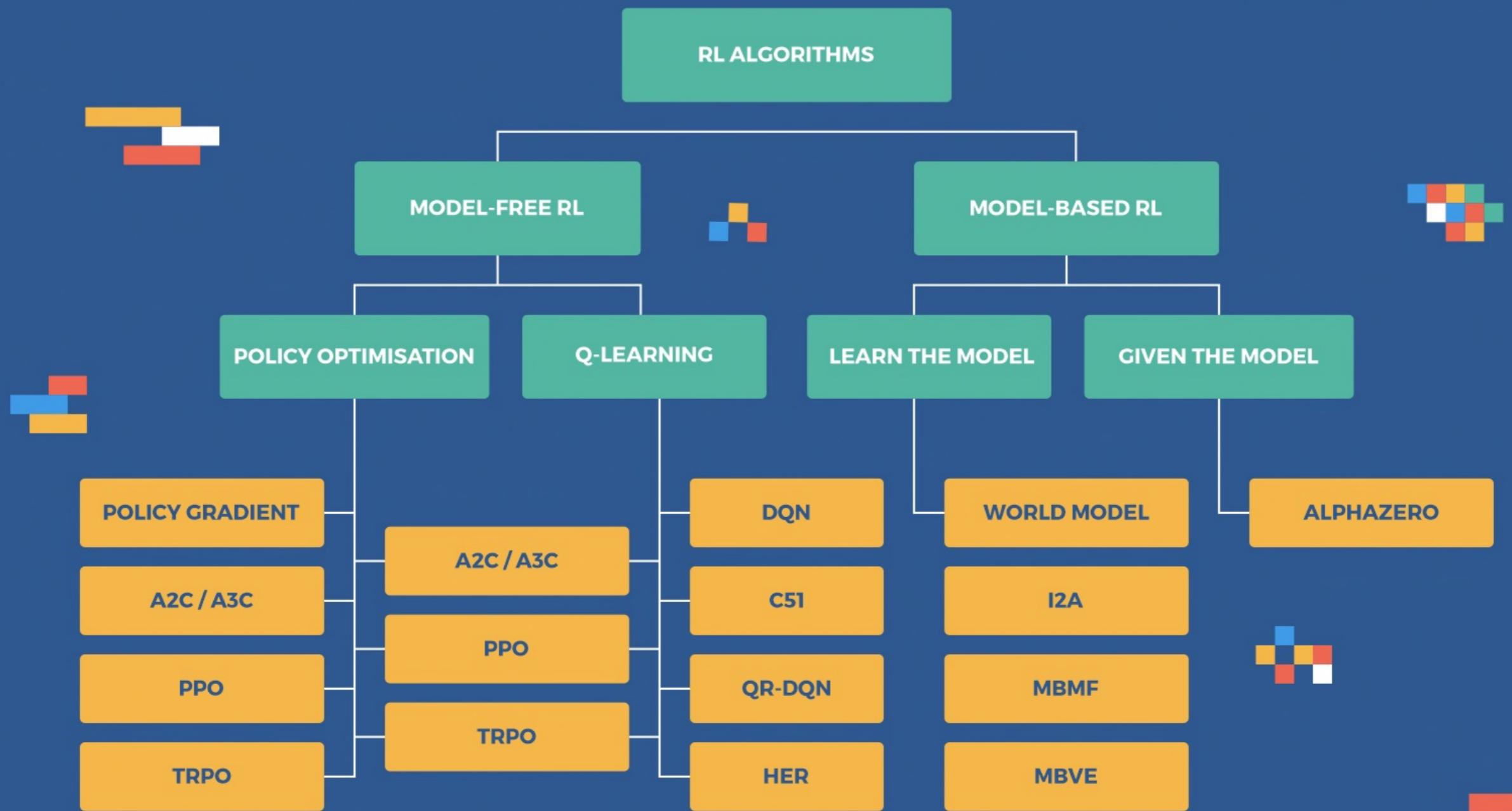


# REINFORCEMENT LEARNING ALGORITHMS MAP



## Fundamentals: Q-Learning

**Initialization (First iteration):**  
For all couples of states  $s$  and actions  $a$ , the Q-values are initialized to 0.

