# NCTU RL HW2-Problem3
# REINFORCE and A2C

廖家鴻 0786009
2021/5/10

# Outline

- Problem3-(a): CartPole-v0 and REINFORCE
  - Network
  - Hyperparameters
  - Results

- Problem3-(b): LunarLander-v2 and A2C
  - Network
  - Hyperparameters
  - Results
  - Remarks

# Problem3-(a):
# CartPole-v0 and REINFORCE

# Prob3-(a): Network

```
45            ########## YOUR CODE HERE (5~10 lines) ##########
46            ### Actor_Net ###
47            self.a_fc0 = nn.Linear(self.state_dim, self.hidden_size)
48            self.a_fc1 = nn.Linear(self.hidden_size, self.hidden_size * 2)
49            self.a_fc2 = nn.Linear(self.hidden_size * 2, self.hidden_size * 2)
50            self.a_fc3 = nn.Linear(self.hidden_size * 2, self.action_dim)
51
52            ### Baseline_Net ### (for estimating value function)
53            self.c_fc0 = nn.Linear(self.state_dim, self.hidden_size)
54            self.c_fc1 = nn.Linear(self.hidden_size, self.hidden_size)
55            self.c_fc2 = nn.Linear(self.hidden_size, self.hidden_size)
56            self.c_fc3 = nn.Linear(self.hidden_size, 1)
57            ########## END OF YOUR CODE ##########
```

1. The hidden_size here is 64

```
63    ▾    def forward(self, state):
64            ########## YOUR CODE HERE (3~5 lines) #######
65            ### Actor_Net ###
66            x = F.relu(self.a_fc0(state))
67            x = F.relu(self.a_fc1(x))
68            x = F.relu(self.a_fc2(x))
69            x = F.softmax(self.a_fc3(x))
70            action_prob = x
71
72            ### Baseline_Net ###
73            y = F.relu(self.c_fc0(state))
74            y = F.relu(self.c_fc1(y))
75            y = F.relu(self.c_fc2(y))
76            y = self.c_fc3(y)
77            baseline_value = y
78            ########## END OF YOUR CODE ##########
79            return action_prob, baseline_value
```

2. I implement the Actor_Net and Baseline_Net separately.

PS: I have tried another setting that is "one state input with two branches output." And this will also work. But here I only present the separated version.

# Prob3-(a): Hyperparameters

```
267  ▼ if __name__ == '__main__':
268        # For reproducibility, fix the random seed
269        random_seed = 20
270        lr = 0.02
271        env = gym.make('CartPole-v0')
272        env.seed(random_seed)
273        torch.manual_seed(random_seed)
274        train(lr)
275        test('CartPole_0.02.pth')
```

1. The learning rate is lr=0.02

2. The discount factor is gamma=0.99

```
109  ▼     def calculate_loss(self, optimizer, gamma=0.99):
135           g_return = self.rewards[t] + gamma*g_return
```
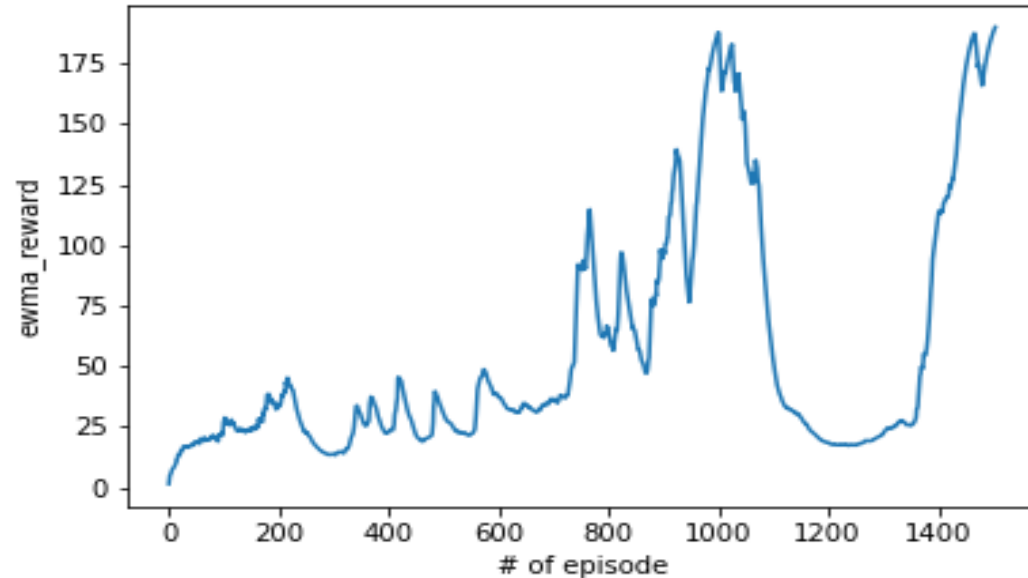
```
170        # Instantiate the policy model and the optimizer
171        model = Policy()
172        optimizer = optim.SGD(model.parameters(), lr=lr)
173        scheduler = Scheduler.StepLR(optimizer, step_size=500, gamma=0.5)
```

