# Reinforcement Learning
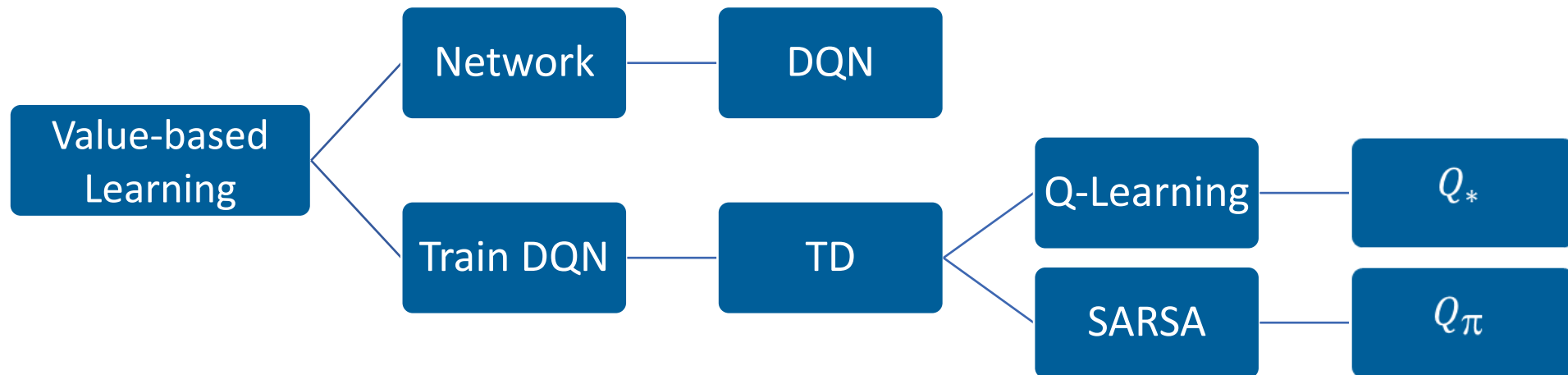
GongYing  2023.11.24

# Contents

# Background

**Goal of RL:** Maximize sum of rewards.

**How:** Take the best <span style="color:red">action</span>, basing on <span style="color:red">Q*(s,a)</span>.

$$a_t = argmax_a Q(s_t, a, \boldsymbol{w})$$

**Problem:** Get to know <span style="color:red">Q*(s,a)</span>.

# DQN

**DQN:** Nueral network named Deep Q Network, denoted as $Q(s,a;\mathbf{w})$.

**Goal:** Use $Q(s,a;\mathbf{w})$ to approximate $Q^*(s,a)$.

state → conv → feature → dense → Q
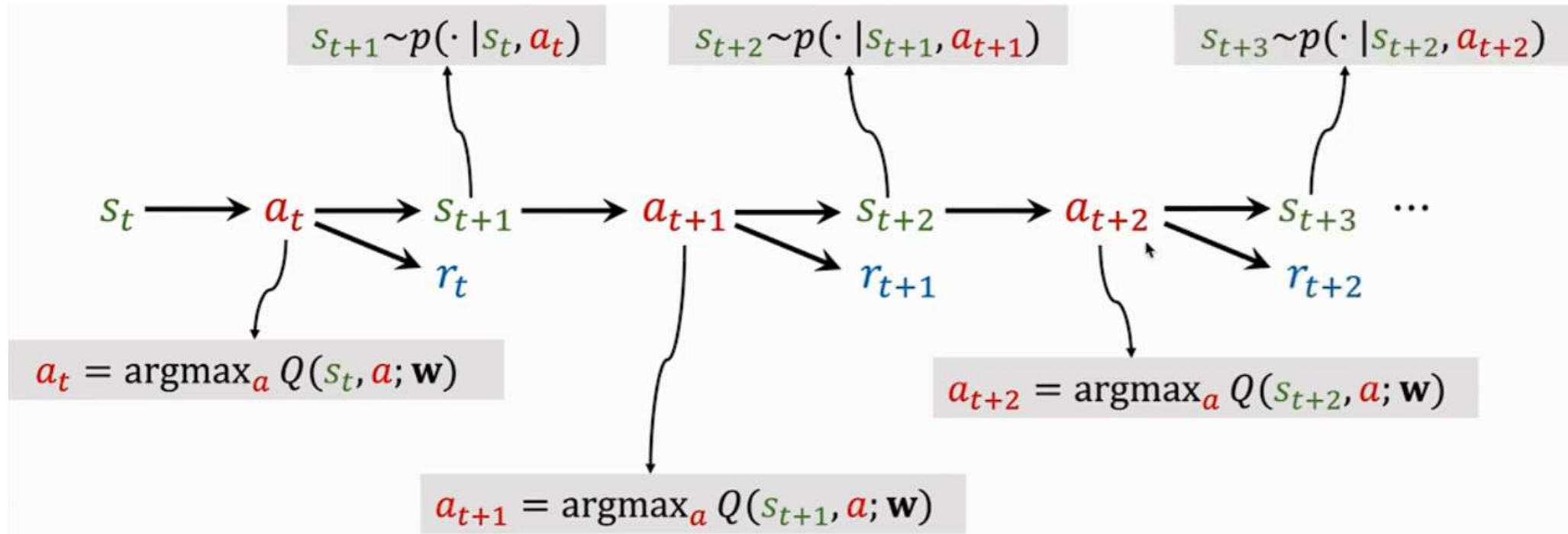
# DQN

