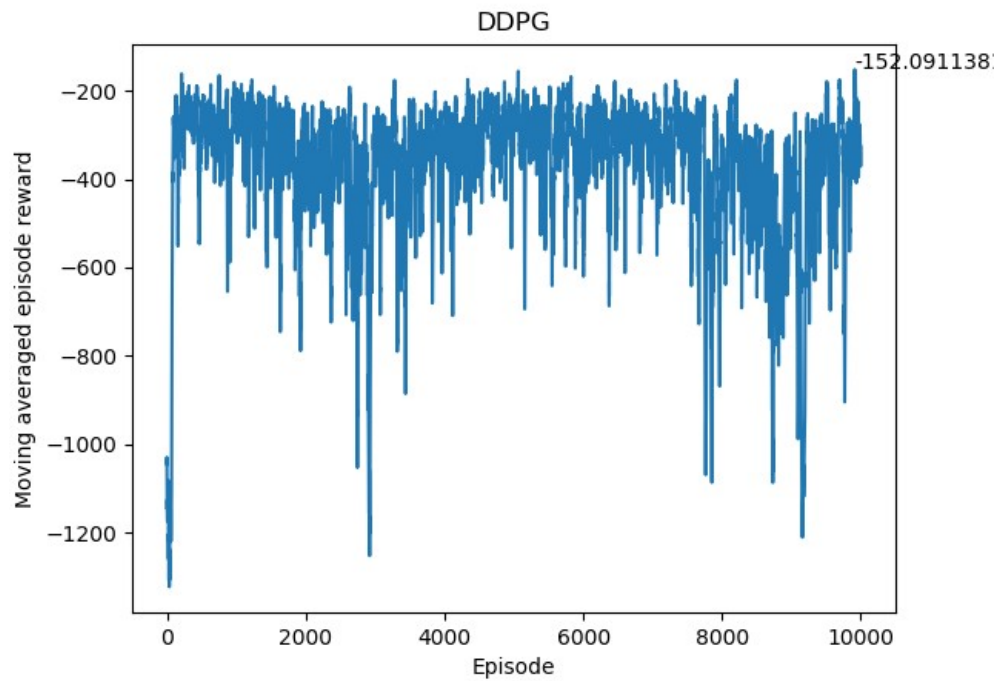