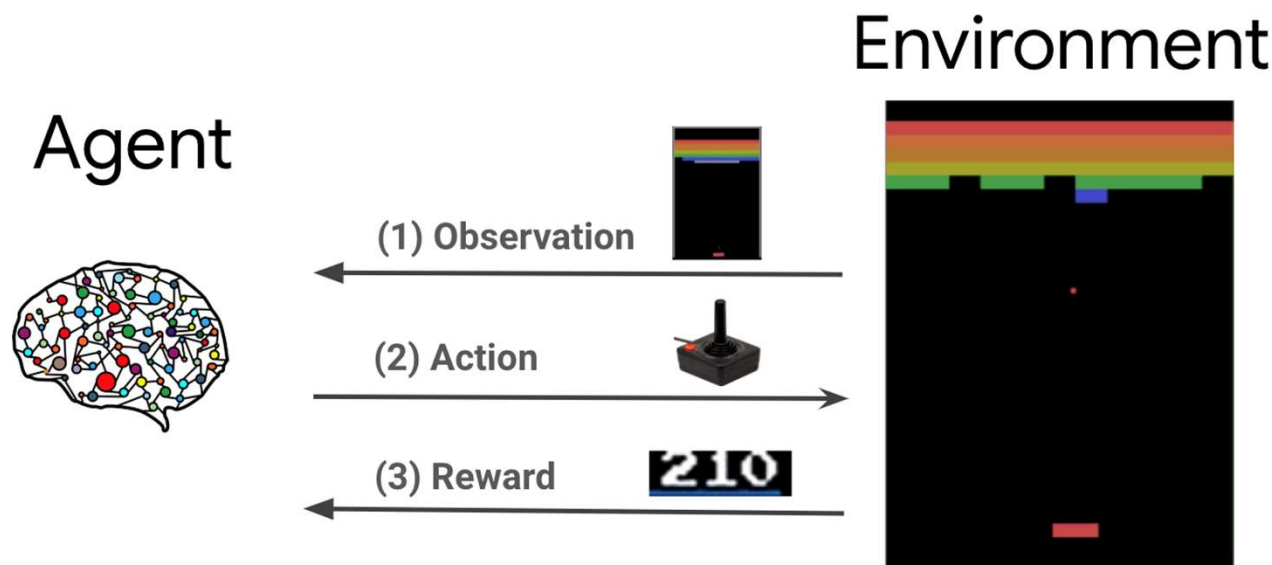


# Deep Q Network

---

# Deep Q Network (DQN)

- A representative Deep Reinforcement Learning (DRL) algorithm
- Proposed in 2013 paper "Playing Atari with Deep Reinforcement Learning".
- The paper showed that a DQN can be trained to play Atari games with human-level accuracy.



# What is DQN?

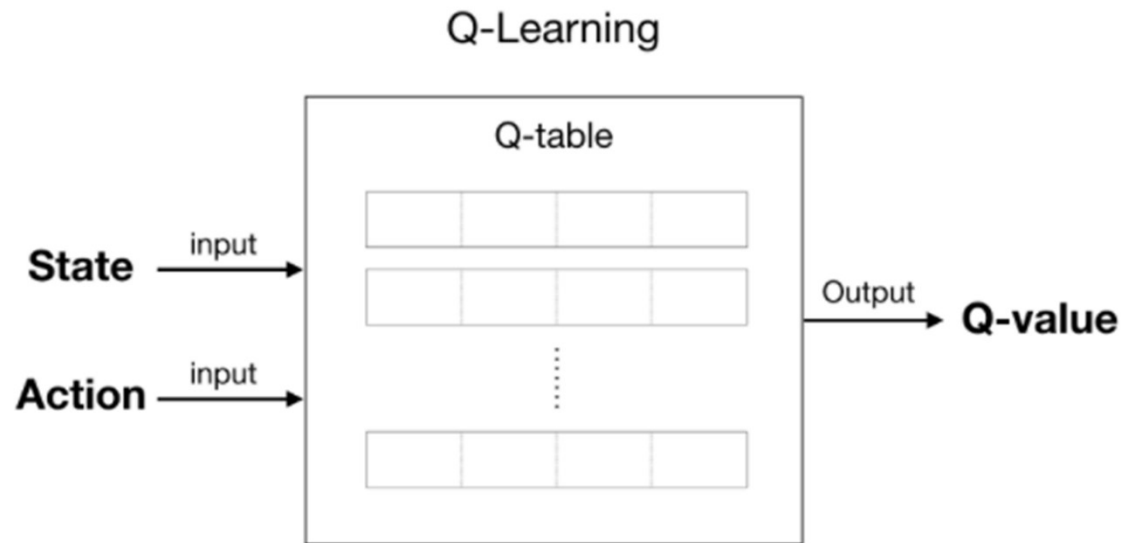
- The objective of reinforcement learning is to find the optimal policy.
- Optimal policy: policy that gives the maximum return.
- To get the policy, we compute the Q function.
- Once we have a Q function, we can extract a policy by choosing actions with maximum Q values.
- For example, if we have a Q table like this:

State	Action	Value
A	up	17
A	down	10
B	up	11
B	down	20

- Our policy is **{A: up, B: down}**.

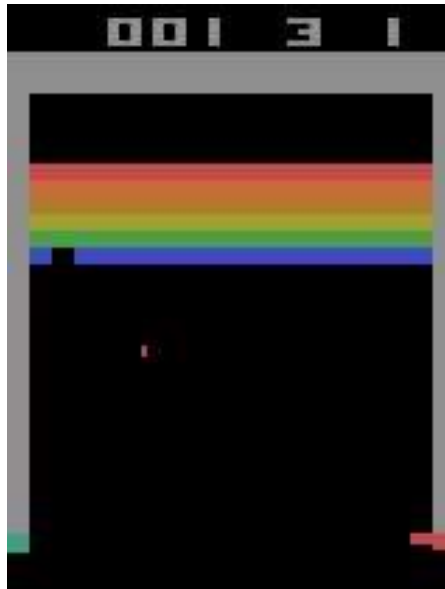
# What is DQN?

