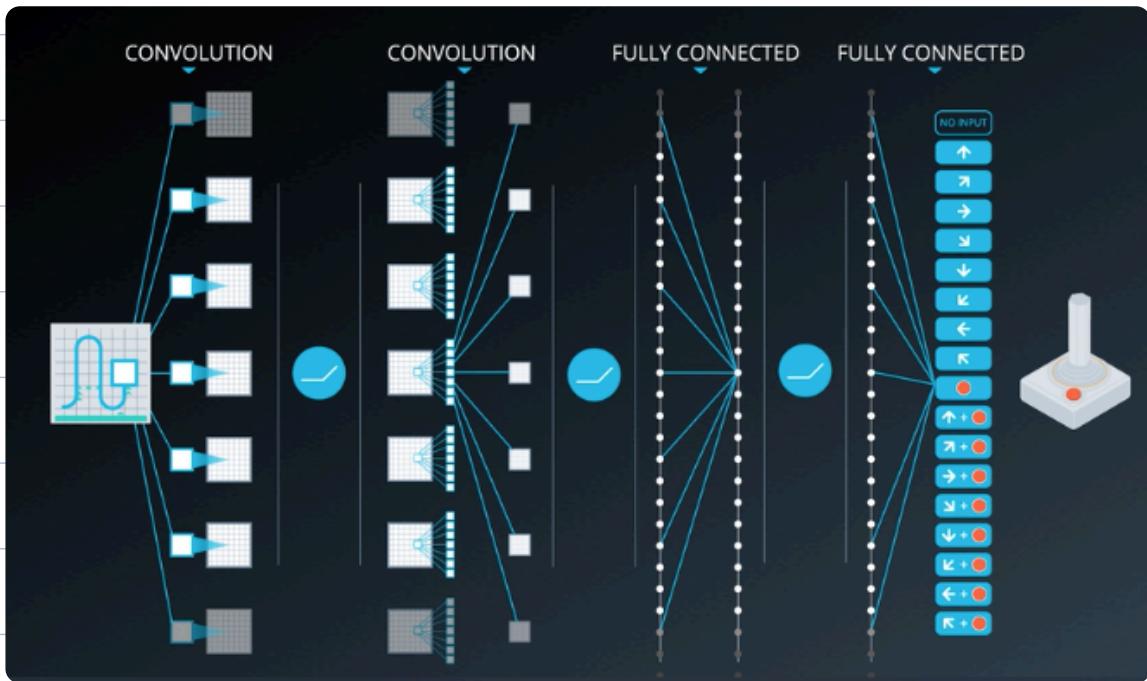




# Deep Q-Networks

---





Deep Q Network Architecture

- ⇒ Deepmind used DQN to play many Atari games.
- ⇒ DQN are able to learn the action value given a state as input.
- ⇒ It used raw image pixels as game state & predicted the action value for each possible game action.
- ⇒ Their proposed approach used 4 consecutive game frames to capture the temporal information.
- ⇒ The network comprised of 2 ReLU activated convolution layers followed by 2 Fully connected layers.
- ⇒ RL is unstable when using neural networks such as DQN. To overcome this issue 2 common methods which were used are:-

## # Experience Replay :-

Reason R1:- When the agent interacts with the environment, the sequence of experience tuples  $(S, A, R, S')$  can be highly correlated.

$(S_t, A_t, R_{t+1}, S_{t+1})$

$(S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2})$

$(S_{t+2}, A_{t+2}, R_{t+3}, S_{t+3})$

The neural network based approaches (DQN) can get

swayed away by these correlations in sequential order

experience tuples.

R<sub>2</sub>:- Some of the (states, actions) can be very rare in an environment & hence throwing them away after a single use (iteration) may not help the network learn.

R<sub>3</sub>:- Some states which are seen initially only can be forgotten by the network if they are not visited again due to progress in the environment.

To overcome these issues a "Replay Buffer" is used which contains a collection of experience tuples  $(S, A, R, S')$  which are added as the agent interacts with the env.

Experience Replay is to then randomly sample from the replay buffer to train the network which prevents action values from oscillating or diverging catastrophically.

It also allows the rare events to be sampled again due to random selection.

based on some policy (random/learned)

Note:- It is good for the agent to first try different actions<sup>↑</sup> on same/different states which are then stored in the buffer & then to train the n/w which allows the agent to explore instead of learning sequentially on the go, which may make it biased towards certain actions.