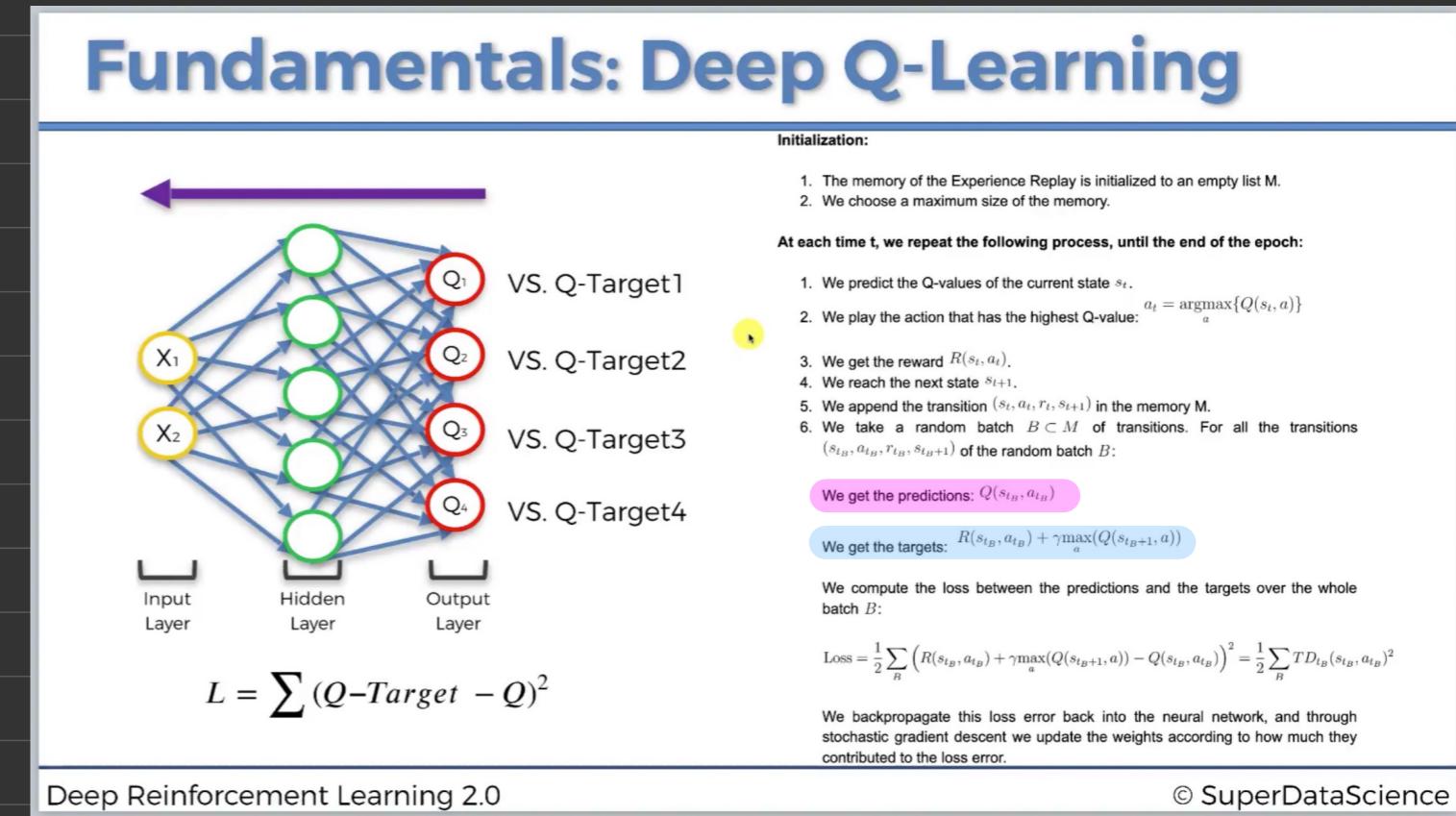
**Next iterations:**  
At each iteration  $t \geq 1$ , we repeat the following steps:

