

Chapter 6. Trust Region Policy Optimization & Proximal Policy Optimization

SAMSUNG AI

Jun 20, 2022

Kihun Kim, Mingyu Park



CORE
Control + Optimization Research Lab

TRPO - Review



TRPO - Review

policy gradient (with importance sampling)

$$\nabla_{\phi} J(\phi) = \mathbb{E}_{s \sim \rho_{\phi_{\text{old}}}, a \sim \pi_{\phi_{\text{old}}}} \left(A^{\pi_{\phi}}(s, a) \frac{\nabla_{\phi} \pi_{\phi}(a|s)}{\pi_{\phi_{\text{old}}}(a|s)} \right)$$

To estimate $A^{\pi_{\phi}}$, we use a separate value function approximator $V(s; \theta)$, and apply **generalized advantage estimation (GAE)**!

Furthermore, we use a stochastic policy $\pi_{\phi}(a|s)$ (Gaussian in our case).

Bad news: in TRPO, we have a lot of extra stuff to implement...
(Hessian-vector product, line search, etc.)

TRPO - Review

GAE?

Given a trajectory $(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T)$ generated by executing the current policy, we first compute

$$\delta_t = r_t + \gamma V(s_{t+1}; \theta) - V(s_t; \theta)$$

Then, GAE is computed as follows:

$$\hat{A}(s_t, a_t) = \sum_{\tau=t}^{T-1} \gamma^{\tau-t} \delta_\tau$$

TRPO - Review

training $V(s; \theta)$?

given a trajectory $(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T)$ generated from π_ϕ , we can compute Monte-Carlo estimates of $V(s_0), V(s_1), \dots, V(s_T)$ as follows:

$$V(s_t) = \sum_{\tau=t}^{T-1} \gamma^{\tau-t} r_\tau + \gamma^T V(s_T)$$

This is the **target** for value function update!

TRPO - Implementation

TRPO - Implementation

What info do we need when we implement **on-policy algorithms**?

1. (s_j, a_j, s'_j, r_j)
2. generalized advantage estimation(GAE)
3. MC estimates of value $V^{\pi_\phi}(s)$
4. probability $\pi_{\phi_{\text{old}}}(a_j|s_j)$

TRPO - Implementation

```
self._obs_mem = np.zeros(shape=(lim, dimS))  
self._act_mem = np.zeros(shape=(lim, dimA))  
self._rew_mem = np.zeros(shape=(lim,))  
self._val_mem = np.zeros(shape=(lim,))  
self._log_prob_mem = np.zeros(shape=(lim,))
```

collected during agent-env interaction

```
# memory of cumulative rewards which are MC-estimates of the current value function  
self._target_v_mem = np.zeros(shape=(lim,))  
# memory of GAE( $\lambda$ )-estimate of the current advantage function  
self._adv_mem = np.zeros(shape=(lim,))
```

