# ▾ Tutorial 5 - DQN and Actor-Critic

Please follow this tutorial to understand the structure (code) of DQNs & get familiar with Actor Critic methods.

References:

Please follow [Human-level control through deep reinforcement learning](https://...) for the original publication as well as the psuedocode. Watch Prof. Ravi's lectures on moodle or nptel for further understanding the core concepts. Contact the TAs for further resources if needed.

## Part 1: DQN

```
'''
Installing packages for rendering the game on Colab
'''

!pip install gym pyvirtualdisplay > /dev/null 2>&1
!apt-get install -y xvfb python-opengl ffmpeg > /dev/null 2>&1
!apt-get update > /dev/null 2>&1
!apt-get install cmake > /dev/null 2>&1
!pip install --upgrade setuptools 2>&1
!pip install ez_setup > /dev/null 2>&1
!pip install gym[atari] > /dev/null 2>&1
!pip install git+https://github.com/tensorflow/docs > /dev/null 2>&1
!pip install gym[classic_control]
```

```
    Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/
    Requirement already satisfied: setuptools in /usr/local/lib/python3.8/dist-packages (
    Collecting setuptools
      Downloading setuptools-67.4.0-py3-none-any.whl (1.1 MB)
                                                 1.1/1.1 MB 42.2 MB/s eta 0:00:00
    Installing collected packages: setuptools
      Attempting uninstall: setuptools
        Found existing installation: setuptools 57.4.0
        Uninstalling setuptools-57.4.0:
          Successfully uninstalled setuptools-57.4.0
    ERROR: pip's dependency resolver does not currently take into account all the package
    ipython 7.9.0 requires jedi>=0.10, which is not installed.
    cvxpy 1.2.3 requires setuptools<=64.0.2, but you have setuptools 67.4.0 which is inco
    Successfully installed setuptools-67.4.0
    Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/
    Requirement already satisfied: gym[classic_control] in /usr/local/lib/python3.8/dist-
    Requirement already satisfied: cloudpickle>=1.2.0 in /usr/local/lib/python3.8/dist-pa
    Requirement already satisfied: numpy>=1.18.0 in /usr/local/lib/python3.8/dist-package
    Requirement already satisfied: gym-notices>=0.0.4 in /usr/local/lib/python3.8/dist-pa
    Requirement already satisfied: importlib-metadata>=4.8.0 in /usr/local/lib/python3.8/
    Collecting pygame==2.1.0
      Downloading pygame-2.1.0-cp38-cp38-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (
                                                 18.3/18.3 MB 78.9 MB/s eta 0:00:00
    Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.8/dist-packages (
```

```
       Installing collected packages: pygame
       Successfully installed pygame-2.1.0
```

```python
'''
A bunch of imports, you don't have to worry about these
'''

import numpy as np
import random
import torch
import torch.nn as nn
import torch.nn.functional as F
from collections import namedtuple, deque
import torch.optim as optim
import datetime
import gym
from gym.wrappers.record_video import RecordVideo
import glob
import io
import base64
import matplotlib.pyplot as plt
from IPython.display import HTML
from pyvirtualdisplay import Display
import tensorflow as tf
from IPython import display as ipythondisplay
from PIL import Image
import tensorflow_probability as tfp
```

```python
'''
Please refer to the first tutorial for more details on the specifics of environments
We've only added important commands you might find useful for experiments.
'''

'''
List of example environments
(Source - https://gym.openai.com/envs/#classic_control)

'Acrobot-v1'
'Cartpole-v1'
'MountainCar-v0'
'''

env = gym.make('CartPole-v1')
env.seed(0)

state_shape = env.observation_space.shape[0]
no_of_actions = env.action_space.n

print(state_shape)
print(no_of_actions)
print(env.action_space.sample())
print("----")
```

```
'''
# Understanding State, Action, Reward Dynamics

The agent decides an action to take depending on the state.

The Environment keeps a variable specifically for the current state.
- Everytime an action is passed to the environment, it calculates the new state and update
- It returns the new current state and reward for the agent to take the next action

'''

state = env.reset()
''' This returns the initial state (when environment is reset) '''

print(state)
print("----")

action = env.action_space.sample()
''' We take a random action now '''

print(action)
print("----")

next_state, reward, done, info = env.step(action)
''' env.step is used to calculate new state and obtain reward based on old state and actio

print(next_state)
print(reward)
print(done)
print(info)
print("----")
```

```
    4
    2
    1
    ----
    [ 0.01369617 -0.02302133 -0.04590265 -0.04834723]
    ----
    1
    ----
    [ 0.01323574  0.17272775 -0.04686959 -0.3551522 ]
    1.0
    False
    {}
    ----
    /usr/local/lib/python3.8/dist-packages/gym/core.py:317: DeprecationWarning: WARN: In:
      deprecation(
    /usr/local/lib/python3.8/dist-packages/gym/wrappers/step_api_compatibility.py:39: Dep
      deprecation(
    /usr/local/lib/python3.8/dist-packages/gym/core.py:256: DeprecationWarning: WARN: Fur
      deprecation(
```

# ▾ DQN

Using NNs as substitutes isn't something new. It has been tried earlier, but the 'human control' paper really popularised using NNs by providing a few stability ideas (Q-Targets, Experience Replay & Truncation). The 'Deep-Q Network' (DQN) Algorithm can be broken down into having the following components.

