. TD learning algorithms:

   - Sarsa algorithm for learning $Q_\pi$.

   - Q-learning algorithm for learning $Q^\star$.

3. Multi-step TD target.

# Derive TD Target

# Discounted Return

Definition of discounted return:

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \gamma^4 \cdot R_{t+4} + \cdots$

$$= \gamma \cdot (R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \gamma^3 \cdot R_{t+4} + \cdots)$$

# Discounted Return

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \gamma^4 \cdot R_{t+4} + \cdots$

$\quad = R_t + \gamma \cdot (R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \gamma^3 \cdot R_{t+4} + \cdots)$

$$= U_{t+1}$$

# Discounted Return

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

- $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot R_{t+3} + \gamma^4 \cdot R_{t+4} + \cdots$

$$= R_t + \gamma \cdot \underbrace{\left( R_{t+1} + \gamma \cdot R_{t+2} + \gamma^2 \cdot R_{t+3} + \gamma^3 \cdot R_{t+4} + \cdots \right)}_{= U_{t+1}}$$

# Derive TD Target

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

# Derive TD Target

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$$= \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \cdot \mathbb{E}[U_{t+1} | s_t, a_t]$$

It is taken w.r.t. $(S_{t+1}, A_{t+1}), (S_{t+2}, A_{t+2}), \cdots$

# Derive TD Target

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$$= \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \boxed{\mathbb{E}[U_{t+1} | s_t, a_t]}$$

$$\mathbb{E}[U_{t+1} | s_t, a_t] \;=\; \mathbb{E}[Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$$

# Derive TD Target

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$$= \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \boxed{\mathbb{E}[U_{t+1} | s_t, a_t]}$$

$$\mathbb{E}[U_{t+1} | s_t, a_t] \;\; = \;\; \mathbb{E}[Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$$

$Q_\pi$ eliminates all the future states and actions from time $t + 1$.

# Derive TD Target

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$$= \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \; \boxed{\mathbb{E}[U_{t+1} | s_t, a_t]}$$

$$\mathbb{E}[U_{t+1} | s_t, a_t] \;=\; \mathbb{E}[Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$$

The expectation is taken w.r.t. only $S_{t+1}$ and $A_{t+1}$.

# Derive TD Target

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$$= \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \cdot \boxed{\mathbb{E}[U_{t+1} | s_t, a_t]}$$

$$= \mathbb{E}[R_t | s_t, a_t] + \gamma \cdot \mathbb{E}[\underline{Q_\pi(S_{t+1}, A_{t+1})} | s_t, a_t].$$

# Derive TD Target

**Identity:** $\boxed{Q_\pi(s_t, a_t)} = \boxed{\mathbb{E}[R_t + \gamma \cdot Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]}$, for all $\pi$.

- Assume $R_t$ depends on $(S_t, A_t, S_{t+1})$.

- $Q_\pi(s_t, a_t) = \mathbb{E}[U_t | s_t, a_t]$

$\qquad = \mathbb{E}[R_t + \gamma \cdot U_{t+1} | s_t, a_t]$

$\qquad = \mathbb{E}[R_t | s_t, a_t] + \gamma \cdot \mathbb{E}[U_{t+1} | s_t, a_t]$

$\qquad = \mathbb{E}[R_t | s_t, a_t] + \gamma \cdot \mathbb{E}[Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t].$

# Derive TD Target

**Identity:** $Q_\pi(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$, for all $\pi$.

- We do not know the expectation.
- Approximate it using Monte Carlo (MC).

# Derive TD Target

**Identity:** $Q_\pi(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$, for all $\pi$.

$y_t$ is its MC approximation.

- Let $(s_{t+1}, r_t)$ be an observation of $(S_{t+1}, R_t)$.

- Sample $a_{t+1} \sim \pi(\cdot | s_{t+1})$.

- TD target: $y_t = r_t + \gamma \cdot Q_\pi(s_{t+1}, a_{t+1})$.

# Derive TD Target

**Identity:** $Q_\pi(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t]$, for all $\pi$.

$y_t$ is its MC approximation.