1. We select a random state  $s_t$  from the possible states.
2. From that state, we play a random action  $a_t$ .
3. We reach the next state  $s_{t+1}$  and we get the reward  $R(s_t, a_t)$ .
4. We compute the Temporal Difference  $TD_t(s_t, a_t)$ :
$$TD_t(s_t, a_t) = R(s_t, a_t) + \gamma \max_a(Q(s_{t+1}, a)) - Q(s_t, a_t)$$
5. We update the Q-value by applying the Bellman equation:
$$Q_t(s_t, a_t) = Q_{t-1}(s_t, a_t) + \alpha TD_t(s_t, a_t)$$

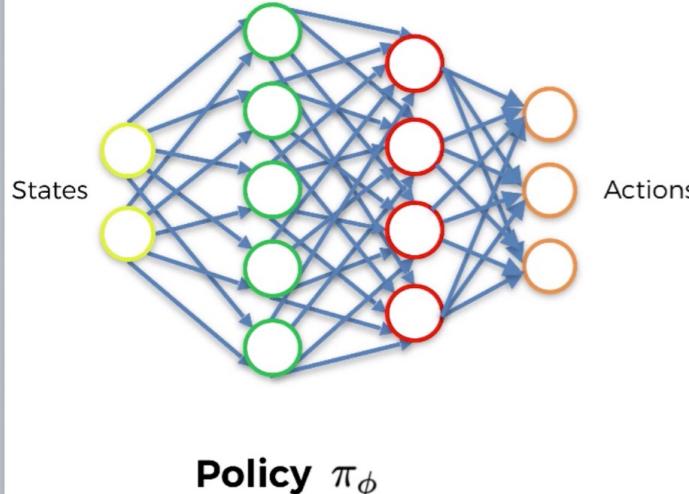
Deep Reinforcement Learning 2.0 © SuperDataScience

It will learn the Q values which is based on Q learning and it will

- TD3 will learn two values at same time
  - The Q-value
  - the parameters to Policy, so that make better choices
- The Q-value represents the quality of being in a state s and take an action a.
- The Q-value of being in state A and take the action Right will be the highest of all others, because it move the agent to the goal.
- In contrast the Q-value of being in state B and took action Right will be the lowest, because the agent will loose the game.



## Fundamentals: Policy Gradient



$$\text{Return: } R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$$

$$\text{Goal: Maximize the expected return } J(\phi) = \mathbb{E}_{s_i \sim p_\pi, a_i \sim \pi} [R_0]$$

### Policy Gradient:

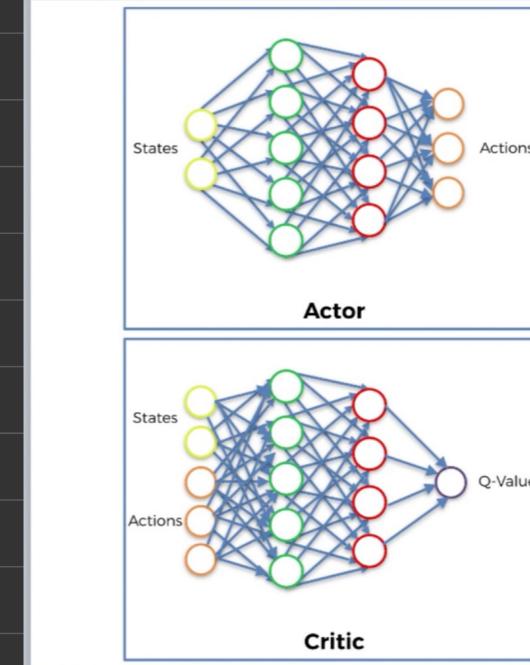
- We compute the gradient of the expected return with respect to the parameters  $\phi$   
 $\nabla_\phi J(\phi)$
- We update the parameters through gradient ascent:

$$\phi_{t+1} = \phi_t + \alpha \nabla_\phi J(\pi_\phi)|_{\phi_t}$$

Deep Reinforcement Learning 2.0

© SuperDataScience

## Fundamentals: Actor-Critic



$$\text{Return: } R_t = \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i)$$

$$\text{Goal: Maximize the expected return } J(\phi) = \mathbb{E}_{s_i \sim p_\pi, a_i \sim \pi} [R_0]$$

