Episode → Change in point.

learned

$$q^{new}(s,a) = (1-\alpha)\underbrace{q(s,a)}_{old\ qvalue} + \alpha\left(R_{t+1} + \gamma\ \underset{a'}{max}\ q(s',a')\right)$$

Episode loop
→ reset env
→ check init onitialize flag for stopping
→ initial reward

Time step loop
→ take action based on exploitation, exploitation        exploitation
                                                          exploitation
→ comp update q table
→ update reward
→ check if episode is compl.
→ add reward of curr epsd to all rewards

Qlearning ≡ Using value iteration to compute optimal Q-value quide.

Deep Q learning ≐ Use neural networks to approximate Q-value function

Experience replay

↓

we store the agent's
experience at each
time step is a DS called
replay memory.

$$C_t = (s_t, a_t, r_{t+1}, s_{t+1})$$

Ⓠ Why we replay memory?

→ To break correlation.

## Summary of Steps

1. Initialize replay memory capacity
2. Initialize the network with random weights (Xavier initialization)
3. For each episode
   1. Initialize the starting state
   2. For each time step
      1. Select an action
         → via exploration or exploitation
      2. Execute selected action in an emulator
      3. Observe reward and next state
      4. Store experience in replay memory.

5. Sample random Batch from replay memory
6. Preprocess states from batch
7. Pass batch of pre-processed states to policy network
8. Calculate loss between target Q-values & Q-values
   • Requires second pass to the network for next state
9. Gradient descent updates weights in the policy network to minimize loss.

## Target network

→ When calculating loss between present Q-values and target Q-values, we used the target network [a clone of the original network] to compute target Q-values.

→ The weights of target network are updated after every x-time steps where x is a hyperparameter

[ Code ] ↘

└→ D Q N class ———┐→ init ()
                                  └→ forward ()

└→ Experience class

└→ Replay memory class ——┐→ init ()
                                          ├→ push ()
                                          └→ sample () & can-provide-sample()

└→ Epsilon greedy strategy ——→ init ()
                                          └→ get-exploration-rate()

└→ RL Agent
                  └→ init ()
                  └→ select-action ()


└→ CartPole Env Manager ()