TD learning: Encourage $Q_\pi(s_t, a_t)$ to approach $y_t$.

# Sarsa

# Tabular Version

- We want to learn $Q_\pi(s, a)$.

- Suppose there are finite number of states and actions.

- Draw a table and learn the entries.

| | Action $a_1$ | Action $a_2$ | Action $a_3$ | Action $a_4$ | $\cdots$ |
|---|---|---|---|---|---|
| State $s_1$ | | | | | |
| State $s_2$ | | | | | |
| State $s_3$ | | | | | |
| $\vdots$ | | | | | |

# Sarsa (tabular version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- Sample $a_{t+1} \sim \pi(\cdot \mid s_{t+1})$, where $\pi$ is a policy function.

- TD target: $y_t = r_t + \gamma \cdot \boxed{Q_\pi(s_{t+1}, a_{t+1})}$.

|  | Action $a_1$ | Action $a_2$ | Action $a_3$ | Action $a_4$ | $\cdots$ |
|---|---|---|---|---|---|
| State $s_1$ |  |  |  |  |  |
| State $s_2$ |  |  |  |  |  |
| State $s_3$ |  |  |  |  |  |
| $\vdots$ |  |  |  |  |  |

# Sarsa (tabular version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- Sample $a_{t+1} \sim \pi(\cdot \mid s_{t+1})$, where $\pi$ is a policy function.

- TD target: $y_t = r_t + \gamma \cdot Q_\pi(s_{t+1}, a_{t+1})$.

- TD error: $\delta_t = Q_\pi(s_t, a_t) - y_t$.

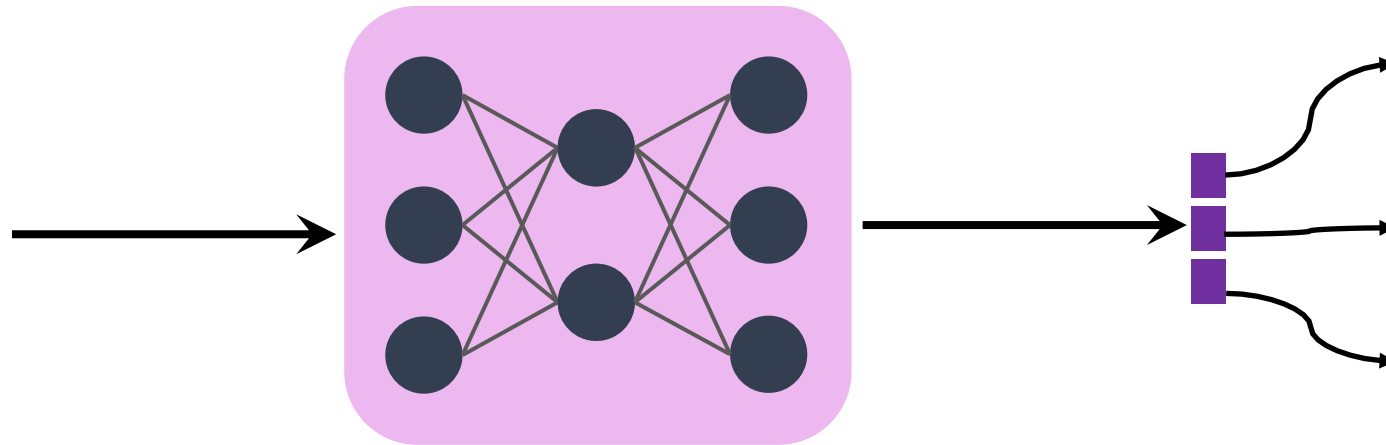- Update: $Q_\pi(s_t, a_t) \leftarrow Q_\pi(s_t, a_t) - \alpha \cdot \delta_t$.

make $Q_\pi(s_t, a_t)$ closer to $y_t$

# Value Network Version

- Approximate $Q_\pi(s, a)$ by a value network, $q(s, a|\mathbf{w})$.



