

Value-Based Reinforcement Learning

Shusen Wang

Value Functions

Action-Value Function $Q(s, a)$

Definition: Discounted return (aka cumulative discounted future reward).

- $R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$ (to infinity.)

Definition: Action-value function.

- $Q_\pi(s, a) = \mathbb{E} [R_t | s, a, \pi].$



- Taken w.r.t. the randomness in the state transition $p(s' | s, a)$.
- The state transition $(s_t, a_t) \mapsto s_{t+1}$ is random.

Action-Value Function $Q(s, a)$

Definition: Discounted return (aka cumulative discounted future reward).

- $R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$ (to infinity.)

Definition: Action-value function.

- $Q_\pi(s, a) = \mathbb{E} [R_t | s, a, \pi].$

Definition: Optimal action-value function.

- $Q^*(s, a) = \max_{\pi} Q_\pi(s, a).$

Deep Q-Network (DQN)

Approximate the Q Function

Goal: Win the game (\approx maximize the total reward.)

Question: If we know $Q^*(s, a)$, what is the best action?

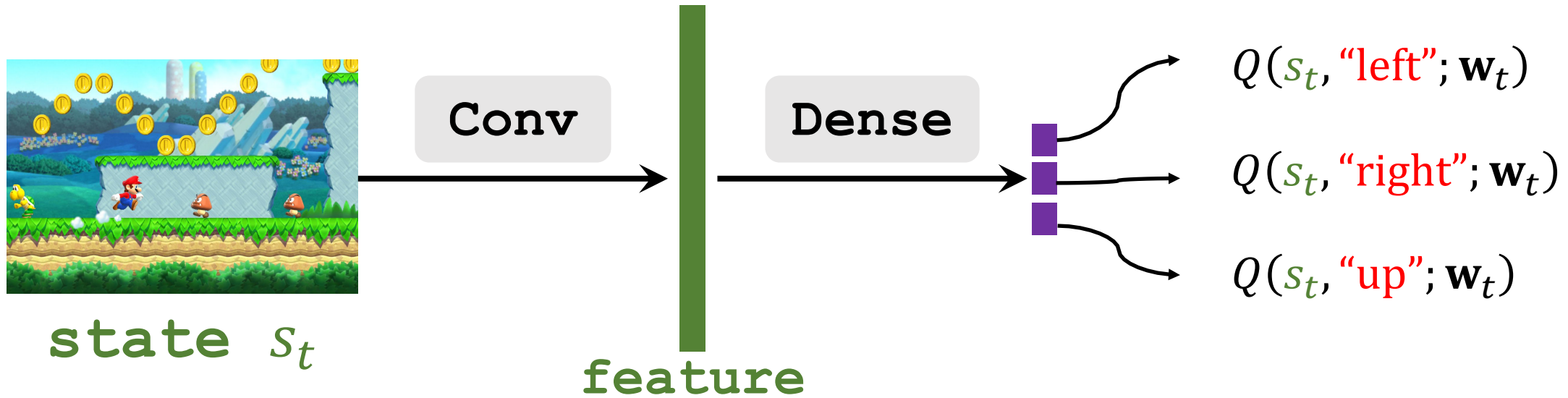
- Obviously, the best action is $a^* = \underset{a}{\operatorname{argmax}} Q(s, a)$.

Challenge: We do not know $Q^*(s, a)$.

- **Solution:** Use a deep neural network to approximate $Q^*(s, a)$.
- Let $Q(s, a; \mathbf{w})$ be a neural network parameterized by \mathbf{w} .
- The inputs are state and action; the output is the approximate Q^* .

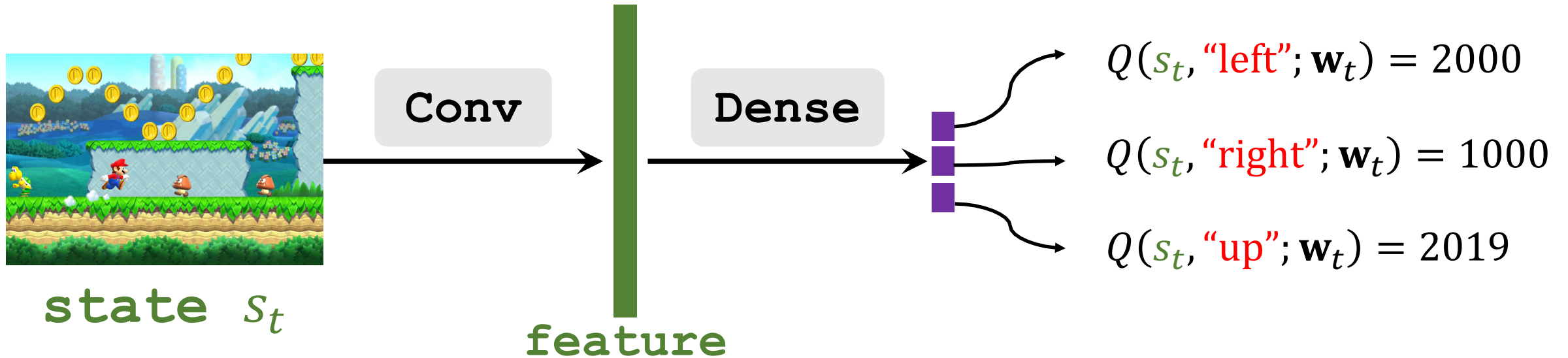
Deep Q Network (DQN)

- Input shape: size of the screenshot.
- Output shape: dimension of action space.



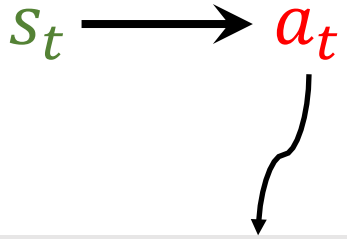
Deep Q Network (DQN)

- Input shape: size of the screenshot.
- Output shape: dimension of action space.



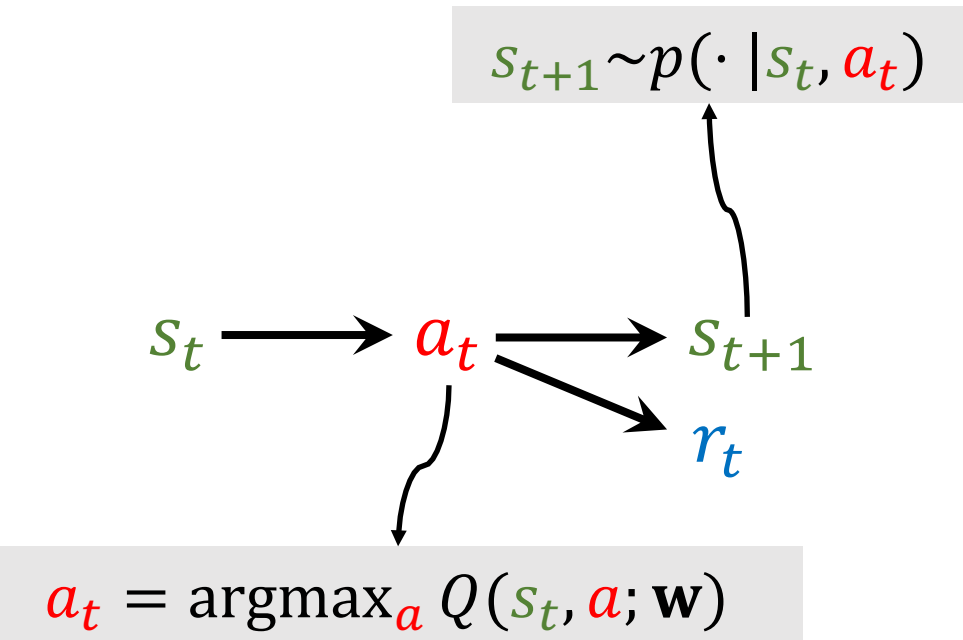
Question: Based on the predictions, what should be the **action**?

