# Lecture-7
# Function approximation

# Function Approximation

## Function approximation

Using function approximation allows us to scale RL to solve realistic problems.

1. We've seen RL finds optimal policies for arbitrary environments, if the value functions $V(s)$ and policies $Q(s, a)$ can be exactly represented in tables

2. But the real world is too large and complex for tables

3. Will RL work for function approximators?

## Example: function approximation

```
Q = np.zeros([n_states, n_actions]) a_p
= Q[s,:]

# action–value table is approximated:
a_p = DeepNeuralNetwork(s)
```
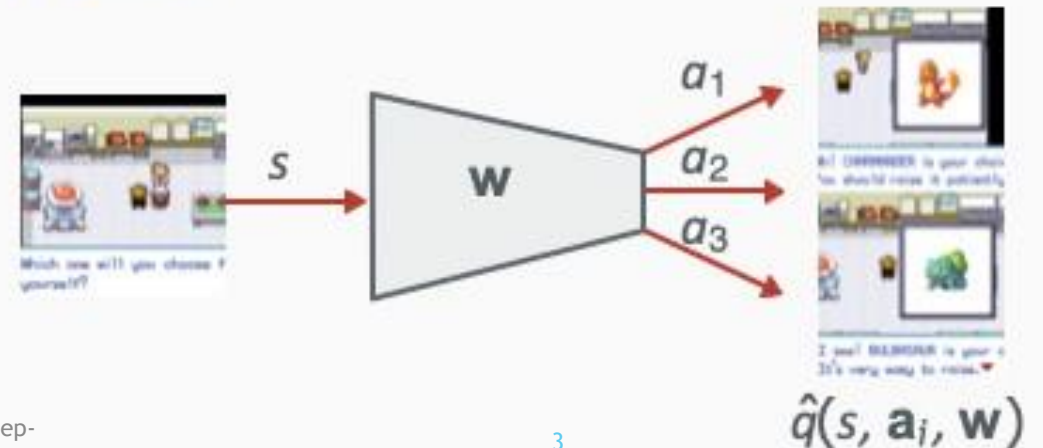
# Function Approximation

**Definition: Function approximation**

There are too many states/actions to fit into memory, which are too slow to process. Therefore we estimate the value function:

$$\hat{v}(S, \mathbf{w}) \approx v_n(S),$$

or for control we'd do:

$$\hat{q}(S, A, \mathbf{w}) \approx q_n(S, A),$$



$$s \longrightarrow \boxed{\mathbf{w}} \longrightarrow \hat{v}(s, \mathbf{w})$$

$$\begin{array}{c} a \\ s \end{array} \longrightarrow \boxed{\mathbf{w}} \longrightarrow \hat{q}(s, a, \mathbf{w})$$

$$s \longrightarrow \boxed{\mathbf{w}} \longrightarrow \begin{array}{c} a_1 \\ a_2 \\ a_3 \end{array} \qquad \hat{q}(s, \mathbf{a}_i, \mathbf{w})$$

# Function Approximation Challenges

**Challenges:** function approximation

However this is not so straightforward:

1. Data is non-stationary

   - As you explore and discover new rewards, the value function can change drastically

2. Data is not i.i.d.

   - When playing a game, all subsequent experiences are highly correlated

**Example:** fog of war

# Incremental methods- stochastic gradient descent for prediction

## Incremental SGD for prediction

Incremental ways to do this, using stochastic gradient descent, to achieve incremental value function approximation.

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_{\mathbf{w}}(v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w}_t))^2$$

$$= \mathbf{w}_t + \alpha(v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w}_t))\nabla_{\mathbf{w}}\hat{v}(S_t, \mathbf{w}_t)$$

How do we compute $v_{\pi}(S_t)$? We substitute it with a target.

# Incremental methods- stochastic gradient descent for prediction

## substituting $V_\pi(S_t)$

For **MC** learning, the target is the return $G_t$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(G_t - \hat{v}(S_t, \mathbf{w}_t))^2$$

For **TD(0)**, the target is $R_{t+1} + \gamma\hat{v}(S_{t+1}, \mathbf{w})$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(R_{t+1} + \gamma\hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}_t))^2$$

For **TD($\lambda$)**, the target is the $\lambda$ return $G_t^\lambda$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(G_t^\lambda - \hat{v}(S_t, \mathbf{w}_t))^2$$

# Incremental methods- stochastic gradient descent for control

**Definition:** action-value function approximation for control

For control, we wish to approximate the action-value function
$$\hat{q}(S, A, \mathbf{w}) \approx q_\pi(S, A)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_{\mathbf{w}}(q_\pi(S_t, A_t) - \hat{q}(S_t, A_t, \mathbf{w}_t))^2$$

$$= \mathbf{w}_t + \alpha(q_\pi(S_t, A_t) - \hat{q}(S_t, A_t, \mathbf{w}_t))\nabla_{\mathbf{w}}\hat{v}(S_t, \mathbf{w}_t)$$

Similarly we substitute $q_\pi(S_t, A_t)$ with a target.

# Incremental methods- stochastic gradient descent for control

**Definition:** substituting $q_\pi(S_t, A_t)$

For **MC** learning, the target is the return $G_t$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(G_t - \hat{q}(S_t, A_t, \mathbf{w}_t))^2$$

For **TD(0)**, the target is $R_{t+1} + \gamma\hat{q}(S_{t+1}, A_{t+1}, \mathbf{w})$:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(R_{t+1} + \gamma\hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w}_t))^2$$

For **TD(λ)**, the target is the λ return:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha\nabla_\mathbf{w}(q_t^\lambda - \hat{q}(S_t, A_t, \mathbf{w}_t))^2$$

# Incremental methods- convergence

In practice, these methods 'chatter' around the near-optimal value function (there's no guarantee when you use function approximation that your improvement step really does improve the policy).

|  | Table | Linear | Non-linear |
|---|---|---|---|
| MC control | ✓ | (✓) | C |
| Sarsa | ✓ | (✓) | C |
| Q-learning | ✓ | C | C |
| Gradient Q-learning | ✓ | ✓ | C |

# Batch learning- experience

Incremental methods have several problems:

1. They are not sample efficient
2. They have strongly correlated updates that break the i.i.d. assumptions of popular SGD algorithms
3. They may rapidly forget rare experiences that would be useful later on

**Experience replay** and prioritized experience replay address these issues by storing experiences and reusing them. These can be sampled uniformly or prioritized to replay important transitions more frequently.