computed when sampling a single episode is done



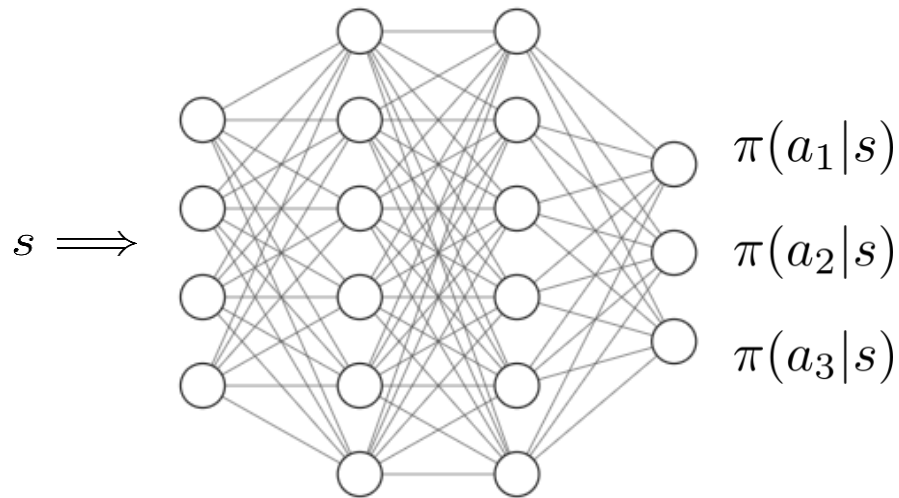
TRPO – Complete Implementation

TRPO – Network Architecture

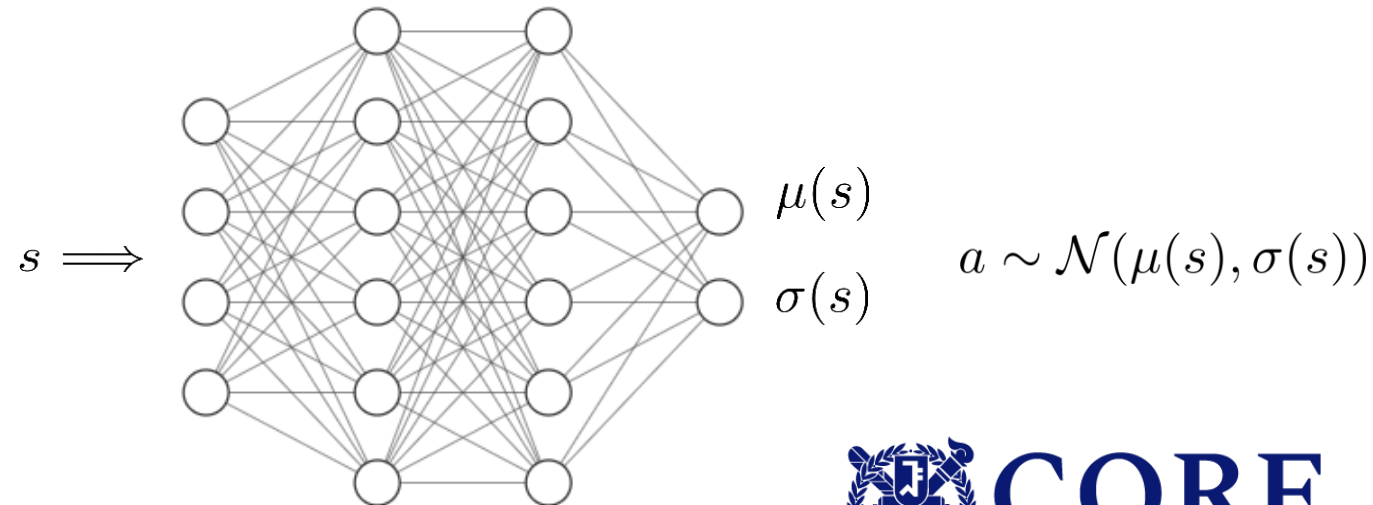
Key components of TRPO agent:

- policy network
 - can handle both discrete/continuous actions
 - network output : **probability distribution** over action space (cf. DDPG)

ex) discrete action space



ex) continuous action space



TRPO : Network Architecture

```
1  class Actor(nn.Module):
2      def __init__(self, obs_dim, act_dim, hidden1, hidden2):
3          super(Actor, self).__init__()
4          self.fc1 = nn.Linear(obs_dim, hidden1)
5          self.fc2 = nn.Linear(hidden1, hidden2)
6          self.fc3 = nn.Linear(hidden2, act_dim) # for \mu
7          self.fc4 = nn.Linear(hidden2, act_dim) # for \sigma
8
9      def forward(self, obs):
10         x = torch.tanh(self.fc1(obs))
11         x = torch.tanh(self.fc2(x))
12         mu = self.fc3(x)
13         log_sigma = self.fc4(x)
14         sigma = torch.exp(log_sigma)
15         return mu, sigma
16
17     def log_prob(self, obs, act):
18         mu, sigma = self.forward(obs)
19         act_distribution = Independent(Normal(mu, sigma), 1)
20         log_prob = act_distribution.log_prob(act)
21         return log_prob
```

