

Variance in Policy Gradients

$$\nabla_\theta J(\theta) = \frac{1}{N} \sum_{\tau} \nabla_\theta \log \pi_\theta(\tau) \times r(\tau)$$

Now, if you add a constant to all rewards, say $+100$ to all rewards, then, the reward for is still practically the same!

But then, the gradients for policy significantly changes!

For example,

say there are 2 trajectories

$$\begin{array}{c} A \\ \hline 500 \end{array} \quad \begin{array}{c} B \\ \hline -500 \end{array}$$

$$\text{Add 1 million to both} \quad \begin{array}{c} 1,000,500 \\ 999,500 \end{array}$$

In this life, one τ looks **BAD!** & one τ looks **GOOD!**

In this life, **BOTH** are looking good!

Also, mathematically since the gradients are in "product" of rewards! if every, $\lambda(\tau)$ gets added a constant b !

Then every gradient is multipled by a constant b !

Say, there are n no's a_1, a_2, \dots, a_n whose mean is μ , & variance is σ^2 . Now, scale

all elements $\boxed{\text{by } b}$ $\rightarrow a_1b, a_2b, \dots, a_nb$

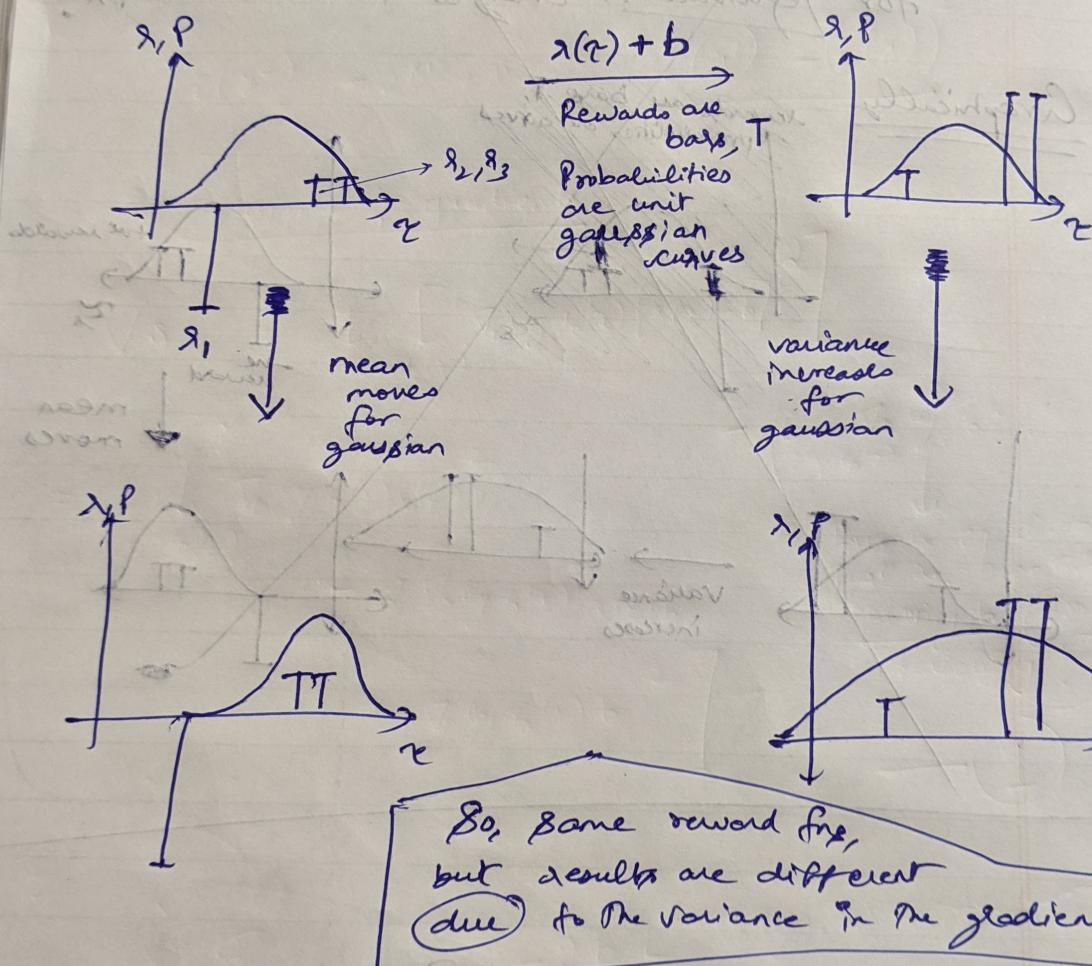
$$\text{mean} = \frac{a_1b + a_2b + \dots + a_nb}{N} = b\mu$$

$$\text{Variance} = \frac{(a_1b - b\mu)^2 + (a_2b - b\mu)^2 + \dots + (a_nb - b\mu)^2}{N} = \underline{\underline{b^2\sigma^2}}$$

The variance of those 'n' numbers increases!