Apply DQN to Play Game

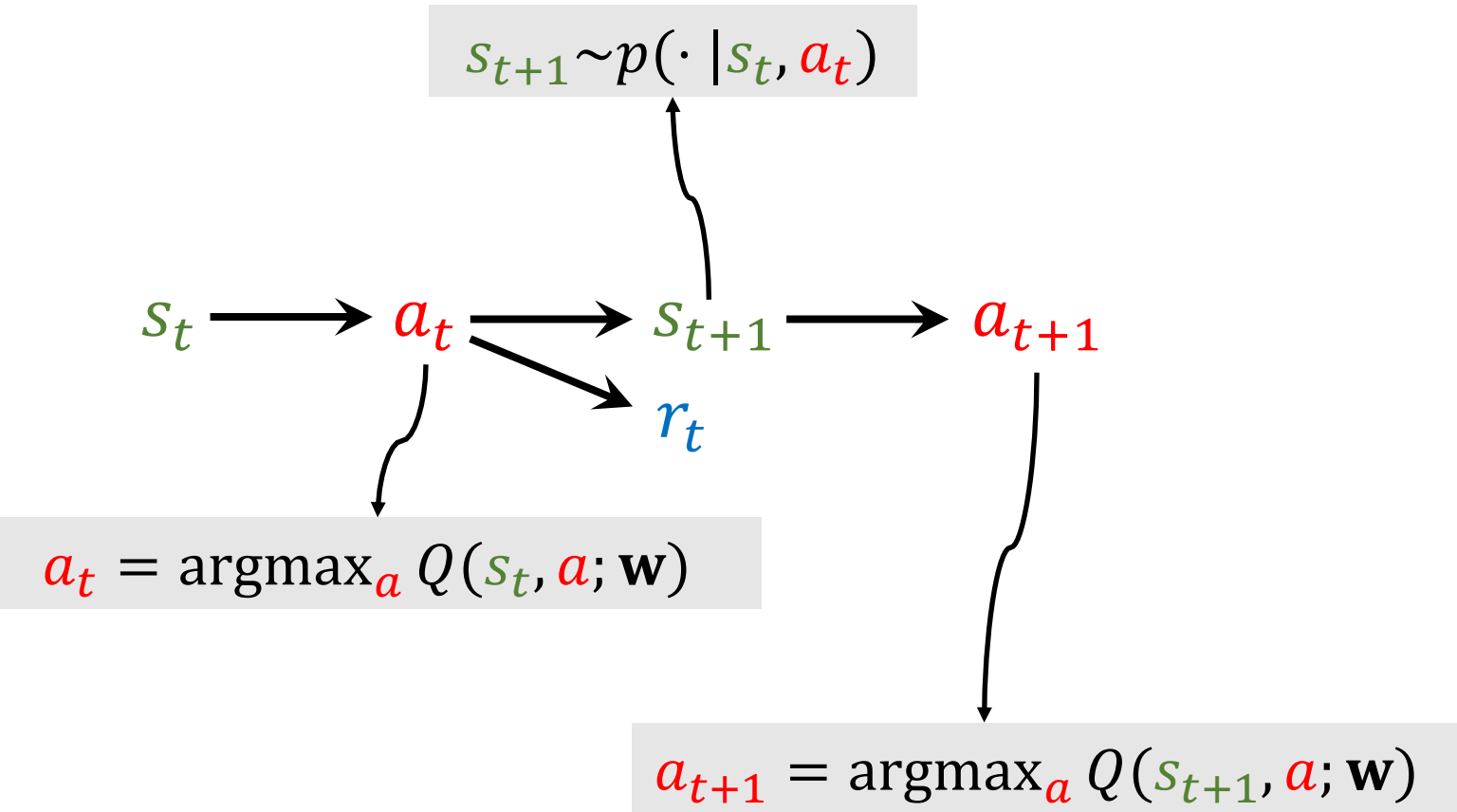


$$a_t = \operatorname{argmax}_a Q(s_t, a; \mathbf{w})$$

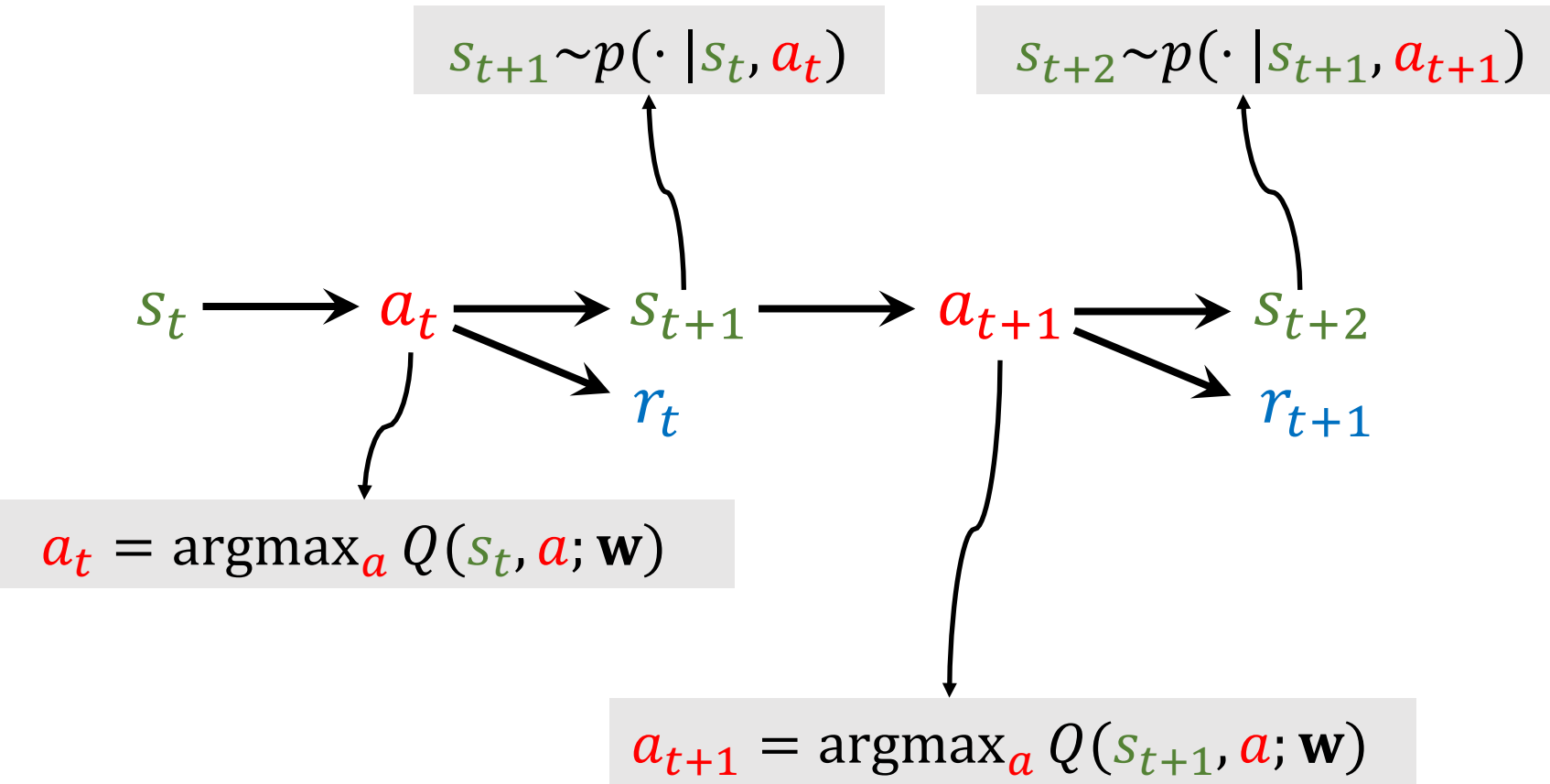
Apply DQN to Play Game



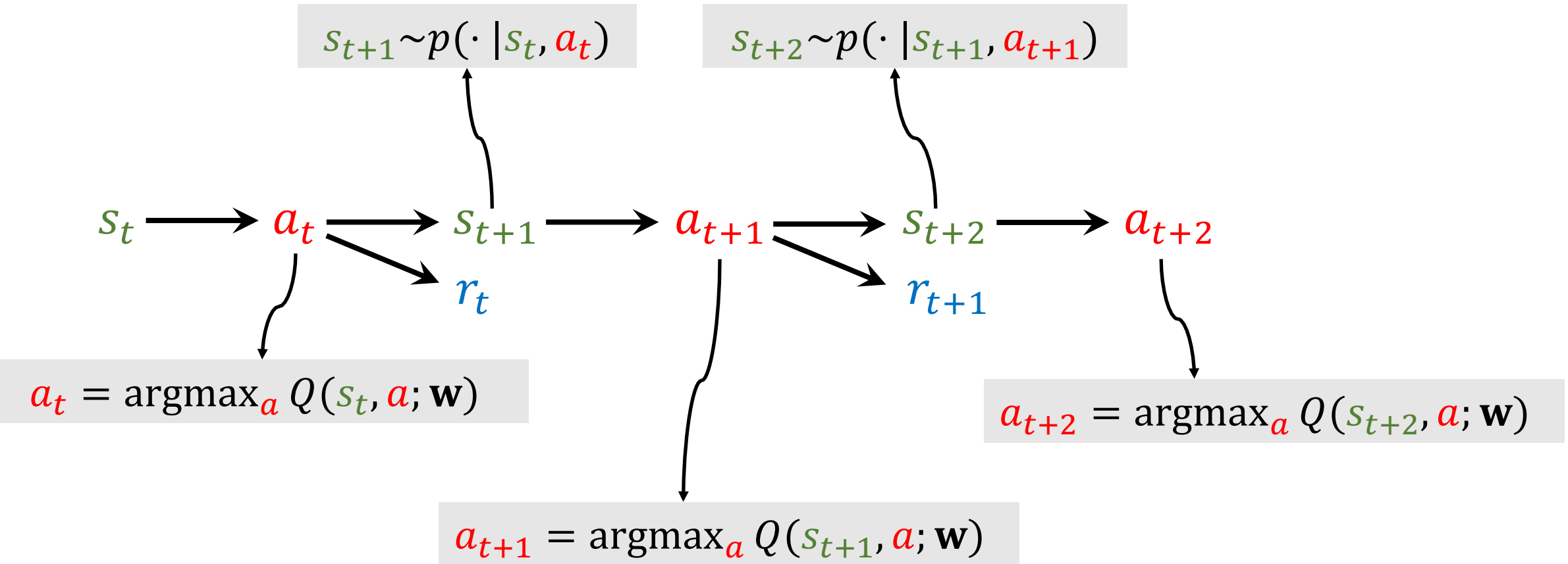
Apply DQN to Play Game



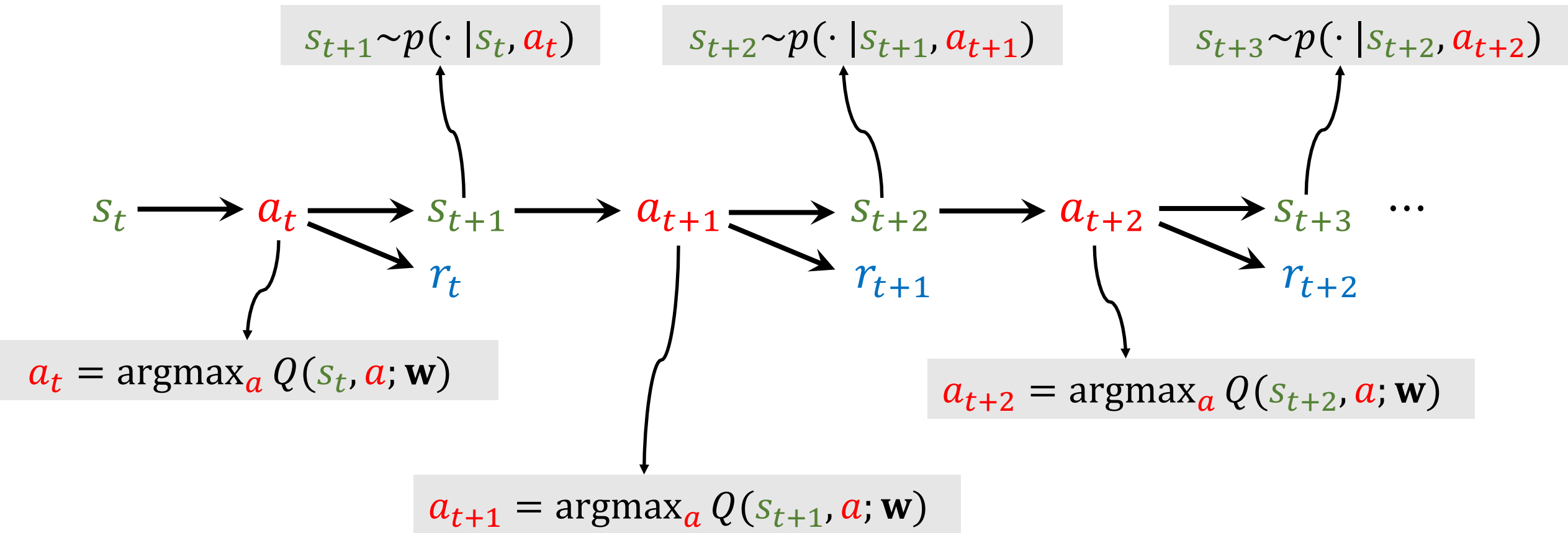