# Batch learning- experience replay with double Q-learning

Putting this all together, we have:

- Sample action from our $\epsilon$-greedy policy or with GLIE

- Store $\langle s_t, a_t, r_{t+1}, s_{t+1} \rangle$ in the experience replay buffer $\mathcal{D}$

- Sample a random minibatch $s, a, r, s' \sim \mathcal{D}$

- Update in the gradient direction from our Q-network towards our Q-targets (with the max).

  - The targets can be frozen [5, 6]

$$\mathcal{L}_i(\mathbf{w}_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \mathbf{w}_i^-) - Q(s, a; \mathbf{w}_i) \right)^2 \right]$$

# Overview of Deep Q Networks

▶ Q Learning builds a Q-table of State-Action values, with dimension *(s, a)*, where *s* is the number of states and *a* is the number of actions. Fundamentally, a Q-table maps state and action pairs to a Q-value.



Q Learning looks up state-action pairs in a Q table

▶ However, in a real-world scenario, the number of states could be huge, making it computationally intractable to build a table.

# Use a Q-Function for real-world problems

▶ To address this limitation we use a Q-function rather than a Q-table, which achieves the same result of mapping state and action pairs to a Q value.



A state-action function is required to handle real-world scenarios with a large state space.

# Neural Networks are the best Function Approximators

▶ Since neural networks are excellent at modeling complex functions, we can use a neural network, which we call a Deep Q Network, to estimate this Q function.

▶ This function maps a state to the Q values of all the actions that can be taken from that state.



State $S$ ⟶ | Deep Q Network $f(S)$ Weights | ⟶ Q-Value $(S, a_1)$
⟶ Q-Value $(S, a_2)$
...
⟶ Q-Value $(S, a_n)$

▶ It learns the network's parameters (weights) such that it can output the Optimal Q values

# DQN Architecture Components

# High-level DQN Workflow

16

# Gather Training Data

Experience Replay selects an ε-greedy action from the current state, executes it in the environment, and gets back a reward and the next state.

# It saves this observation as a sample of training data.
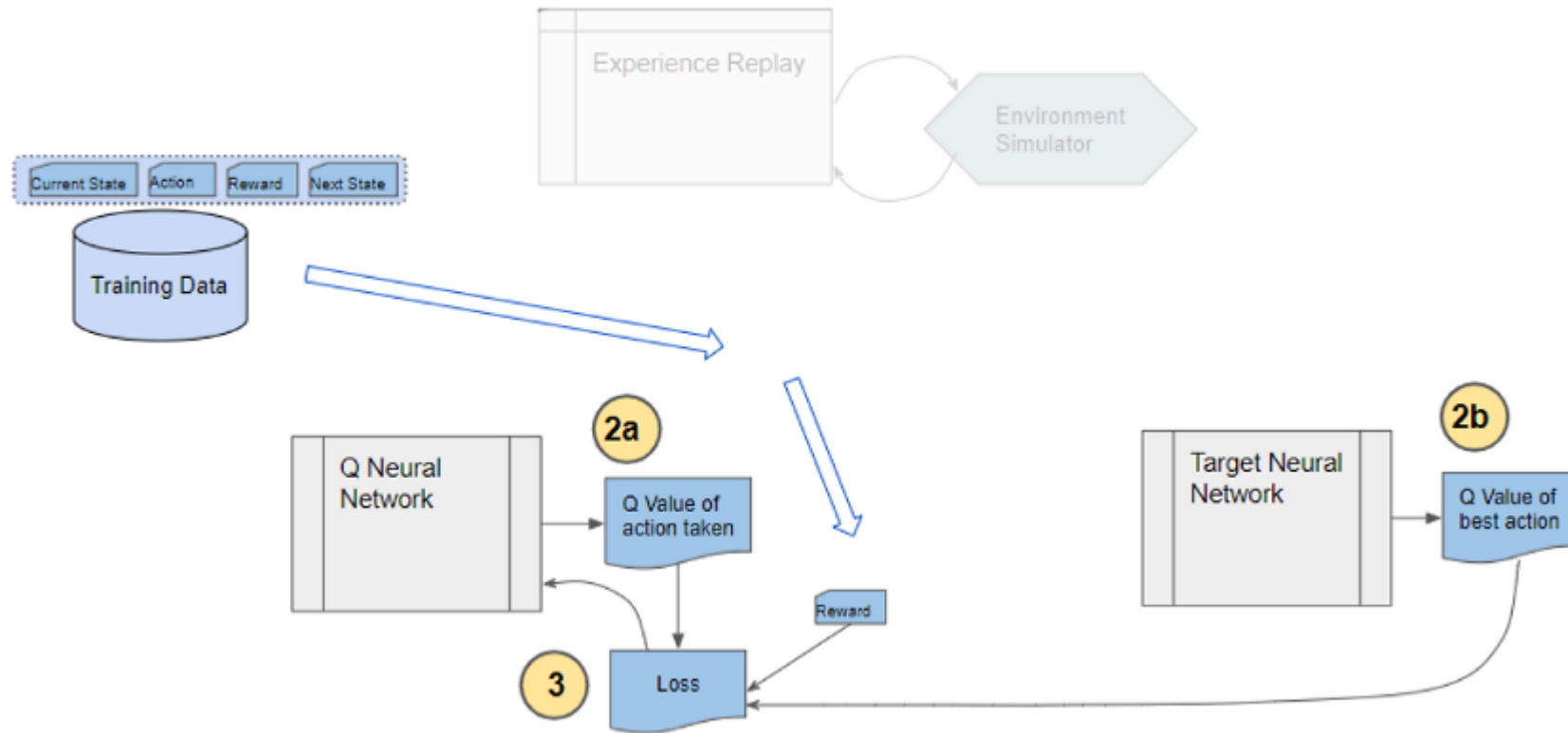
# Let's see the next phase of the flow.