## Q-Network:

The neural network used as a function approximator is defined below

```
'''
### Q Network & Some 'hyperparameters'

QNetwork1:
Input Layer - 4 nodes (State Shape) \
Hidden Layer 1 - 64 nodes \
Hidden Layer 2 - 64 nodes \
Output Layer - 2 nodes (Action Space) \
Optimizer - zero_grad()

QNetwork2: Feel free to experiment more
'''

import torch
import torch.nn as nn
import torch.nn.functional as F


'''
Bunch of Hyper parameters (Which you might have to tune later **wink wink**)
'''
BUFFER_SIZE = int(1e5)  # replay buffer size
BATCH_SIZE = 64         # minibatch size
GAMMA = 0.99            # discount factor
LR = 5e-4               # learning rate
UPDATE_EVERY = 20       # how often to update the network (When Q target is present)


class QNetwork1(nn.Module):

    def __init__(self, state_size, action_size, seed, fc1_units=128, fc2_units=64):
        """Initialize parameters and build model.
        Params
        ======
            state_size (int): Dimension of each state
            action_size (int): Dimension of each action
            seed (int): Random seed
            fc1_units (int): Number of nodes in first hidden layer
            fc2_units (int): Number of nodes in second hidden layer
        """
```

```
        super(QNetwork1, self).__init__()
        self.seed = torch.manual_seed(seed)
        self.fc1 = nn.Linear(state_size, fc1_units)
        self.fc2 = nn.Linear(fc1_units, fc2_units)
        self.fc3 = nn.Linear(fc2_units, action_size)

    def forward(self, state):
        """Build a network that maps state -> action values."""
        x = F.relu(self.fc1(state))
        x = F.relu(self.fc2(x))
        return self.fc3(x)
```

▼ Replay Buffer:

This is a 'deque' that helps us store experiences. Recall why we use such a technique.

```
import random
import torch
import numpy as np
from collections import deque, namedtuple

device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

class ReplayBuffer:
    """Fixed-size buffer to store experience tuples."""

    def __init__(self, action_size, buffer_size, batch_size, seed):
        """Initialize a ReplayBuffer object.

        Params
        ======
            action_size (int): dimension of each action
            buffer_size (int): maximum size of buffer
            batch_size (int): size of each training batch
            seed (int): random seed
        """
        self.action_size = action_size
        self.memory = deque(maxlen=buffer_size)
        self.batch_size = batch_size
        self.experience = namedtuple("Experience", field_names=["state", "action", "reward
        self.seed = random.seed(seed)

    def add(self, state, action, reward, next_state, done):
        """Add a new experience to memory."""
        e = self.experience(state, action, reward, next_state, done)
        self.memory.append(e)

    def sample(self):
        """Randomly sample a batch of experiences from memory."""
        experiences = random.sample(self.memory, k=self.batch_size)

        states = torch.from_numpy(np.vstack([e.state for e in experiences if e is not None
        actions = torch.from_numpy(np.vstack([e.action for e in experiences if e is not No
```

```
        rewards = torch.from_numpy(np.vstack([e.reward for e in experiences if e is not No
        next_states = torch.from_numpy(np.vstack([e.next_state for e in experiences if e i
        dones = torch.from_numpy(np.vstack([e.done for e in experiences if e is not None])

        return (states, actions, rewards, next_states, dones)

    def __len__(self):
        """Return the current size of internal memory."""
        return len(self.memory)
```

## Truncation:

We add a line (optionally) in the code to truncate the gradient in hopes that it would help with the stability of the learning process.