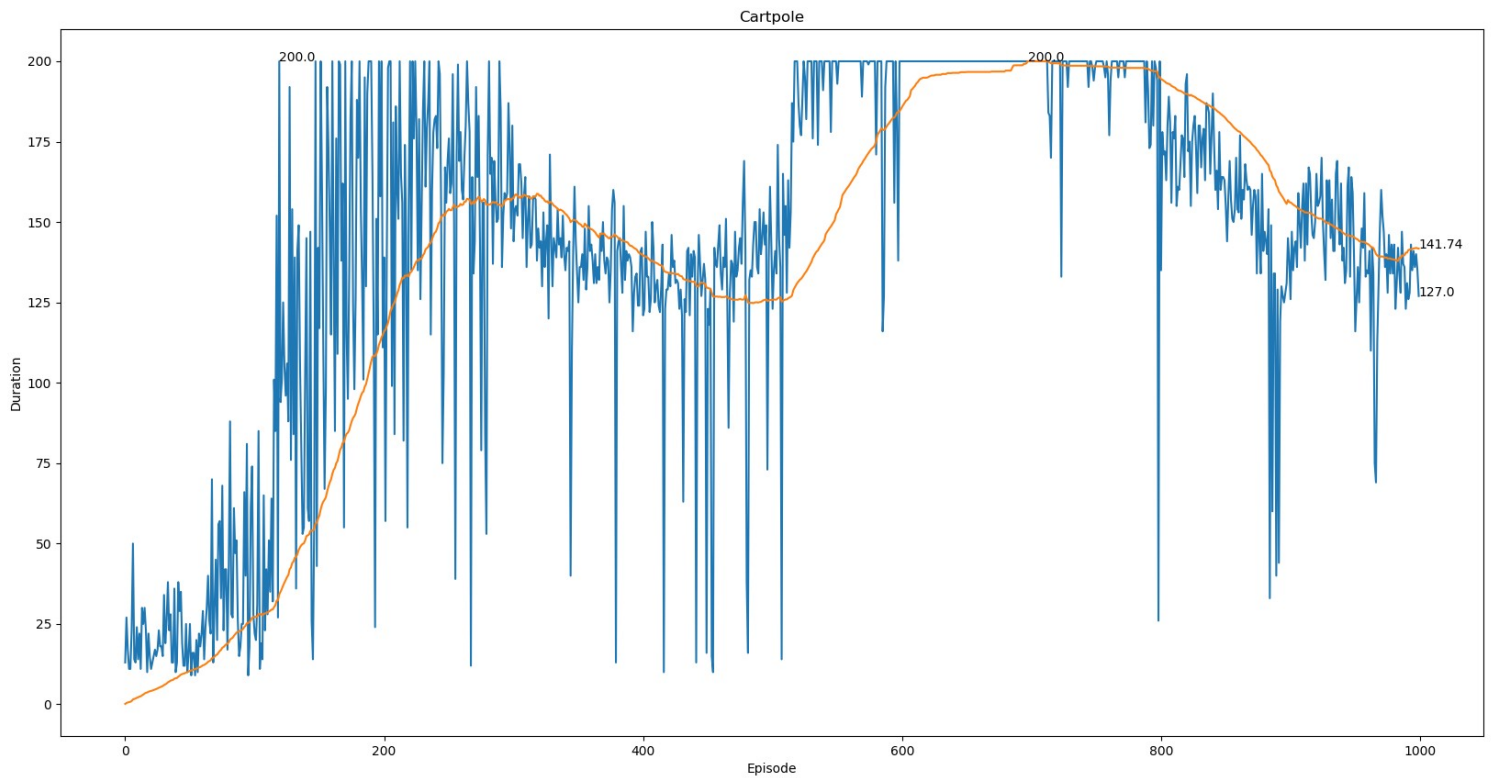# 一、Plot at CartPole-v0 and Pendulum