1. Actor Network (By actor we just mean policy)

Why $\log \sigma(s)$ instead of $\sigma(s)$?
 \implies We have a constraint $\sigma(s) > 0$!

gives $\log \pi(a^i | s^i)$'s for (s^i, a^i) pairs

TRPO : Network Architecture

```
1  class Critic(nn.Module):
2      # critic  $V(s; \theta)$ 
3      def __init__(self, obs_dim, hidden1, hidden2):
4          super(Critic, self).__init__()
5          self.fc1 = nn.Linear(obs_dim, hidden1)
6          self.fc2 = nn.Linear(hidden1, hidden2)
7          self.fc3 = nn.Linear(hidden2, 1)
8
9      def forward(self, obs):
10         x = torch.tanh(self.fc1(obs))
11         x = torch.tanh(self.fc2(x))
12
13         return self.fc3(x)
```

2. Critic Network : super-easy!



TRPO : Action Selection by Agent

```
1  class TRPOAgent:
2      def __init__(self, obs_dim, act_dim, hidden1=64, hidden2=32):
3          ...
4
5      def act(self, obs):
6          obs = torch.tensor(obs, dtype=torch.float).to(device)
7          with torch.no_grad():
8              mu, sigma = self.pi(obs)
9              act_distribution = Independent(Normal(mu, sigma), 1)
10             action = act_distribution.sample()
11             log_prob = act_distribution.log_prob(action)
12             val = self.V(obs)
13         action = action.cpu().numpy()
14         log_prob = log_prob.cpu().numpy()
15         val = val.cpu().numpy()
16
17     return action, log_prob, val
```

Gaussian distribution $\mathcal{N}(\mu(s^i), \sigma(s^i))$

$a^i \sim \pi_\phi(s^i) := \mathcal{N}(\mu(s^i), \sigma(s^i))$

$\log \pi_\phi(a^i | s^i)$

$V(s^i; \theta)$



TRPO : Policy & Value Update

Unfortunately, the parameter updates in TRPO...

1. Fisher-vector product
2. conjugate gradient descent
3. backtracking line search

TRPO – Complete Implementation

1. Fisher-vector product

```
1 def fisher_vector_product(v, actor, obs_batch, cg_damping=1e-2):
2     v.detach()
3     kl = torch.mean(kl_div(actor=actor, old_actor=actor, obs_batch=obs_batch))
4     kl_grads = torch.autograd.grad(kl, actor.parameters(), create_graph=True)
5     kl_grad = torch.cat([grad.view(-1) for grad in kl_grads])
6     kl_grad_p = torch.sum(kl_grad * v)
7     Iv = torch.autograd.grad(kl_grad_p, actor.parameters()) # product of Fisher information I and v
8     Iv = flatten(Iv)
9     return Iv + v * cg_damping
```

Basically, Fish information matrix is **Hessian** of a function.
 \implies costly (requires n^2 numbers of 2nd-order differentiation)

However, computation of Ax can easily be done using automatic differentiation:
use

$$(\nabla_{\theta}^2 f)v = \nabla_{\theta}(\underbrace{\nabla_{\theta} f \cdot v}_{\text{scalar}}).$$

TRPO – Complete Implementation

2. conjugate gradient method

```
1  def cg(f_Ax, b, actor, obs_batch, cg_iters=10, residual_tol=1e-10):
2      p = b.clone()
3      r = b.clone()
4      x = torch.zeros_like(b)
5      rdotr = r.dot(r)
6      for i in range(cg_iters):
7          z = f_Ax(p, actor, obs_batch)
8          v = rdotr / p.dot(z)
9          x += v*p
10         r -= v*z
11         newrdotr = r.dot(r)
12         mu = newrdotr/rdotr
13         p = r + mu*p
14         rdotr = newrdotr
15         if rdotr < residual_tol:
16             break
17     return x
```

works if computation of Ax is possible for any x (avoids taking A^{-1})