## Tutorial Agent Code:

```python
class TutorialAgent():

    def __init__(self, state_size, action_size, seed):

        ''' Agent Environment Interaction '''
        self.state_size = state_size
        self.action_size = action_size
        self.seed = random.seed(seed)

        ''' Q-Network '''
        self.qnetwork_local = QNetwork1(state_size, action_size, seed).to(device)
        self.qnetwork_target = QNetwork1(state_size, action_size, seed).to(device)
        self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)

        ''' Replay memory '''
        self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)

        ''' Initialize time step (for updating every UPDATE_EVERY steps)        -Needed
        self.t_step = 0

    def step(self, state, action, reward, next_state, done):

        ''' Save experience in replay memory '''
        self.memory.add(state, action, reward, next_state, done)

        ''' If enough samples are available in memory, get random subset and learn '''
        if len(self.memory) >= BATCH_SIZE:
            experiences = self.memory.sample()
            self.learn(experiences, GAMMA)

        """ +Q TARGETS PRESENT """
        ''' Updating the Network every 'UPDATE_EVERY' steps taken '''
        self.t_step = (self.t_step + 1) % UPDATE_EVERY
        if self.t_step == 0:
```

```python
                    self.qnetwork_target.load_state_dict(self.qnetwork_local.state_dict())

    def act(self, state, eps=0.):

        state = torch.from_numpy(state).float().unsqueeze(0).to(device)
        self.qnetwork_local.eval()
        with torch.no_grad():
            action_values = self.qnetwork_local(state)
        self.qnetwork_local.train()

        ''' Epsilon-greedy action selection (Already Present) '''
        if random.random() > eps:
            return np.argmax(action_values.cpu().data.numpy())
        else:
            return random.choice(np.arange(self.action_size))

    def learn(self, experiences, gamma):
        """ +E EXPERIENCE REPLAY PRESENT """
        states, actions, rewards, next_states, dones = experiences

        ''' Get max predicted Q values (for next states) from target model'''
        Q_targets_next = self.qnetwork_target(next_states).detach().max(1)[0].unsqueeze(1)

        ''' Compute Q targets for current states '''
        Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))

        ''' Get expected Q values from local model '''
        Q_expected = self.qnetwork_local(states).gather(1, actions)

        ''' Compute loss '''
        loss = F.mse_loss(Q_expected, Q_targets)

        ''' Minimize the loss '''
        self.optimizer.zero_grad()
        loss.backward()

        ''' Gradiant Clipping '''
        """ +T TRUNCATION PRESENT """
        for param in self.qnetwork_local.parameters():
            param.grad.data.clamp_(-1, 1)

        self.optimizer.step()
```

## ▾ Here, we present the DQN algorithm code.

```python
''' Defining DQN Algorithm '''

state_shape = env.observation_space.shape[0]
action_shape = env.action_space.n

def dqn(n_episodes=10000, max_t=1000, eps_start=1.0, eps_end=0.01, eps_decay=0.995):

    scores = []
```

```python
    ''' list containing scores from each episode '''

    scores_window_printing = deque(maxlen=10)
    ''' For printing in the graph '''

    scores_window= deque(maxlen=100)
    ''' last 100 scores for checking if the avg is more than 195 '''

    eps = eps_start
    ''' initialize epsilon '''

    for i_episode in range(1, n_episodes+1):
        state = env.reset()
        score = 0
        for t in range(max_t):
            action = agent.act(state, eps)
            next_state, reward, done, _ = env.step(action)
            agent.step(state, action, reward, next_state, done)
            state = next_state
            score += reward
            if done:
                break

        scores_window.append(score)
        scores_window_printing.append(score)
        ''' save most recent score '''

        eps = max(eps_end, eps_decay*eps)
        ''' decrease epsilon '''

        print('\rEpisode {}\tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_windo
        if i_episode % 10 == 0:
            scores.append(np.mean(scores_window_printing))
        if i_episode % 100 == 0:
            print('\rEpisode {}\tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_wi
        if np.mean(scores_window)>=195.0:
            print('\nEnvironment solved in {:d} episodes!\tAverage Score: {:.2f}'.format(i_
            break
    return [np.array(scores),i_episode-100]

''' Trial run to check if algorithm runs and saves the data '''

begin_time = datetime.datetime.now()
agent = TutorialAgent(state_size=state_shape,action_size = action_shape,seed = 0)



dqn()


time_taken = datetime.datetime.now() - begin_time

print(time_taken)
```

```
    Episode 100      Average Score: 38.24
    Episode 200      Average Score: 144.32
```

```
Episode 231      Average Score: 195.80
Environment solved in 131 episodes!      Average Score: 195.80
0:01:58.300658
```

## ▾ Task 1a

Understand the core of the algorithm, follow the flow of data. Identify the exploration strategy used.

### Task 1b

Out of the two exploration strategies discussed in class ($\epsilon$-greedy & Softmax). Implement the strategy that's not used here.

### Task 1c

How fast does the agent 'solve' the environment in terms of the number of episodes? (Cartpole-v1 defines "solving" as getting average reward of 195.0 over 100 consecutive trials)

How 'well' does the agent learn? (reward plot?) The above two are some 'evaluation metrics' you can use to comment on the performance of an algorithm.

Please compare DQN (using $\epsilon$-greedy) with DQN (using softmax). Think along the lines of 'no. of episodes', 'reward plots', 'compute time', etc. and add a few comments.