# 二、 Implement/adjustment of the network structure & each loss function

- CartPole-v0：Deep Q-learning with experience replay

```
12  class Net(nn.Module):
13      def __init__(self, input_size, hidden_size, output_size):
14          super().__init__()
15          self.linear1 = nn.Linear(input_size, hidden_size)
16          self.linear2 = nn.Linear(hidden_size, output_size)
17      def forward(self, x):
18          x = F.relu(self.linear1(x))
19          x = self.linear2(x)
20          return x
21
22  class Agent(object):
23      def __init__(self, **kwargs):
24          for key, value in kwargs.items():
25              setattr(self, key, value)
26          self.eval_net = Net(self.state_space_dim, 32, self.action_space_dim).to(device)
27          self.target_net = Net(self.state_space_dim, 32, self.action_space_dim).to(device)
28          self.target_net.load_state_dict(self.eval_net.state_dict())
```

在邊很單純的就是兩層 fully connect，中間夾了一個 relu activation 而已，latent vector 按照 spec 設成 32

比較特別的點是分成兩個 network，一個做主要的決定 action，一個來當作要 regression 的目標，而每過 K 個 iteration（這邊我設 50），就更新 regression 的目標成第一個 net 的 weight

如此就不會讓 regression 目標漂浮不定，比較好 train

而這邊 Loss 則是用 Mean Square Error Loss

- Pendulum-v0：DDPG

這邊則是用了叫做 Actor Critic 的技巧，Actor 負責決定動作，Critic 負責產生 Q 值

如此可以解決當 Action 是一個連續的值時，要找 maximum Q 的情況

而兩者的 model 也是很簡單，就是單純的 fully connect

Actor 的 fully connected 中間利用 relu 來連結，最後用 tanh 來輸出

Critic 中間用 relu，最後就不經過 activation function 了，直接輸出

而 Loss 一樣是用 Mean Square Error Loss

```
37  class ActorNet(nn.Module):
38      def __init__(self):
39          super(ActorNet, self).__init__()
40          self.fc = nn.Linear(3, 400)
41          self.fc2 = nn.Linear(400, 300)
42          self.mu_head = nn.Linear(300, 1)
43
44      def forward(self, s):
45          x = torch.relu(self.fc(s))
46          x = torch.relu(self.fc2(x))
47          u = torch.tanh(self.mu_head(x))
48          return u
49
50
51  class CriticNet(nn.Module):
52      def __init__(self):
53          super(CriticNet, self).__init__()
54          self.fc = nn.Linear(3, 400)
55          self.fc2 = nn.Linear(401, 300)
56          self.v_head = nn.Linear(300, 1)
57
58      def forward(self, s, a):
59          x = F.relu(self.fc(s))
60          x = F.relu(self.fc2(torch.cat([x, a], dim=1)))
61          state_value = self.v_head(x)
62          return state_value
```

# 三、Implement the training process of deep Q-learning

```python
def learn(self):
    if (len(self.buffer)) < self.batch_size:
        return

    samples = random.sample(self.buffer, self.batch_size)
    s0, a0, r1, s1 = zip(*samples)
    s0 = torch.tensor(s0, dtype=torch.float).to(device)
    a0 = torch.tensor(a0, dtype=torch.long).view(self.batch_size, -1).to(device)
    r1 = torch.tensor(r1, dtype=torch.float).view(self.batch_size, -1).to(device)
    s1 = torch.tensor(s1, dtype=torch.float).to(device)

    y_true = r1 + self.gamma * torch.max(self.target_net(s1).detach(), dim=1)[0].view(self.batch_size, -1)
    y_pred = self.eval_net(s0).gather(1, a0)

    loss_fn = nn.MSELoss()
    loss = loss_fn(y_pred, y_true)

    self.optimizer.zero_grad()
    loss.backward()
    self.optimizer.step()
```

這邊一樣是用 TD learning 的技巧，不同的是 Q 值的計算變成了用 neural network 來取代
一開始先從歷史紀錄（buffer）sample 出 batch_size 份量的 episodes
predict 值會是經過 self.eval_net(s0)後，取這個 action(a0)後的值
而 true value 則是利用 target_net(s1)後最大的值，加上這步的 reward
下面就把這兩個值做 Mean Square Error Loss，再 back propogate 就可以了

# 四、Implement of epsilon-greedy action select method

```python
34        def update_iter(self):
35            self.steps += 1
36            self.epsi *= self.decay
37            if self.steps % self.copy_weight == 0:
38                self.target_net.load_state_dict(self.eval_net.state_dict())
39
40        def act(self, s0):
41            if random.random() < self.epsi:
42                a0 = random.randrange(self.action_space_dim)
43            else:
44                s0 = torch.tensor(s0, dtype=torch.float).view(1,-1).to(device)
45                a0 = torch.argmax(self.eval_net(s0)).item()
46            return a0
```

這邊有兩個 function 來完成 epsilon-greedy，一個是 update_iter，每次 iteration 完後就會 call 這個
function，來更新目前 iteration 次數，以及 epsi threshold 乘上一個 decay 值
而在實際決定動作的時候，會先 random 一個數值，如果小於 epsi threshold 的話，就會隨機挑一個 action
做，如此就可以完成 epsilon-greedy action 了

# 五、The mechanism of critic updating

Sample random minibatch of $N$ transitions $(s_j, a_j, r_j, s_{j+1})$ from R

Set $y_i = r_i + \gamma Q'(s_{t+1}, \mu'(s_{t+1}|\theta^{\mu'})|\theta^{Q'})$

Update critic by minimizing the loss: $L = \frac{1}{N}\sum_i(y_i - Q(s_i, a_i|\theta^Q))^2$