### Deterministic Policy Gradient:

- In actor-critic methods, the policy, known as the actor, can be updated through the deterministic policy gradient algorithm:  
 $\nabla_\phi J(\phi) = \mathbb{E}_{s \sim p_\pi} [\nabla_a Q^\pi(s, a)|_{a=\pi(s)} \nabla_\phi \pi_\phi(s)]$
- We update the policy parameters through gradient ascent:

$$\phi_{t+1} = \phi_t + \alpha \nabla_\phi J(\pi_\phi)|_{\phi_t}$$

Deep Reinforcement Learning 2.0

© SuperDataScience

- In the previous model we try to estimate the Q-value.
- Given the Q-value the Agent apply some policy to choose the best action.
- The policy may be pick the action that gives the higher reward  
 $\text{argmax}_a [Q(s, a)]$
- Now this kind of algo try to learn the best policy, not the best Q-value.

- In this model we have two networks
  1. Actor "Policy" get the states and output the possible actions
  2. The critic that get States, Actions and produces the Q-value
- Actor try to find the best action.
- Critic try to find the best Q-values.

# Fundamentals: Taxonomy of AI models

Model-Free vs. Model-Based

Value-Based vs. Policy Based

Off-Policy vs. On-Policy

Deep Reinforcement Learning 2.0

© SuperDataScience

## Model Free

- Do not try to learn how the env work.
- Just take the env output, do some process and choose the action.
- Q-Learning, Actor-Critic and most of the Algo.

## Model Based

- Try to understand the env and create a model to represent it
- The model is trying to capture the transition func and the reward func.

If, after learning the Agent can make predictions about what the next state/reward is [before] it do the action than it's a Model-Based, otherwise if it can't predict what next state/reward will, it's a Model Free

## Value Based

- Agent calculates the Q-value and choose the action that results in the highest reward.
- the agent doesn't learn the action, it learn the value, and the policy will be choose the highest action value.  $\text{Argmax}_a [Q(s,a)]$

## Policy Based

- Instead of predict the Q-value to later choose the action, It predict the Action directly from observed state.

## Off-Policy

- use a experience Replay Memory to store past transitions so that it can sample and learn from the history

