

# 强化学习 HW5

— 人 饶翔云 520030910366

## Problem 1

Coding:

TRPO:

```
1  """ ----- Programming 1: implement the linear search to find best
2  parameter for actor (you may refer to original TRPO paper,
3  Appendix C) ----- """
4  """ YOUR CODE HERE """
5  # load new parameter to new actor
6  torch.nn.utils.convert_parameters.vector_to_parameters(new_para,
7                                                         new_actor.parameters())
8  # compute kl divergence and new objective
9  new_action_dists = torch.distributions.Categorical(new_actor(states))
10 kl = torch.mean(torch.distributions.kl.kl_divergence(old_action_dists,
11                                                         new_action_dists))
12 new_obj = self.compute_surrogate_obj(states, actions, advantage,
13                                     old_log_probs, new_actor)
14 improve = new_obj - old_obj
15 # check if kl is less than constraint and improve is large enough
16 if kl < self.kl_constraint and improve > 0:
17     self.actor.load_state_dict(new_actor.state_dict())
18     break
19 """ ----- Programming 1 ----- """
```

```
1  """ ----- Programming 2: implement the conjugate_gradient function,
2  the linear search function, to update actor parameter (you may refer to original
3  TRPO paper, Section 6) ----- """
4  """ YOUR CODE HERE """
5  # compute gradient vector
6  grads = torch.cat([grad.view(-1) for grad in grads])
7  # compute search direction
8  step_dir = self.conjugate_gradient(grads.data, states, old_action_dists)
9  # get hessian vector product
10 hd = self.hessian_matrix_vector_product(states, old_action_dists, step_dir)
11 # get max coef
12 max_coef = torch.sqrt(2 * self.kl_constraint / (torch.dot(step_dir, hd) + 1e-8))
13 # get max vec
14 max_vec = max_coef * step_dir
15 # line search
16 self.line_search(states, actions, advantage, old_log_probs,
17                 old_action_dists, max_vec)
18 """ ----- Programming 2 ----- """
```

```

1  """ ----- Programming 3: Compute GAE and update the parameter of actor
2  and critic ----- """
3  """ YOUR CODE HERE """
4  # compute GAE
5  values = self.critic(states)
6  next_values = self.critic(next_states)
7  td_target = rewards + self.gamma * next_values * (1 - dones)
8  td_delta = td_target - values
9  advantage = compute_advantage(self.gamma, self.lmbda, td_delta)
10 # normalize advantage if needed
11 advantage = (advantage - advantage.mean()) / (advantage.std() + 1e-10)
12 # update critic
13 critic_loss = torch.nn.functional.mse_loss(values, td_target.detach())
14 self.critic_optimizer.zero_grad()
15 critic_loss.backward()
16 self.critic_optimizer.step()
17 # get gradient of loss and hessian vector product of kl divergence
18 old_action_dists = torch.distributions.Categorical(self.actor(states).detach())
19 old_log_probs = torch.log(old_action_dists.probs.gather(1, actions)).detach()
20 # update actor
21 self.policy_learn(states, actions, old_action_dists, old_log_probs, advantage)
22 """ ----- Programming 3 ----- """