### Submission Steps

Task 1: Add a text cell with the answer.

Task 2: Add a code cell below task 1 solution and use 'Tutorial Agent Code' to build your new agent (with a different exploration strategy).

Task 3: Add a code cell below task 2 solution running both the agents to solve the CartPole v-1 environment and add a new text cell below it with your inferences.

## ▾ Solution for Task 1

Task 1a): The exploration strategy used is $\epsilon$-greedy.

Task 1b)

```
class TutorialAgent():

    def __init__(self, state_size, action_size, seed):
```

```python
        ''' Agent Environment Interaction '''
        self.state_size = state_size
        self.action_size = action_size
        self.seed = random.seed(seed)

        ''' Q-Network '''
        self.qnetwork_local = QNetwork1(state_size, action_size, seed).to(device)
        self.qnetwork_target = QNetwork1(state_size, action_size, seed).to(device)
        self.optimizer = optim.Adam(self.qnetwork_local.parameters(), lr=LR)

        ''' Replay memory '''
        self.memory = ReplayBuffer(action_size, BUFFER_SIZE, BATCH_SIZE, seed)

        ''' Initialize time step (for updating every UPDATE_EVERY steps)          -Needed
        self.t_step = 0

    def step(self, state, action, reward, next_state, done):

        ''' Save experience in replay memory '''
        self.memory.add(state, action, reward, next_state, done)

        ''' If enough samples are available in memory, get random subset and learn '''
        if len(self.memory) >= BATCH_SIZE:
            experiences = self.memory.sample()
            self.learn(experiences, GAMMA)

        """ +Q TARGETS PRESENT """
        ''' Updating the Network every 'UPDATE_EVERY' steps taken '''
        self.t_step = (self.t_step + 1) % UPDATE_EVERY
        if self.t_step == 0:

            self.qnetwork_target.load_state_dict(self.qnetwork_local.state_dict())

    def act(self, state, tau=0.1):

        state = torch.from_numpy(state).float().unsqueeze(0).to(device)
        self.qnetwork_local.eval()
        with torch.no_grad():
            action_values = self.qnetwork_local(state)
        self.qnetwork_local.train()

        ''' Softmax action selection (Implemented) '''
        action_values=action_values.cpu().data.numpy()

        x=np.array([action_values[0][id]/tau for id in range(self.action_size)])
        pi=(np.exp(x - np.max(x)) / np.exp(x - np.max(x)).sum())

        return np.random.choice(np.arange(self.action_size),p=pi)

    def learn(self, experiences, gamma):
        """ +E EXPERIENCE REPLAY PRESENT """
        states, actions, rewards, next_states, dones = experiences

        ''' Get max predicted Q values (for next states) from target model'''
```

```python
        Q_targets_next = self.qnetwork_target(next_states).detach().max(1)[0].unsqueeze(1)

        ''' Compute Q targets for current states '''
        Q_targets = rewards + (gamma * Q_targets_next * (1 - dones))

        ''' Get expected Q values from local model '''
        Q_expected = self.qnetwork_local(states).gather(1, actions)

        ''' Compute loss '''
        loss = F.mse_loss(Q_expected, Q_targets)

        ''' Minimize the loss '''
        self.optimizer.zero_grad()
        loss.backward()

        ''' Gradiant Clipping '''
        """ +T TRUNCATION PRESENT """
        for param in self.qnetwork_local.parameters():
            param.grad.data.clamp_(-1, 1)

        self.optimizer.step()
```

## Task 1c

```python
state_shape = env.observation_space.shape[0]
action_shape = env.action_space.n

def dqn(n_episodes=10000, max_t=1000, tau=1):

    scores = []
    ''' list containing scores from each episode '''

    scores_window_printing = deque(maxlen=10)
    ''' For printing in the graph '''

    scores_window= deque(maxlen=100)
    ''' last 100 scores for checking if the avg is more than 195 '''


    for i_episode in range(1, n_episodes+1):
        state = env.reset()
        score = 0
        for t in range(max_t):
            action = agent.act(state,tau)
            next_state, reward, done, _ = env.step(action)
            agent.step(state, action, reward, next_state, done)
            state = next_state
            score += reward
            if done:
                break

        scores_window.append(score)
```

```
            scores_window_printing.append(score)
            ''' save most recent score '''


        print('\rEpisode {}\tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_windo
        if i_episode % 10 == 0:
            scores.append(np.mean(scores_window_printing))
        if i_episode % 100 == 0:
            print('\rEpisode {}\tAverage Score: {:.2f}'.format(i_episode, np.mean(scores_wi
        if np.mean(scores_window)>=195.0:
            print('\nEnvironment solved in {:d} episodes!\tAverage Score: {:.2f}'.format(i_
            break
    return [np.array(scores),i_episode-100]

''' Trial run to check if algorithm runs and saves the data '''

begin_time = datetime.datetime.now()
agent = TutorialAgent(state_size=state_shape,action_size = action_shape,seed = 0)


dqn()


time_taken = datetime.datetime.now() - begin_time

print(time_taken)
```
```
    Episode 100     Average Score: 144.54
    Episode 112     Average Score: 196.87
    Environment solved in 12 episodes!     Average Score: 196.87
    0:01:18.291599
```

# ▾ Task 1c

We need to add +100 to episode count as when priniting we have used i_episode-100 when environment is solved.