What is conjugate gradient for?
 \implies To efficiently solve $As = g$

TRPO – Complete Implementation

3. backtracking line search

```
1 def backtracking_line_search(...):
2     backtrac_coef = 1.0
3     alpha = 0.5
4     beta = 0.5
5     flag = False
6     expected_improve = (actor_loss_grad * maximal_step).sum(0, keepdim=True)
7     for i in range(10):
8         new_params = params + backtrac_coef * maximal_step
9         update_model(actor, new_params)
10        new_actor_loss = surrogate_loss(actor, adv, states, old_policy.detach(), actions)
11        loss_improve = new_actor_loss - actor_loss
12        expected_improve *= backtrac_coef
13        improve_condition = loss_improve / expected_improve
14        kl = kl_div(actor=actor, old_actor=old_actor, obs_batch=states)
15        kl = kl.mean()
16        if kl < max_kl and improve_condition > alpha:
17            flag = True
18            break
19        backtrac_coef *= beta
20    if not flag:
21        params = flat_params(old_actor)
22        update_model(actor, params)
```

$\theta \leftarrow \theta + \beta s$ where s : search direction found through CG

check if new policy stays within **trust region**

decrease β and repeat the verification

TRPO – Complete Implementation

```
1  def train(env, agent, max_iter, gamma=0.99, lr=3e-4, lam=0.95, delta=1e-3, steps_per_epoch=4000, eval_interval=4000):
2      ...
3      for epoch in range(num_epochs):
4          state = env.reset()
5          step_count = 0
6          ep_reward = 0
7          for t in range(steps_per_epoch):
8              action, log_prob, v = agent.act(state)
9              next_state, reward, done, _ = env.step(action)
10             memory.append(state, action, reward, v, log_prob)
11             ep_reward += reward
12             step_count += 1
13             if (step_count == max_ep_len) or (t == steps_per_epoch - 1):
14                 s_last = torch.tensor(next_state, dtype=torch.float).to(device)
15                 v_last = agent.V(s_last).item()
16                 memory.compute_values(v_last)
17             elif done:
18                 v_last = 0.0
19                 memory.compute_values(v_last)
```

