

Reinforcement Learning

Chapter 5: RL via Policy Gradient

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

Department of Electrical and Computer Engineering
University of Toronto

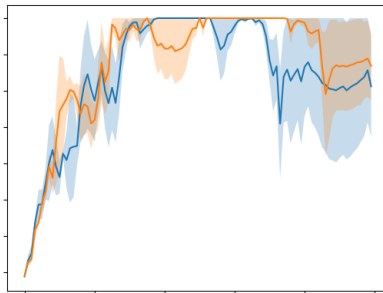
Fall 2025

Observing Basic PGM

Let's break the problem down: even by using *baseline*, we still observe *instability* while we use PGM

If we plot the *average reward* we collect through time

- We might see it getting improved up to some point
- It then could drop drastically at some other point



Main Reason: *Estimate Variance*

The main reason for this behavior is *high variance of the gradient estimator*:
recall that we estimate the gradient of *average value* by

$$\hat{\nabla} \mathcal{J}(\pi_{\theta}) = \sum_{t=0}^{T-1} U_t \nabla \log \pi_{\theta}(A_t | S_t)$$

This is true that

$$\mathbb{E}_{\pi_{\theta}} \left\{ \hat{\nabla} \mathcal{J}(\pi_{\theta}) \right\} = \nabla \mathcal{J}(\pi_{\theta})$$

However, its variance, i.e.,

$$\text{Var} \left\{ \hat{\nabla} \mathcal{J}(\pi_{\theta}) \right\} = \mathbb{E}_{\pi_{\theta}} \left\{ \left(\hat{\nabla} \mathcal{J}(\pi_{\theta}) - \nabla \mathcal{J}(\pi_{\theta}) \right)^2 \right\}$$

could be very large: one given sample take us far away from the true direction!

Reducing Variance: Using Mini-Batches

- + But, isn't that always the case in SGD?! We assume that those errors cancel each other out! Right?!
- That's right! But apparently, it's not working here!

To find out an explanation to this behavior, let's try to **reduce the variance** by using larger mini-batches

- Collect B sample trajectories by policy π_{θ}
- Compute the gradient estimator for each trajectory

$$\hat{\nabla}_b \mathcal{J}(\pi_{\theta}) = \sum_{t=0}^{T-1} U_t \nabla \log \pi_{\theta}(A_t | S_t)$$

- Average the estimators to get a better estimator

$$\hat{\nabla} \mathcal{J}(\pi_{\theta}) = \frac{1}{B} \sum_{b=1}^B \hat{\nabla}_b \mathcal{J}(\pi_{\theta})$$

Advantage PGM: Mini-Batch Version

AdvantagePGM() :

1: Initiate with θ and learning rate α

2: **while** *interacting* **do**

3: **for** *mini-batch* $b = 1 : B$ **do**

4: Sample the environment with policy π_θ

$$S_0, A_0 \xrightarrow{R_1} S_1, A_1 \xrightarrow{R_2} \dots \xrightarrow{R_{T-1}} S_{T-1}, A_{T-1} \xrightarrow{R_T} S_T$$

5: **for** $t = 0 : T - 1$ **do**

6: Compute *sample advantage* $U_t = R_{t+1} + \gamma v_{\pi_\theta}(S_{t+1}) - v_{\pi_\theta}(S_t)$

7: **end for**

8: Compute sample gradient $\hat{\nabla}_b = \sum_{t=0}^{T-1} U_t \nabla \log \pi_\theta(A_t | S_t)$

9: **end for**

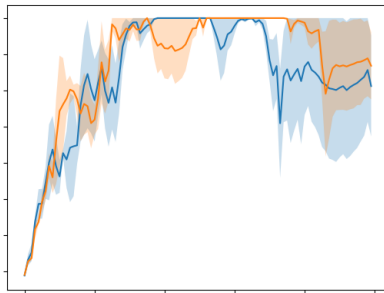
10: Update policy network

$$\theta \leftarrow \theta + \frac{\alpha}{B} \sum_{b=1}^B \hat{\nabla}_b$$

11: **end while**

Observing Mini-Batch

After trying mini-batch PGM: we see that the variance of the curve slightly improves; however, we still see that problem