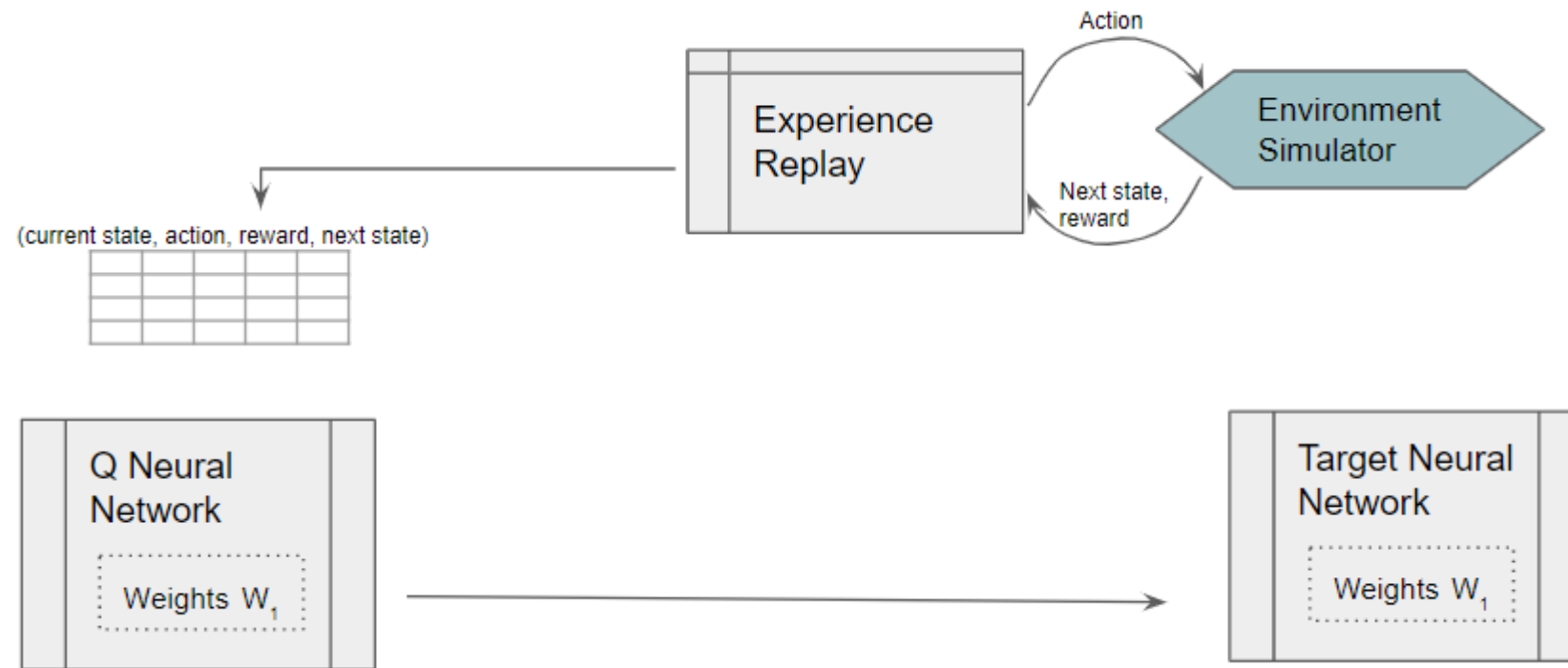# Q Network predicts Q-value

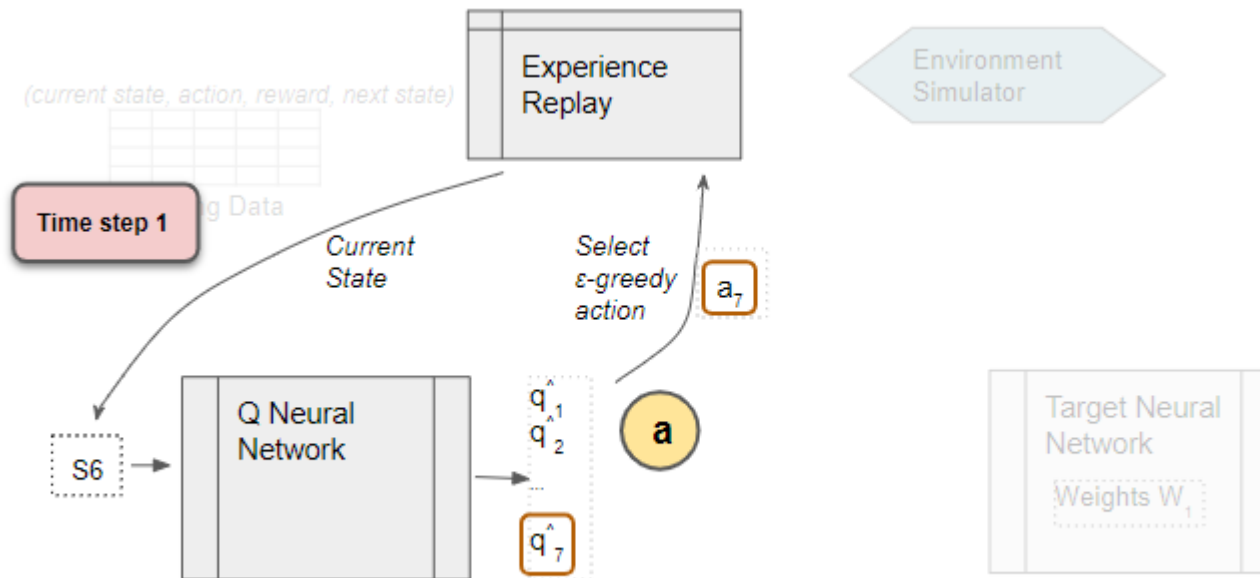# Target Network predicts Target Q-value

# Compute Loss and Train Q Network

# Why do we need Experience Replay?

▶ All of the actions and observations that the agent has taken from the beginning (limited by the capacity of the memory, of course) are stored.

