

# DQN (Deep Q Networks)

Deep Q Learning

algorithm that learned to play many Atari video games better than humans.

raw pixel data (exactly what a human would see on screen)

has a neural network operating like a function approximator

pass an image of the game state, and o/p is a vector of action values, with the max value indicating the action to be taken.

As reinforcement signal, fed back the change in game score at each time step.

NN initialized with random values, actions taken are all over the place, but overtime it begins to associate situations & sequences in the game with apt actions & learns to actually play the game well.

DQN is designed to produce q value for each every possible action in a forward pass

(action for the state can then be chosen stochastically or greedily)

architecture: 3 layered CNN (to exploit spatial relationships)

(Regularized linear units)

ReLU activation

1 fully connected HL (ReLU)

1 fully connected o/p layer that produced vector of action values

statespace for atari (210x160 colored)

(84x84 monochrome)

deepmind performed some minimal processing

eq. images

allowed them to

use more optimized NN operations on a few

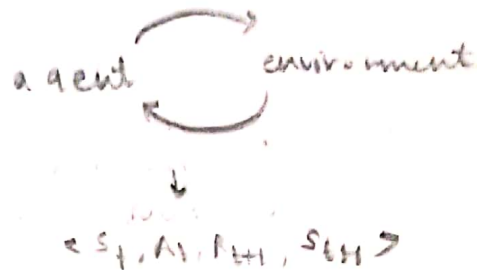
## 56) Experience replay

↓

originally proposed to make more efficient use of observed experiences.

↓

for each time step, we get the tuple, learn from it, and then discard it, moving onto the next tuple in the following time step.



a little wasteful → we could maybe learn a little bit more, if we stored the tuples somewhere

↓

some states are pretty rare, and some actions are pretty costly.

↓

nice to be able to recall such experiences

↓

replay buffer

store each experienced tuple in this buffer as we are interacting with the environment, and then sample a small batch to learn from them)

↓

As a result, we are able to learn from individual tuples, multiple times, recall rare occurrences, and make better use of our experience.

every action  $A_t$  affects the next state  $S_t$  in some way, which can mean that a sequence of tuples can be highly co-related.

↓

A naive & learning approach that learns from these experiences in sequential order runs the risk of being swayed by the effects of this co-relation.

↓

(randomized)  
since sampling will lead to random sequences, the sequential correlation based learning can be avoided. (avoid action values from oscillating or diverging catastrophically)

↓

example: tennis: rallying against the wall,

↓

more confident with my forehand, than my backhand. and I can hit the ball fairly straight, hence the ball keeps coming back in the same area.

i.e. to the right, and hence we keep hitting forehand shots

↓

but rest of the state space is not being explored,

so I try different combinations of states and actions and sometimes mistakes are made but eventually best policy is learnt

↓

for discrete space, this seems okay, but for continuous space, things may start to fall apart.

using epsilon + greedy method



every action  $A_t$  affects the next state  $S_t$  in some way, which can mean that a sequence of tuples can be highly co-related.

↓

A naive & learning approach that learns from these experiences in sequential order runs the risk of being swayed by the effects of this co-relation.

↓

(randomized)

Since sampling will lead to random sequences, the sequential correlation based learning can be avoided. (avoid action values from oscillating or diverging catastrophically)

↓

example: tennis: rallying against the wall,

↓

more confident with my forehand, than my backhand. and I can hit the ball fairly straight, hence the ball keeps coming back in the same area.

i.e. to the right, and hence we keep hitting forehand shots

↓

but rest of the state space is not being explored,

so I try different combination of states and actions and sometimes mistakes are made but eventually best policy is learnt

↓

for discrete space, this seems okay, but for continuous space, things may start to fall apart.

using epsilon + greedy method

(58)

\* The ball can come anywhere between the extremes left & extreme right



discretization of the state space may lead to no learning for specific holes i.e. (not visited during exploration)  
margin for deviation is not much,



Better to use a function approximator using RBF kernels or a  $Q$ -network that can generalize learning across the space.

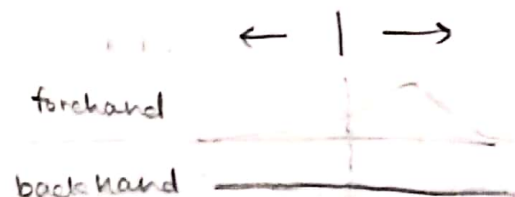


when the ball comes to the right and we successfully employ a ~~right~~ forehand, state value function changes slightly, i.e. it becomes more true around the exact region where the ball came, but also raises the value for the forehand shot in general across the state space.

of course the effect is less pronounced away from the exact spot, but over time it can add up. and that's exactly what happens when we try to learn while playing, processing each exp. tuple in order.



state correlation problem here leads to the agent playing forehand everywhere



i) stop learning while ~~playing~~ practicing



thus time is best spent in trying  
out different shots, playing randomly  
and exploring the state space

ii) hence it becomes imp. to memorize the experience, then  
later analysing what went wrong & what went  
right. , when I am at home and resting or I  
take a break.