3. The learning rate scheduler is set with step_size=500 and gamma=0.5

# Prob3-(a): Results

```
231            # check if we have "solved" the cart pole problem
232  ▼         if ewma_reward >= 190.:
233                torch.save(model.state_dict(), './CartPole_0.02.pth')
```



```
Ep 1503 Length: 200 reward: 200.0000    ewma reward: 188.9586
Ep 1504 Length: 200 reward: 200.0000    ewma reward: 189.5107
Ep 1505 Length: 200 reward: 200.0000    ewma reward: 190.0352
Solved! Running reward is now 190.03516076169097 and the last episode runs to 200 time steps!
Episode 1    Reward: 200.0
Episode 2    Reward: 200.0
Episode 3    Reward: 200.0
Episode 4    Reward: 200.0
Episode 5    Reward: 200.0
Episode 6    Reward: 200.0
Episode 7    Reward: 200.0
Episode 8    Reward: 200.0
Episode 9    Reward: 200.0
Episode 10   Reward: 200.0
```

1. I save the best agent when the ewma_reward achieve 190.

2. The left figure shows the curve of ewma_reward for all training episodes.

3. At this setting, the total training episodes are 1505. And the last episode runs to 200 steps.

4. I follow the sample code to test the trained agent for 10 episodes and all get the maximum reward.

# Problem3-(b):
# LunarLander-v2 and A2C

# Prob3-(b): Network

```
43          self.hidden_size = 64
44          ########## YOUR CODE HERE (5~10 lines) ##########
45          self.fc0 = nn.Linear(self.state_dim, self.hidden_size)
46          ### Actor_Net ###
47          self.a_fc0 = nn.Linear(self.state_dim, self.hidden_size)
48          self.a_fc1 = nn.Linear(self.hidden_size, self.hidden_size * 2)
49          self.a_fc2 = nn.Linear(self.hidden_size * 2, self.hidden_size * 2)
50          self.a_fc3 = nn.Linear(self.hidden_size * 2, self.action_dim)
51
52          ### Baseline_Net ### (for estimating value function)
53          self.c_fc0 = nn.Linear(self.state_dim, self.hidden_size)
54          self.c_fc1 = nn.Linear(self.hidden_size, self.hidden_size)
55          self.c_fc2 = nn.Linear(self.hidden_size, self.hidden_size)
56          self.c_fc3 = nn.Linear(self.hidden_size, 1)
57          ########## END OF YOUR CODE ##########
```

1. The hidden_size here is 64.

```
63      def forward(self, state):
64          ########## YOUR CODE HERE (3~5 lines) ##########
65          # s = F.relu(self.fc0(state))
66          ### Actor_branch ###
67          x = F.relu(self.a_fc0(state))
68          x = F.relu(self.a_fc1(x))
69          x = F.relu(self.a_fc2(x))
70          x = F.sigmoid(self.a_fc3(x))
71          action_prob = x
72
73          ### Baseline_branch ### (for estimating value function)
74          y = F.relu(self.c_fc0(state))
75          y = F.relu(self.c_fc1(y))
76          y = F.relu(self.c_fc2(y))
77          y = self.c_fc3(y)
78          baseline_value = y
79          ########## END OF YOUR CODE ##########
80          return action_prob, baseline_value
```

2. Again, I implement the Actor_Net and Baseline_Net separately.

PS: I also tried another setting that is "one state input with two branches output", but it is hard work.

# Prob3-(a): Hyperparameters

```
261  ▼ if __name__ == '__main__':
262        # For reproducibility, fix the random seed
263        random_seed = 20
264        lr = 0.01
265        env = gym.make('LunarLander-v2')
266        env.seed(random_seed)
267        torch.manual_seed(random_seed)
268        train(lr)
269        test('LunarLander_0.01.pth')
```