# Fixed Q-Targets:- The Q-learning update step in neural n/w is as follows:-

$$J(w) = E_{\pi} \left[ (q_{\pi}(s, a) - \hat{q}(s, a, w))^2 \right]$$

$$\nabla_w J(w) = -2 (q_{\pi}(s, a) - \hat{q}(s, a, w)) \nabla_w \hat{q}(s, a, w)$$

$$\Delta w = -\alpha \frac{1}{2} \nabla_w J(w)$$

$$= -\alpha (q_{\pi}(s, a) - \hat{q}(s, a, w)) \nabla_w \hat{q}(s, a, w) \quad \text{--- ①}$$

To update  $w$ ,  
same  $w$  used as target q-value estimator

$$\Delta w = \alpha (R + \gamma \max_a \hat{q}(s', a, w) - \hat{q}(s, a, w)) \nabla_w \hat{q}(s, a, w) \quad (2)$$

The problem in moving from eq. (1) to eq.(2) is that  $q_\pi(s, a)$  is the target action value that we hope to achieve. But, instead we are again estimating it with model parameters ' $w$ ' which is mathematically incorrect.

Here, we are trying to update a guess with a guess which can potentially lead to harmful correlations as we are trying to update ' $w$ ' with the same ' $w$ ' parameters in the Target Q-value estimation, leading to unstable training.

Instead, we can use the following update rule:-

decoupled

$$\Delta w = \alpha (R + \gamma \max_a \hat{q}(s', a, w^-) - \hat{q}(s, a, w)) \nabla_w \hat{q}(s, a, w)$$

where  $w^-$  are the weights of a separate target network that are not changed during the learning step.

After a certain number of learning steps, updated ' $w$ ' weights are copied into ' $w^-$ ' and the process is repeated.

## Algorithm: Deep Q-Learning

- Initialize replay memory  $D$  with capacity  $N$
- Initialize action-value function  $\hat{q}$  with random weights  $w$
- Initialize target action-value weights  $w^- \leftarrow w$
- **for** the episode  $e \leftarrow 1$  to  $M$ :
  - Initial input frame  $x_i$
  - Prepare initial state:  $s \leftarrow \phi(x_i)$
  - **for** time step  $t \leftarrow 1$  to  $T$ :
    - Choose action  $A$  from state  $S$  using policy  $\pi \leftarrow \epsilon\text{-Greedy}(\hat{q}(S, A, w))$
    - Take action  $A$ , observe reward  $R$ , and next input frame  $x_{t+1}$
    - Prepare next state:  $s' \leftarrow \phi(x_{t-2}, x_{t-1}, x_t, x_{t+1})$
    - Store experience tuple  $(s, A, R, s')$  in replay memory  $D$
    - $S \leftarrow s'$

$\phi(x_i)$  is a pre-processing function that takes the initial game frame to convert it into required size by stack 4 frames together

SAMPLE	Take action $A$ , observe reward $R$ , and next input frame $x_{t+1}$ Prepare next state: $s' \leftarrow \phi(x_{t-2}, x_{t-1}, x_t, x_{t+1})$ Store experience tuple $(s, A, R, s')$ in replay memory $D$ $S \leftarrow s'$
LEARN	Obtain random minibatch of tuples $(s_j, a_j, r_j, s_{j+1})$ from $D$ Set target $y_j = r_j + \gamma \max_a \hat{q}(s_{j+1}, a, w^-)$ Update: $\Delta w = \alpha (y_j - \hat{q}(s_j, a_j, w)) \nabla_w \hat{q}(s_j, a_j, w)$ Every $C$ steps, reset: $w^- \leftarrow w$

④ You are encouraged to read the original DQN paper by Deepmind which gives many other optimization techniques.

## Updating Fixed Target Q Network:-

While copying the weights of  $w$  to  $w^-$ , the weights are not completely overwritten by  $w$ , but instead weights of both  $w$  &  $w^-$  are combined together using a parameter  $\tau_{\text{au}}(\tau)$  which governs the division as follows:-

$$w_{\text{new}} \leftarrow \tau * w + (1-\tau) * w_{\text{old}} \quad \text{where } 0 \leq \tau \leq 1$$

## Improvements over DQN

### # Double DQN:-

DQNs have a problem of "Overestimation of Q-values".

Recall, that in DQN:-

$$\Delta w = \alpha \left( R + \gamma \max_a \hat{q}(s', a, w^-) - \hat{q}(s, A, w) \right) \nabla_w \hat{q}(s, A, w)$$

$\downarrow \text{rewrite as}$

$$R + \gamma \hat{q}(s', \underset{\text{state}}{\underset{|}{\arg\max_a}} \hat{q}(s', a, w^-), w^-)$$

$\underset{\text{Q-value}}{|} \quad \underset{\text{action with max Q-value}}{|}$

When written in this manner, it is clear that the DQN can do mistakes especially in the initial training steps when it has very less knowledge & hasn't tried out all possible actions to know which is a good action.

This leads to an overestimation of Q-values since we are taking the maximum among a set of noisy numbers.

To reduce this overestimation, Double Q-Learning can be used, where you estimate the best action using one set of parameters & get the q-value for that action using other set of parameters to reduce the overestimation.