Advantage: more comprehensive set of examples,  
we can generalize patterns from across these,  
recalling them in whatever order we please



after a round of learning, I can go back  
getting a round of experience, and the  
process is repeated.

(learning more robust policy)



building a database, and learning  
a mapping from them, which in a way  
reduces the RL problem or atleast value learning  
portion to a supervised learning scenario.



multiple passes  
over the same  
experience

and this allows us to apply  
machine learning techniques, models  
and algorithms prevalent in  
supervised learning literature

(60)

Fixed  $Q$  targets

↓

$Q$  learning → form of TD learning

(desired value)

→ target value

$$\Delta w = \alpha (R + \underbrace{\gamma \max_a \hat{q}(s', a, w)}_{\text{target value}} - \underbrace{\hat{q}(s, a, w)}_{\text{current estimate}}) \nabla_w \hat{q}(s, a, w)$$

Difference : TD error

↓

goal is to reduce the difference between these values

↓

The TD target here is supposed to be a replacement for the true value function  $q_\pi(s, A)$ , which is unknown to us.

↓

We defined the squared error loss function

wrt  $q_\pi(s, A)$ , later differentiated

with respect to  $w$  to get gradient descent

update rule.  $q_\pi(s, A)$  not dependent

on approximation parameters,

resulting in a simple update rule.

↓

But our TD target is dependent on these

parameters, which means simply

replacing the  $q_\pi(s, A)$  with our TD

target, <sup>is</sup> ~~might not be the best~~ method to approach this. (mathematically incorrect)



(\*) we might get away with that in practice because the changes are pretty small and change generally occurs in the right direction. (61)

↓

but if  $\alpha = 1$ , and changes being made are significant, then issues creep in.

↓

less of a concern when lookup table or dictionary is used, since  $(s, a)$  pair values are stored separately, but here the values are correlated and change at some pair impacts nearby values as well, as values are intrinsically tied together through function parameters.

↓

more like chasing a moving target (not optimal)

↓

better to set a short target, reach it & then change it. and repeat.

↓

decoupling the target position from the donkey's action. giving it a more stable environment and stops it from oscillating and all.

↓

(fix the function parameters used to generate our target)

↓

$w' \leftarrow w$

↓

i.e. move towards one  $w'$  for certain timesteps and then update  $w'$  with the current  $w$ , and so on. ← changing  $w$  for next steps, then update  $w'$  with latest  $w$  and so on.



(62)

i.e. we employ

→ immediate return for choosing a certain action

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

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

$x \rightarrow$  feature vector for environment state repr.

initialization methods.

\* Deep Q Learning Algorithm : part 1 : sampling  
↓ (building experience)

a) choose action  $A$  from states using policy

$$\pi \leftarrow \epsilon - \text{greedy}(\hat{q}(s, A, w))$$

b) Take action  $A$ , observe reward  $R$ , next input frame  $x_{t+1}$

c) prepare ~~any~~ next state

$$s' \leftarrow \phi((x_{t-2}, x_{t-1}, x, x_{t+1}, x_{t+2}))$$

d) Store experience tuple  $(s, A, R, s')$  in replay memory D

$$s \leftarrow s'$$

part 2: learn phase

a) obtain random minibatch of tuples  $(s_j, a_j, r_j, s_{j+1})$  from  $D$ .

b) set target  $y_j = r_j + \gamma \max_a \hat{q}(s_{j+1}, a, w^-)$

\* c) update :  $\Delta w = \alpha (y_j - \hat{q}(s_j, a_j, w)) \nabla_w \hat{q}(s_j, a_j, w)$

$$w = w + \Delta w$$

d) every  $c$  steps, reset :  $w^- \leftarrow w$

(6.3)

These tasks are not directly dependent on each other. Therefore multiple sampling ~~beats~~ steps, then one learning steps or even multiple learning steps with different random tuple batches.

↓

Other part of the algorithm is meant to support this core basis of the DQN.

↓

a) initialize empty memory buffer  $D$  (capacity:  $N$ )

(since this buffer is finite in size, we can use a circular queue, that retains the most recent experience tuples)

b) initialize action-value function  $\hat{q}$  with random weights

like the ones available in pytorch. (network weights)

c)  $w^- \leftarrow w$  (initial)

\* Now for each episode & each time step, we observe a raw screen image or i/p frame  $x_t$ , which needs to be preprocessed, (to capture temporal features)

↓

Stack i/p frames to build state vector

↓

$$S \leftarrow \phi(x_t)$$

\* memory is not cleared after each episode, this enables us to build and learn <sup>from</sup> experience across episodes.