- In previous chapters, we used a Q table.
- "Learning" was done by updating the Q table iteratively.
- In a Q table, an entry consists of a (state, action) pair and its Q value.
- Basically, a Q table is a **function**.
- The **input** to the table is a (state, action).
- The **output** from the table is a **Q value**.



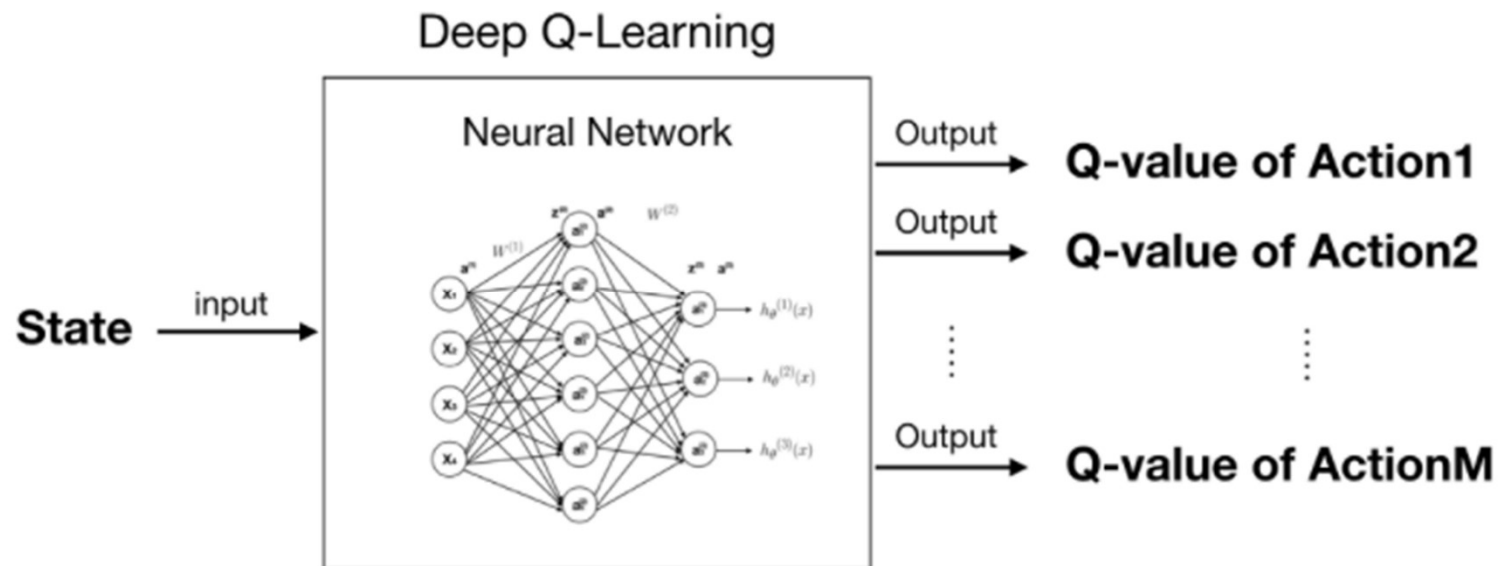
# What is DQN?

- The problem with using a Q table is when the number of states is too large.
- Say we have an environment with 100,000 states and 10 possible actions. Then we have 1,000,000 entries in our Q table.
- It will be very expensive to compute the Q values of all state-actions pairs.
- e.g.) The 'Breakout-v0' environment has  $210 \times 160 \times 3$  different states.



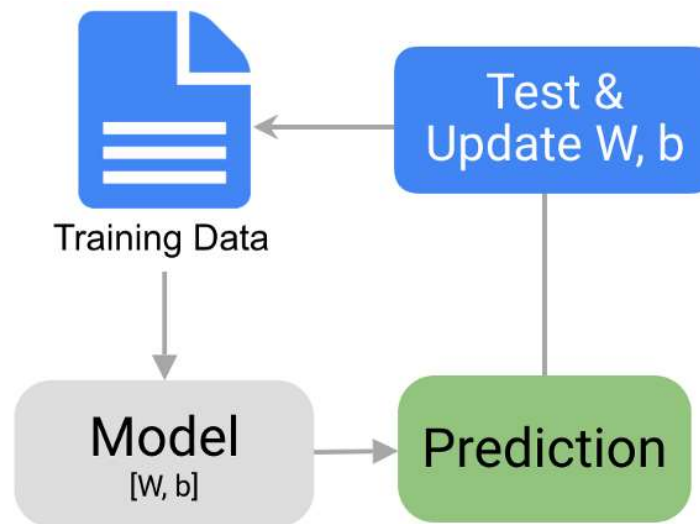
# What is DQN?

- Instead of using a Q-table, we use a neural network.
- The neural network is basically a **function approximator** for the Q function.
- The input to neural network is **a vector representing the state**.
- The output of the network is **the Q-values for each action in the action space**.
- The neural network can estimate Q values for unvisited states, based on the experience from similar states.



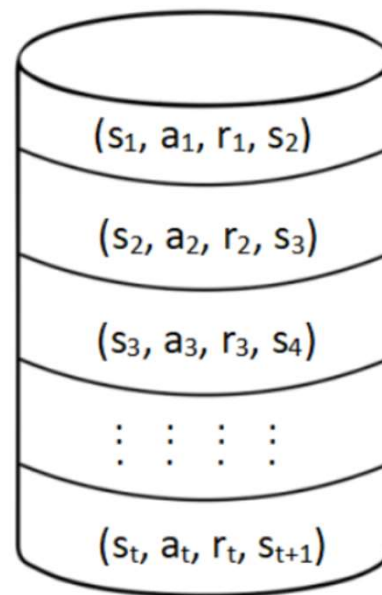
# Understanding DQN

- As Q tables need updates, a DQN must be trained.
- The DQN is trained using **supervised learning**, where we provide **training samples and their labels (target value)**.
- In supervised learning, **stochastic gradient descent** is typically used, where we provide a **batch of samples** to the network.