Sample outputs: For tau-0.1 Environment solved in 819 episodes! Average Score: 196.80 0:05:08.328808

For tau=1 Environment solved in 12 episodes! Average Score: 196.87 0:01:18.291599

For tau=10 Environment solved in -16 episodes! Average Score: 197.21 0:01:02.774652

Inferences:

1. For increasing values of tau we see that the algorithm converges faster than epsilon-greedy algorithm.
2. Computation time- Softmax takes more time to choose action since at each step we require to calculate the probabilities, and as when the value function explodes or tends to zero, it can become quite become computationally intensive.

3. The reward based plots indicate that both agents could learn to play the game. As tau increases there is more exploration and the entropy of the system increases, leading to faster learning of the agent.

4. Q learning does not define a policy. We could do a softmax over the state action values with some dampening/sharpening, but it is probably difficult to tune this parameter such that it works for multiple environments with different reward distributions. This could also be said for entropy regularization, but since policy-based methods already define a policy, an entropy term might be more obvious and easier-to-interpret in that case. Thus epsilon-greedy is more commonly used in off policy algos.

## Part 2: One-Step Actor-Critic Algorithm

**Actor-Critic methods** learn both a policy $\pi(a|s;\theta)$ and a state-value function $v(s;w)$ simultaneously. The policy is referred to as the actor that suggests actions given a state. The estimated value function is referred to as the critic. It evaluates actions taken by the actor based on the given policy. In this exercise, both functions are approximated by feedforward neural networks.

- The policy network is parametrized by $\theta$ - it takes a state $s$ as input and outputs the probabilities $\pi(a|s;\theta) \ \forall \ a$
- The value network is parametrized by $w$ - it takes a state $s$ as input and outputs a scalar value associated with the state, i.e., $v(s;w)$
- The single step TD error can be defined as follows:
$$\delta_t = R_{t+1} + \gamma v(s_{t+1};w) - v(s_t;w)$$
- The loss function to be minimized at every step ($L_{tot}^{(t)}$) is a summation of two terms, as follows:
$$L_{tot}^{(t)} = L_{actor}^{(t)} + L_{critic}^{(t)}$$
  where,
$$L_{actor}^{(t)} = -\log \pi(a_t|s_t;\theta)\delta_t$$
$$L_{critic}^{(t)} = \delta_t^2$$
- **NOTE: Here, weights of the first two hidden layers are shared by the policy and the value network**

  - First two hidden layer sizes: [1024, 512]
  - Output size of policy network: 2 (Softmax activation)
  - Output size of value network: 1 (Linear activation)

## Initializing Actor-Critic Network

```python
class ActorCriticModel(tf.keras.Model):
    """
    Defining policy and value networkss
    """
    def __init__(self, action_size, n_hidden1=1024, n_hidden2=512):
        super(ActorCriticModel, self).__init__()

        #Hidden Layer 1
        self.fc1 = tf.keras.layers.Dense(n_hidden1, activation='relu')
        #Hidden Layer 2
        self.fc2 = tf.keras.layers.Dense(n_hidden2, activation='relu')

        #Output Layer for policy
        self.pi_out = tf.keras.layers.Dense(action_size, activation='softmax')
        #Output Layer for state-value
        self.v_out = tf.keras.layers.Dense(1)

    def call(self, state):
        """
        Computes policy distribution and state-value for a given state
        """
        layer1 = self.fc1(state)
        layer2 = self.fc2(layer1)

        pi = self.pi_out(layer2)
        v = self.v_out(layer2)

        return pi, v
```