Apply DQN to Play Game



Apply DQN to Play Game



Apply DQN to Play Game

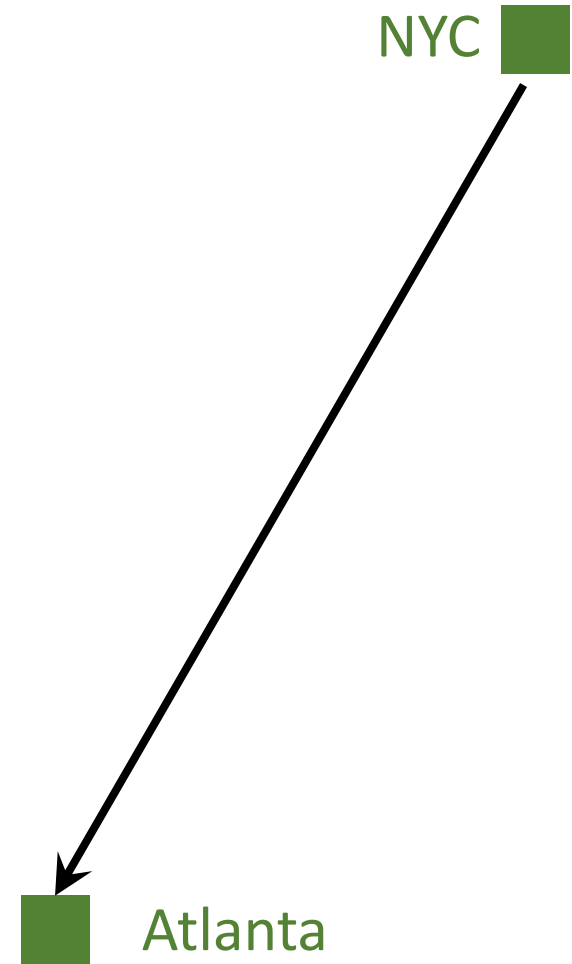


Temporal Difference (TD) Learning

Temporal Difference (TD) Learning

- I want to drive from NYC to Atlanta.
- Model $Q(\mathbf{w})$ estimate the time cost, e.g., 1000 minutes.