TRPO – Complete Implementation

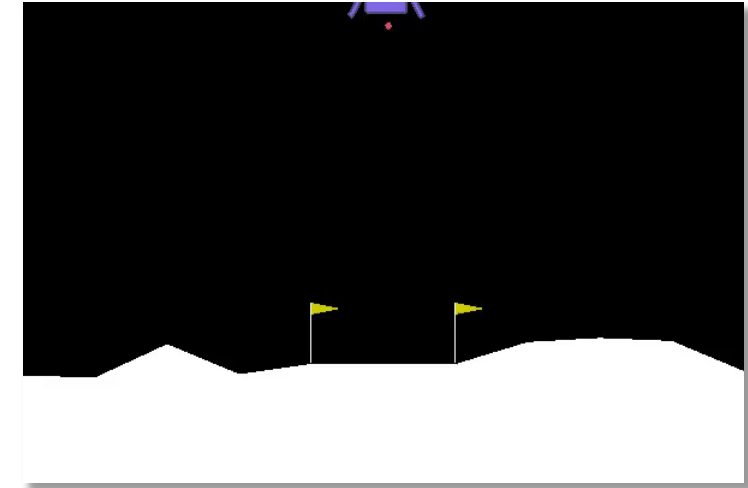
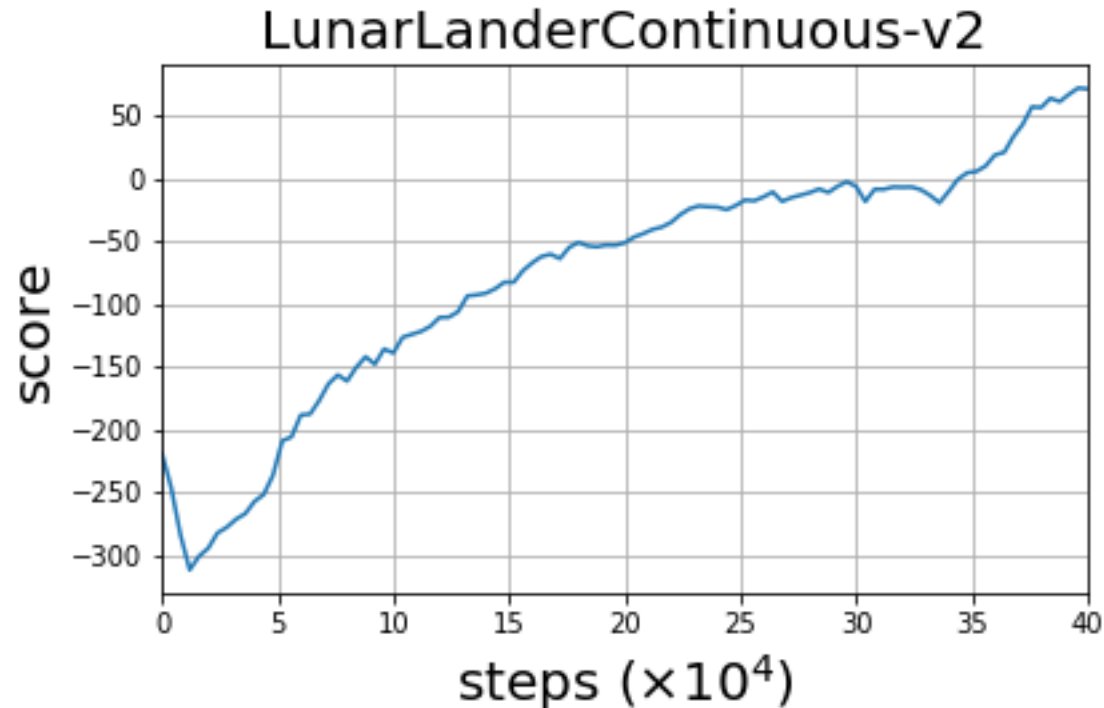
```
21     state = next_state
22     if done:
23         state = env.reset()
24         step_count = 0
25         ep_reward = 0
26     total_t += 1
27     update(agent, memory, critic_optim, delta, num_updates=1)
28     return
```

TRPO - Experiment

OpenAI Gym LunarLanderContinuous-v2

Training DDPG on the task was extremely unstable...

What about TRPO?



Proximal Policy Optimization(PPO) : What's different?

```
1  for i in range(num_iter):
2      log_probs, ent = self.pi.compute_log_prob(states, actions)
3      r = torch.exp(log_probs - old_log_probs)
4      clipped_r = torch.clamp(r, 1 - self.epsilon, 1 + self.epsilon)
5      single_step_obj = torch.min(r * A, clipped_r * A)
6      pi_loss = -torch.mean(single_step_obj)
7      v = self.V(states)
8      V_loss = torch.mean((v - target_v) ** 2)
9      ent_bonus = torch.mean(ent)
10     loss = pi_loss + 0.5 * V_loss - 0.01 * ent_bonus
11     self.optimizer.zero_grad()
12     loss.backward()
13     self.optimizer.step()
14     if i == num_iter - 1:
15         kl = torch.mean(old_log_probs - log_probs).item()
16         if kl > 1.5 * self.kl_threshold:
17             break
18     return
```

simple, surrogate loss function for policy network parameters!
⇒ 1st-order optimization (no Hessian needed)