# Tutorial #13

## Deep Q-learning

- Recall the Bellman equations for

(i) Value function:
$$V(s) = \max_{a \in \mathcal{A}} \mathbb{E}\left[ \pi(a,s,s') + \gamma V(s') \right]$$
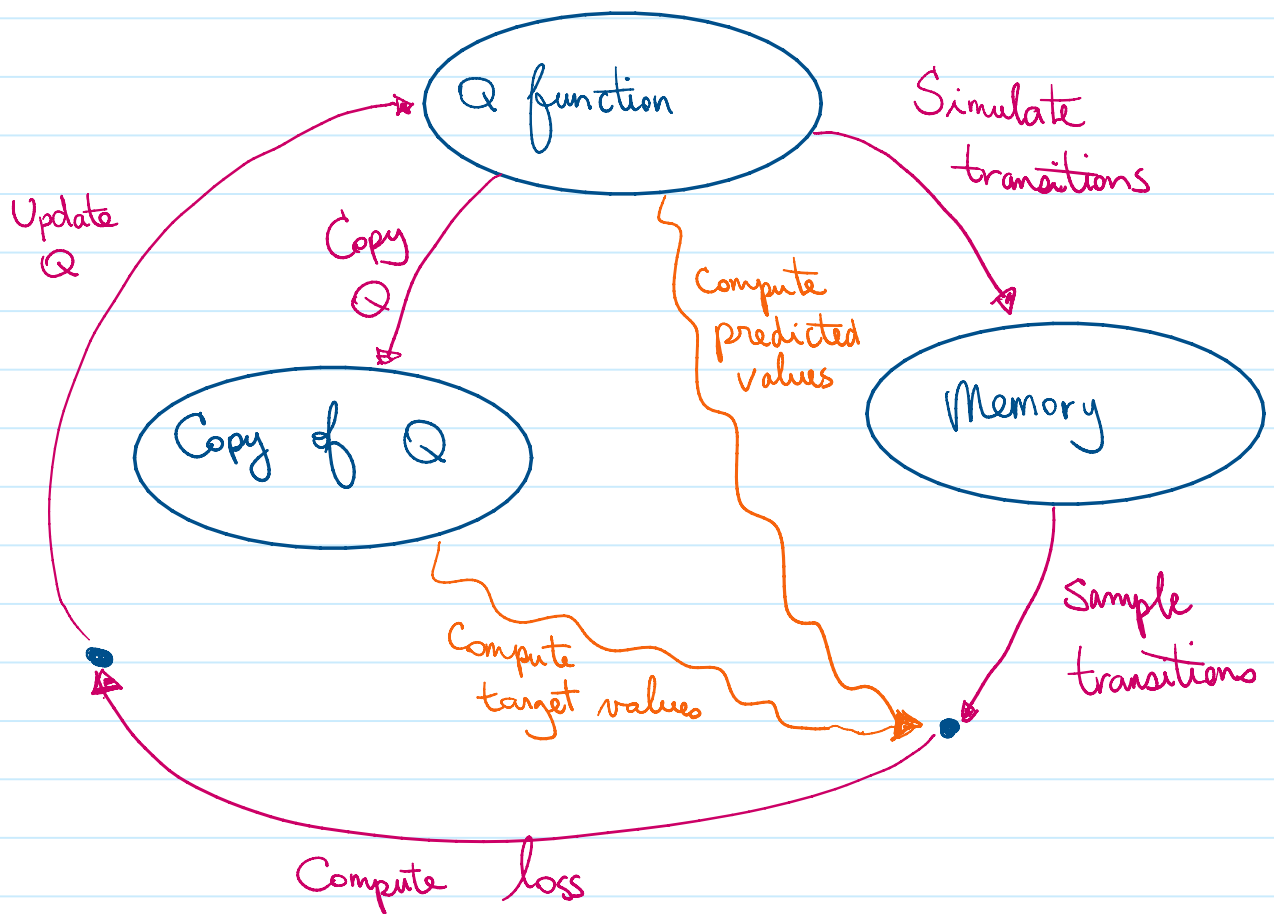
(ii) Q function:
$$Q(s,a) = \mathbb{E}\left[ \pi(a,s,s') + \gamma V(s') \right]$$
$$= \mathbb{E}\left[ \pi(a,s,s') + \gamma \max_{a' \in \mathcal{A}'} Q(s',a') \right].$$

- We want to use the Q-function (as a neural network) to find the best policy.

$\boxed{\text{Representation of the algorithm}}$



$\boxed{\text{More details on deep Q-learning}}$

- Q function (main and target networks):

  ANN with state dimension as inputs and number of actions as outputs, to give

  $$Q(s, a^{(1)}), Q(s, a^{(2)}), \ldots, Q(s, a^{(|A|)}).$$

- Memory: a replay buffer with a fixed size to store & sample transitions $(S_t, a_t, r_t, S_{t+1})$

- Simulate transitions:

  $\varepsilon$-greedy policy, where we select the best action (with respect to the Q function) with probability $1-\varepsilon$, and a random action with probability $\varepsilon$. It allows exploration & exploitation.

- Predicted Q value = $Q(S_t, a_t)$

- Target Q value = $r_t + \gamma \max\limits_{a' \in A} Q(S_{t+1}, a')$

- Update loss: Mean squared error between the predicted and target values.

- Other implementation details:

  (i) frequence of the copy of Q as the target network K.

  (ii) decay the $\varepsilon$ in the $\varepsilon$-greedy.

  (iii) mini-batch size when sampling transitions from the memory.

  (iv) select only valid actions.