Question: How do I update the model?

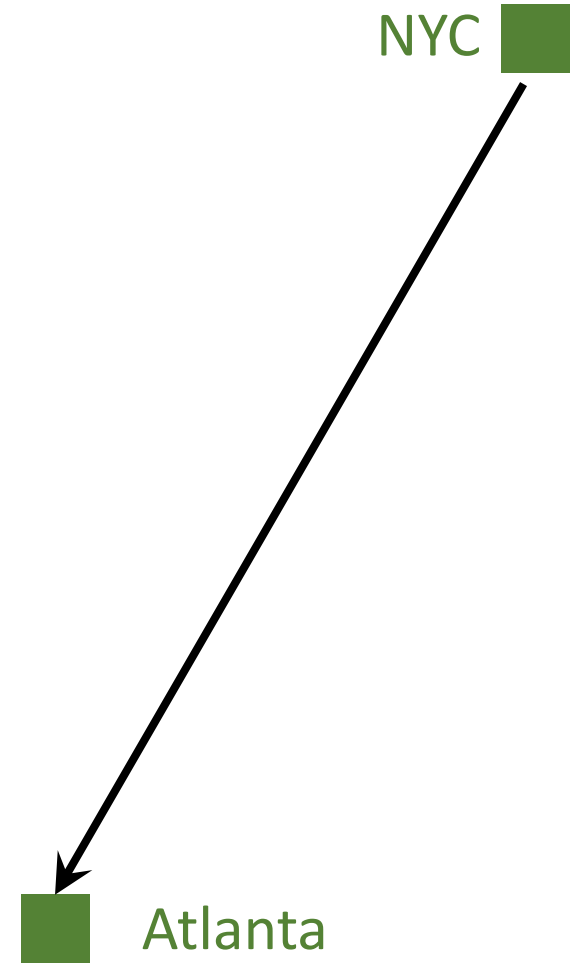


Temporal Difference (TD) Learning

- I want to drive from NYC to Atlanta.
- Model $Q(\mathbf{w})$ estimate the time cost, e.g., 1000 minutes.

Question: How do I update the model?

- Make a prediction: $q = Q(\mathbf{w})$, e.g., $q = 1000$.
- Finish the trip and get the target y , e.g., $y = 860$.
- Loss: $L = \frac{1}{2}(q - y)^2$.
- Gradient: $\frac{\partial L}{\partial \mathbf{w}} = \frac{\partial q}{\partial \mathbf{w}} \cdot \frac{\partial L}{\partial q} = (q - y) \cdot \frac{\partial Q(\mathbf{w})}{\partial \mathbf{w}}$.
- Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{\partial L}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.

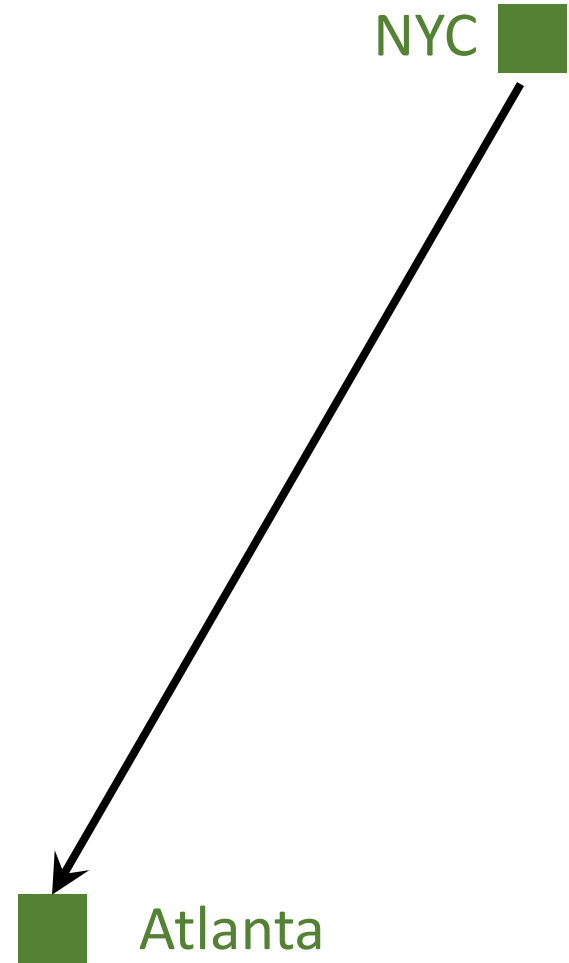


Temporal Difference (TD) Learning

- I want to drive from NYC to Atlanta.
- Model $Q(\mathbf{w})$ estimate the time cost, e.g., 1000 minutes.

Question: How do I update the model?

- Can I update the model **before finishing the trip**?

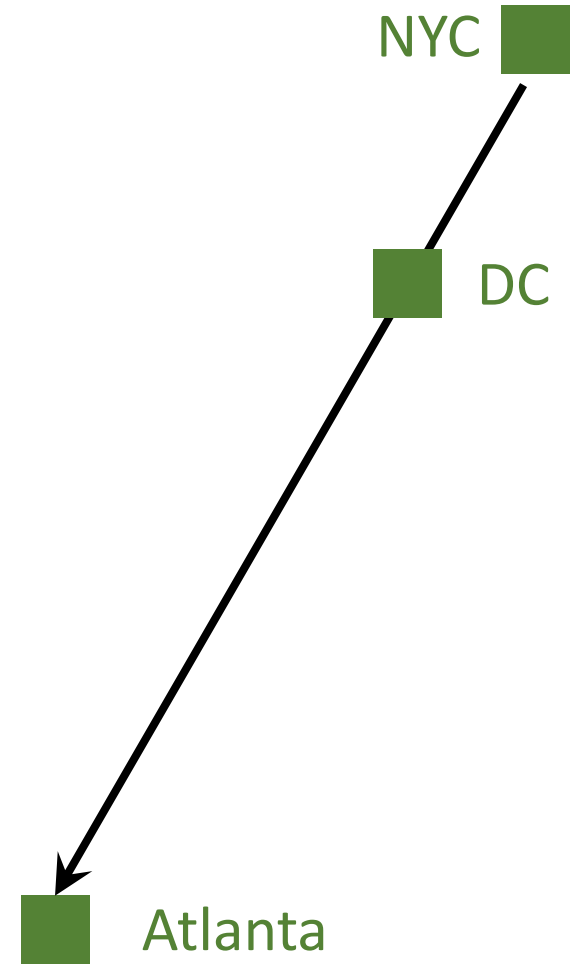


Temporal Difference (TD) Learning

- I want to drive from NYC to Atlanta (via DC).
- Model $Q(\mathbf{w})$ estimate the time cost, e.g., 1000 minutes.