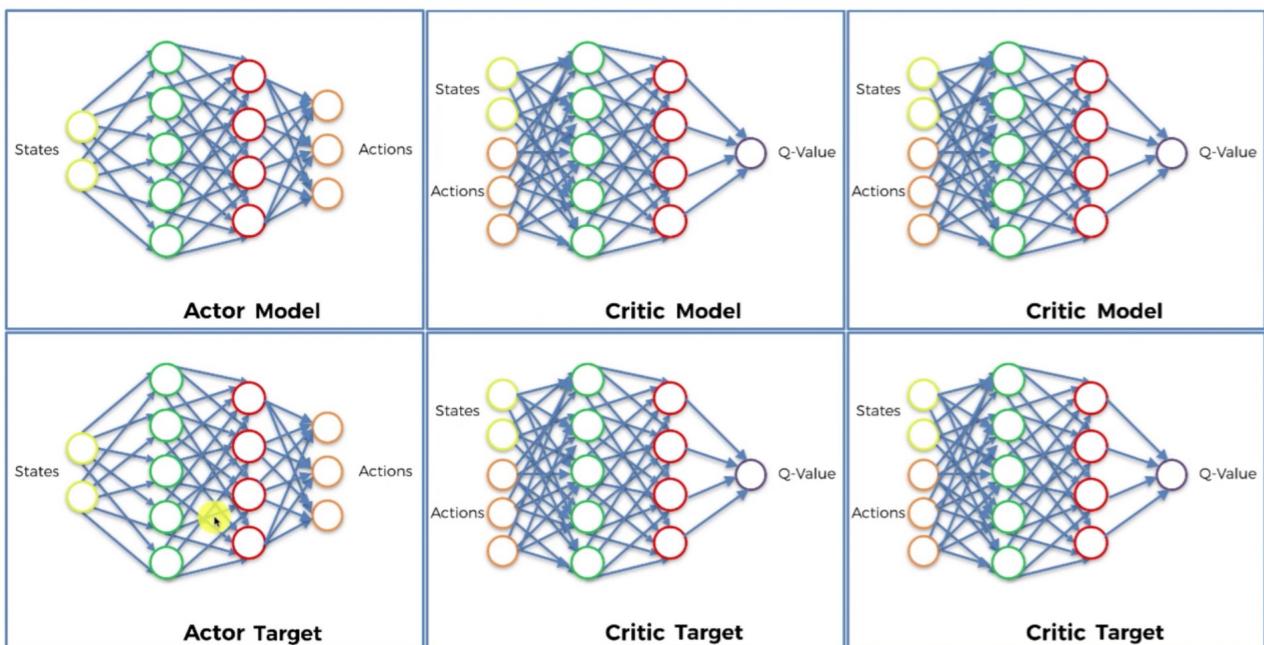
## On-Policy

- Do not store past transitions. Learn just from new observations.

	Deep Q-Learning	Policy Gradient	Actor-Critic	World Models
Model-Free vs. Model-Based	Model-Free	Model-Free	Model-Free	Model-Based
Value-Based vs. Policy-Based	Value-Based	Policy-Based	Both	Policy-Based
Off-Policy vs. On-Policy	Off-Policy	On-Policy	Off-Policy	On-Policy

# Q-learning Step

## Twin-Delayed DDPG (TD3)



Deep Reinforcement Learning 2.0

© SuperDataScience

- This algo produces continuous outputs, instead of distinct output actions as Q-learning
- Remember the max [Q(s,a)] in Q-learning. This only works in discrete number of actions  $a$
- The Actor predicts the Policy
- The Critic predicts the Q-values

### Initialization:

Step 1: We initialize the Experience Replay memory, with a size of 20000. We will populate it with each new transition.  
 Step 2: We build one neural network for the Actor model and one neural network for the Actor target.  
 Step 3: We build two neural networks for the two Critic models and two neural networks for the two Critic targets.