- + What does this say then?
- It says that the problem is simply coming from **high variance**

Alternative Look at Advantage Optimization

In the latest version of PGM: we saw that the gradient of the *average value* can be computed as

$$\nabla \mathcal{J}(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} \{ u_{\pi_{\theta}}(S, A) \nabla \log \pi_{\theta}(A|S) \}$$

Let's assume that we can compute it exactly: then, we will use gradient descent to find optimal θ , i.e.,

AdvantageGD():

- 1: Start with some initial θ_0
- 2: **for** $k = 1 : K$ **do**
- 3: Compute the *exact gradient*

$$\nabla \mathcal{J}(\pi_{\theta_k}) = \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \nabla \log \pi_{\theta}(A|S) \mid \theta = \theta_k \right\}$$

- 4: Update the parameters as $\theta_{k+1} = \theta_k + \alpha \nabla \mathcal{J}(\pi_{\theta_k})$
- 5: **end for**

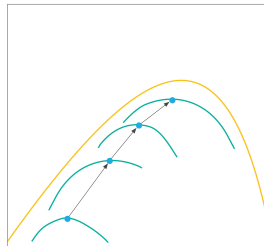
Alternative Look: Surrogate Function

Let us now define the following **surrogate** function at point θ_k

$$\mathcal{L}_k(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \frac{\pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)} \right\}$$

We should pay attention that

- ① We have a **sequence of surrogate** functions
 - ↳ Each function is defined locally at point θ_k
- ② $\pi_{\theta}(A|S)$ is the **only term that belongs to θ**
 - ↳ Everything else depends on θ_k which is a **constant**



Alternative Look: *Surrogate Function*

$$\mathcal{L}_k(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \frac{\pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)} \right\}$$

Let's compute the gradient of this surrogate function at $\theta = \theta_k$

$$\begin{aligned} \nabla \mathcal{L}_k(\pi_{\theta_k}) &= \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \frac{\nabla \pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)} \right\} \\ &= \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \nabla \log \pi_{\theta}(A|S) \mid_{\theta=\theta_k} \right\} \\ &= \nabla \mathcal{J}(\pi_{\theta_k}) \end{aligned}$$

This is exactly the gradient that we update our policy with!

Alternative Look: GD with Surrogate Function

So, we could re-write the gradient descent loop as

AdvantageGD() :

- 1: Start with some initial θ_0
- 2: **for** $k = 1 : K$ **do**
- 3: Compute the **exact gradient** $\nabla \mathcal{L}_k(\pi_{\theta_k})$
- 4: Update the parameters as $\theta_{k+1} = \theta_k + \alpha \nabla \mathcal{L}_k(\pi_{\theta_k})$
- 5: **end for**

- + Now, what's the point?! It's still same problem!
- Sure! But, let's see what we are doing now

In each iteration we add gradient scaled with α to the previous parameters

- Why we do that?
- + We want to increase $\mathcal{L}_k(\pi_{\theta})$ maximally
- Exactly! So, why don't we simply replace θ_{k+1} with its maximizer?!

Alternative Look: GD with Surrogate Function

We re-write the gradient descent loop as follows

AdvantageGD() :

- 1: Start with some initial θ_0
- 2: **for** $k = 1 : K$ **do**
- 3: Compute the **surrogate function** $\mathcal{L}_k(\pi_\theta)$
- 4: Update the parameters as $\theta_{k+1} = \operatorname{argmax}_\theta \mathcal{L}_k(\pi_\theta)$
- 5: **end for**

This algorithm does the **learning rate** tuning by itself

It **gets us rid of** specifying the **learning rate** α

- + Nice job! But, how are we su[supposed] to find the **surrogate function**?
- You could guess! By **sampling**!