**state $s$**

**Value Network**
**(parameterized by $\mathbf{w}$)**

$q(s, \text{"left"}; \mathbf{w})$

$q(s, \text{"right"}; \mathbf{w})$

$q(s, \text{"up"}; \mathbf{w})$

# Value Network Version

- Approximate $Q_\pi(s, a)$ by a value network, $q(s, a|\mathbf{w})$.

- Note that the $Q_\pi(s, a)$ and $q(s, a|\mathbf{w})$ depends on $\pi$.

- $q$ is used as the critic who evaluates the actor. (Actor-Critic Method.)

- We seek to learn the parameter, $\mathbf{w}$.

# Sarsa (Value Network Version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- Sample $a_{t+1} \sim \pi(\cdot \mid s_{t+1})$, where $\pi$ is a policy function.

- TD target: $y_t = r_t + \gamma \cdot q(s_{t+1}, a_{t+1} \mid \mathbf{w})$.

- TD error: $\delta_t = q(s_t, a_t \mid \mathbf{w}) - y_t$.

- SGD: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot \delta_t \cdot \dfrac{\partial \, q(s_t, a_t \mid \mathbf{w})}{\partial \, \mathbf{w}}$.

# Sarsa: Recap

- **Goal:** Learn the action-value function $Q_\pi$.

- **Tabular version** (directly learn $Q_\pi$).

  - There are finite state and actions.

  - Draw a table, and update the table using Sarsa.

- **Value network version** (function approximation).

  - Approximate $Q_\pi$ by the value network $q(s, a|\mathbf{w})$.

  - Update the parameter, $\mathbf{w}$, using Sarsa.

  - Application: actor-critic method.

# Q-Learning

# Sarsa VS Q-Learning

- Sarsa is for training action-value function, $Q_\pi(s, a)$.

- TD target: $y_t = r_t + \gamma \cdot Q_\pi(s_{t+1}, a_{t+1})$.

- We used Sarsa for value network (critic).

# Sarsa VS Q-Learning

- Sarsa is for training action-value function, $Q_\pi(s, a)$.

- TD target: $y_t = r_t + \gamma \cdot Q_\pi(s_{t+1}, a_{t+1})$.

- We used Sarsa for value network (critic).


- Q-learning is for training the optimal action-value function, $Q^\star(s, a)$.

- TD target: $y_t = r_t + \gamma \cdot \max_a Q^\star(s_{t+1}, a)$.

- We used Q-learning for DQN.

# Derive TD Target

- We have proved that for all $\pi$,

$$Q_\pi(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_\pi(S_{t+1}, A_{t+1}) | s_t, a_t].$$

- If $\pi$ is the optimal policy $\pi^\star$, then

$$Q_{\pi^\star}(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q_{\pi^\star}(S_{t+1}, A_{t+1}) | s_t, a_t].$$

- We denote $Q_{\pi^\star}$ by $Q^\star$.

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^\star(S_{t+1}, A_{t+1}) | s_t, a_t].$

# Derive TD Target

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^\star(S_{t+1}, A_{t+1}) | s_t, a_t]$.

- The action $A_{t+1}$ is computed by

$$A_{t+1} = \operatorname*{argmax}_{a} Q^\star(S_{t+1}, a).$$

- Thus $Q^\star(S_{t+1}, A_{t+1}) = \max_{a} Q^\star(S_{t+1}, a)$.

# Derive TD Target

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}[R_t + \gamma \cdot Q^\star(S_{t+1}, A_{t+1}) | s_t, a_t]$.

- The action $A_{t+1}$ is computed by

$$A_{t+1} = \underset{a}{\mathrm{argmax}}\, Q^\star(S_{t+1}, a).$$

- Thus $Q^\star(S_{t+1}, A_{t+1}) = \underset{a}{\max}\, Q^\star(S_{t+1}, a)$.

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}\left[R_t + \gamma \cdot \underset{a}{\max}\, Q^\star(S_{t+1}, a) \,\middle|\, s_t, a_t\right]$.