So, if the reward fn is modified equivalently by adding a constant, Then the gradients are multiplied by a constant factor increasing its variance QUADRATICALLY

Graphically, let's say there are 2 reward fns off by a constant factor



So, if we take the mean of rewards & subtract it to all rewards, this would make both reward fns in this example equivalent!

This mean subtraction would nullify any constants that may have been added to the rewards!

Expectation
gradient
mean

This is same for

r_1
 $J(\theta)$

difference

mean

B

OBVIOUSLY

Also,

By adding a constant to all rewards, we [only] reduce - The variance & gradients, keeping its mean

[Same]

Proof

$$\nabla_\theta J(\theta)$$

$$= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\tau) [a(\tau) - b]]$$

$$= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\tau) a(\tau)] - \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\tau) \times b]$$

Unlike $a(\tau)$,
b doesn't change for
each τ , b is
constant for all τ s,
so, it can be taken
common

$$= -b \mathbb{E}_{\tau \sim \pi_\theta(\tau)} [\nabla_\theta \log \pi_\theta(\tau)]$$

$$= -b \sum_{\text{all } \tau} \pi_\theta(\tau) \nabla_\theta \log \pi_\theta(\tau)$$

$$= -b \sum_{\text{all } \tau} V_\theta \pi_\theta(\tau)$$

$$= -b \nabla_\theta \left(\sum_{\text{all } \tau} \pi_\theta(\tau) \right)$$

$$= -b \nabla_\theta (1) \rightarrow 0$$

This is because, adding
same constant to all

π_θ (Just) increases

$J(\theta)$ by a constant, so,

differentiation of $J(\theta)$ would

not differentiate the constant +
make it 0!

But The variance of it is

not zero!

Proof
PTO

Variance of \bar{x} has whose mean is μ !

$$\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_n - \mu)^2}{n}$$

denoted as σ^2 standard deviation of sample = $\sigma^2 = \frac{1}{n} \sum (x_i - \mu)^2$

where the process N consists of

sample

$$= \frac{x_1^2 + x_2^2 + \dots + x_n^2}{N} - 2\mu \left(\frac{x_1 + x_2 + \dots + x_n}{N} \right) + \frac{N\mu^2}{N}$$

$\left[\frac{d - (x_1 + x_2 + \dots + x_n)}{N} \right] \downarrow \mu$

$$= E(x^2) - 2\mu^2 + \mu^2 = E(x^2) - (E(x))^2$$

So, variance of gradients

$$\text{Var} = E \left[\frac{\partial \log \pi_\theta(x)}{\partial \theta} (\hat{x}(x) - b) \right]^2$$

$$- \left(E \frac{\partial \log \pi_\theta(x)}{\partial \theta} \right)^2$$

This is equal to
Just the expectation,
ignoring b , (Proven
Previously!)

①

$$- 2 \left(E \frac{\partial \log \pi_\theta(x)}{\partial \theta} \right) \hat{x}(x) b$$

$$+ \left(E \frac{\partial \log \pi_\theta(x)}{\partial \theta} b \right)^2$$

old
 σ
(if b was
never
added)

$$E \left[\left(\frac{\partial \log \pi_\theta(x)}{\partial \theta} \hat{x}(x) \right) \right]^2$$
$$- \left(E \left[\frac{\partial \log \pi_\theta(x)}{\partial \theta} \hat{x}(x) \right] \right)^2$$

$$- 2b E \left(\frac{\partial \log \pi_\theta(x)}{\partial \theta} \right) \hat{x}(x)$$

$$+ b^2 E \left(\frac{\partial \log \pi_\theta(x)}{\partial \theta} \right)^2$$

So, old σ + non-zero
terms.

Hence Proved, the variance has CHANGED

So, what should b , be, so, variance is
@ minimum?

$$\frac{\partial \text{Var}}{\partial b} = 0 \quad \text{Eq from ①}$$

$$\begin{aligned} \text{Var} &= E \left[\left(V_0 \log \pi_b(r) \lambda(r) \right)^2 \right] \\ &\quad - \left(E \left[V_0 \log \pi_b(r) \lambda(r) \right] \right)^2 \\ \frac{\partial \text{Var}}{\partial b} &= \underbrace{-2b E \left(V_0 \log \pi_b(r) \right)^2 \lambda(r)}_{\text{O}} \\ &\quad + b^2 E \left(V_0 \log \pi_b(r) \right)^2 \\ &\quad - 2 E \left(V_0 \log \pi_b(r) \right)^2 \lambda(r) \\ &\quad + 2b E \left(V_0 \log \pi_b(r) \right)^2 \end{aligned}$$

$$0 = \cancel{2b E \left(V_0 \log \pi_b(r) \right)^2} - \cancel{2 E \left(V_0 \log \pi_b(r) \right)^2 \lambda(r)}$$

$$\frac{E \left(\left(V_0 \log \pi_b(r) \right)^2 \lambda(r) \right)}{E \left(\left(V_0 \log \pi_b(r) \right)^2 \right)} = b$$

So, expectation of rewards weighted by

(No log $\pi_b(r)$)

~~These reward terms involving the
 So, the approximation, etc just~~