Understanding Surrogate Function

Let's get back to the definition of the **surrogate** function: *it is easy to interpret it as importance sampling*

$$\mathcal{L}_k(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \frac{\pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)} \right\}$$

We sample trajectories by policy π_{θ_k}

- We collect **advantage samples** $u_{\pi_{\theta_k}}(S, A)$
- We do know the policy samples $\pi_{\theta_k}(A|S)$

But we compute the average with respect to π_{θ}

This gives us a **very good explanation** of why PGM is **unstable**

Review: Importance Sampling Trade-off

Consider the basic setting in *importance sampling*:

we have samples from $X \sim p(x)$ but we want to estimate $X \sim q(x)$

If we could sample from $q(x)$

$$\mu = \mathbb{E}_q \{X\}$$

The variance of the estimate is

$$\begin{aligned} \text{Var} \{X\} &= \mathbb{E}_q \left\{ (X - \mu)^2 \right\} \\ &= \mathbb{E}_q \{X^2\} - \mu^2 = \sigma^2 \end{aligned}$$

Review: Importance Sampling Trade-off

Consider the basic setting in *importance sampling*:

we have samples from $X \sim p(x)$ but we want to estimate $X \sim q(x)$

Now that we sample $p(x)$

$$\bar{\mu} = \mathbb{E}_p \left\{ X \frac{q(X)}{p(X)} \right\} = \mathbb{E}_q \{X\} = \mu$$

The variance of this estimate is

$$\begin{aligned} \text{Var} \{X\} &= \mathbb{E}_p \left\{ \left(X \frac{q(X)}{p(X)} \right)^2 \right\} - \mu^2 = \int x^2 \frac{q^2(x)}{p^2(x)} p(x) - \mu^2 \\ &= \int x^2 \frac{q(x)}{p(x)} q(x) - \mu^2 = \mathbb{E}_q \left\{ X^2 \frac{q(X)}{p(X)} \right\} - \mu^2 \neq \sigma^2 \end{aligned}$$

Review: Importance Sampling Trade-off

Consider the basic setting in *importance sampling*:

we have samples from $X \sim p(x)$ but we want to estimate $X \sim q(x)$

If we could sample from $q(x)$

$$\mu = \mathbb{E}_q \{X\}$$

and the estimate variance is

$$\sigma^2 = \mathbb{E}_q \{X^2\} - \mu^2$$

Now that we sample $p(x)$

$$\mu = \mathbb{E}_q \{X\}$$

and the estimate variance is

$$\bar{\sigma}^2 = \mathbb{E}_q \left\{ X^2 \frac{q(X)}{p(X)} \right\} - \mu^2$$

With *importance sampling*, *variance* scales with *ratio of the distributions*

Back to Surrogate: Root of High Variance

Gradient descent based on *surrogate functions* optimizes

$$\mathcal{L}_k(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta_k}} \left\{ u_{\pi_{\theta_k}}(S, A) \frac{\pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)} \right\}$$

We can assume that by sampling the environment, we are estimating this *surrogate function* by *importance sampling* from samples of policy π_{θ_k} : let $\hat{\mathcal{L}}_k(\pi_{\theta})$ be our estimate; then we could say

$$\text{Var} \left\{ \hat{\mathcal{L}}_k(\pi_{\theta}) \right\} \propto \frac{\pi_{\theta}(A|S)}{\pi_{\theta_k}(A|S)}$$

- This explains why we see unreliable estimates in PGM!
- + How exactly?!
- Let's break it down!

Back to Surrogate: Root of High Variance

$$\text{Var} \left\{ \hat{\mathcal{L}}_k (\pi_{\theta}) \right\} \propto \frac{\pi_{\theta} (A|S)}{\pi_{\theta_k} (A|S)}$$

Looking at this variance, we could say: we should not **naively** update θ_k by maximizing its **surrogate function**. We should update it such that

- ① **surrogate function** is maximized, and
- ② the policy specified by $\pi_{\theta_{k+1}}$ is rather close to π_{θ_k}

-
- + But aren't we doing that?! We just change the policy parameters slightly, i.e., $\|\theta_{k+1} - \theta_k\|^2$ is typically small
 - Well! That doesn't say anything about **difference between $\pi_{\theta_{k+1}}$ and π_{θ_k}**

Observation: Sensitivity of Policy Network

Sensitivity of Policy

Even if we change parameters θ *slightly*, π_θ can change *hugely*!