In Double DQN we already have 2 parameter sets ( $w$  &  $w^-$ ) for local & target Q-networks which are used as follows to reduce overestimation:-

$$R + \gamma \hat{q}(s', \underset{a}{\operatorname{argmax}} \hat{q}(s', a, w), w^-)$$

select best action  
evaluate that action

This way both the models should agree on that action otherwise the q-value returned won't be high which solves the problem of overestimation.

## # Prioritized Experience Replay :-

Another issue in DQN is that the Replay Buffer samples experiences with a uniform probability. This can be a problem since some experience might be more useful than the rest & some may be very rare.

This motivates us to have a Prioritized Experience Replay Buffer.

To assign a probability to an experience, we can use the amount of TD Error the experience generates as a measure of importance of that experience :-

$$\frac{\delta_t = R_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, w^-) - \hat{q}(s_t, A, w)}{\frac{\text{TD Target}}{\text{TD Expected}}}$$

then, Priority ( $p_t$ ) of a sample can be :-

$$p_t = |s_t|$$

& can be stored in the replay buffer, with the experience as :-  $\langle s_t, A_t, R_{t+1}, s_{t+1}, p_t \rangle$

then, the Sampling Probability of each sample can be expressed as :-

$$P(i) = \frac{p_i}{\sum p_k}$$

When an experience is picked, its priority is updated using the newly computed TD error.

Some improvements :-

1. If the TD error of a sample is zero, the sample may never be picked even though it may be useful in future. To avoid leaving out an experience, we can add a small constant ( $\epsilon$ ) to every priority value:-

$$\text{Now, Priority } p_t = |\delta_t| + \epsilon$$

2. Greedily Sampling using the probabilities may lead to overfitting over some samples which have high probability. To avoid this we can reintroduce some sort of uniform sampling as follows:-

$$\text{Sampling Probability } P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}, \text{ where } 0 \leq \alpha \leq 1$$

$\alpha = 0$ , leads to complete uniform sampling

$\alpha = 1$ , leads to complete priority sampling

3. We only want to sample experiences from the buffer using priorities but such a sampling also introduces a bias in the learned Q-values.

To avoid this, we can modify the update rule :-

$$\text{Modified Update Rule: } \Delta w = \alpha \left( \frac{1}{N} \cdot \frac{1}{P(i)} \right)^b \delta_i \Delta_w \hat{q}(s_i, a_i, w)$$

Importance - sampling weight

here,  $N$  = Number of samples in the replay buffer

$$b = \text{constant such that } 0 < b \leq 1$$

value of ' $b$ ' can be gradually increased to 1 as the q-values converge overtime with training.

## # Dueling Network:-

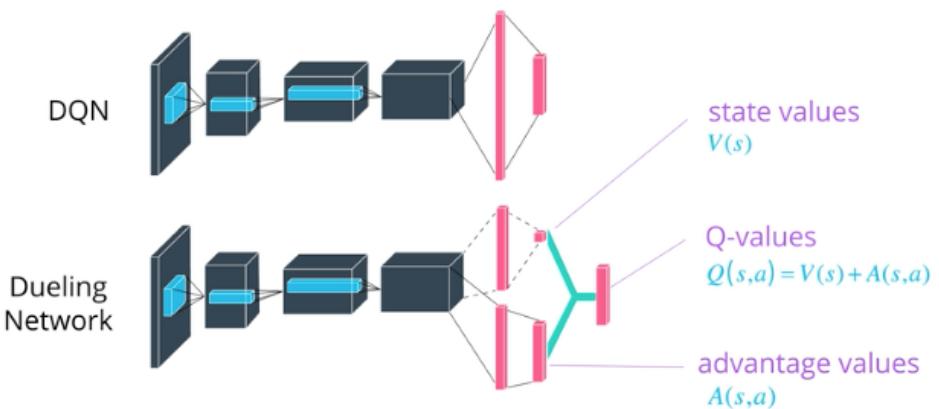
The motivation behind Dueling Network

N/W architecture is that mostly

State values remain the same

where as the action values

may vary.



To fulfill this, the FC layers are divided into 2 parts, State Values & Advantage Values.

State values gives out a scalar value ( $V(s)$ ) corresponding to the state

Advantage values gives a vector of values which corresponds to the values of each action.

The outputs are then combined as follows to get the Q-value :-

$$Q(s,a) = V(s) + A(s,a)$$

which is calculated as follows for stable training :-

$$Q(s,a) = V(s) + \left( A(s,a) - \frac{1}{|A|} \sum_{a'} A(s,a') \right)$$

$\downarrow$   
No. of Actions

## Rainbow DQN:-

Deepmind combined six different extensions of DQN together & named the algorithm as Rainbow DQN. The 6 optimizations were :-

- Double DQN
- Multi-step bootstrap targets (revisited in A3C)
- Prioritized Experience Replay
- Distributional DQN
- Dueling DQN
- Noisy DQN

These techniques combined together gave state-of-the-art results.