# Training Samples for DQN

- What are **training samples** we use to train a DQN?
- In Q-learning, we learn experience by going through episodes.
- In a state, the agent performs an action according to the policy (such as an epsilon-greedy policy), and moves to the next state.
- From this we get a **transition**, which is  $(s, a, r, s')$ . This transition becomes a training sample.
- We save the transitions in a buffer called **replay buffer**.





# Training Samples for DQN

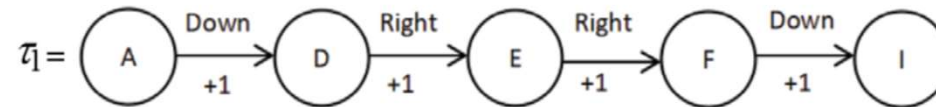
- Process for constructing a replay buffer

1. Initialize the replay buffer  $\mathcal{D}$ .
2. For each episode perform *step 3*.
3. For each step in the episode:
  1. Make a transition, that is, perform an action  $a$  in the state  $s$ , move to the next state  $s'$ , and receive the reward  $r$ .
  2. Store the transition information  $(s, a, r, s')$  in the replay buffer  $\mathcal{D}$ .

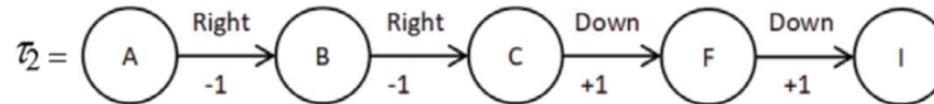
# Training Samples for DQN

- Example: running episodes in the Grid World

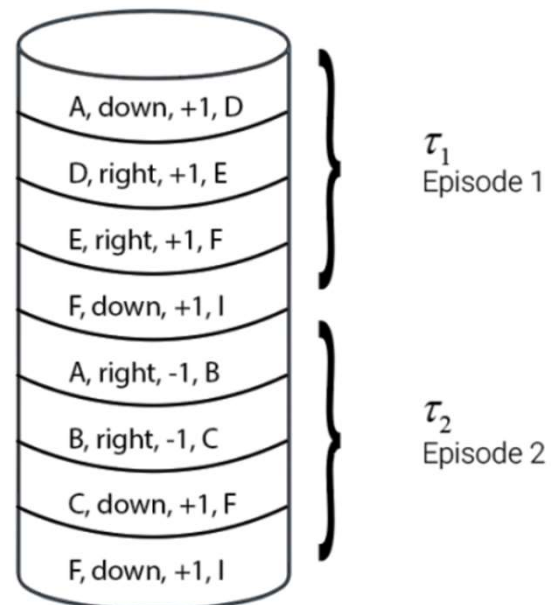
Episode 1:



Episode 2:



replay buffer



# Training Samples for DQN

- We collect transitions from many episodes and store them in the replay buffer.
- When training a model, we select **minibatches** from the replay buffer.
- In order to avoid correlation, we select **random samples** to create a minibatch.
- This process is called **Experience Replay**.
- Since we have a limited size for the replay buffer, old transition samples are replaced with new transition samples.

# Loss Function

- Our objective of using a neural network is to estimate Q values of state-action pairs. This is a **regression** task.
- For regression task, we generally use the mean squared error (MSE) as the loss function.

$$\text{MSE} = \frac{1}{K} \sum_{i=1}^K (y_i - \hat{y}_i)^2$$

- In the equation,  $y$  is the target value,  $\hat{y}$  is the predicted value, and  $K$  is the number of training samples in the minibatch.

# Loss Function

- The predicted value  $\hat{y}$  is the outcome of the model.
- What is the target value  $y$ ?
- In the Bellman optimality equation, the optimal Q value can be obtained using the following equation.

$$Q^*(s, a) = \mathbb{E}_{s' \sim p}[R(s, a, s') + \gamma \max_{a'} Q^*(s', a')]$$

- We remove the expectation from the equation. We will approximate the expectation by sampling K number of transitions from the replay buffer and taking the average value.

$$Q^*(s, a) = r + \gamma \max_{a'} Q^*(s', a')$$

- We denote  $R(s, a, s')$  as  $r$ .

# Loss Function

- Since we want our Q values to become optimal, we set the target Q value as:
  - $Q^*(s, a) = r + \gamma \max_{a'} Q(s', a')$
- Also, we denote the predicted value as:
  - $Q_\theta(s, a)$
- Thus, the difference between target value and predicted value is:
  - $Q^*(s, a) - Q_\theta(s, a) = r + \gamma \max_{a'} Q(s', a') - Q_\theta(s, a)$
- This is the temporal difference error used in the update rule of Q-learning.

# Loss Function

- We use the MSE loss. Thus, our loss function can be expressed as:

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \max_{a'} Q_{\theta}(s'_i, a') - Q_{\theta}(s_i, a_i))^2$$

- $y_i = r_i + \gamma \max_{a'} Q_{\theta}(s'_i, a')$ 
  - Since we only have a model  $\theta$  we extract the maximum Q value from the model.
- $\hat{y}_i = Q_{\theta}(s_i, a_i)$

- If the next state  $s'$  is a terminal state, we cannot compute the Q value because we do not take any action in the terminal state.
- Thus, if  $s'$  is terminal, we define  $y_i = r_i$ .
- In summary, our loss function will be:

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (y_i - Q_{\theta}(s_i, a_i))^2 \quad y_i = \begin{cases} r_i & \text{if } s' \text{ is terminal} \\ r_i + \gamma \max_{a'} Q_{\theta}(s'_i, a') & \text{if } s' \text{ is not terminal} \end{cases}$$

# The Target Network

- In our loss function, both the target value and the predicted value come from the same model.