▶ Then a batch of samples is randomly selected from this memory.

▶ This ensures that the batch is 'shuffled' and contains enough diversity from older and newer samples (eg. from several regions of the factory floor and under different conditions) to allow the network to learn weights that generalize to all the scenarios that it will be required to handle.
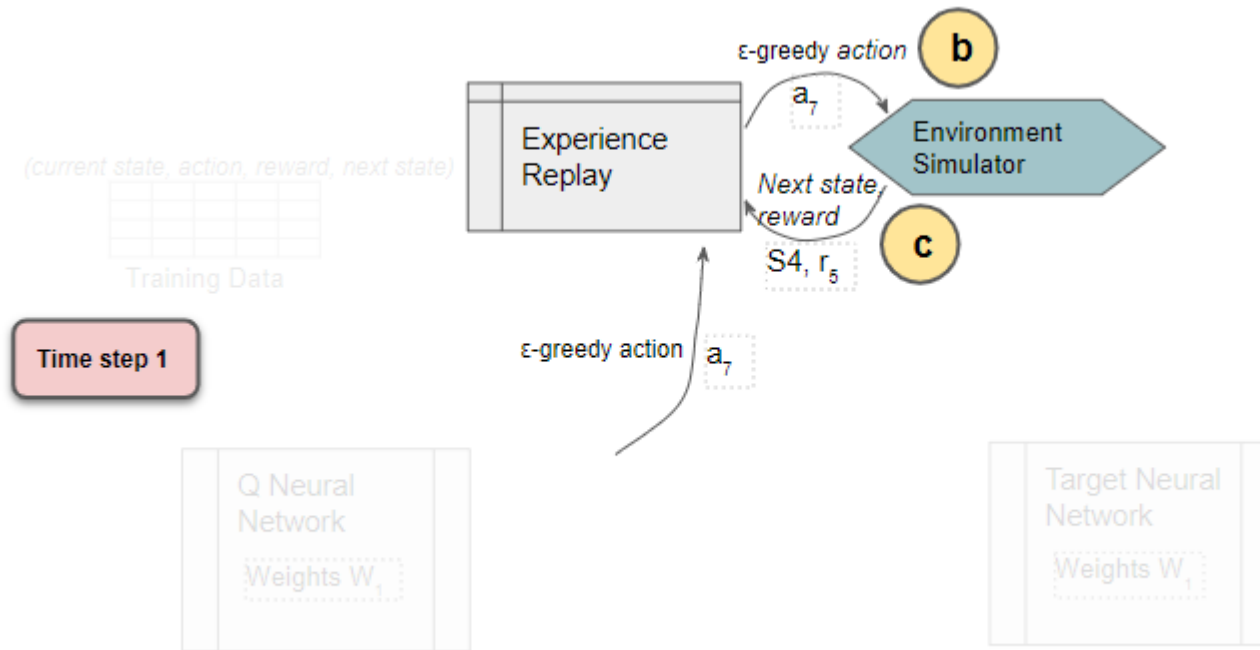
https://towardsdatascience.com/reinforcement-learning-explained-visually-part-5-deep-q-networks-step-by-step-5a5317197f4b

24

# Overall Flow--1. Initialization



Action

Experience Replay

Environment Simulator

Next state, reward

(current state, action, reward, next state)

