

1. Learning algorithm

1.1. Overview of the technique

The agent is implemented using the Deep Q-Learning method (along with Replay Buffer and Fixed Q-Target technique). At each time step, the agent will take an action, store the transition in its replay buffer (Sample) and then sample a minibatch from the buffer to update the primary Q-network with the target Q-network. The target network is updated either after a number of steps by copying the primary network's weights or slowly update by τ (a hyperparameter):

$$\text{target} = \tau * \text{primary} + (1 - \tau) * \text{target}$$

The algorithm in detail:

Algorithm: Deep Q-Learning

- Initialize replay memory D with capacity N
- Initialize action-value function \hat{q} with random weights \mathbf{w}
- Initialize target action-value weights $\mathbf{w}^- \leftarrow \mathbf{w}$
- **for** the episode $e \leftarrow 1$ to M :
 - Initial input frame x_1
 - Prepare initial state: $S \leftarrow \phi(\langle x_1 \rangle)$
 - **for** time step $t \leftarrow 1$ to T :

SAMPLE

Choose action A from state S using policy $\pi \leftarrow \epsilon\text{-Greedy}(\hat{q}(S, A, \mathbf{w}))$
Take action A , observe reward R , and next input frame x_{t+1}
Prepare next state: $S' \leftarrow \phi(\langle x_{t-2}, x_{t-1}, x_t, x_{t+1} \rangle)$
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 $y_j = r_j + \gamma \max_a \hat{q}(s_{j+1}, a, \mathbf{w}^-)$
Update: $\Delta \mathbf{w} = \alpha (y_j - \hat{q}(s_j, a_j, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s_j, a_j, \mathbf{w})$
Every C steps, reset: $\mathbf{w}^- \leftarrow \mathbf{w}$

Besides, in order to help the agent learn better, I have implemented the algorithm with 2 improvements: Duelling DQN and Double DQN to avoid overestimating problem.

*) Double DQN:

When learning:

- Vanilla DQN estimate the true value for computing loss by the maximum action-value at the next state, which return from the target network. This can lead to overestimation of action-value.

Formula for this true value estimation can be written as:

$$R + \gamma \hat{q}(S', \underset{a}{\operatorname{argmax}} \hat{q}(S', a, \mathbf{w}), \mathbf{w})$$

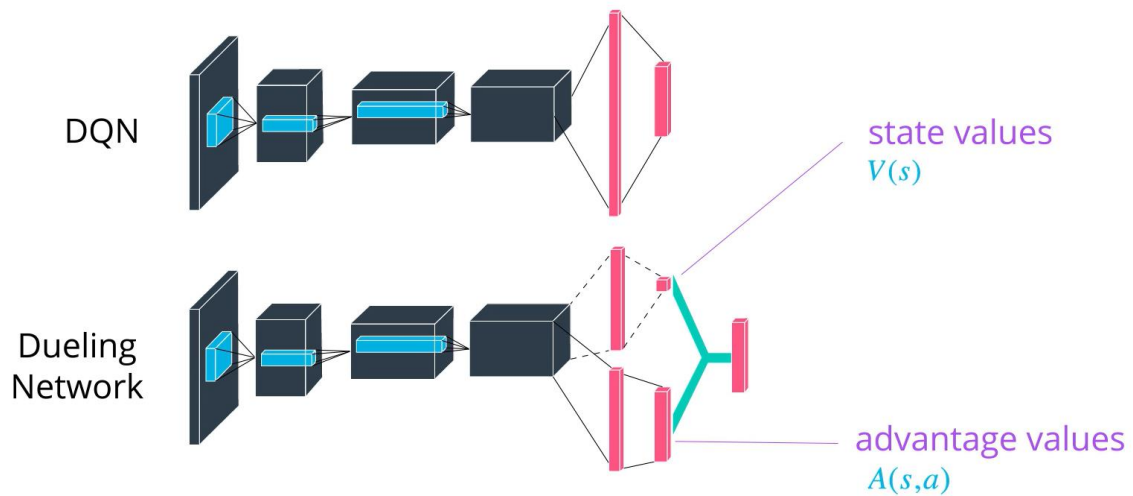
(Select best action from the target network and evaluate that action on the target network)

- Double DQN solve this by select the best action from the primary network and evaluate that action on the target network.
- Formula for this estimation

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

*) Dueling DQN

Its network architecture is different from vanilla DQN. It uses 2 streams, 1 to compute the value of state and 1 to compute the advantage of each action in that state.