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \max_{a'} \underbrace{Q_{\theta}(s'_i, a')}_{\text{Compute using } \theta} - \underbrace{Q_{\theta}(s_i, a_i)}_{\text{Compute using } \theta})^2$$

- This causes instability in the MSE and the network learns poorly. It also causes a lot of divergence during training.
  - When we update the network parameters  $\theta$ , both the target and the predicted value changes.
  - The predicted value keeps on trying to be the same as the target value, but the target value keeps on changing due to the update on the network parameter  $\theta$ .



# The Target Network

- It helps if we **freeze the target value** for a while and **compute only the predicted value** so that our predicted value matches the target value.
- In order to apply this idea, we prepare two models represented by parameters  $\theta$  and  $\theta'$ . They have exactly the same architecture.
- We freeze the target Q-network  $\theta'$  for a while, and update the main Q-network  $\theta$ .

$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \underbrace{\max_{a'} Q_{\theta'}(s'_i, a')}_{\substack{\text{Compute} \\ \text{using } \theta'}} - \underbrace{Q_{\theta}(s_i, a_i)}_{\substack{\text{Compute} \\ \text{using } \theta}})^2$$

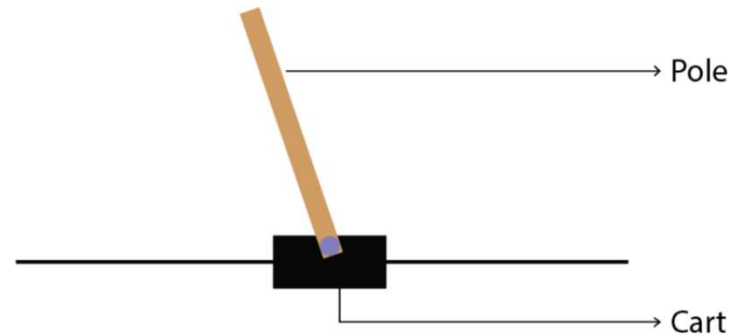
- Every once in a while, the parameters  $\theta$  is copied to  $\theta'$ .

# The DQN algorithm

1. Initialize the main network parameter  $\theta$  with random values
2. Initialize the target network parameter  $\theta'$  by copying the main network parameter  $\theta$
3. Initialize the replay buffer  $\mathcal{D}$
4. For  $N$  number of episodes, perform *step 5*
5. For each step in the episode, that is, for  $t = 0, \dots, T-1$ :
  1. Observe the state  $s$  and select an action using the epsilon-greedy policy, that is, with probability epsilon, select random action  $a$  and with probability 1-epsilon, select the action  $a = \arg \max_a Q_\theta(s, a)$
  2. Perform the selected action and move to the next state  $s'$  and obtain the reward  $r$
  3. Store the transition information in the replay buffer  $\mathcal{D}$
  4. Randomly sample a minibatch of  $K$  transitions from the replay buffer  $\mathcal{D}$
  5. Compute the target value, that is,  $y_i = r_i + \gamma \max_{a'} Q_{\theta'}(s'_i, a')$
  6. Compute the loss,  $L(\theta) = \frac{1}{K} \sum_{i=1}^K (y_i - Q_\theta(s_i, a_i))^2$
  7. Compute the gradients of the loss and update the main network parameter  $\theta$  using gradient descent:  $\theta = \theta - \alpha \nabla_\theta L(\theta)$
  8. Freeze the target network parameter  $\theta'$  for several time steps and then update it by just copying the main network parameter  $\theta$

# Cart-Pole Balancing using DQN

- State space: 4 continuous values
  - cart position
  - cart velocity
  - pole angle
  - pole velocity at the tip
- Action space: 2 discrete actions
  - push cart to the right
  - push cart to the left
- Reward
  - The agent acquires +1 reward for every timestep until the termination
- Terminating condition
  - The pole is more than 15 degrees from vertical
  - The cart moves more than 2.4 units from the center
- Since we have a continuous state space, it is difficult to use a Q-table for learning.



# Cart-Pole Balancing using DQN [ex016]

- libraries
  - we use the 'collections' library to manage the replay buffer.
    - we are going to use a **double-ended queue** (deque) for the buffer.
  - we use the 'random' library to sample a random subset from a list

```
# libraries
import gym
import collections
import random
```

- we use the **Pytorch** library to train a neural network.
  - other possibility is to use Keras/Tensorflow.

```
# pytorch library is used for deep learning
import torch
import torch.nn as nn
import torch.optim as optim
import torch.nn.functional as F
```

# Cart-Pole Balancing using DQN

- hyperparameters
  - learning rate: 0.0005
  - discount factor ( $\gamma$ ): 0.98
  - batch size: 32
  - size of the replay buffer: 50000

```
# hyperparameters
learning_rate = 0.0005
gamma = 0.98
buffer_limit = 50000      # size of replay buffer
batch_size = 32
```

# Cart-Pole Balancing using DQN

- class ReplayBuffer
  - '.\_\_init\_\_' initializes the replay buffer
  - 'put' adds a new transition to the buffer
  - 'sample' creates a batch by randomly selecting transitions from the buffer
    - It makes lists of **s**, **a**, **r**, **s'**, and **done**.
    - It also converts the list into tensors.
  - 'size' returns the number of transitions stored in the buffer