由於 critic 作用是 Q 值的產生，所以我們更新也只要考慮這部份就好了
我們會先對 buffer sample 一切 tranitions 出來做更新用
第二行是在計算 target value，分別是 reward 和 target network 產生的值
第三行就是把 target value 和 evaluation network 產生的 Q 值做 Mean Square Error Loss，然後再更新就好了
要注意的是第二行在計算 target value 時，會用到 actor 的 target network 來產生 action，在實做時這部份要小心不要更新到了

# 六、The mechanism of actor updating

Update the actor policy using the sampled gradient:

$$\nabla_{\theta^\mu}\mu|_{s_i} \approx \frac{1}{N}\sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)}\nabla_{\theta^\mu}\mu(s|\theta^\mu)|_{s_i}$$

他的 loss 會是 sampled 的 Q 對 a 的 gradient，用 code 寫成的話會變成下面樣子

```
a_loss = -self.eval_cnet(s, self.eval_anet(s)).mean()
a_loss.backward()
```

actor 負責的是 action 的選擇，所以會希望他所選的 action 能盡量讓 critic 產生的 Q 值越大越好，而換成 loss 要越小越好的話就多加一個負號就可以了

# 七、How to calculate the gradients
對於 actor 的 gradient，前項 Q 對 a 的 gradient 可以看做是將 actor predict 的結果丟到 critic 後，對結果最 back propagate
後項則是 actor 產生的 action 結果做 gradient
所以整體來看可以看做是一連串下來做 back propagation 就可以了

# 八、How the code work

- CartPole-v0

首先先設定好 environment 以及 model 參數

下方是主程式，每次 environment 都會 reset，然後因為我是在遠端機器，所以就不處理 X server 的問題了，直接把 render 關掉

每次都會讓 agent 對這個 state 做動作，丟回 env 後就記錄下來
然後讓 agent 去 learn 他

```python
env = gym.make('CartPole-v0')
params = {
    'gamma': 0.95,
    'epsi_high': 1,
    'epsi_low': 0.001,
    'decay': 0.995,
    'lr': 0.0005,
    'capacity': 10000,
    'batch_size': 128,
    'copy_weight': 50,
    'state_space_dim': env.observation_space.shape[0],
    'action_space_dim': env.action_space.n
}
agent = Agent(**params)
agent.load_weight()
score = []
mean = []
finish = False
```

後面就是一些更新次數或 epsilon 以及紀錄結果的 code，就不詳細解釋了

```python
for episode in range(100):
    s0 = env.reset()
    total_reward = 1
    while True:
        # env.render()
        a0 = agent.act(s0)
        s1, r1, done, _ = env.step(a0)

        if done:
            r1 = -1

        agent.put(s0, a0, r1, s1)

        if done:
            break

        total_reward += r1
        s0 = s1
        agent.learn()
    agent.update_iter()
    # if total_reward==200 and not finish:
    #     agent.save_weight()
    #     print('save')
    # if sum(score[-100:])/100 == 200:
    #     finish = True
    #     print('avg score reach 200')
    score.append(total_reward)
    print(f'Episode {episode}: score {total_reward}, avg {sum(score[-100:])/100}')
    mean.append(sum(score[-100:])/100)
agent.plot(score, mean)
```

右圖是 Net 的架構以及用一個 Agent class 把 model 和訓練等 function 都包在一起

```python
10    device = torch.device('cuda:2')
11    class Net(nn.Module):
12        def __init__(self, input_size, hidden_size, output_size):
13            super().__init__()
14            self.linear1 = nn.Linear(input_size, hidden_size)
15            self.linear2 = nn.Linear(hidden_size, output_size)
16        def forward(self, x):
17            x = F.relu(self.linear1(x))
18            x = self.linear2(x)
19            return x
20
21    class Agent(object):
22        def __init__(self, **kwargs):
23            for key, value in kwargs.items():
24                setattr(self, key, value)
25            self.eval_net = Net(self.state_space_dim, 32, self.action_space_dim).to(device)
26            self.target_net = Net(self.state_space_dim, 32, self.action_space_dim).to(device)
27            self.target_net.load_state_dict(self.eval_net.state_dict())
28
29            self.optimizer = optim.Adam(self.eval_net.parameters(), lr=self.lr)
30            self.buffer = []
31            self.steps = 0
32            self.epsi = self.epsi_high
33            self.training = True
```