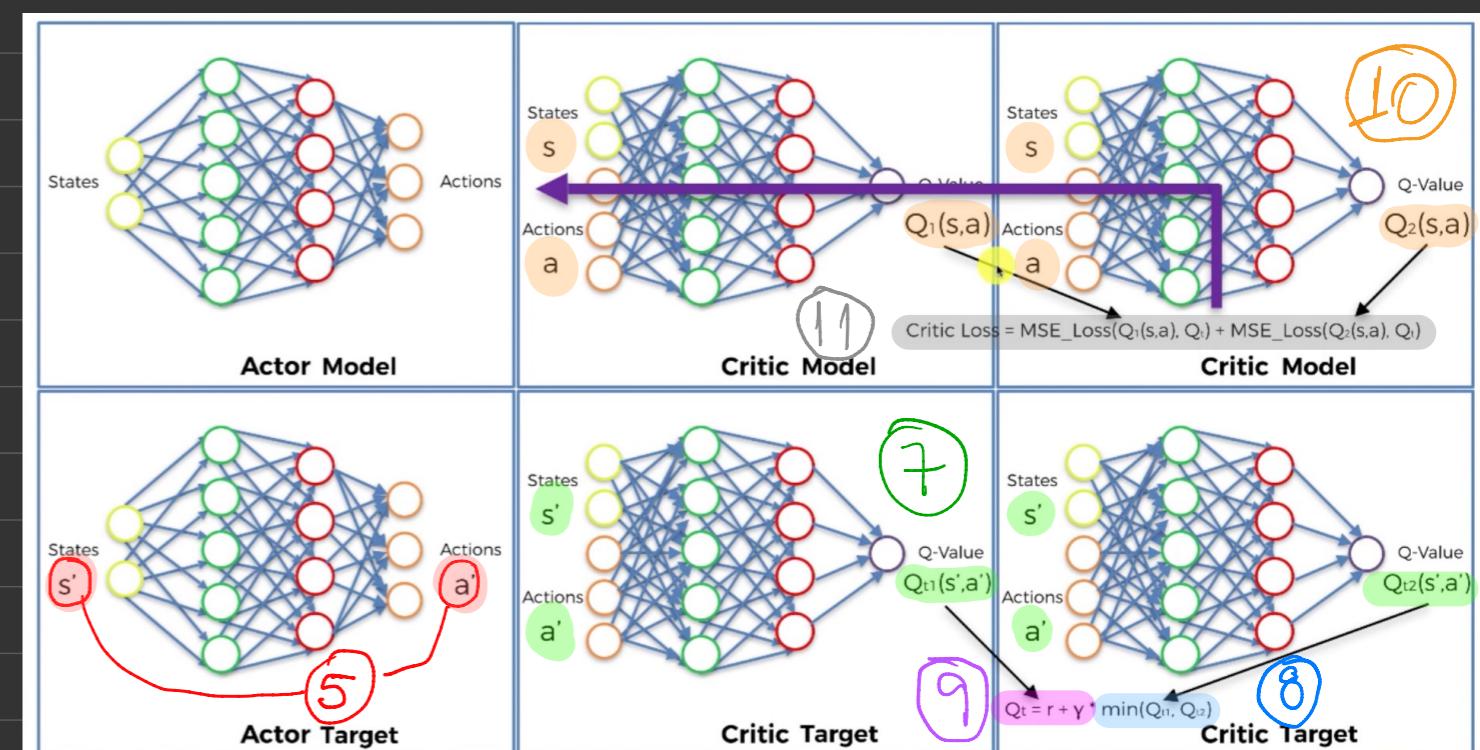
## Twin-Delayed DDPG (TD3)

### Initialization:

Step 1: We initialize the Experience Replay memory, with a size of 20000. We will populate it with each new transition.  
 Step 2: We build one neural network for the Actor model and one neural network for the Actor target.  
 Step 3: We build two neural networks for the two Critic models and two neural networks for the two Critic targets.

**Training Process - We run a full episode with first 10,000 actions played randomly, and then with actions played by the Actor model. Then we repeat the following steps:**

- We sample a batch of transitions  $(s, s', a, r)$  from the memory. Then for each element of the batch:
- From the next state  $s'$ , the Actor target plays the next action  $a'$ .
- We add Gaussian noise to this next action  $a'$  and we clamp it in a range of values supported by the environment.
- The two Critic targets take each the couple  $(s', a')$  as input and return two Q-values  $Q_{t1}(s', a')$  and  $Q_{t2}(s', a')$  as outputs.
- We keep the minimum of these two Q-values:  $\min(Q_{t1}, Q_{t2})$ . It represents the approximated value of the next state.
- We get the final target of the two Critic models, which is:  $Q_t = r + \gamma * \min(Q_{t1}, Q_{t2})$ , where  $\gamma$  is the discount factor.
- The two Critic models take each the couple  $(s, a)$  as input and return two Q-values  $Q_1(s, a)$  and  $Q_2(s, a)$  as outputs.
- We compute the loss coming from the two Critic models:  $\text{Critic Loss} = \text{MSE\_Loss}(Q_1(s, a), Q_t) + \text{MSE\_Loss}(Q_2(s, a), Q_t)$ .
- We backpropagate this Critic loss and update the parameters of the two Critic models with a SGD optimizer.