# Derive TD Target

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}\left[R_t + \gamma \cdot \max_a Q^\star(S_{t+1}, a) \,\middle|\, s_t, a_t\right].$

# Derive TD Target

**Identity:** $Q^\star(s_t, a_t) = \mathbb{E}\left[ R_t + \gamma \cdot \max_a Q^\star(S_{t+1}, a) \,\middle|\, s_t, a_t \right].$

$$\approx y_t$$

- Let $(s_{t+1}, r_t)$ be an observation of $(S_{t+1}, R_t)$.

- TD target: $y_t = r_t + \gamma \cdot \max_a Q^\star(s_{t+1}, a).$

# Q-Learning (tabular version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- TD target: $y_t = r_t + \gamma \cdot \max_a Q^\star(s_{t+1}, a)$.

# Q-Learning (tabular version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- TD target: $y_t = r_t + \gamma \cdot \boxed{\max_a Q^\star(s_{t+1}, a)}$.

| | Action $a_1$ | Action $a_2$ | Action $a_3$ | Action $a_4$ | $\cdots$ |
|---|---|---|---|---|---|
| State $s_1$ | | | | | |
| State $s_2$ | | | | | |
| State $s_3$ | | | | | |
| $\vdots$ | | | | | |

# Q-Learning (tabular version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- TD target: $y_t = r_t + \gamma \cdot \max_a Q^\star(s_{t+1}, a)$.

- TD error: $\delta_t = Q^\star(s_t, a_t) - y_t$.

- Update:   $Q^\star(s_t, a_t) \leftarrow Q^\star(s_t, a_t) - \alpha \cdot \delta_t$.

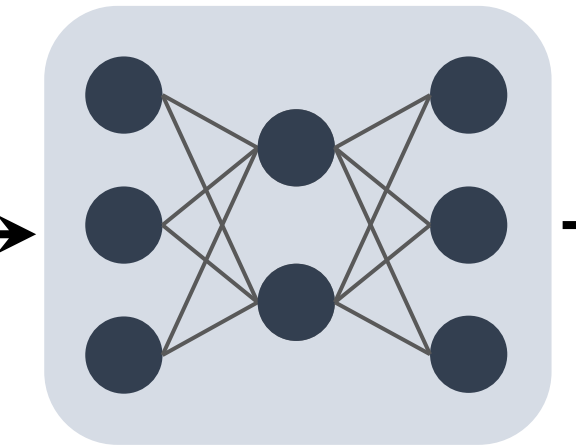make $Q^\star(s_t, a_t)$ closer to $y_t$

# DQN Version

- Approximate $Q^\star(s, a)$ by a value network, $Q(s, a | \mathbf{w})$.

- The network is called deep Q network (DQN).



**state $s$**

**DQN**
**(parameterized by $\mathbf{w}$)**

$Q(s, \text{"left"}; \mathbf{w})$

$Q(s, \text{"right"}; \mathbf{w})$

$Q(s, \text{"up"}; \mathbf{w})$

# DQN Version

- Approximate $Q^\star(s, a)$ by a value network, $Q(s, a|\mathbf{w})$.

- The network is called deep Q network (DQN).

- It controls the agent by: $a_t = \max\limits_{a} Q(s_t, a|\mathbf{w})$

- We seek to learn the parameter, $\mathbf{w}$.

# Q-Learning (DQN Version)

- Observe $(s_t, a_t, r_t, s_{t+1})$.

- TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a \mid \mathbf{w})$.

- TD error: $\delta_t = Q(s_t, a_t \mid \mathbf{w}) - y_t$.

- Update: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot \delta_t \cdot \dfrac{\partial\, q(s_t, a_t \mid \mathbf{w})}{\partial\, \mathbf{w}}$.

# Q-Learning: Recap

- **Goal:** Learn the optimal action-value function $Q^\star$.

- **Tabular version** (directly learn $Q^\star$).

  - There are finite state and actions.

  - Draw a table, and update the table using Q-learning.

- **DQN version** (function approximation).

  - Approximate $Q^\star$ by the value network $Q(s, a | \mathbf{w})$.