再來是每次 iteration 結束後，會更新 step 次數以及 epsilon 的數值，另外次數如果到達一定程度的話就更新 target network 的參數

```python
44        def update_iter(self):
45            self.steps += 1
46            self.epsi *= self.decay
47            if self.steps % self.copy_weight == 0:
48                self.target_net.load_state_dict(self.eval_net.state_dict())
49
50        def act(self, s0):
51            if random.random() < self.epsi:
52                a0 = random.randrange(self.action_space_dim)
53            else:
54                s0 = torch.tensor(s0, dtype=torch.float).view(1,-1).to(device)
55                a0 = torch.argmax(self.eval_net(s0)).item()
56            return a0
57
58        def put(self, *transition):
59            if len( self.buffer)==self.capacity:
60                self.buffer.pop(0)
61            self.buffer.append(transition)
```

```python
63        def learn(self):
64            if (len(self.buffer)) < self.batch_size:
65                return
66
67            if not self.training:
68                return
69
70            samples = random.sample(self.buffer, self.batch_size)
71            s0, a0, r1, s1 = zip(*samples)
72            s0 = torch.tensor(s0, dtype=torch.float).to(device)
73            a0 = torch.tensor(a0, dtype=torch.long).view(self.batch_size, -1).to(device)
74            r1 = torch.tensor(r1, dtype=torch.float).view(self.batch_size, -1).to(device)
75            s1 = torch.tensor(s1, dtype=torch.float).to(device)
76
77            y_true = r1 + self.gamma * torch.max(self.target_net(s1).detach(), dim=1)[0].view(self.batch_size, -1)
78            y_pred = self.eval_net(s0).gather(1, a0)
79
80            loss_fn = nn.MSELoss()
81            loss = loss_fn(y_pred, y_true)
82
83            self.optimizer.zero_grad()
84            loss.backward()
85            self.optimizer.step()
```

最後是 learn，首先會先檢查 buffer 夠不夠 batch_size 的數量，再來就是 sample，然後做各種前處理，再做上面有說過的更新就好了

- Pendulum-v0

主程式如右，一樣先設定好 environment 環境，然後就開始跑 iteration
每次 iteration 都限制最多只有 200 回合，然後也是跑和存到 buffer，如果 buffer 是滿的話就更新參數

```
150  def main():
151      env = gym.make('Pendulum-v0')
152      env.seed(args.seed)
153
154      agent = Agent()
155
156      training_records = []
157      running_reward, running_q = -1000, 0
158      for i_ep in range(10000):
159          score = 0
160          state = env.reset()
161
162          for t in range(200):
163              action = agent.select_action(state)
164              state_, reward, done, _ = env.step(action)
165              score += reward
166              if args.render:
167                  env.render()
168              agent.store_transition(Transition(state, action, (reward + 8) / 8, state_))
169              state = state_
170              if agent.memory.isfull:
171                  q = agent.update()
172                  running_q = 0.99 * running_q + 0.01 * q
173
174          running_reward = running_reward * 0.9 + score * 0.1
175          training_records.append(TrainingRecord(i_ep, running_reward))
```

```
84   class Agent():
85       max_grad_norm = 0.5
86
87       def __init__(self):
88           self.training_step = 0
89           self.var = 1.
90           self.eval_cnet, self.target_cnet = CriticNet().float().to(device), CriticNet().float().to(device)
91           self.eval_anet, self.target_anet = ActorNet().float().to(device), ActorNet().float().to(device)
92           self.memory = Memory(10000)
93           self.optimizer_c = optim.Adam(self.eval_cnet.parameters(), lr=0.001)
94           self.optimizer_a = optim.Adam(self.eval_anet.parameters(), lr=0.0001)
95           self.Loss = nn.MSELoss()
96
97       def select_action(self, state):
98           state = torch.from_numpy(state).float().to(device).unsqueeze(0)
99           mu = self.eval_anet(state)
100          dist = Normal(mu, torch.tensor(self.var, dtype=torch.float).to(device))
101          action = dist.sample()
102          action.clamp(-2.0, 2.0)
103          return (action.item(),)
104
105      def save_param(self):
106          torch.save(self.eval_anet.state_dict(), 'ddpg_anet_params.pkl')
107          torch.save(self.eval_cnet.state_dict(), 'ddpg_cnet_params.pkl')
108
109      def store_transition(self, transition):
110          self.memory.update(transition)
```