```

## PPO:

```

1  """ ----- Programming 4: Compute Advantage Function ----- """
2  """ YOUR CODE HERE """
3  # compute td delta by using Bellman equation
4  td_target = rewards + self.gamma * self.critic(next_states) * (1 - dones)
5  td_delta = td_target - self.critic(states)
6  # compute advantage
7  advantage = compute_advantage(self.gamma, self.lmbda, td_delta)
8  # normalize advantage
9  # advantage = (advantage - advantage.mean()) / (advantage.std() + 1e-12)
10 # initialize old_log_probs for computing ratio
11 old_probs = self.actor(states)
12 dist = torch.distributions.Categorical(old_probs)
13 old_log_probs = dist.log_prob(actions.squeeze(1)).detach()
14 old_probs = old_probs.view((-1, self.action_dim)).detach()
15 """ ----- Programming 4 ----- """

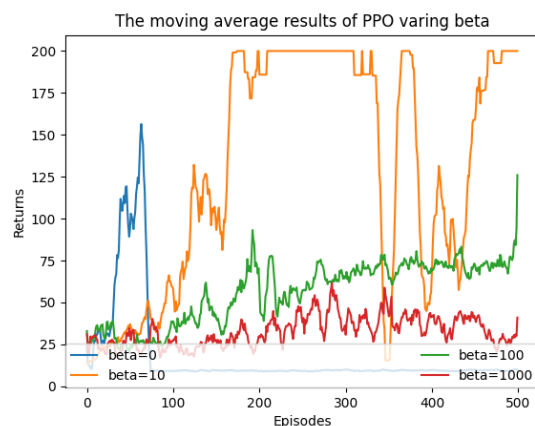
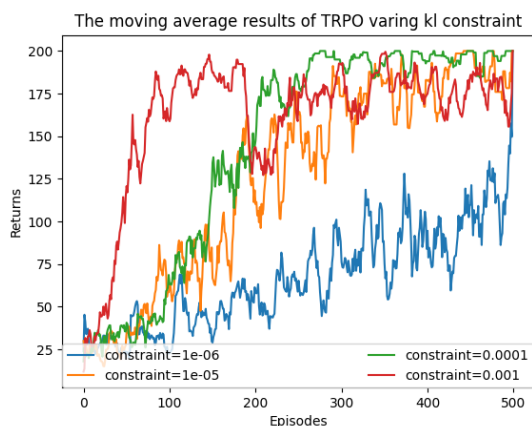
```

```

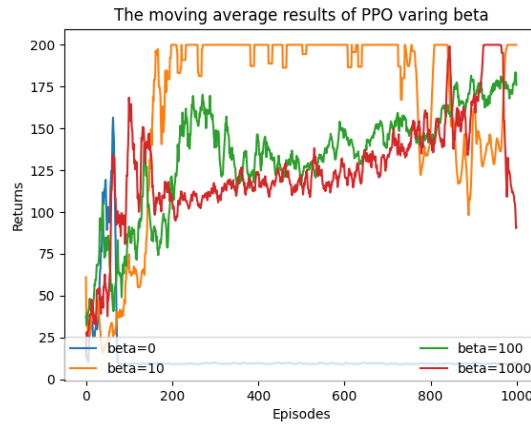
1  for _ in range(self.epochs):
2      """ ----- Programming 5: Update the parameter of actor and critic
3      (you may refer to original PPO paper) ----- """
4      """ YOUR CODE HERE """
5      from torch.distributions import Categorical
6      from torch.nn.functional import kl_div
7      # compute new_log_probs
8      new_probs = self.actor(states)
9      dist = Categorical(new_probs)
10     new_log_probs = dist.log_prob(actions.squeeze(1))
11     # compute ratio
12     ratio = torch.exp(new_log_probs - old_log_probs.detach())
13     # compute surrogate loss
14     surrogate_loss = ratio * advantage
15     # compute kl divergence (penalty)
16     penalty = kl_div(input = new_probs.log(), target = old_probs.detach(),
17                     reduction='batchmean')
18     # update actor
19     self.actor_optimizer.zero_grad()
20     actor_loss = - (surrogate_loss - self.beta * penalty).mean()
21     actor_loss.backward()
22     self.actor_optimizer.step()
23     # update critic
24     self.critic_optimizer.zero_grad()
25     critic_loss = F.mse_loss(self.critic(states), td_target.detach())
26     critic_loss.backward()
27     self.critic_optimizer.step()
28     # I want to update beta but it's not required.
29     # dl = kl_div(input=self.actor(states).log(),
30     #             target=old_probs.detach(), reduction='sum')
31     # print(dl)
32     #if abs(dl) >= 1.5 * 1e-2:
33     #    self.beta *= 2
34     #if abs(dl) <= 1e-2 / 1.5:
35     #    self.beta *= 0.5
36     # print(self.beta)
37     """ ----- Programming 5 ----- """

```

## Results:



下图是自适应beta。



## Answers:

(a)

从图一我们可以看到，随着kl-constraint增大，每次进行游戏得到的总返回值是增大的。且kl-constraint越大，总返回值随着episode增大而增大的速度也越大。这是因为kl-constraint通过限制新的分布与旧的分布之间的KL散度，限制了actor梯度更新的范围。当kl-constraint增大，actor梯度更新的范围增大，随着训练轮数（episode）的增加，actor更新也就越快，所得到的总返回值也就越大。

(b)

从图二中我们可以看到，随着 $\beta$ 的增大，总返回值随着episode增大的趋势开始减小。这是因为 $\beta$ 和PPO中计算loss的KL-penalty惩罚项相关。当 $\beta$ 大的时候，这个惩罚项就变得很大，减小了loss的值，从而减小了梯度更新范围，让更新速度变慢。当 $\beta$ 为0时，我们从图中可以看到，PPO算法由于没有kl-penalty的约束而迅速发生了崩坏现象，导致总返回值变得很低。而 $\beta$ 不为0的时候，可以看到PPO算法迅速收敛。除此之外，通过比较图二图三，我们可以发现， $\beta$ 基本用于防止初始时通过限制新分布和旧分布之间的距离来防止，PPO算法的崩坏，而通过观察输出，之后 $\beta$ 基本降到一个接近于0的数。这意味着，崩坏现象通常在训练初期发生，这可以通过初始化一个相对大的 $\beta$ 来避免。在训练中后期可以通过调整 $\beta$ 的值来加速梯度更新。

(c)

他们在某种程度上是相似的，比如都和梯度更新范围有关。因为kl-constraint限制了新的分布和旧的分布之间的距离从而限制了梯度更新范围，而在PPO中，继承了TRPO的思想，使用 $\beta * KLD$ 的惩罚项来约束新旧分布之间的距离。所以他们是相似的。