Given this observation, we could modify our gradient descent loop as

AdvantageGD():

- 1: Start with some initial θ_0
- 2: **for** $k = 1 : K$ **do**
- 3: Compute the *surrogate function* $\mathcal{L}_k(\pi_\theta)$
- 4: Update the parameters as

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}_k(\pi_\theta) \quad \text{subject to} \quad \pi_\theta \text{ and } \pi_{\theta_k} \text{ are close}$$

5: **end for**

+ How can we quantify “ π_θ and π_{θ_k} being close”?

Review: Kullback-Leibler Divergence

KL Divergence

Kullback-Leibler divergence between two distributions p and q is defined as

$$D_{\text{KL}}(p\|q) = \mathbb{E}_p \left\{ \log \left(\frac{p(X)}{q(X)} \right) \right\} = \int_x \log \left(\frac{p(x)}{q(x)} \right) p(x)$$

*We can use this definition to find the **divergence** between π_{θ} and π_{θ_k}*

$$\bar{D}_{\text{KL}}(\pi_{\theta}\|\pi_{\theta_K}) = \mathbb{E}_{S \sim d_{\pi_{\theta_k}}} \left\{ \mathbb{E}_{\pi_{\theta}} \left\{ \log \left(\frac{\pi_{\theta}(\textcolor{red}{A}|\textcolor{green}{S})}{\pi_{\theta_k}(\textcolor{red}{A}|\textcolor{green}{S})} \right) \right\} \right\}$$

Review: Properties of KL Divergence

KL divergence shows interesting properties

- It is **zero** when distributions are the **same**

$$D_{\text{KL}}(p\|p) = 0$$

and **increases** when they get **more different**

- It is always **non-negative**, i.e., for any p and q

$$D_{\text{KL}}(p\|q) \geq 0$$

↳ This property is often called **Gibbs' inequality**

But remember that KL divergence is asymmetric, i.e.,

$$D_{\text{KL}}(\textcolor{red}{p}\|\textcolor{green}{q}) \neq D_{\text{KL}}(\textcolor{green}{q}\|\textcolor{red}{p})$$

Trust Region Policy Gradient Method

Back to our modified gradient descent loop, we could write

AdvantageGD() :

- 1: Start with some initial θ_0
- 2: **for** $k = 1 : K$ **do**
- 3: Compute the *surrogate function* $\mathcal{L}_k(\pi_\theta)$
- 4: Update the parameters as

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmax}} \mathcal{L}_k(\pi_\theta) \quad \text{subject to} \quad \bar{D}_{\text{KL}}(\pi_\theta \| \pi_{\theta_k}) \leq d_{\max}$$

5: **end for**

This modified approach is called

Trust Region PGM

*since it computes the best policy gradient within a **trusted resion***

Surrogate Optimization: *Exact Solution*

- + *How can we solve the optimization in the loop then?*
- *As you could guess, we are going to find a way around it!*

The concrete way to solve it is to use **regularization**: we want to solve

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \mathcal{L}_k(\pi_{\theta}) \quad \text{subject to} \quad \bar{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k}) \leq d_{\max}$$

We solve instead

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \mathcal{L}_k(\pi_{\theta}) - \beta (\bar{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k}) - d_{\max})$$

for some β that potentially minimizes the **regularized** objective

*This is going to be computationally very **expensive**!*

Surrogate Optimization: Natural Policy Gradient

We instead use *Taylor expansion* to approximate both *surrogate* and *constraint*

Taylor Expansion

An analytic function $f(x)$ can be expanded around point x_0 as

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + \dots$$

Let's start with the *surrogate function*

$$\mathcal{L}_k(\pi_{\theta}) = \mathcal{L}_k(\pi_{\theta_k}) + \nabla \mathcal{L}_k(\pi_{\theta_k})^T (\theta - \theta_k) + \varepsilon$$

In Assignment 3, we show $\mathcal{L}_k(\pi_{\theta_k}) = 0$: so, setting $\nabla_k = \nabla \mathcal{L}_k(\pi_{\theta_k})$ we get

$$\mathcal{L}_k(\pi_{\theta}) \approx \nabla_k^T (\theta - \theta_k)$$

Surrogate Optimization: *Natural Policy Gradient*

Next, we go for *constraint term*

$$\begin{aligned}\bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) &= \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}_k} \parallel \pi_{\boldsymbol{\theta}_k}) + \nabla \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) \big|_{\boldsymbol{\theta}_k}^{\top} (\boldsymbol{\theta} - \boldsymbol{\theta}_k) \\ &\quad + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_k)^{\top} \nabla^2 \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) \big|_{\boldsymbol{\theta}_k} (\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \varepsilon\end{aligned}$$