上圖是 Agent，將訓練的東西都包在一起
對 Actor 和 Critic 上面的問題已經有看過了，這邊就直接省略說明
對於 Actor 和 Critic 一樣都會有 Evaluation net 和 Target net
Memory 則是用來存 episode history 的類別

selection action 裡面，為了讓 agent 能有一點變化，會是讓 actor 出來一個 mu，然後對此做一個 normal distribution 的抽樣，而其 variance 會逐步減小，最後為了不讓他衝過頭，會在對兩邊做 crop，超過或低於都直接變成最大最小值

```
112    def update(self):
113        self.training_step += 1
114
115        transitions = self.memory.sample(32)
116        s = torch.tensor([t.s for t in transitions], dtype=torch.float).to(device)
117        a = torch.tensor([t.a for t in transitions], dtype=torch.float).view(-1, 1).to(device)
118        r = torch.tensor([t.r for t in transitions], dtype=torch.float).view(-1, 1).to(device)
119        s_ = torch.tensor([t.s_ for t in transitions], dtype=torch.float).to(device)
120
121        with torch.no_grad():
122            q_target = r + args.gamma * self.target_cnet(s_, self.target_anet(s_))
123        q_eval = self.eval_cnet(s, a)
124
125        # update critic net
126        self.optimizer_c.zero_grad()
127        c_loss = self.Loss(q_eval, q_target).mean()
128        c_loss.backward()
129        nn.utils.clip_grad_norm_(self.eval_cnet.parameters(), self.max_grad_norm)
130        self.optimizer_c.step()
131
132        # update actor net
133        self.optimizer_a.zero_grad()
134        a_loss = -self.eval_cnet(s, self.eval_anet(s)).mean()
135        a_loss.backward()
136        nn.utils.clip_grad_norm_(self.eval_anet.parameters(), self.max_grad_norm)
137        self.optimizer_a.step()
138
139        if self.training_step % 300 == 0:
140            self.target_cnet.load_state_dict(self.eval_cnet.state_dict())
141        if self.training_step % 301 == 0:
142            self.target_anet.load_state_dict(self.eval_anet.state_dict())
143
144        self.var = max(self.var * 0.999, 0.01)
145
146        return q_eval.mean().item()
```

最後是整個 model 的更新，我們會從 buffer 裡 sample batch_size 大小的份數，然後做完各
種前處理之後就是實做上面解釋過的，對 actor 和 critic 的更新了
而這邊為了方便，沒有讓 target network 做 soft update，而是跟 DQN 一樣到一定 episode
就覆蓋過去，結果也顯示還是不錯的

右圖則是存 history 的 class
更新方式是用一個 data_pointer 紀錄
這串 list 的頭是哪裡，滿了後要更新
就把他蓋過去

sample 也就直接用 numpy 的 random
choice 來取出一定數量就完成了

```
65    class Memory():
66        data_pointer = 0
67        isfull = False
68
69        def __init__(self, capacity):
70            self.memory = np.empty(capacity, dtype=object)
71            self.capacity = capacity
72
73        def update(self, transition):
74            self.memory[self.data_pointer] = transition
75            self.data_pointer += 1
76            if self.data_pointer == self.capacity:
77                self.data_pointer = 0
78                self.isfull = True
79
80        def sample(self, batch_size):
81            return np.random.choice(self.memory, batch_size)
```