Question: How do I update the model?

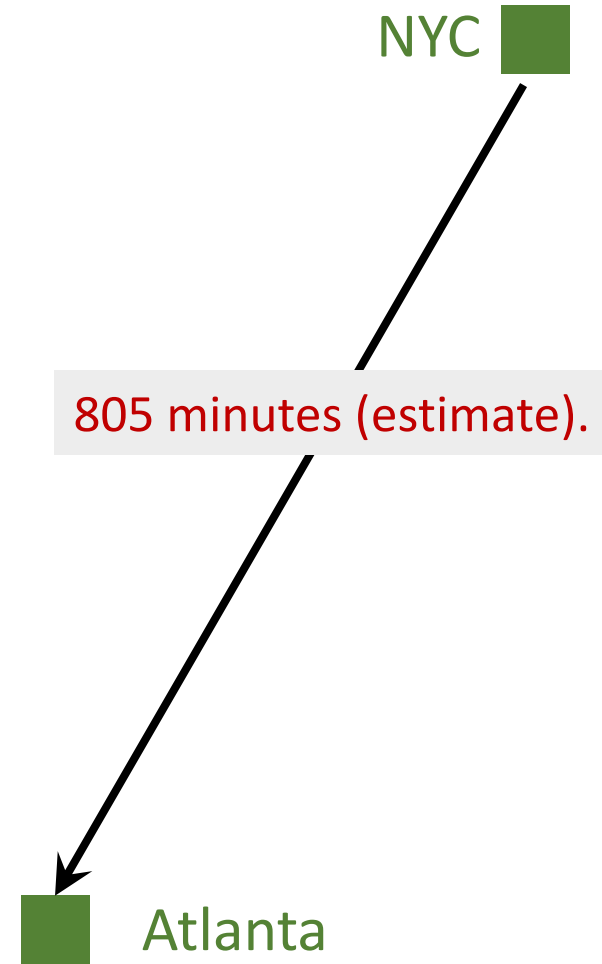
- Can I update the model before finishing the trip?
- Can I get a better \mathbf{w} as soon as I arrived DC?



Temporal Difference (TD) Learning

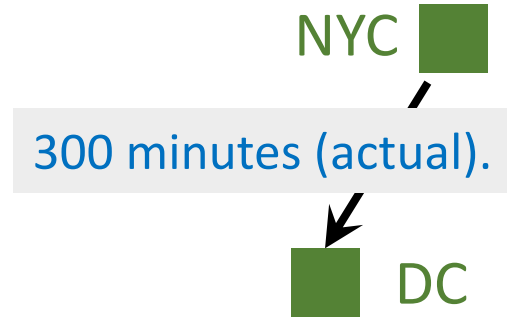
- Model's estimate:

NYC to Atlanta: 805 minutes (estimate).



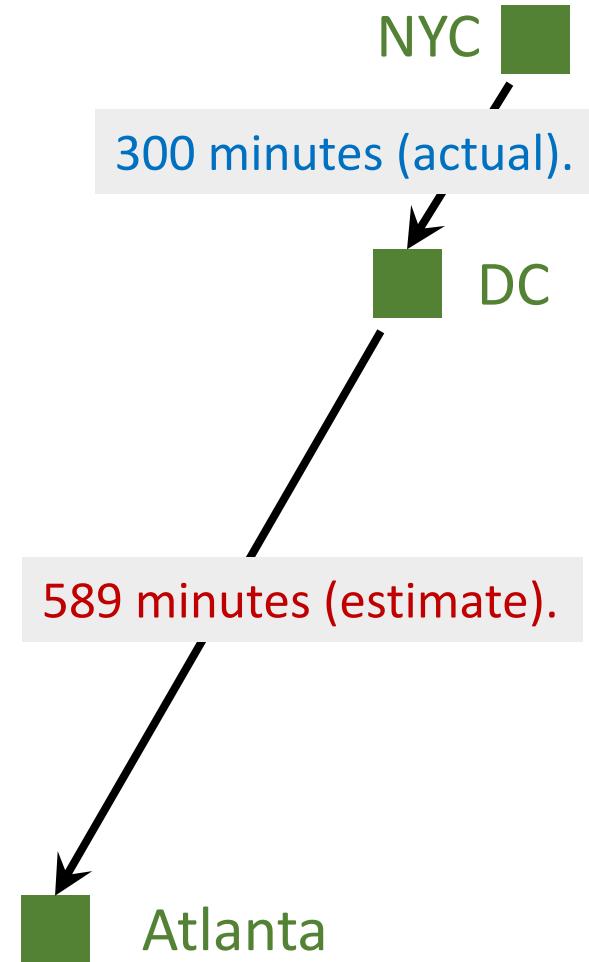
Temporal Difference (TD) Learning

- Model's estimate:
 NYC to Atlanta: 805 minutes (estimate).
- I arrived at DC; actual time cost:
 NYC to DC: 300 minutes (actual).



Temporal Difference (TD) Learning

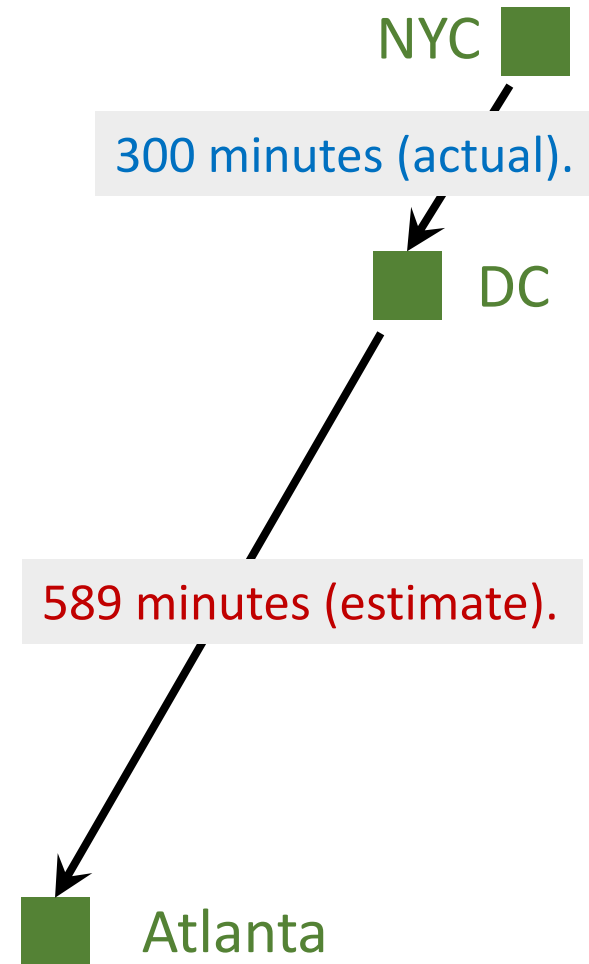
- Model's estimate:
NYC to Atlanta: 805 minutes (estimate).
- I arrived at DC; actual time cost:
NYC to DC: 300 minutes (actual).
- Model now updates its estimate:
DC to Atlanta: 589 minutes (estimate).