  - Update the parameter, $\mathbf{w}$, using Q-learning.

# Multi-Step Target

# Multi-Step Return

**Identity:** $U_t = R_t + \gamma \cdot \boxed{U_{t+1}}$

Replace $U_{t+1}$ by $R_{t+1} + \gamma \cdot U_{t+2}$

- It follows that $\boxed{U_{t+1} = R_{t+1} + \gamma \cdot U_{t+2}}$

# Multi-Step Return

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

Replace $U_{t+1}$ by $R_{t+1} + \gamma \cdot U_{t+2}$

**Identity:** $U_t = R_t + \gamma \cdot (R_{t+1} + \gamma \cdot U_{t+2}).$

# Multi-Step Return

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

Replace $U_{t+1}$ by $R_{t+1} + \gamma \cdot U_{t+2}$

**Identity:** $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot U_{t+2}.$

# Multi-Step Return

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

**Identity:** $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot U_{t+2}.$

**Identity:** $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot U_{t+3}.$

# Multi-Step Return

**Identity:** $U_t = R_t + \gamma \cdot U_{t+1}.$

**Identity:** $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot U_{t+2}.$

**Identity:** $U_t = R_t + \gamma \cdot R_{t+1} + \gamma^2 \cdot R_{t+2} + \gamma^3 \cdot U_{t+3}.$

**Identity:** $U_t = \sum_{i=0}^{m-1} \gamma^i \cdot R_{t+i} + \gamma^m \cdot U_{t+m}.$

# Multi-Step TD Targets

**Identity:** $U_t = \sum_{i=0}^{m-1} \gamma^i \boxed{R_{t+i}} + \gamma^m \cdot \boxed{U_{t+m}}$

- $m$-step TD target for learning $Q_\pi$ (i.e., Sarsa):

$$y_t = \sum_{i=0}^{m-1} \gamma^i \cdot \boxed{r_{t+i}} + \gamma^m \cdot \boxed{Q_\pi(s_{t+m}, a_{t+m})}.$$

# Multi-Step TD Targets

**Identity:** $U_t = \sum_{i=0}^{m-1} \gamma^i \boxed{R_{t+i}} + \gamma^m \cdot \boxed{U_{t+m}}$

- $m$-step TD target for learning $Q_\pi$ (i.e., Sarsa):

$$y_t = \sum_{i=0}^{m-1} \gamma^i \cdot r_{t+i} + \gamma^m \cdot Q_\pi(s_{t+m}, a_{t+m}).$$

- $m$-step TD target for learning $Q^\star$ (i.e., Q-learning):

$$y_t = \sum_{i=0}^{m-1} \gamma^i \cdot \boxed{r_{t+i}} + \gamma^m \cdot \boxed{\max_a Q^\star(s_{t+m}, a)}$$

# Comparison

- One-step TD target for learning $Q^\star$ (or DQN):

$$y_t = r_t + \gamma \cdot \max_a Q^\star(s_{t+1}, a).$$

- $m$-step TD target for learning $Q^\star$ (or DQN):

$$y_t = \sum_{i=0}^{m-1} \gamma^i \cdot r_{t+i} + \gamma^m \cdot \max_a Q^\star(s_{t+m}, a).$$

- If $m$ is suitably tuned, $m$-step target works better than one-step target [1].

**Reference:**

1. Hossel et al. Rainbow: combining improvements in deep reinforcement learning. In *AAAI*, 2018.

# Summary

# Summary

1. Mathematically derived the TD target.

2. TD learning algorithms:

   - Sarsa algorithm for learning $Q_\pi$ and value network (critic).

   - Q-learning algorithm for learning $Q^\star$ and DQN.

3. Multi-step TD target.

# Improvements for TD Learning

1. Experience replay (ER) and its variants (e.g., prioritized ER.)

    - Reuse experience.

    - Eliminate correlation.

2. Target network and double DQN.

    - Address the overestimation issue of Q-learning.

3. Multi-step TD target.

# Thank you!