Q Neural Network

Weights $W_1$

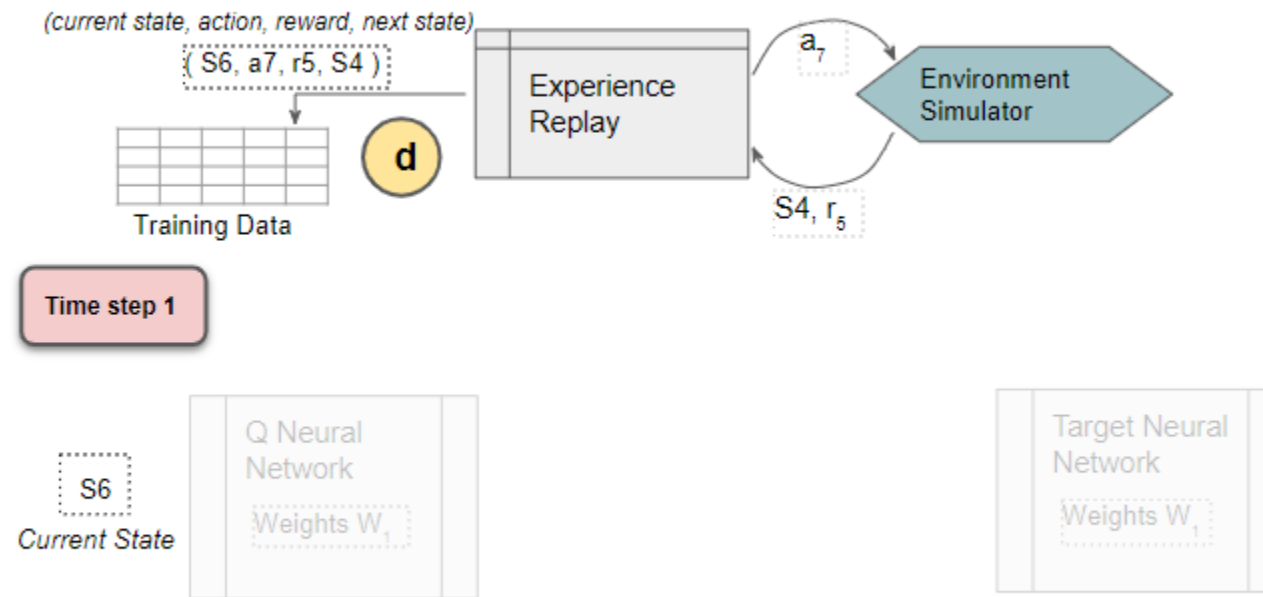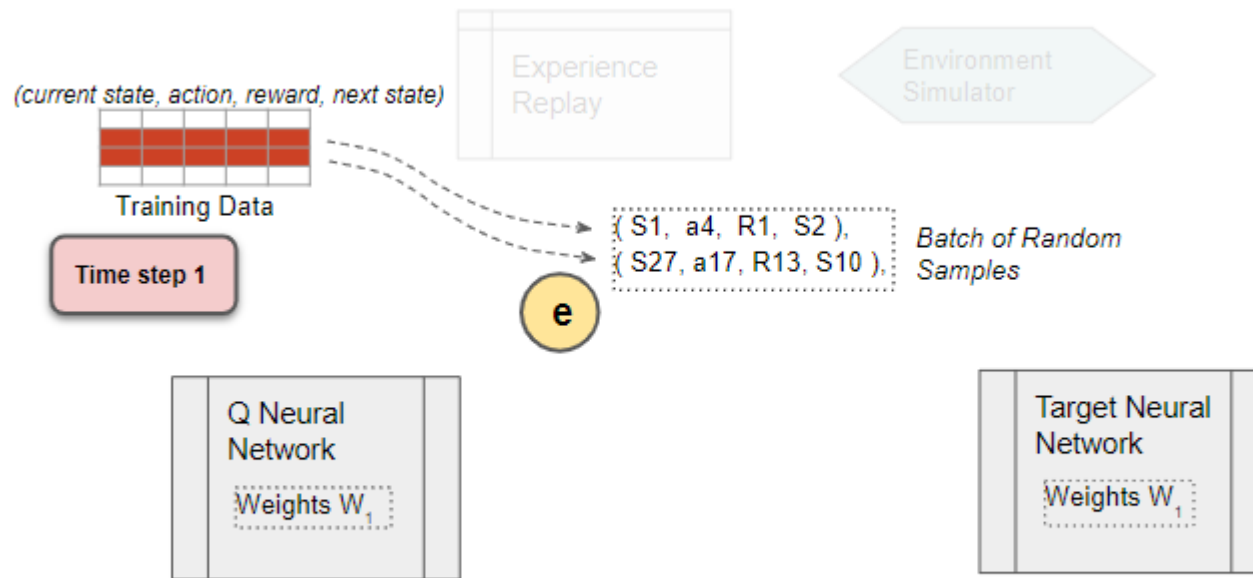Target Neural Network

Weights $W_1$

# 2. Experience Replay

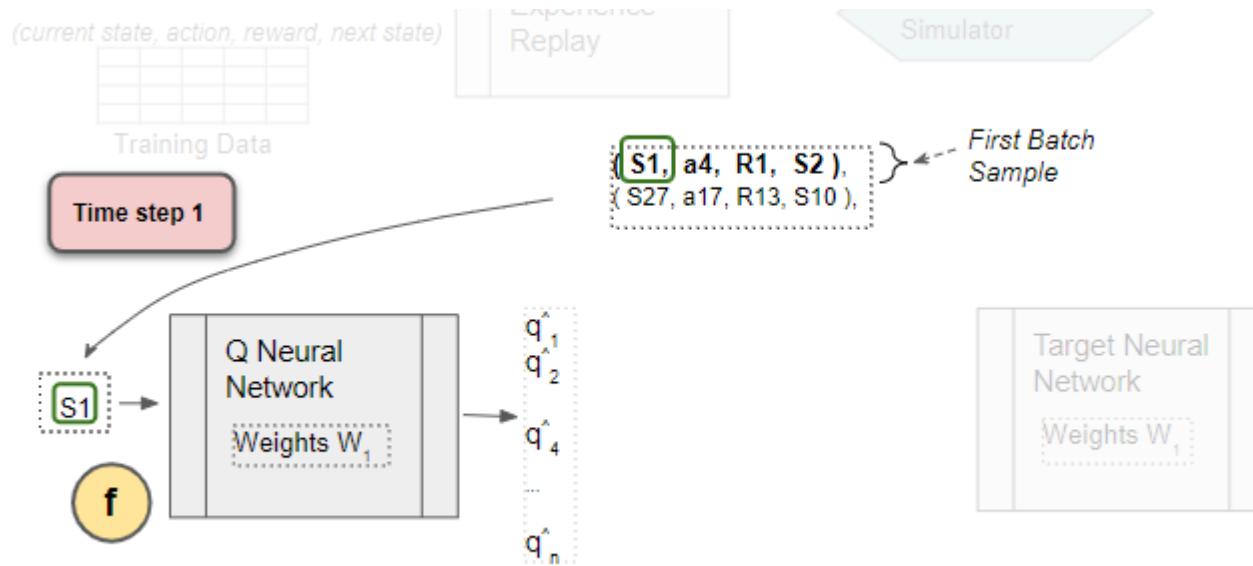# Experience Replay executes the ε-greedy action and receives the next state and reward