## Agent Class

**Task 2a:** Write code to compute $\delta_t$ inside the Agent.learn() function

```python
class Agent:
    """
    Agent class
    """
    def __init__(self, action_size, lr=0.001, gamma=0.99, seed = 85):
        self.gamma = gamma
        self.ac_model = ActorCriticModel(action_size=action_size)
        self.ac_model.compile(tf.keras.optimizers.Adam(learning_rate=lr))
        np.random.seed(seed)

    def sample_action(self, state):
        """
        Given a state, compute the policy distribution over all actions and sample one act
        """
        pi,_ = self.ac_model(state)

        action_probabilities = tfp.distributions.Categorical(probs=pi)
        sample = action_probabilities.sample()
```

```python
            return int(sample.numpy()[0])

    def actor_loss(self, action, pi, delta):
        """
        Compute Actor Loss
        """
        return -tf.math.log(pi[0,action]) * delta

    def critic_loss(self,delta):
        """
        Critic loss aims to minimize TD error
        """
        return delta**2

    @tf.function
    def learn(self, state, action, reward, next_state, done):
        """
        For a given transition (s,a,s',r) update the paramters by computing the
        gradient of the total loss
        """
        with tf.GradientTape(persistent=True) as tape:
            pi, V_s = self.ac_model(state)
            _, V_s_next = self.ac_model(next_state)

            V_s = tf.squeeze(V_s)
            V_s_next = tf.squeeze(V_s_next)


            #### TO DO: Write the equation for delta (TD error)
            ## Write code below
            delta =reward+self.gamma*V_s_next-V_s
            loss_a = self.actor_loss(action, pi, delta)
            loss_c =self.critic_loss(delta)
            loss_total = loss_a + loss_c

        gradient = tape.gradient(loss_total, self.ac_model.trainable_variables)
        self.ac_model.optimizer.apply_gradients(zip(gradient, self.ac_model.trainable_vari
```

## Train the Network

```python
env = gym.make('CartPole-v1')

#Initializing Agent
agent = Agent(lr=1e-4, action_size=env.action_space.n)
#Number of episodes
episodes = 1800
tf.compat.v1.reset_default_graph()

reward_list = []
average_reward_list = []
begin_time = datetime.datetime.now()

for ep in range(1, episodes + 1):
```

```python
        state = env.reset().reshape(1,-1)
        done = False
        ep_rew = 0
        while not done:
            action = agent.sample_action(state) ##Sample Action
            next_state, reward, done, info = env.step(action) ##Take action
            next_state = next_state.reshape(1,-1)
            ep_rew += reward  ##Updating episode reward
            agent.learn(state, action, reward, next_state, done) ##Update Parameters

            state = next_state ##Updating State
        reward_list.append(ep_rew)

        if ep % 10 == 0:
            avg_rew = np.mean(reward_list[-10:])
            print('Episode ', ep, 'Reward %f' % ep_rew, 'Average Reward %f' % avg_rew)

        if ep % 100:
            avg_100 =  np.mean(reward_list[-100:])
            if avg_100 > 195.0:
                print('Stopped at Episode ',ep-100)
                break

time_taken = datetime.datetime.now() - begin_time
print(time_taken)
```

```
    /usr/local/lib/python3.8/dist-packages/gym/core.py:317: DeprecationWarning: WARN: In:
      deprecation(
    /usr/local/lib/python3.8/dist-packages/gym/wrappers/step_api_compatibility.py:39: Dep
      deprecation(
    Episode  10 Reward 21.000000 Average Reward 34.700000
    Episode  20 Reward 44.000000 Average Reward 36.700000
    Episode  30 Reward 76.000000 Average Reward 58.400000
    Episode  40 Reward 62.000000 Average Reward 65.700000
    Episode  50 Reward 49.000000 Average Reward 71.200000
    Episode  60 Reward 54.000000 Average Reward 51.400000
    Episode  70 Reward 45.000000 Average Reward 44.500000
    Episode  80 Reward 90.000000 Average Reward 75.800000
    Episode  90 Reward 121.000000 Average Reward 94.300000
    Episode  100 Reward 136.000000 Average Reward 120.800000
    Episode  110 Reward 104.000000 Average Reward 135.500000
    Episode  120 Reward 134.000000 Average Reward 111.900000
    Episode  130 Reward 111.000000 Average Reward 125.200000
    Episode  140 Reward 51.000000 Average Reward 59.200000
    Episode  150 Reward 64.000000 Average Reward 68.300000
    Episode  160 Reward 35.000000 Average Reward 55.500000
    Episode  170 Reward 70.000000 Average Reward 59.100000
    Episode  180 Reward 197.000000 Average Reward 105.400000
    Episode  190 Reward 157.000000 Average Reward 189.700000
    Episode  200 Reward 148.000000 Average Reward 116.700000
    Episode  210 Reward 129.000000 Average Reward 145.400000
    Episode  220 Reward 55.000000 Average Reward 105.700000
    Episode  230 Reward 181.000000 Average Reward 131.100000
    Episode  240 Reward 133.000000 Average Reward 214.100000
    Episode  250 Reward 201.000000 Average Reward 177.700000
    Episode  260 Reward 197.000000 Average Reward 304.500000
    Episode  270 Reward 493.000000 Average Reward 275.700000
    Episode  280 Reward 425.000000 Average Reward 306.300000
```
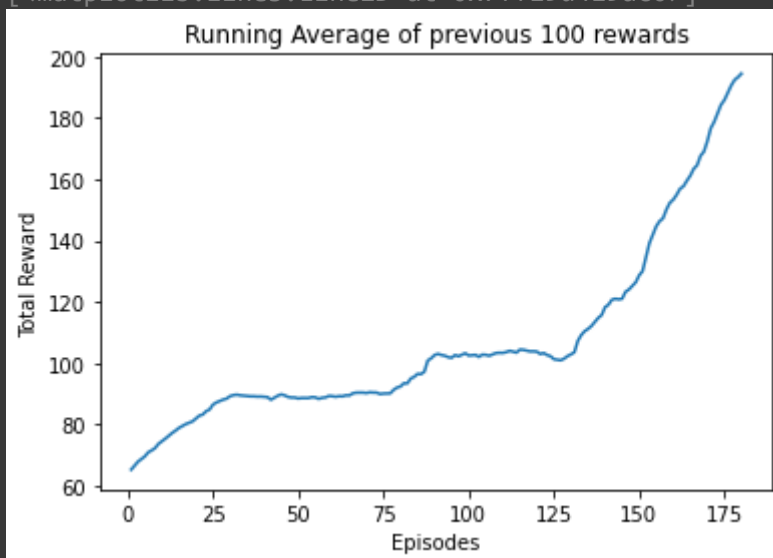
```
Stopped at Episode  180
0:06:31.474817
```