In Assignment 3, we show that

$$\begin{aligned}\bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}_k} \parallel \pi_{\boldsymbol{\theta}_k}) &= 0 \\ \nabla \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) \big|_{\boldsymbol{\theta}_k} &= \mathbf{0}\end{aligned}$$

So, by defining $\mathbf{H}_k = \nabla^2 \bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) \big|_{\boldsymbol{\theta}_k}$ we have

$$\bar{D}_{\text{KL}}(\pi_{\boldsymbol{\theta}} \parallel \pi_{\boldsymbol{\theta}_k}) \approx \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_k)^{\top} \mathbf{H}_k (\boldsymbol{\theta} - \boldsymbol{\theta}_k)$$

Surrogate Optimization: Natural Policy Gradient

Now, let's replace these approximations

$$\begin{aligned}\mathcal{L}_k(\pi_{\theta}) &\approx \nabla_k^{\top} (\theta - \theta_k) \\ \bar{D}_{\text{KL}}(\pi_{\theta} \parallel \pi_{\theta_k}) &\approx \frac{1}{2} (\theta - \theta_k)^{\top} \mathbf{H}_k (\theta - \theta_k)\end{aligned}$$

into the optimization problem

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmax}} \nabla_k^{\top} (\theta - \theta_k) \quad \text{subject to} \quad \frac{1}{2} (\theta - \theta_k)^{\top} \mathbf{H}_k (\theta - \theta_k) \leq d_{\max}$$

This is a classic **linear programming** whose solution is given by

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2d_{\max}}{\nabla_k^{\top} \mathbf{H}_k^{-1} \nabla_k}} \mathbf{H}_k^{-1} \nabla_k$$

Surrogate Optimization: Natural Policy Gradient

This is like classic update with *linear correction* and *tuned learning rate*

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \sqrt{\frac{2d_{\max}}{\nabla_k^T \mathbf{H}_k^{-1} \nabla_k}} \mathbf{H}_k^{-1} \nabla_k$$

This is often referred to as

natural policy gradient

It gives a better direction for update; however,

- It could still *not* increase *surrogate* or *deviate constraint*
 - ↳ This is due to the *inaccuracy* of approximations
- It requires estimate of \mathbf{H}_k which is *computationally expensive*
- It also needs to invert estimate of \mathbf{H}_k which is again *expensive*

Natural PGM

NaturalPGM():

- 1: Start with some initial θ
- 2: **while** *interacting* **do**
- 3: **for** *mini-batch* $b = 1 : B$ **do**
- 4: Sample the environment with policy π_θ

$$S_0, A_0 \xrightarrow{R_1} S_1, A_1 \xrightarrow{R_2} \dots \xrightarrow{R_{T-1}} S_{T-1}, A_{T-1} \xrightarrow{R_T} S_T$$

- 5: **end for**
- 6: Estimate $\hat{\nabla}$ and $\hat{\mathbf{H}}$ from samples
- 7: Update the parameters as

$$\theta \leftarrow \theta + \sqrt{\frac{2d_{\max}}{\hat{\nabla}^\top \hat{\mathbf{H}} \hat{\nabla}}} \hat{\mathbf{H}}^{-1} \hat{\nabla}$$

- 8: **end while**

TRPO and PPO Algorithms

There are two sets of solutions to overcome the challenges in *natural* PGM

- Trust Region Policy Optimization

- ↳ Regularize learning rate via backtracking line
 - ↳ This guarantees the satisfaction of constraint
- ↳ Use sampling to find the estimate \hat{H}
- ↳ Use gradient conjugate to compute the inverse

- Proximal Policy Optimization

- ↳ Skip all these steps by computationally-efficient clipping
 - ↳ The clipping guarantees the satisfaction of constraint
- ↳ We do not need to find estimate \hat{H} anymore