8)

- Taking the  $\min(Q_{t1}, Q_{t2})$  is the same as the  $\max[Q(s,a)]$  of Q-learning

$$\tilde{a} \leftarrow \pi_{\phi'}(s') + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$$

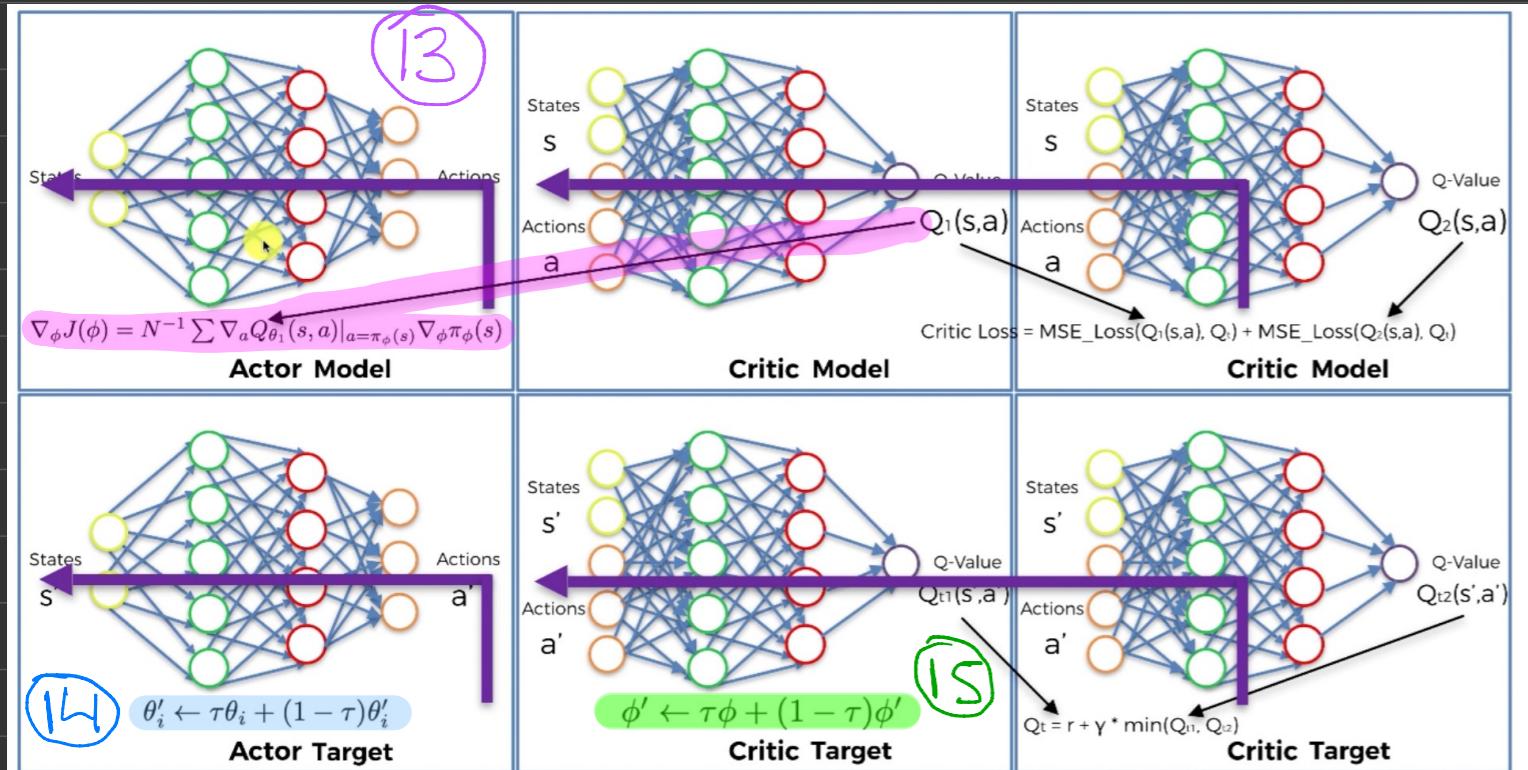
$$\tilde{a} \leftarrow \text{clip}(\tilde{a})$$

6

- Also taking the  $\min()$  helps to avoid overestimate the target Q-value.