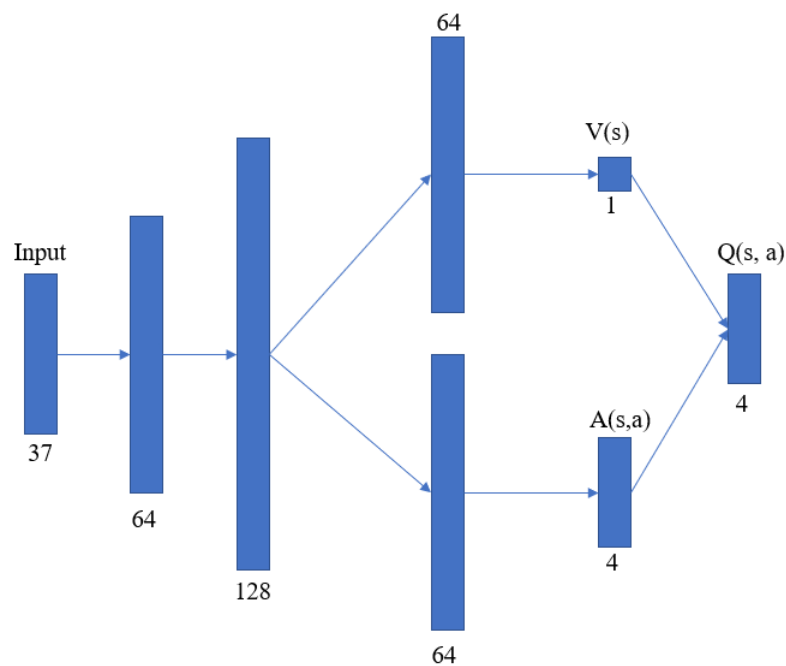


The output of the dueling network is computed as follow (forcing the advantage function estimator to have 0 advantage at the chosen action):

$$Q(s, a) = V(s) + (A(s, a) - \frac{1}{Na} \sum_{a'} A(s, a')) \quad (Na: \# \text{ actions})$$

1.2. Network architecture

The network architecture for both primary and target Q-network is described below (the number is # neurons at each hidden layers)



1.3. Hyperparameters select

```
class Config:
    # seed for randomness
    SEED = 3737

    # name of the unity environment
    UNITY_ENV = "Banana_Windows_x86_64\Banana.exe"

    # training hyperparams
    BATCH_SIZE = 64
    BUFFER_SIZE = 100000
    TAU = 1e-3
    NETWORK_ARCHITECTURE = [64, 128, 64]
    LEARN EVERY = 4
    EPSILON_START = 1.0
    EPSILON_END = 0.05
    EPSILON_DECAY = 0.9995
    EPSILON_FOR_INFER = 0.05
    GAMMA = 0.99
    LR = 0.001
    DEVICE = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

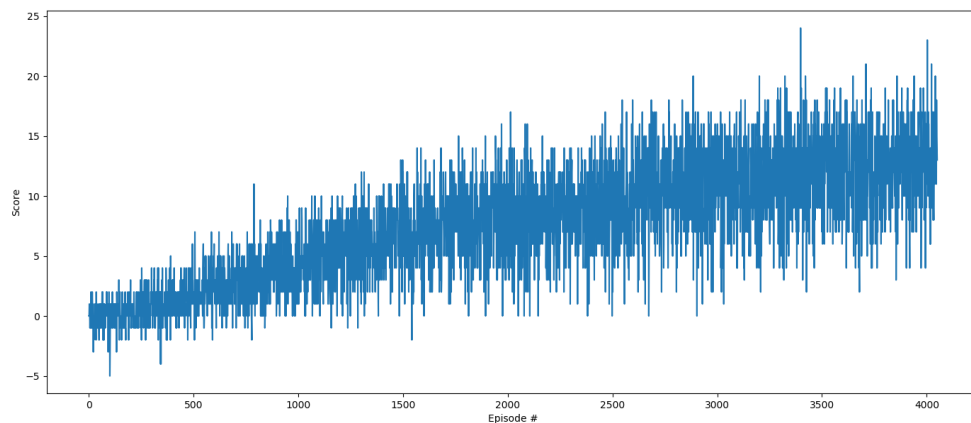
    # batch size for learning
    # buffer max capacity
    # hyperparam update the target network
    # hidden layers in the network
    # agent will sample and learn after LEARN EVERY steps
    # epsilon start at first episode
    # epsilon min
    # epsilon decay after each episode
    # epsilon for inferencing
    # discount factor
    # learning rate

    NUM_EPISODE = 5000
    NUM_EPISODE_INFER = 100
    MAX_TIMESTEP = 300

    # num of training episodes
    # num of inferencing episodes
    # max timestep for each episode

    CHECKPOINT = "checkpoint.pth"
    # checkpoint path of trained weight
```

2. Plot of rewards



3. Ideas for future works

The result I have shown is not too good, which solve the problem in about 4k episodes 😞 I think some modifications will make the result better:

- Tuning hyperparameters such as epsilon (start, end, decay) and the architecture of the Q-network probably enhance the learning ability of the agent
- Moreover, many other techniques can be used in order to help the agent learn better such as Prioritized Replay Buffer, Distributing DQN, ... or even combine all – Rainbow.