# 九、Other study or improvement for the project

其實老師在上這邊時非常的快速，我是自己額外看完 youtube 上李宏毅影片後才寫這份作業的

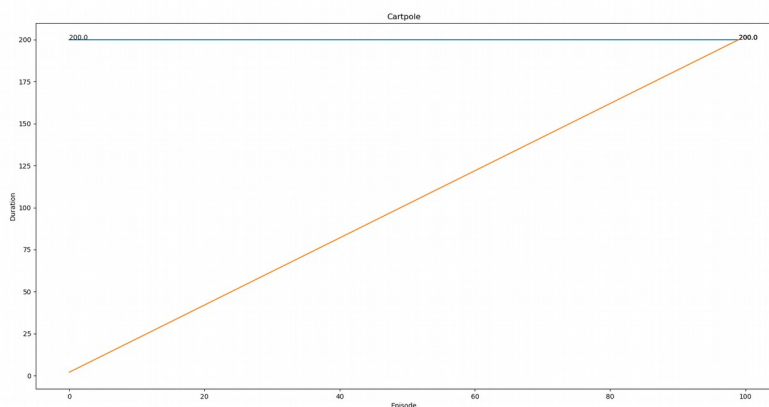但其實自己寫的部份也沒有說非常多，在網路上有找到別人寫好的版本，參考後修改一些地方就下去跑了

其中也發現別人包成 class 的寫法和對於 sample history data 的寫法，甚至是 model 參數的設定方式都不一樣，收穫十分良多
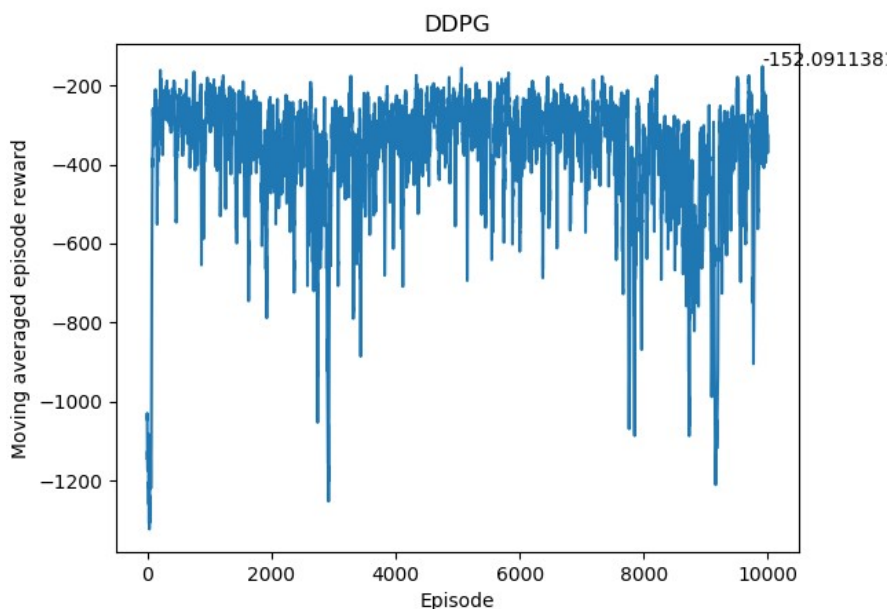
pendulum 會要跑好幾個小時，但幸好沒太多問題，很快就 train 起來了

cartpole 則是 train 到很好的時候再多 train 就會不小心壞掉，應該是跟那個 epsilon 有關係，所以我在 training 階段，只要連續到 200 一定次數就不再 store weight，結果就可以像下面一樣都到 200 了

# 十、Performance

- CartPole: 200



- Pendulum：-152.091138