# Policy learning Step

- Step 13: Once every two iterations, we update our Actor model by performing gradient ascent on the output of the first Critic model:  $\nabla_{\phi} J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s, a)|_{a=\pi_{\phi}(s)} \nabla_{\phi} \pi_{\phi}(s)$ , where  $\phi$  and  $\theta_1$  are resp. the weights of the Actor and the Critic.
- Step 14: Still once every two iterations, we update the weights of the Actor target by Polyak averaging:  $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$
- Step 15: Still once every two iterations, we update the weights of the Critic target by Polyak averaging:  $\phi' \leftarrow \tau \phi + (1 - \tau) \phi'$



Polyak uses a small number  $\tau$  0.8. This means we update the target network with a very small portion of [Model] weights, instead of Replace the entire target weights by the Target.

Gradient Ascent to maximize the value

In Step 13 we do:

$$a = \text{Action}(state)$$

$$\text{critic\_out} = \text{Critic-1}(state, a)$$

$$\text{critic\_out} = -\text{critic\_out}.mean()$$

$\text{critic\_out.backward()}$

Reduce to a single value

Perform Grad. Asc. to update the actor weights so that it produces an Action that Results in a higher  $\text{Q}_1(s, a)$

```

def train(self, iterations, batch_size=100, gamma=0.99, tau=0.005,
          policy_noise=0.2, noise_clip=0.5, policy_freq=2):
    if self.replay_buffer.mem_cntr < batch_size:
        return

    for i in range(iterations):
        # 4. Sample from memory
        batch_state, batch_action, batch_reward, batch_state_, batch_done = self.replay_buffer.sample(batch_size)
        state = torch.tensor(batch_state, dtype=torch.float, device=device)
        action = torch.tensor(batch_action, dtype=torch.float, device=device)
        reward = torch.tensor(batch_reward, dtype=torch.float, device=device).view(-1, 1)
        state_ = torch.tensor(batch_state_, dtype=torch.float, device=device)
        done = torch.tensor(batch_done, dtype=torch.bool, device=device).view(-1, 1)

        # 5. Actor target play state_
        action_ = self.actor_target(state_)

        # 6. Add gaussian noise the same shape as action (bs, n_actions)
        noise = torch.Tensor(batch_action).data.normal_(0, policy_noise).clamp(-noise_clip, noise_clip).to(device)
        action_ = (action_ + noise).clamp(-self.max_action, self.max_action)

        # 7. Critic target play with state_, action_
        q1_target, q2_target = self.critic_target(state_, action_)

        # 8. Min of q1_target and q2_target
        min_q = torch.min(q1_target, q2_target) # remove from backpropagation

        # 9. Calculate bellman equation
        q_target = reward + (gamma * min_q).detach()
        q_target[done] = 0.0

        # 10. Critic play with state, action
        q1, q2 = self.critic(state, action)

        # 11. calculate critic loss of each Critic and the q_target
        critic_loss_1 = torch.nn.functional.mse_loss(q1, q_target)
        critic_loss_2 = torch.nn.functional.mse_loss(q2, q_target)
        critic_loss = critic_loss_1 + critic_loss_2

        # 12. backpropagation
        self.critic_optimizer.zero_grad()
        critic_loss.backward()
        self.critic_optimizer.step()

        # 13. Delayed update of Actor model by perform gradient ascent trying to
        if i % policy_freq == 0:
            actor_loss = - self.critic.q1(state, self.actor(state)).mean()
            self.actor_optimizer.zero_grad()
            actor_loss.backward()
            self.actor_optimizer.step()

        # 14. Update actor_target using Polyak Average
        for param, target_param in zip(self.actor.parameters(), self.actor_target.parameters()):
            target_param.data.copy_(tau * param.data + (1 - tau) * target_param.data)

        # 15. Update critic_target using Polyak Average
        for param, target_param in zip(self.critic.parameters(), self.critic_target.parameters()):
            target_param.data.copy_(tau * param.data + (1 - tau) * target_param.data)
    
```