# Each such result is a sample observation which will later be used as training data.
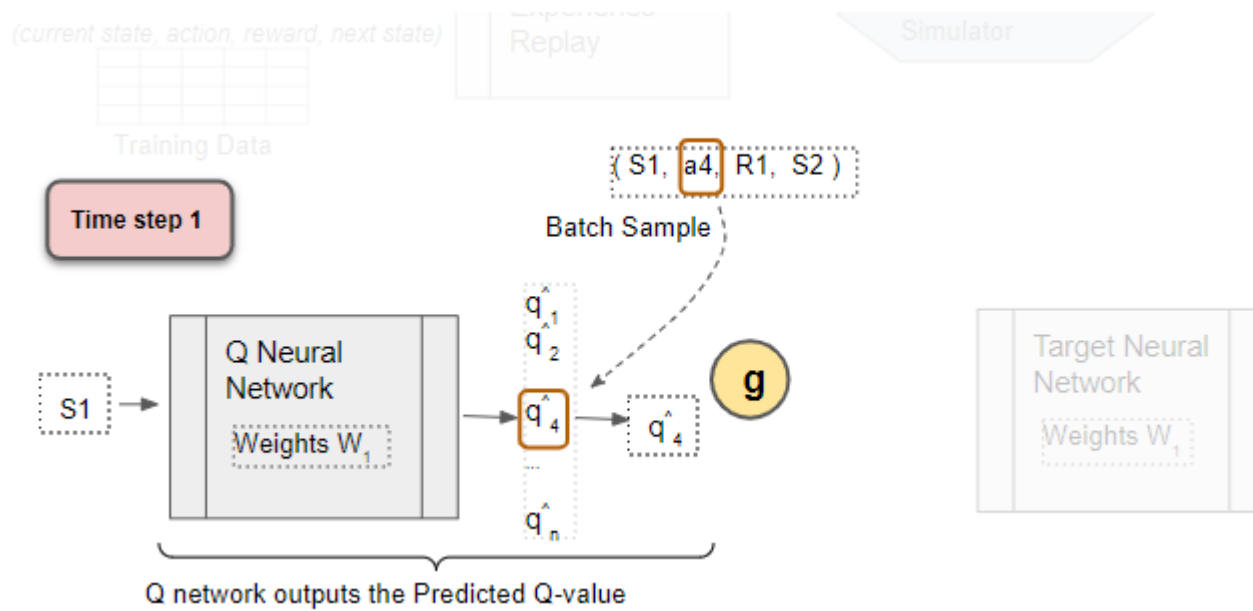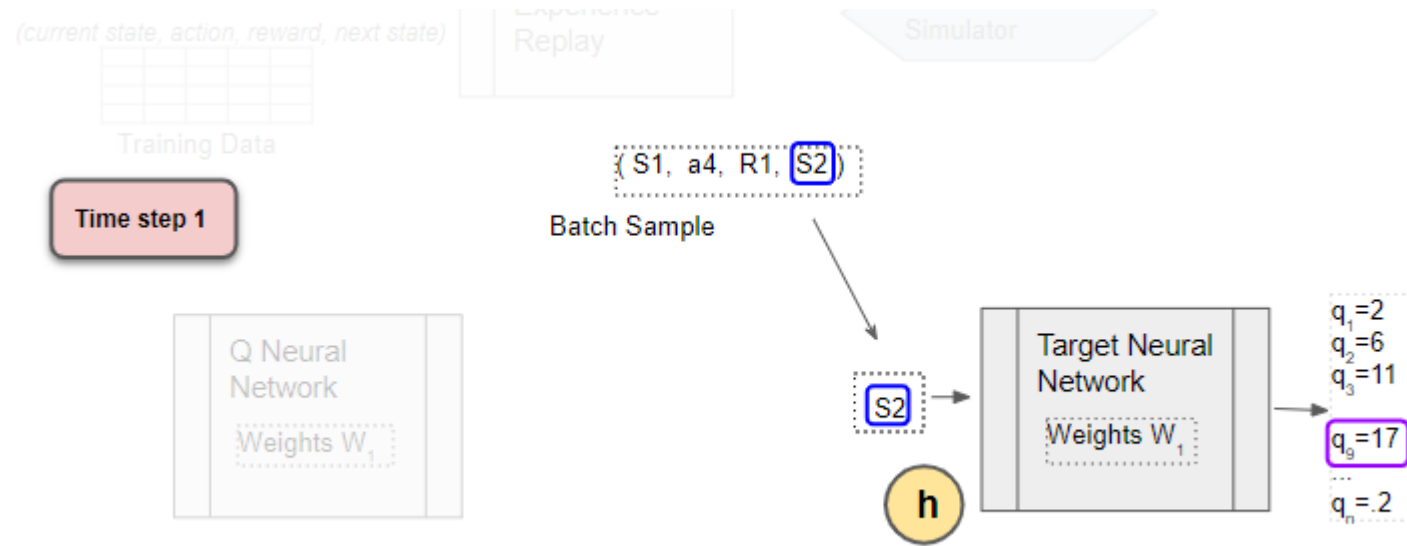
# 3. Select random training batch

(current state, action, reward, next state)