## Task 2b: Plot total reward curve

In the cell below, write code to plot the total reward averaged over 100 episodes (moving average)

```python
plt.xlabel("Episodes")
plt.ylabel("Total Reward")
plt.title("Running Average of previous 100 rewards")
for i in range(100,ep):
  average_reward_list.append(np.mean(reward_list[i-100:i]))

plt.plot([i+1 for i in range(len(average_reward_list))],average_reward_list)
```

```
[<matplotlib.lines.Line2D at 0x7ff19d419ac0>]
```



## Code for rendering (source)

```python
# Render an episode and save as a GIF file

display = Display(visible=0, size=(400, 300))
display.start()


def render_episode(env: gym.Env, model: tf.keras.Model, max_steps: int):
  screen = env.render(mode='rgb_array')
  im = Image.fromarray(screen)

  images = [im]
```

```python
    state = tf.constant(env.reset(), dtype=tf.float32)
    for i in range(1, max_steps + 1):
      state = tf.expand_dims(state, 0)
      action_probs, _ = model(state)
      action = np.argmax(np.squeeze(action_probs))
      state, _, done, _ = env.step(action)
      state = tf.constant(state, dtype=tf.float32)

      # Render screen every 10 steps
      if i % 10 == 0:
        screen = env.render(mode='rgb_array')
        images.append(Image.fromarray(screen))

      if done:
        break

  return images


# Save GIF image
images = render_episode(env, agent.ac_model, 200)
image_file = 'cartpole-v1.gif'
# loop=0: loop forever, duration=1: play each frame for 1ms
images[0].save(
    image_file, save_all=True, append_images=images[1:], loop=0, duration=1)
```
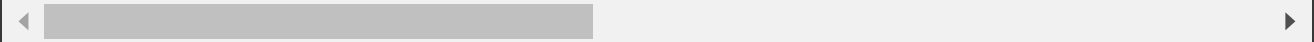
```
/usr/local/lib/python3.8/dist-packages/gym/core.py:43: DeprecationWarning: WARN: The
See here for more information: https://www.gymlibrary.ml/content/api/
  deprecation(
```

```python
import tensorflow_docs.vis.embed as embed
embed.embed_file(image_file)
```

```
!jupyter nbconvert --to html "/content/drive/MyDrive/Colab Notebooks/ME19B118_Tutorial_5_D
```

```
[NbConvertApp] WARNING | pattern '/content/drive/MyDrive/Colab Notebooks/ME19B118_
This application is used to convert notebook files (*.ipynb)
        to various other formats.

        WARNING: THE COMMANDLINE INTERFACE MAY CHANGE IN FUTURE RELEASES.