1. The learning rate is lr=0.01

```
103        def calculate_loss(self, gamma=0.95):
132            baseline_loss = nn.MSELoss()
133            value_loss = baseline_loss(value, reward+(gamma*next_value) )
134
135            advantage = (reward + gamma*next_value)  - value
136            policy_loss = -log_act_prob * advantage.detach()
```
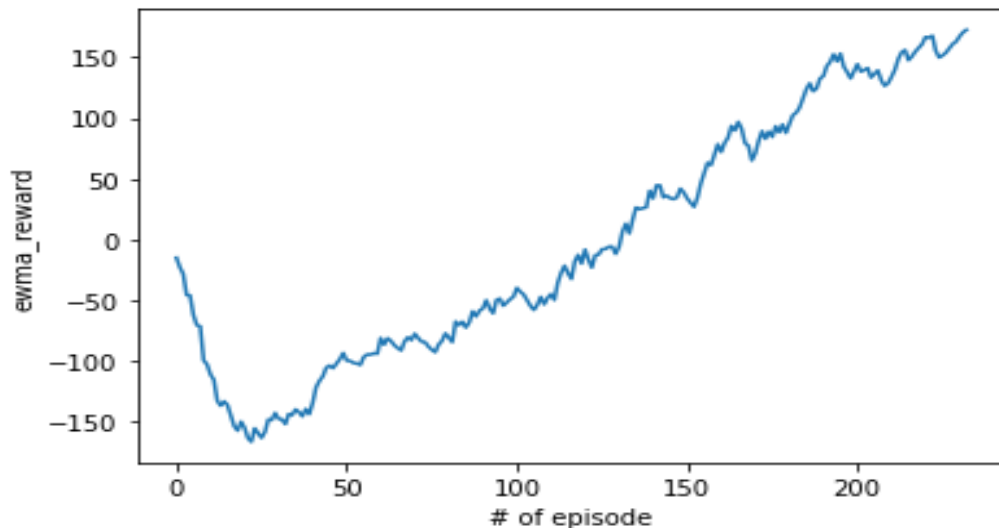
2. The discount factor is gamma=0.95

```
160        # Instantiate the policy model and the optimizer
161        model = Policy()
162        optimizer = optim.SGD(model.parameters(), lr=lr)
163        # optimizer = optim.Adam(model.parameters(), lr=lr)
164        scheduler = Scheduler.StepLR(optimizer, step_size=100, gamma=0.8)
```

3. The learning rate scheduler is set with step_size=100 and gamma=0.8

# Prob3-(b): Results

```
222          # check if we have "solved" the cart pole problem
223    ▼     if ewma_reward >= 172.:
224             torch.save(model.state_dict(), './LunarLander_0.01.pth')
225    ▼        print("Solved! Running reward is now {} and "
226                   "the last episode runs to {} time steps!".format(ewma_reward, t))
227             break
```



1. I save the best agent when the ewma_reward achieve 172.

2. The left figure shows the curve of ewma_reward for all training episodes.

```
Ep 230   Length: 516 R: 200.5952   ewma reward: 162.8868
Ep 231   Length: 267 R: 247.0709   ewma reward: 167.0960
Ep 232   Length: 395 R: 233.3892   ewma reward: 170.4107
Ep 233   Length: 870 R: 212.3251   ewma reward: 172.5064
Solved! Running reward is now 172.506412644235 and the last episode runs to 870 time steps!
Episode 1    Reward: 247.8139565405427
Episode 2    Reward: 277.3608585658985
Episode 3    Reward: -1.5739138103989347
Episode 4    Reward: 202.5240771279636
Episode 5    Reward: 244.8683823638374
Episode 6    Reward: 235.36878862993836
Episode 7    Reward: 141.4510078744156
Episode 8    Reward: 239.24533874901292
Episode 9    Reward: 74.15252161222428
Episode 10   Reward: -191.0576380074957
```

3. At this setting, the total training episodes are 233.

4. I follow the sample code to test the trained agent for 10 episodes, only episode 3 and 10 are failedg.

# Prob3-(b): Remarks

```
########## YOUR CODE HERE (10-15 lines) ##########
for t in itertools.count(start=1):
    action = model.select_action(state)
    next_state, reward, done, _ = env.step(action)

    ep_reward += reward
    state = next_state
    model.rewards.append(reward)

    loss, policy_loss, value_loss = model.calculate_loss(reward, next_state)
    # optimizer.zero_grad()
    # loss.backward()
    # nn.utils.clip_grad_norm_(model.parameters(), 3)
    # optimizer.step()

    ### Update Actor_Net ###
    optimizer.zero_grad()
    policy_loss.backward(retain_graph=True)
    nn.utils.clip_grad_norm_(model.parameters(), 3)
    optimizer.step()

    ### Update Behavior_Net ###
    optimizer.zero_grad()
    value_loss.backward()
    nn.utils.clip_grad_norm_(model.parameters(), 3)
    optimizer.step()

    if done:
        break

model.clear_memory()
########## END OF YOUR CODE ##########
```

The key point that I finally success is that the network update frequency.

Originally, I update the network per episode as 3-(a), but it always fail no matter I tune the hyperparameters.

After I change the way to update the network by every step, the agent finally works!!!

# The End