| | | | | |
|---|---|---|---|---|
| | | | | |

Training Data

Time step 1

Experience Replay

Environment Simulator

( S1, a4, R1, S2 ),
( S27, a17, R13, S10 ),

*Batch of Random Samples*

e

Q Neural Network

Weights $W_1$

Target Neural Network

Weights $W_1$

# Use the current state from the sample as input to predict the Q values for all actions

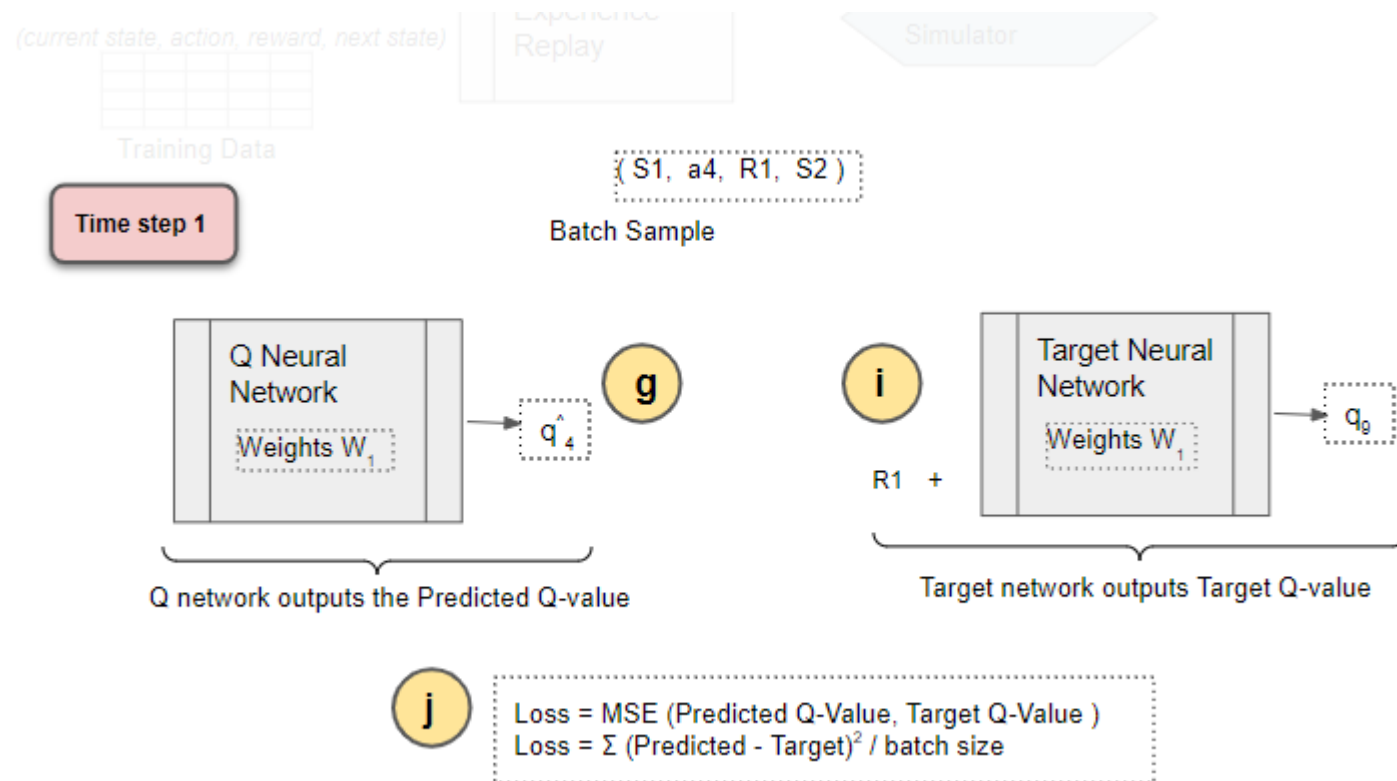# 4. Select the Predicted Q-value

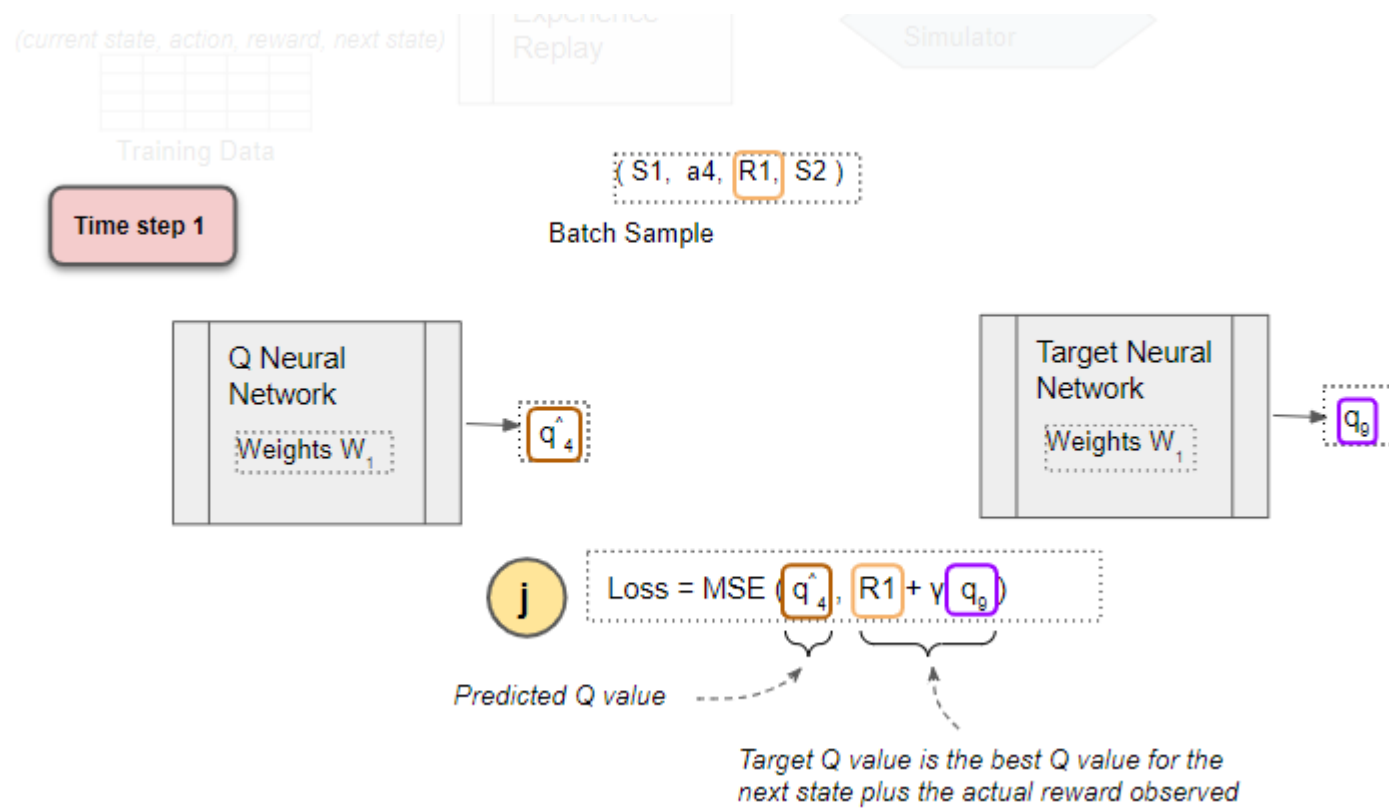# Use the next state from the sample as input to the Target network