Temporal Difference (TD) Learning

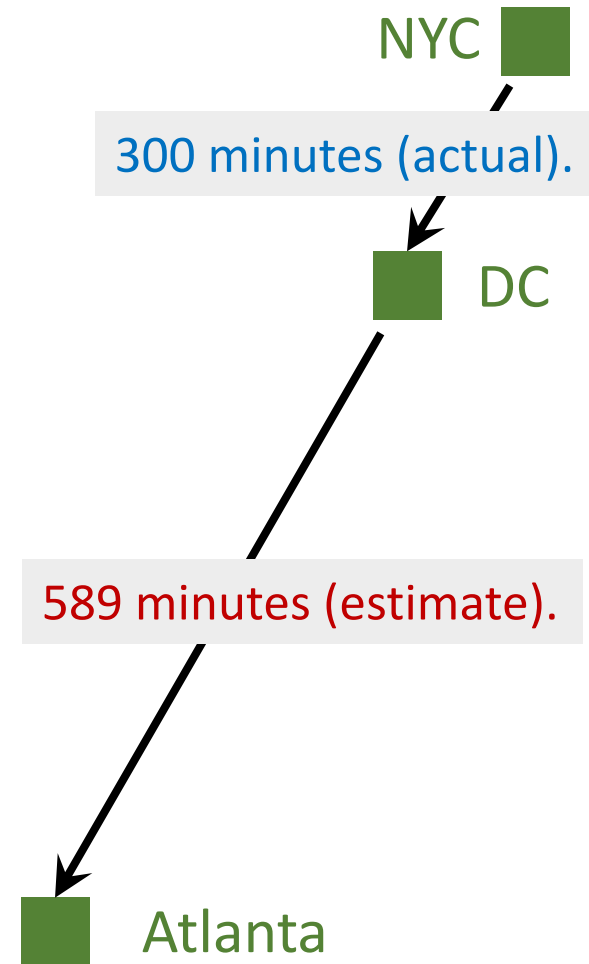
- Model's estimate: $Q(\mathbf{w}) = 805$ minutes .
- Updated estimate: $300 + 589 = 889$ minutes.

TD target.



Temporal Difference (TD) Learning

- Model's estimate: $Q(\mathbf{w}) = 805$ minutes .
- Updated estimate: $300 + 589 = 889$ minutes.
↓
TD target.
- TD target $y = 889$ is a more reliable estimate than 805 .

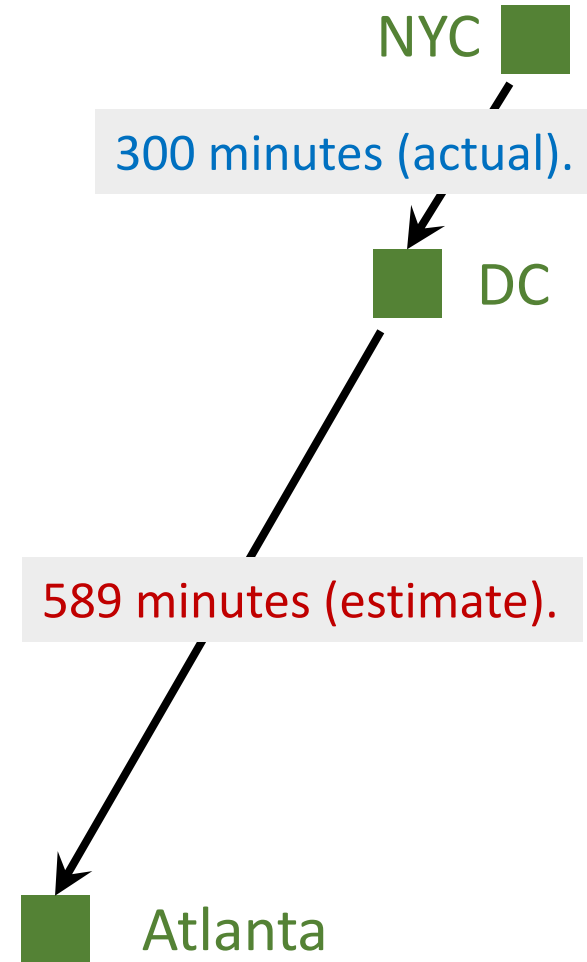


Temporal Difference (TD) Learning

- Model's estimate: $Q(\mathbf{w}) = 805$ minutes .
- Updated estimate: $300 + 589 = 889$ minutes.

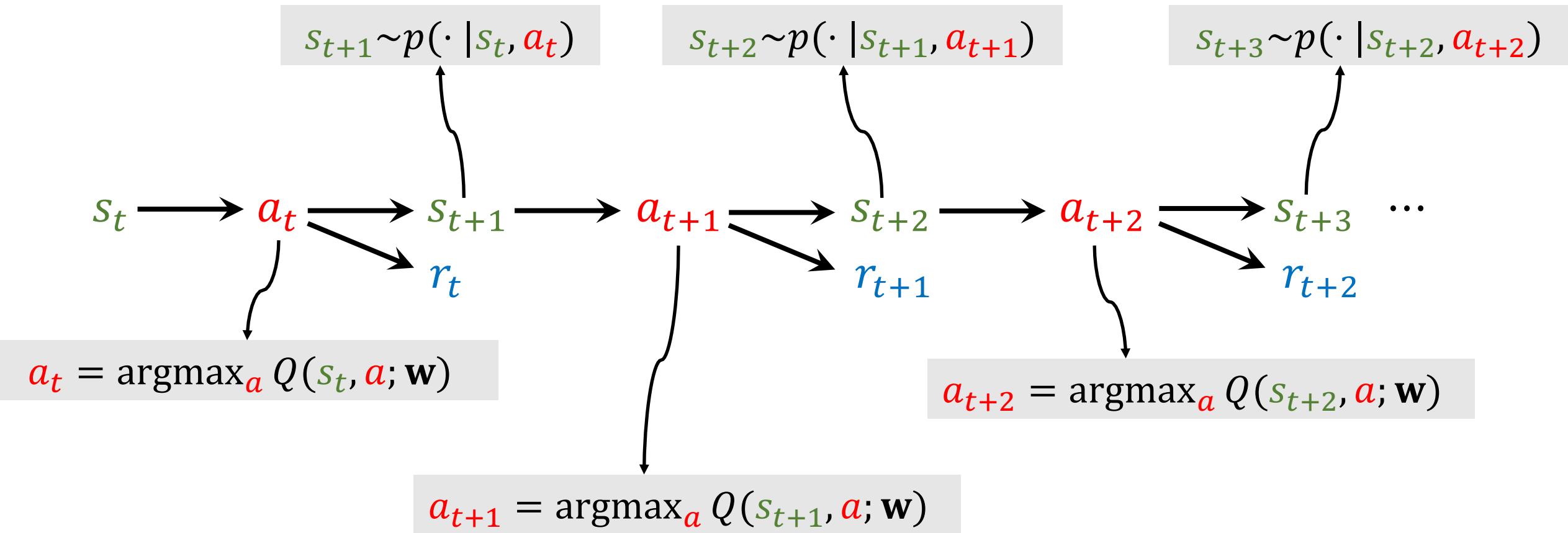
↙
TD target.

- TD target $y = 889$ is a more reliable estimate than 805 .
- Loss: $L = \frac{1}{2} (Q(\mathbf{w}) - y)^2$.
- Gradient: $\frac{\partial L}{\partial \mathbf{w}} = (805 - 889) \cdot \frac{\partial Q(\mathbf{w})}{\partial \mathbf{w}}$.
- Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{\partial L}{\partial \mathbf{w}} \bigg|_{\mathbf{w}=\mathbf{w}_t}$.