Options
=======
The options below are convenience aliases to configurable class-options,
as listed in the "Equivalent to" description-line of the aliases.
To see all configurable class-options for some <cmd>, use:
    <cmd> --help-all

--debug
    set log level to logging.DEBUG (maximize logging output)
    Equivalent to: [--Application.log_level=10]
--show-config
    Show the application's configuration (human-readable format)
    Equivalent to: [--Application.show_config=True]
--show-config-json
    Show the application's configuration (json format)
    Equivalent to: [--Application.show_config_json=True]
--generate-config
    generate default config file
    Equivalent to: [--JupyterApp.generate_config=True]
-y
    Answer yes to any questions instead of prompting.
    Equivalent to: [--JupyterApp.answer_yes=True]
--execute
    Execute the notebook prior to export.
    Equivalent to: [--ExecutePreprocessor.enabled=True]
--allow-errors
    Continue notebook execution even if one of the cells throws an error and inclu
    Equivalent to: [--ExecutePreprocessor.allow_errors=True]
--stdin
    read a single notebook file from stdin. Write the resulting notebook with defa
    Equivalent to: [--NbConvertApp.from_stdin=True]
--stdout
    Write notebook output to stdout instead of files.
    Equivalent to: [--NbConvertApp.writer_class=StdoutWriter]
--inplace
    Run nbconvert in place, overwriting the existing notebook (only
            relevant when converting to notebook format)
    Equivalent to: [--NbConvertApp.use_output_suffix=False --NbConvertApp.export_f
--clear-output
    Clear output of current file and save in place,
            overwriting the existing notebook.
    Equivalent to: [--NbConvertApp.use_output_suffix=False --NbConvertApp.export_f
--no-prompt
    Exclude input and output prompts from converted document.
    Equivalent to: [--TemplateExporter.exclude_input_prompt=True --TemplateExporte
--no-input
    Exclude input cells and output prompts from converted document.
            This mode is ideal for generating code-free reports.
    Equivalent to: [--TemplateExporter.exclude_output_prompt=True --TemplateExport
```
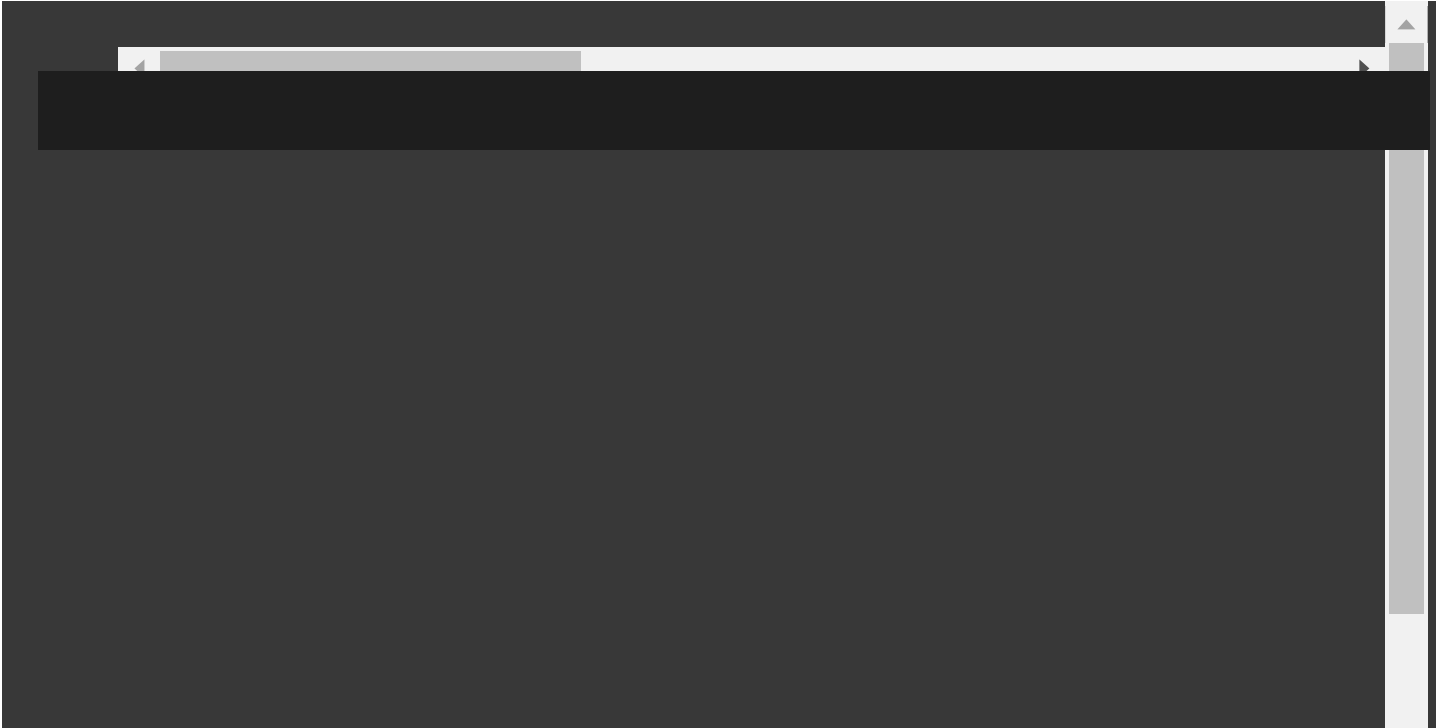
✓  0s      completed at 1:52 AM