# DQNs

## Basic Structure - Idea:

Agent | D.N.N.



States

Input layer

Hidden layers

Output layer

$Q(s,a)$

max $Q(s,a)$

Optimal action $(a_t)$

Environment

Reward $(r_t)$

$r_{t+1}$

State $(s_t)$

$s_{t+1}$

Q-learning function: $Q(s,a) = r + \gamma \max_{a'} Q(s',a')$

Πρόβλημα! Πρέπει να ξέρω $Q(s',a')$, $\forall (s',a') \rightarrow$ inefficiency.

Solution: Deep Q-Network.

Updating Q-learning function - Bellman Equation!

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot \ell \left[ \overbrace{r + \gamma \max_{a'} Q(s',a')}^{\text{immediate reward}} - Q(s,a) \right]$$

learning rate

Discount (what it should be factor (How much we value the future")

target $Q(s,a)$

current $Q(s,a)$

Καταλαβαίνουμε οτι νευρωνικό μας πρέπει να ελαχιστοποιήσει την

$$Q(s,a) - \left[ r + \gamma \max_{a'} Q(s',a') \right] \quad \leftarrow \text{loss function}$$

## Target network

- Τα $Q(s,a)$, $Q(s',a')$ δεν πρέπει να υπολογίζονται με τον ίδιο τρόπο
  - Δεν έχει νόημα!

Φτιάχνουμε δυο όμοια νευρωνικά: $\begin{cases} \text{policy-net} \leftarrow \{ \text{απόφαση} / \text{training} \} \\ \text{target-net} \leftarrow \{ Q(s',a') \} \end{cases}$

Τα βάρη του policy-net ανακινώνται συνεχώς, ενώ του target-net περιοδικά.

⚠ Notation: Βάρη policy: $\theta$ , Βάρη target: $\theta'$

DQN Bellman Equation: $Q(s,a) \leftarrow Q(s,a) + \ell \left[ \underbrace{r + \gamma \max_{a'} Q(s',a'; \theta') - Q(s,a; \theta)}_{\text{loss function}} \right]$

\* 

Μεχρι τώρα: $\begin{cases} \text{Αρχιτεκτονική} \checkmark \\ \text{Loss function για υπολογισμό των βάρων} \checkmark \end{cases}$

\* Το target_net θα ψιμνεύεται smoothly (soft update) ως εξής:

$$\theta' \leftarrow \tau \theta + (1-\tau)\theta' \text{, όπου } \tau \ll 1.$$

## Πως κάνουμε train:

- Selecting action: φτιάχνουμε ένα epsilon
  οπότε παιρνουμε τυχαίο αριθμό $x$ στο $[0,1)$



Αν $x <$ epsilon τότε: random-action ←Exploration

Αλλιώς: χρησιμοποίησε το policy-net ←Exploitation

└ This is called epsilon greedy policy.

- **Calculating the reward**

  Depends on the environment,

  General Rules
  - small and dense rewards > Large and sparse ones
  - give penalties to control learning.
  - be creative, there is no single best solution

- **Experience Replay**

  Αποθηκεύουμε στω μνήμα το (state, action, next state, reward) για κάθε action που κάνουμε. Έτσι μπορούμε να υποβάλουμε το loss σε batches των experiences που έχουμε μαζέψει και τελικά να ανανεώνουμε τα βάρη.