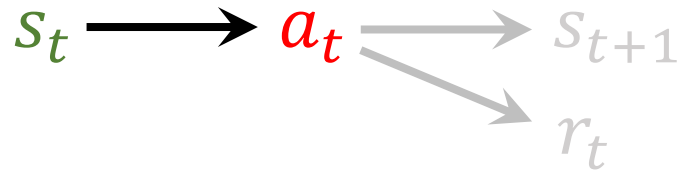


TD Learning for DQN

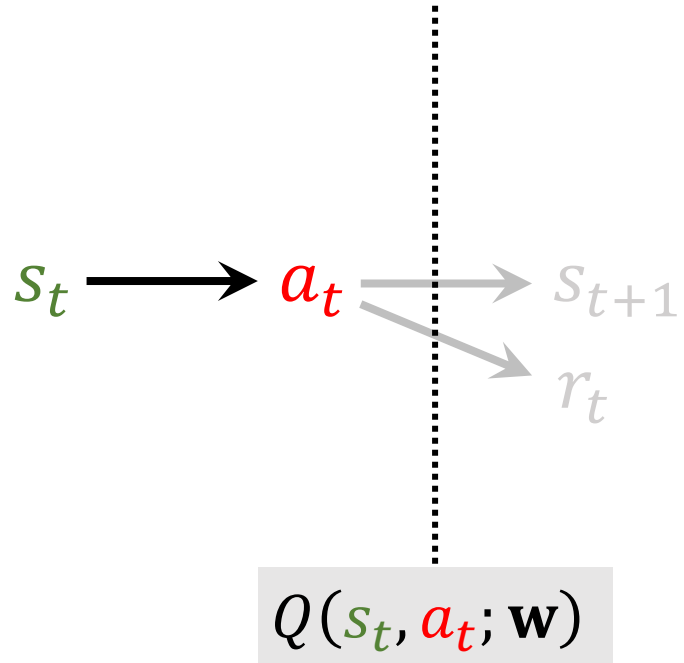
Apply DQN to Play Game



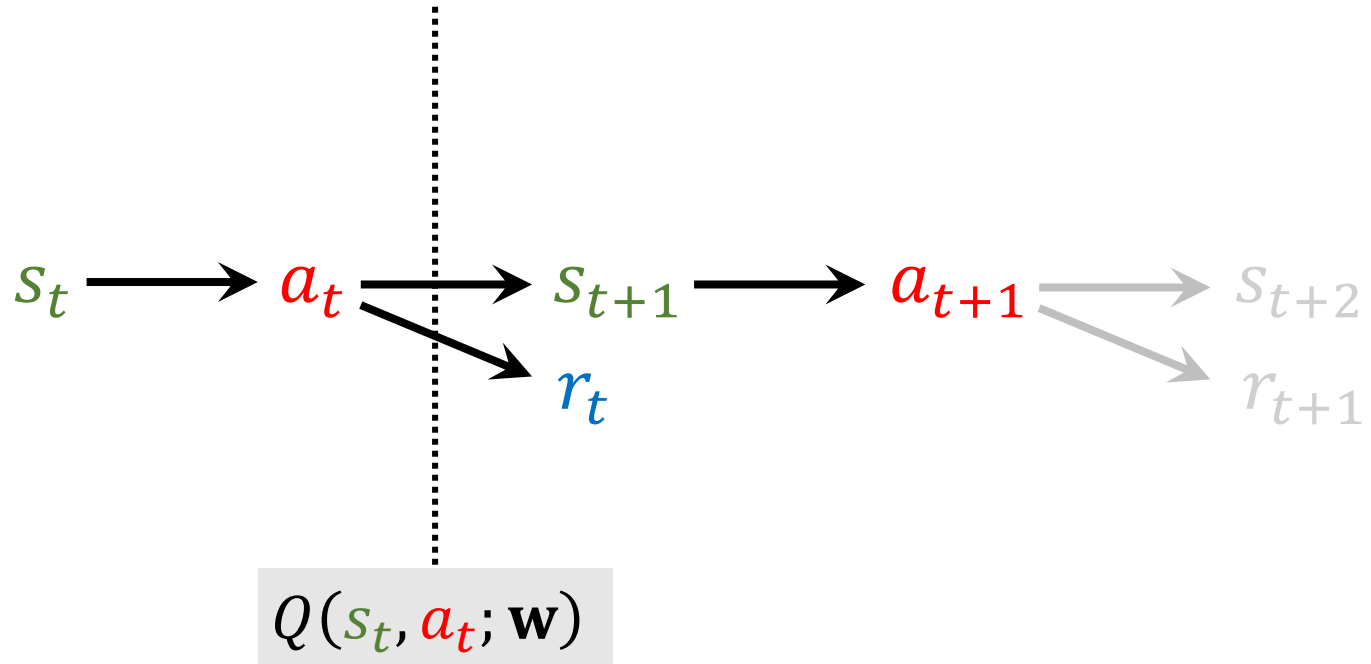
Apply DQN to Play Game



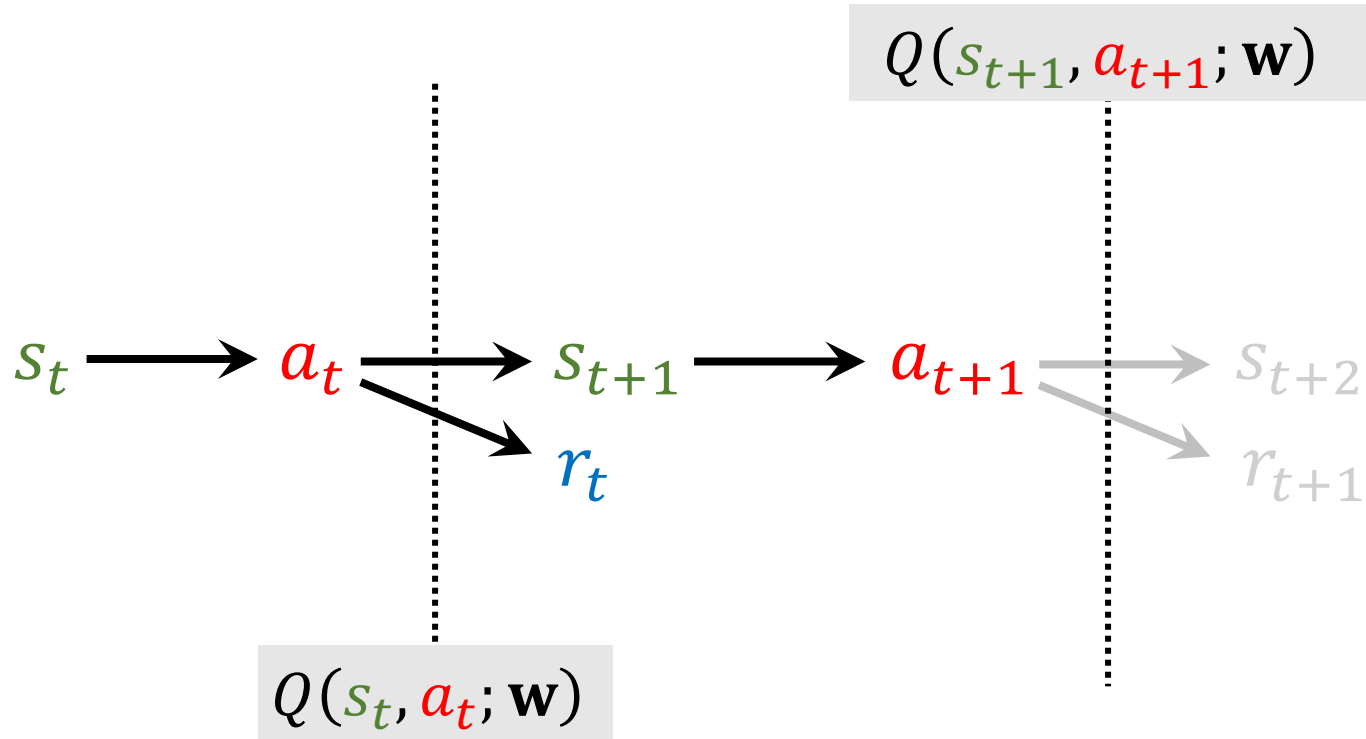
Apply DQN to Play Game



Apply DQN to Play Game

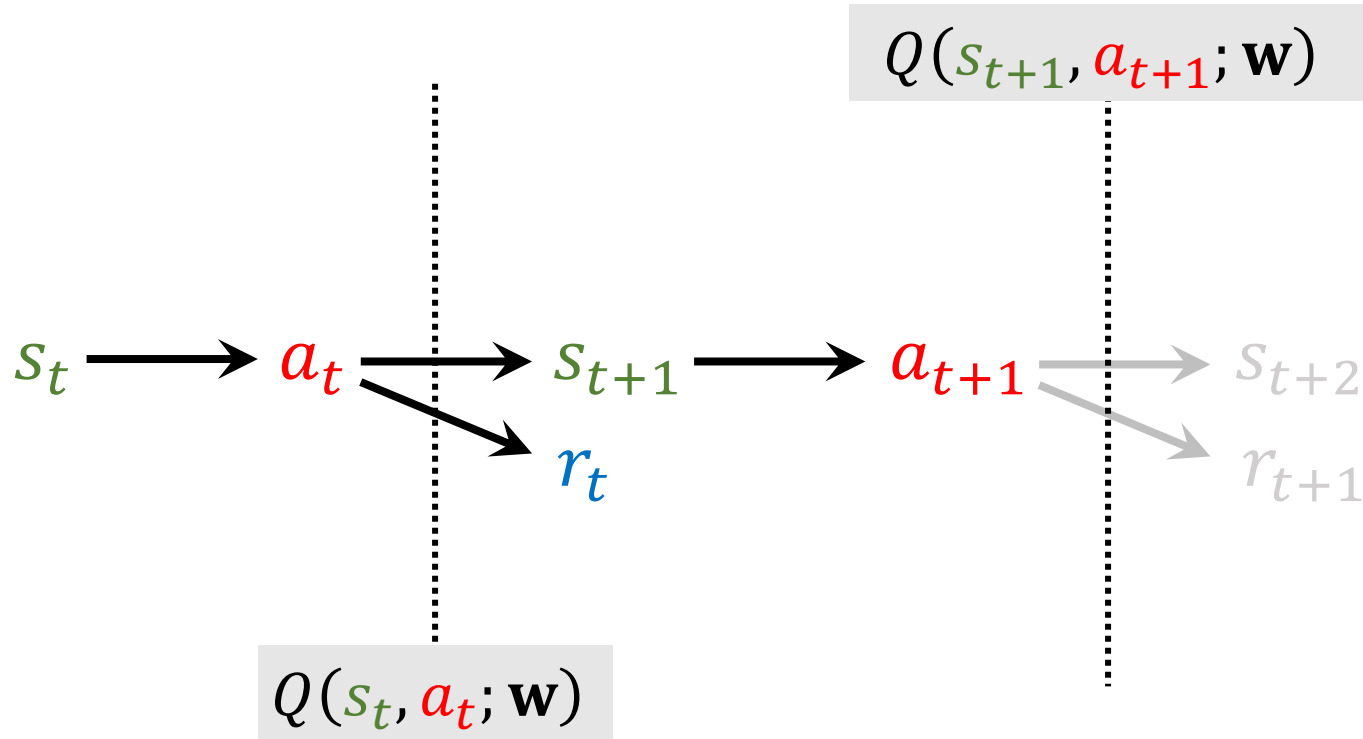


Apply DQN to Play Game



Train DQN using TD learning

If it is accurate estimate, then $Q(s_{t+1}, a_{t+1}; \mathbf{w}) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots]$



If it is accurate estimate, then $Q(s_t, a_t; \mathbf{w}) = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots]$

Train DQN using TD learning

If it is accurate estimate, then $Q(s_{t+1}, a_{t+1}; \mathbf{w}) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots]$



If DQN is accurate estimate, then

$$\begin{aligned} Q(s_t, a_t; \mathbf{w}) &= r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w}) \\ &= r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}) \end{aligned}$$



If it is accurate estimate, then $Q(s_t, a_t; \mathbf{w}) = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots]$

Train DQN using TD learning

If DQN is accurate estimate, then

$$\begin{aligned} Q(s_t, a_t; \mathbf{w}) &= r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w}) \\ &= r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}) \end{aligned}$$

Train DQN using TD learning

If DQN is accurate estimate, then

$$\begin{aligned} Q(s_t, a_t; \mathbf{w}) &= r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w}) \\ &= r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}) \end{aligned}$$

Old estimate (less reliable)

TD target (more reliable estimate of the value)

Train DQN using TD learning

If DQN is accurate estimate, then

$$\begin{aligned} Q(s_t, a_t; \mathbf{w}) &= r_t + \gamma \cdot Q(s_{t+1}, a_{t+1}; \mathbf{w}) \\ &= r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}) \end{aligned}$$

- TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}_t)$.
- Loss: $L_t = \frac{1}{2} [Q(s_t, a_t; \mathbf{w}) - y_t]^2$.
- Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot \frac{\partial L_t}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$

Summary

Action-Value Function Approximation

Definition: Optimal action-value function.

- $Q^*(s, a) = \max_{\pi} \mathbb{E} [R_t | s, a, \pi].$

DQN: Approximate $Q^*(s, a)$ using a neural network.

- $Q(s, a; \mathbf{w})$ is a neural network parameterized by \mathbf{w} .
- Input: observed state s (e.g., a screenshot of game.)
- Output: a vector, each entry of which corresponds to an action a .

Temporal Difference (TD) Learning

Algorithm: One iteration of TD learning.

1. Observe state s_t and action a_t .
2. Predict the value: $q_t = Q(s_t, a_t; \mathbf{w}_t)$.
3. Differentiate the value network: $\mathbf{d}_t = \frac{\partial Q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.

Temporal Difference (TD) Learning

Algorithm: One iteration of TD learning.

1. Observe state s_t and action a_t .
2. Predict the value: $q_t = Q(s_t, a_t; \mathbf{w}_t)$.
3. Differentiate the value network: $\mathbf{d}_t = \frac{\partial Q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}_t}$.
4. Environment provides new state s_{t+1} and reward r_t .
5. Compute TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}_t)$.

Temporal Difference (TD) Learning

Algorithm: One iteration of TD learning.

1. Observe state s_t and action a_t .
2. Predict the value: $q_t = Q(s_t, a_t; \mathbf{w}_t)$.
3. Differentiate the value network: $\mathbf{d}_t = \left. \frac{\partial Q(s_t, a_t; \mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}_t}$.
4. Environment provides new state s_{t+1} and reward r_t .
5. Compute TD target: $y_t = r_t + \gamma \cdot \max_a Q(s_{t+1}, a; \mathbf{w}_t)$.
6. Gradient descent: $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \cdot (q_t - y_t) \cdot \mathbf{d}_t$.