state → conv → feature → dense → Q

**Example:** If action space is {"left", "up", "right"}, then **Q** is a

3-dimensional vector $\begin{pmatrix} Q(s_t,"left";\boldsymbol{w}_t) \\ Q(s_t,"up";\boldsymbol{w}_t) \\ Q(s_t,"right";\boldsymbol{w}_t) \end{pmatrix}$

# DQN

## Apply DQN in a game



$$s_{t+1} \sim p(\cdot \mid s_t, a_t)$$

$$s_{t+2} \sim p(\cdot \mid s_{t+1}, a_{t+1})$$

$$s_{t+3} \sim p(\cdot \mid s_{t+2}, a_{t+2})$$

$$s_t \longrightarrow a_t \longrightarrow s_{t+1} \longrightarrow a_{t+1} \longrightarrow s_{t+2} \longrightarrow a_{t+2} \longrightarrow s_{t+3} \quad \cdots$$

$$r_t \qquad r_{t+1} \qquad r_{t+2}$$

$$a_t = \mathrm{argmax}_a \, Q(s_t, a; \mathbf{w})$$

$$a_{t+2} = \mathrm{argmax}_a \, Q(s_{t+2}, a; \mathbf{w})$$

$$a_{t+1} = \mathrm{argmax}_a \, Q(s_{t+1}, a; \mathbf{w})$$

# TD – Q-Learning

**TD:** Temporal difference, used to train DQN, including

Q-Learning & SARSA.

**Q-Learning** aims to approximate Q*, while **SARSA** aims

to approximate $Q_\pi$.

# TD – Q-Learning

Observe transition: $(s_t, a_t, r_t, s_{t+1})$

$$\hat{q} = Q(s_t, a_t; \boldsymbol{w})$$

TD target: $y_t = r_t + \gamma * \max_a Q(s_{t+1}, a; \boldsymbol{w}) \rightarrow$ *Bellman Optimality Equation*

TD error: $\delta_t = Q(s_t, a_t; \boldsymbol{w}) - y_t$

Loss: $\delta_t^2/2$

Gradient: $\dfrac{\partial \delta_t^2/2}{\partial \boldsymbol{w}} = \delta_t * \dfrac{\partial Q(s_t, a_t; \boldsymbol{w})}{\partial \boldsymbol{w}}$

Update: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \alpha * \delta_t * \dfrac{\partial Q(s_t, a_t; \boldsymbol{w})}{\partial \boldsymbol{w}}$

# TD – SARSA

SARSA: StateActionRewardStateAction

$$(s_t, a_t, r_t, s_{t+1}, a_{t+1})$$

SARSA is used to approximate $Q_\pi$.

# TD – SARSA

Observe transition: $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$

$$\hat{q} = q(s_t, a_t; \boldsymbol{w})$$

TD target: $y_t = r_t + \gamma * q(s_{t+1}, a_{t+1}; \boldsymbol{w})$

TD error: $\delta_t = q(s_t, a_t; \boldsymbol{w}) - y_t$

Loss: $\delta_t^2 / 2$

Gradient: $\dfrac{\partial \delta_t^2 / 2}{\partial \boldsymbol{w}} = \delta_t * \dfrac{\partial q(s_t, a_t; \boldsymbol{w})}{\partial \boldsymbol{w}}$

Update: $\boldsymbol{w} \leftarrow \boldsymbol{w} - \alpha * \delta_t * \dfrac{\partial q(s_t, a_t; \boldsymbol{w})}{\partial \boldsymbol{w}}$

# On-policy & off-policy

**On-policy**: Behavior policy is the same with target policy.

**Off-policy**: Behavior policy is **not** the same with target policy.

**Behavior policy**: Used to control agent and gain experience.

**Target policy**: A certain policy that RL aims to find.

For example, Q-Learning is off-policy while SARSA is on-policy.

# Experience replay

Meaning: Save interaction records (experience) into arrays and reuse them later to train the agent.

Advantages: Elimination of relevance and faster convergence speed.

Only applies to off-policy solutions.

# Thank you.