(64)

## \* Improvements over DQN

(most prominent ones)



a) Double DQN

b) Prioritized Experience Replay

c) Dueling DQN

think budget  
over estimates



a kind of bias where we have the error (value is more than

• Overestimation of action values that Q-learning is prone to (the actual value)

prone to



TD target



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

$$\gamma W \hat{q}(s, a, W)$$

TD target : Q value for the state  $s'$ , and the action that results in the maximum Q-value among all possible actions from that state



Choosing based on max value of action value function might not be the best strategy, as the

Since Q-values are still evolving, the agents won't have gathered enough information to figure out which action is best.

accuracy of our Q-values depend a lot on what actions have been tried, and what neighbouring states have been explored



over estimation of Q-values.

more robustness ✓

(max among set of noisy numbers)



Double Q-learning



selected best action as per one set of parameters  $w$ ,  
but evaluate using  $w'$

↓

this is like 2 separate function approximators  
agreeing upon the best possible action.

↓

2 value functions  
↓  
randomized pick &  
drop.

$R + \gamma \hat{q}(s', \arg \max_a \hat{q}(s', a, w), w')$

selected best action (from local  
two  
networks)  
↓  
evaluation  
basis of that  
action.

↓

if  $w$  picks an action that is not best according to  $w'$ ,  
 $Q$ -value returned is not that high.

↓

In the long run, this prevents the algorithm from  
propagating incidental high rewards that may have  
been obtained by chance & do not reflect long-term  
returns.

•) prioritized experience replay:

but some experiences might be  
more important than others

↓

these exp. might be rare as well

↓

uniform sampling might result in very  
small chance of these experiences  
getting selected.

↓

also, since finite memory, these exp.  
might get lost.

basic principle behind  
experience replay

↓

sample state tuples,  
store in buffer, and  
randomly pick sample  
of tuples from buffer to  
learn

↓

breaks correlation b/w  
consecutive experiences &  
stabilizes the algorithm

(66)

hence the requirement for prioritized experience replay comes in.

↓

criteria to assign priorities to each tuple?

↓

a) TD error delta

(the bigger the error, the more we expect to learn from that tuple)

$$\delta_i = R_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a, w) - \hat{q}(s_t, A, w)$$

↓

$p_i = |\delta_i|$  → store with each corresponding tuple in the replay buffer

↓

and when creating batches, we can use these values to compute sampling probability

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

↓

and as tuple is picked, corresponding values can be updated using a newly updated TD error using the latest  $q$  values.

a) if  $TD_{error} = 0$ , probability of being picked  $= 0$

↓

doesn't necessarily mean, we have nothing to learn from the experience.

↓

might be the case of limited learning, i.e.

it is possible our estimate was close due to the limited samples we visited till that point.

↓

to prevent from starvation,

$$P_i = |\delta_i| + \epsilon$$

b) greedily using these max. values might cause only certain states to be explored. (replayed over & over)

↓

resulting in a sort of overfitting on that subset.

↓

to avoid this, introduce concept of uniform sampling, (random)

↓

add another <sup>hyper</sup>parameter  $A$ , which is used to redefine the sampling probability as,

$$P(i) = \frac{P_i^A}{\sum_k P_k^A}$$

↓

$A$  can be varied to adjust the contribution of greedy & random sampling.

$A = 0$ : pure uniform randomness

$A = 1$ : greedy priority choosing



68

modified update rule:

$$\Delta w = \alpha \left( \frac{1}{N} \cdot \frac{1}{P(c_i)} \right)^b \delta_i \nabla_m \hat{q}(s_t, a_t, w)$$

size of

buffer

↓

b to control,

how much these weights affect learning.

important sampling weight

to compensate for introduced bias.

→ more imp. towards the end, when the Q-values begin to converge.

↓

b can be moved from a low value to 1, over time.

• Dueling Networks.

↓

considered is to use 2 streams

↓

conventional architecture for this is a series of conv layers (for feature extraction), and then fully connected dense layer to come up with probabilities.

i/p state (feature vector)

↓

get action function values

(Q values)

→ 1 stream estimates state value function :  $v(c)$

↓

and other estimates

the advantage of each action

↓

\* branch off with their respective fully connected layers.

Finally the desired  $Q$  values obtained by combining the state & advantage values.

(69)

The ~~Q-values~~ values of most states don't vary a lot across actions, so it makes sense to try & directly estimate them, but we still need to capture

$Q(s)$  the difference different actions make in each state

↓  
advantage function.

(OpenAI gym)

↓  
some modifications are required to be made on the normal model architecture.

↓  
other extensions: a) <sup>\*</sup> learning from multi-step.

bootstrap targets

b) Distributional DQN.

c) noisy DQN.

↓  
each of the six extensions address a different issue with the original DQN implementation

↓  
The algorithm that tested performance using all the ~~three~~ six extensions together: rainbow (deepmind)

↓  
outperforms individual modifications & achieves state of the art on Atari 2600 games.