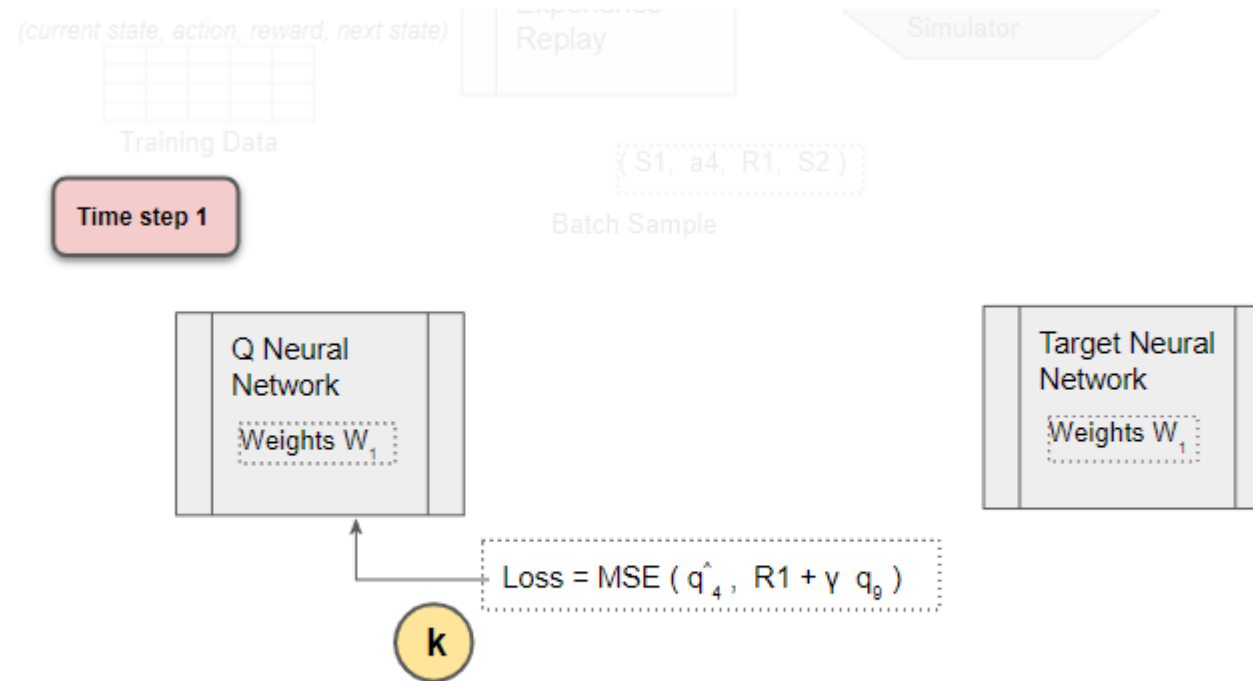
# 5. Get the Target Q Value
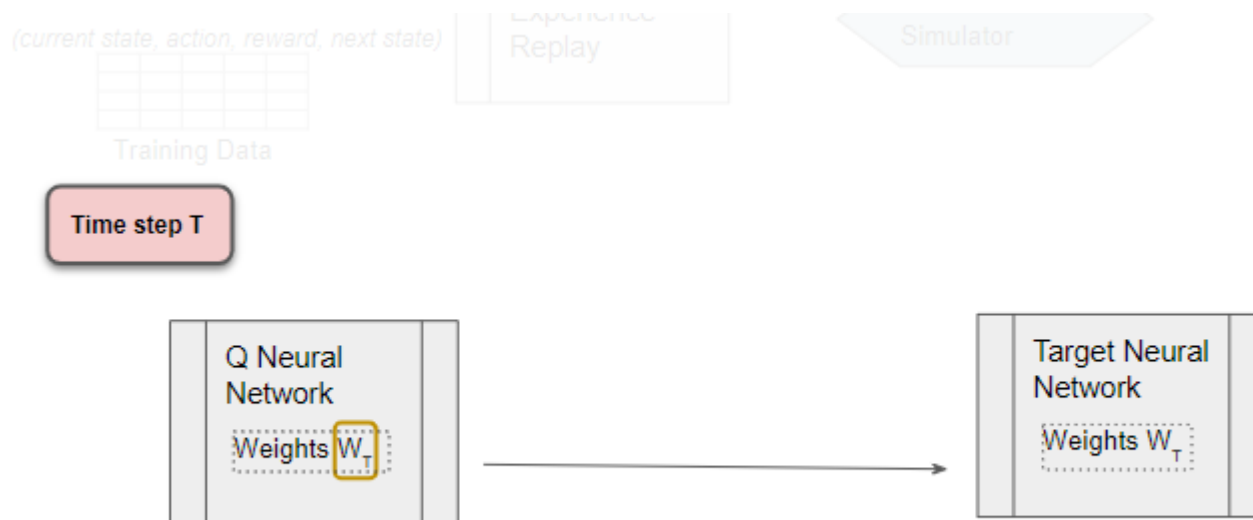
# 6. Compute Loss

# 7. Back-propagate Loss to Q-Network

# 8. After T time-steps, copy Q Network weights to Target Network

# Summary

- The DQN works in a similar way as Q-Learning.

- Since it is a neural network, it uses a Loss function rather than an equation.

- It also uses the Predicted (ie. Current) Q Value, Target Q Value, and observed reward to compute the Loss to train the network, and thus improve its predictions.