# Cart-Pole Balancing using DQN

```
class ReplayBuffer():
    def __init__(self):
        self.buffer = collections.deque(maxlen=buffer_limit)    # double-ended queue

    def put(self, transition):
        self.buffer.append(transition)

    def sample(self, n):
        mini_batch = random.sample(self.buffer, n)
        s_lst, a_lst, r_lst, s_prime_lst, done_mask_lst = [], [], [], [], []

        for transition in mini_batch:
            s, a, r, s_prime, done_mask = transition
            s_lst.append(s)
            a_lst.append([a])
            r_lst.append([r])
            s_prime_lst.append(s_prime)
            done_mask_lst.append([done_mask])

        return torch.tensor(s_lst, dtype=torch.float), torch.tensor(a_lst), \
            torch.tensor(r_lst), torch.tensor(s_prime_lst, dtype=torch.float), \
            torch.tensor(done_mask_lst)

    def size(self):
        return len(self.buffer)
```

# Cart-Pole Balancing using DQN

- class Qnet
  - '.\_\_init\_\_' defines layers of the model
    - nn.Linear(4, 128) is a fully connected layer with 4 inputs and 128 outputs
    - The model used here has one hidden layer with 128 neurons
  - 'forward' is called when an input is passed to the Qnet object
    - it takes a single argument, x, which is the input vectors
    - F.relu is applies ReLU to the input
  - 'sample\_action' selects an action according to the epsilon-greedy policy.
    - If 'coin' is smaller than epsilon, either 0 or 1 is chosen randomly.
    - If 'coin' is larger than epsilon, the action with the maximum Q value is chosen.



# Cart-Pole Balancing using DQN

```
class Qnet(nn.Module):
    def __init__(self):
        super(Qnet, self).__init__()
        self.fc1 = nn.Linear(4, 128)
        self.fc2 = nn.Linear(128, 128)
        self.fc3 = nn.Linear(128, 2)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x

    def sample_action(self, obs, epsilon):
        out = self.forward(obs)
        coin = random.random()
        if coin < epsilon:
            return random.randint(0,1)
        else :
            return out.argmax().item()
```

# Cart-Pole Balancing using DQN

- def train: trains the Q network using minibatches in the replay buffer
  - First, sample a batch from the replay buffer.
  - Pass the states  $s$  as the input to the model  $q$  and get the output  $q\_out$ .
  - From  $q\_out$ , select values for actions taken in each sample and assign to  $q\_a$ .
  - Calculate  $\max\_q\_prime$  which is  $\max_{a'} Q_{\theta}(s'_i, a')$  for all samples in the batch.
  - Calculate the target value  $target$ .
    - If  $done\_mask$  is 1, the target value is equal to  $r$ .
    - If  $done\_mask$  is 0, the target value is equal to  $r + \gamma \max_{a'} Q_{\theta}(s'_i, a')$ .
  - Calculate the MSE loss.
$$L(\theta) = \frac{1}{K} \sum_{i=1}^K (r_i + \gamma \max_{a'} Q_{\theta'}(s'_i, a') - Q_{\theta}(s_i, a_i))^2$$
  - Calculate the gradient.
$$\theta = \theta - \alpha \nabla_{\theta} L(\theta)$$
  - Update the parameters.

# Cart-Pole Balancing using DQN

```
def train(q, q_target, memory, optimizer):  
    for i in range(10):  
        s,a,r,s_prime,done_mask = memory.sample(batch_size)  
  
        q_out = q(s)  
        q_a = q_out.gather(1,a)  
        max_q_prime = q_target(s_prime).max(1)[0].unsqueeze(1)  
        target = r + gamma * max_q_prime * done_mask  
        loss = F.mse_loss(q_a, target)  
  
        optimizer.zero_grad()  
        loss.backward()  
        optimizer.step()
```

# Cart-Pole Balancing using DQN

- The main function
- Create the environment
- Create two Q networks: q, and q\_target.
- Create and initialize a replay buffer

```
def main():  
    env = gym.make('CartPole-v1')  
    q = Qnet()  
    q_target = Qnet()  
    q_target.load_state_dict(q.state_dict())  
    memory = ReplayBuffer()
```

# Cart-Pole Balancing using DQN

- Initialize variables
  - print\_interval: the interval for printing out the progress
    - After print\_interval, we also copy parameters from q to q\_target.
  - score: average score during a duration of print\_interval.
  - optimizer: the gradient descent algorithm
    - We use the Adam optimizer here.
    - Learning rate is given as an argument to the optimizer.

```
print_interval = 20  
score = 0.0  
optimizer = optim.Adam(q.parameters(), lr=learning_rate)
```

# Cart-Pole Balancing using DQN

- For each episode,
- Set the epsilon for the episode
  - epsilon is used for selecting actions using the epsilon-greedy policy.
  - In the first episode, epsilon is set to 0.08.
  - In the later episodes, epsilon is decreased linearly until it reaches 0.01.
  - We promote more exploration in the initial stage of training, but reduce exploration as we go.
- Reset the environment to the initial state.
- Set variable 'done' to False. The episode will end when 'done' becomes True.

```
for n_epi in range(2000):  
    epsilon = max(0.01, 0.08 - 0.01*(n_epi/200)) #Linear annealing from 8% to 1%  
    s = env.reset()  
    done = False
```

# Cart-Pole Balancing using DQN

- In the episode,
  - Sample an action using the epsilon-greedy policy.
  - Perform action and get the transition result ( $s'$ ,  $r$ , done).
  - Insert the transition into the replay buffer.
  - Move on to the next state.
  - Add reward to 'score'.
  - If the new state is a terminal state, end the episode.

```
while not done:
    a = q.sample_action(torch.from_numpy(s).float(), epsilon)
    s_prime, r, done, info = env.step(a)
    done_mask = 0.0 if done else 1.0
    memory.put((s,a,r/100.0,s_prime, done_mask))
    s = s_prime

    score += r
    if done:
        break
```

# Cart-Pole Balancing using DQN

- After one episode is over,
- Train model  $q$  using samples from the replay buffer.
  - We only do training if there are more than 2000 samples in the replay buffer.
- For each `print_interval`,
- Copy the parameters from  $q$  to  $q\_target$ .
- Print the average score, memory size, and the current epsilon value.
- Reset score to 0.

```
if memory.size()>2000:
    train(q, q_target, memory, optimizer)

if n_epi%print_interval==0 and n_epi!=0:
    q_target.load_state_dict(q.state_dict())
    print("n_episode : {}, score : {:.1f}, n_buffer : {}, eps : {:.1f}%".format(
        n_epi, score/print_interval, memory.size(), epsilon*100))
    score = 0.0
```



# Cart-Pole Balancing using DQN

- When we are done running the predefined number of episodes, we close the environment.
- Let the function main() be called when the file is executed.

```
env.close()

if __name__ == '__main__':
    main()
```

- Try running the file
  - It is considered successful if the average score exceeds 200 constantly.

# Double DQN

- “Deep Reinforcement Learning with Double Q-Learning”, *David Silver*, 2015
- Vanilla DQN has a problem of "overestimation".
- Overestimation
  - In DQN, the target value is computed as:
    - $y = r + \gamma \max_{a'} Q_{\theta'}(s', a')$ ,  $\theta'$  is the target network.
  - In other words, the value of (s, a) is the reward of (s, a) plus the maximum action value in the next state s'.
  - In Q-learning, we use  $\max_{a'}(Q(s', a'))$  to estimate this value.
  - There is an assumption here:  $\max_{a'}(Q(s', a'))$  is assumed to be the value of the best action in state s'.
  - However, especially in the early stage of learning, this may not be true.
  - $E(\max_{a'}(Q(s', a'))) \geq \max_{a'}(E(Q(s', a')))$ , which means  $\max_{a'}(Q(s', a'))$  is larger than any action value in average.
  - Because of this, the Q value is overestimated. In other words, the Q value will be calculated larger than what it should be.

# Double DQN

- What are the consequences of overestimation?
  - overestimation can slow down learning
  - non-best actions are overestimated and takes long time to drop
- Solution: Double DQN
  - Instead of selecting target value as:
    - $y = r + \gamma \max_{a'} Q_{\theta'}(s', a')$
  - We change the target value computation as:
    - $y = r + \gamma Q_{\theta'}(s', \arg \max_{a'} Q_{\theta}(s', a'))$
    - the main network parameter  $\theta$  is used for action selection.
    - the target network parameter  $\theta'$  is used for Q value computation.

# Double DQN

- Action ( $a'$ ) selection
  - We compute the Q values of all the next state-action pairs using the main network parameterized by  $\theta$ , and then we select action  $a'$  which has the maximum Q value.
  - $y = r + \gamma Q_{\theta'}(s', \arg \max_{a'} Q_{\theta}(s', a'))$
- Q value computation
  - Once we have selected action  $a'$ , then we compute the Q value using the target network parameterized by  $\theta'$  for the selected action  $a'$ .
  - $y = r + \gamma Q_{\theta'}(s', a')$

# Double DQN

- Implementation
  - From the code implementing DQN, we only need to change the part where we calculate the target value.
- DQN

```
q_out = q(s)
q_a = q_out.gather(1,a)

max_q_prime = q_target(s_prime).max(1)[0].unsqueeze(1)

target = r + gamma * max_q_prime * done_mask
loss = F.mse_loss(q_a, target)
```

- Double DQN

```
q_out = q(s)
q_a = q_out.gather(1,a)

argmax_Q = q(s_prime).max(1)[1].unsqueeze(1)
max_q_prime = q_target(s_prime).gather(1, argmax_Q)

target = r + gamma * max_q_prime * done_mask
loss = F.mse_loss(q_a, target)
```

# Dueling DQN

- "Dueling Network Architectures for Deep Reinforcement Learning", Z. Wang et al.
- **Advantage function**
  - The difference between the Q function and the value function
  - $A(s, a) = Q(s, a) - V(s)$
  - Q function: the expected return an agent would obtain starting from state  $s$ , performing action  $a$ , and following the policy  $\pi$ .
  - Value function: the expected return an agent would obtain starting from state  $s$  and following the policy  $\pi$ .
  - The advantage function tells us **how good the action  $a$  is compared to the average actions in state  $s$ .**

## Dueling DQN

- In DQN, we feed the state as input, and our network computes the Q value for all actions in the state.
- As an alternative approach, we can compute the Q values using the advantage function.
- We can rewrite the preceding equation as:
  - $Q(s, a) = V(s) + A(s, a)$
- We can compute the Q value by adding the value function and the advantage function.

# Dueling DQN

- Motivation
  - For some environments, it is unnecessary to know the value of each action at every timestep.
  - Example: Atari game Enduro, where it is not necessary to know which action to take until collision is imminent.
  - By separating two estimators (one for  $V$  and one for  $A$ ), the model can learn which states are valuable, without having to learn the effect of each action for each state.

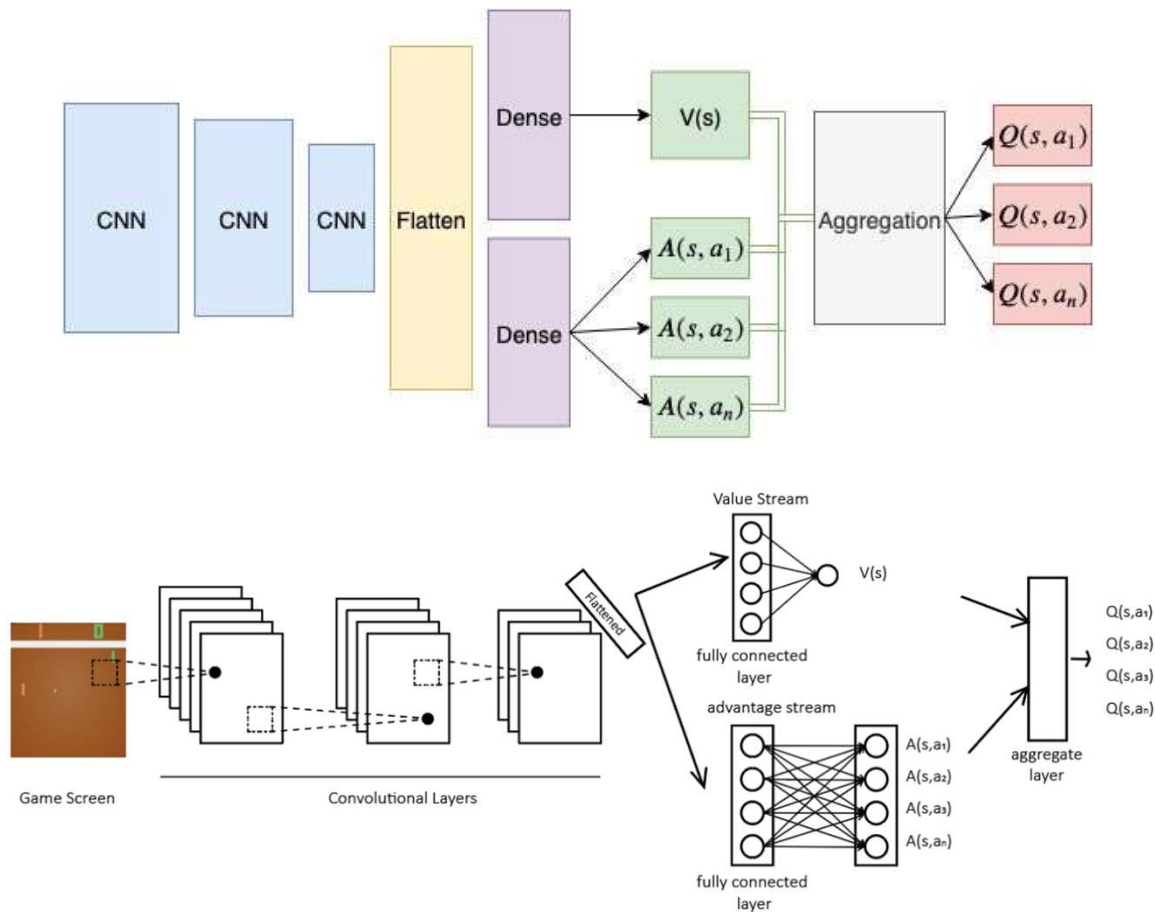




# Dueling DQN

- Model architecture

- The neural network model is separated into two models:  $V$  and  $A$ .
- Then, the model is aggregated by summing  $V$  and  $A$  to produce  $Q$ .



# Dueling DQN

- Benefit of using dueling architecture
  - Efficient learning of state-value function
  - With every update of the Q values, V is updated.
  - In a single-stream architecture, only the value for one of the actions is updated.
  - More frequent updating of the value in the dueling approach allocates more resources to V, and thus allows for better approximation of the state values, which in turn improves the accuracy of Q-learning.

# Dueling DQN

- Problem of identifiability
  - Since we are adding  $V$  and  $A$  to obtain  $Q$ , we cannot tell the values of  $V$  and  $A$  by looking at  $Q$ .
  - When trying to learn some optimal  $Q$  value for a state-action pair  $(s, a)$ , it is difficult to separately learn  $V(s)$  and  $A(s, a)$ .
  - Example
    - Suppose the "true" value of state  $s$  is 50, and the "true" advantage in state  $s$  for action  $a$  is 10. Then,  $Q(s, a)$  will be 60.
    - However,  $Q(s, a)$  is also 60 when  $V(s) = -1000$  and  $A(s, a) = 1060$ .
    - There is no guarantee that the "true" values of  $V(s)$  and  $A(s, a)$  are being learned separately and uniquely from each other.

# Dueling DQN

- Solution to the Problem of identifiability
  - Normalize the advantage function
  - We subtract the average advantage of all actions from  $A(s, a)$ .
  - $\mathcal{A}$  is the size of the action space.

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

# Dueling DQN

- Implementation
  - We can implement DuelingQnet class, and use it instead of Qnet.

```
class DuelingQnet(nn.Module):
    def __init__(self):
        super(DuelingQnet, self).__init__()
        self.fc1 = nn.Linear(4, 128)
        self.fc_value = nn.Linear(128, 128)
        self.fc_adv = nn.Linear(128, 128)
        self.value = nn.Linear(128, 1)
        self.adv = nn.Linear(128, 2)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        v = F.relu(self.fc_value(x))
        a = F.relu(self.fc_adv(x))
        v = self.value(v)
        a = self.adv(a)
        a_avg = torch.mean(a)
        q = v + a - a_avg
        return q
```

```
def sample_action(self, obs, epsilon):
    out = self.forward(obs)
    coin = random.random()
    if coin < epsilon:
        return random.randint(0,1)
    else :
        return out.argmax().item()
```

## Additional Topics: Prioritized Experience Replay

- In DQN, we randomly sample a minibatch of  $K$  transitions from the replay buffer and train the network.
- However, **some transitions may be more important to learn than others.**
- If so, we can assign priority to transitions and sample transitions based on their priorities.
- So which transitions can we consider as more important than others?

## Additional Topics: Prioritized Experience Replay

- For a transition, we calculate the TD error to calculate the loss function.
  - $\delta = r + \gamma \max_{a'} Q_{\theta'}(s', a') - Q_{\theta}(s, a)$
- A transition that has a higher TD error implies that the transition is not correct, so we need to learn more about that transition.
- A transition that has a low TD error implies that the transition is already good.
- We can learn more from our mistakes rather than something we are already good at.
- Similarly, we can assign higher priority to transitions with high TD error.
- Two methods for prioritizing transitions
  - proportional prioritization
  - rank-based prioritization

# Additional Topics: Prioritized Experience Replay

- Proportional Prioritization
  - We set the priority value  $p$  to be the TD error.
    - $p_i = |\delta_i|$
  - Using this equation,  $p_i$  will be 0 if the TD error is 0.
  - If a transition has 0 priority, it may never be used in training.
  - So we add a small number to the priority to avoid 0.
    - $p_i = |\delta_i| + \epsilon$
  - Instead of having the priority as a raw number, we can convert it into a probability value ranging from 0 to 1.
    - $P(i) = \frac{p_i}{\sum_k p_k}$
  - In this equation, we add a control parameter  $\alpha$  which has a value in the range of 0 and 1. The  $\alpha$  parameter decides how much probability we are going to assign to the transitions with high probability.
    - If  $\alpha$  is close to 1, we pick high priority transitions very frequently.
    - If  $\alpha$  is close to 0, we are close to choosing random transitions.
    - $P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$



# Additional Topics: Prioritized Experience Replay

- Rank-based prioritization
  - Rank of a transition  $i$ 
    - location of the transition in the replay buffer where the transitions are sorted from high TD error to low TD error.
  - We define the priority of transition  $i$  using rank as:
    - $p_i = \frac{1}{\text{Rank}(i)}$
  - We convert the priority into probability
    - $P(i) = \frac{p_i}{\sum_k p_k}$
  - We add control parameter  $\alpha$ 
    - $P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$

# Additional Topics: Prioritized Experience Replay

- Correcting the bias
  - When we use prioritized experience replay, it is highly biased towards samples with high priority. It can lead to overfitting.
  - To combat the bias, we use **important sampling**.
  - For a transition, an important weight is calculated as follows:
    - $w_i = \left(\frac{1}{N} \cdot \frac{1}{P(i)}\right)^\beta$
    - $N$  is the length of the replay buffer
    - $P(i)$  is the probability of the transition  $i$
    - $\beta$  is a parameter for controlling the importance weight
    - Typically, we start with small value (such as 0.4) and increase it towards 1.

## Additional Topics: Smooth L1 Loss

- MSE loss is also called L2 loss

$$Loss(y, f(x)) = \sum_{i=1}^N (y_i - f(x_i))^2$$

- We also have an L1 loss

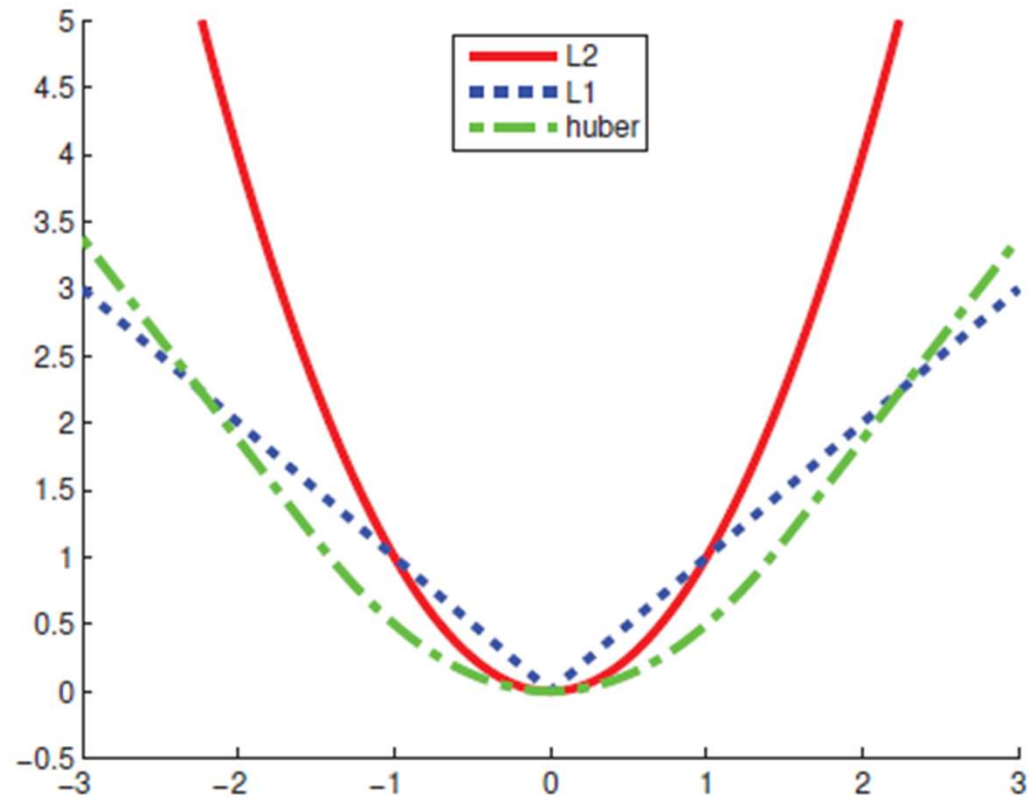
$$Loss(y, f(x)) = \sum_{i=1}^N |y_i - f(x_i)|$$

- The L1 loss has non-differentiable points but are more robust to outliers.
- The L2 loss is differentiable in all points but are weaker to outliers.
- Smooth L1 loss (also called Huber Loss) is a combination of the two.

$$Loss_{\delta}(y, f(x)) = \begin{cases} \frac{1}{2}((y_i - f(x_i))^2 & \text{for } |y_i - f(x_i)| \leq \delta, \\ \delta |y_i - f(x_i)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

# Additional Topics: Smooth L1 Loss

- Comparison of L2, L1, and Huber loss



End of Class

---

Questions?

Email: [jso1@sogang.ac.kr](mailto:jso1@sogang.ac.kr)