

Trust Region methods

Konpat Preechakul
Chulalongkorn University
November 2019

Recap policy gradient

- Policy gradient theorem

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

- Variance problem
- Advantage

$$A(s, a) = Q(s, a) - V(s)$$

- Variance reduction
 - A2C

$$\nabla J(\theta) = \mathbb{E}_{s,a} [A(s, a) \nabla \log \pi_\theta(a|s)]$$

Recap policy gradient

- Off-policy gradient
 - Off-policy critic
 - Off-policy actor

$$\nabla_{\theta} J(\theta) \approx \sum_s d^b(s) \sum_a Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(a|s)$$

- Deterministic policy gradient
 - DDPG

$$\nabla_{\theta} J(\theta) = \sum_s d^{\pi}(s) \nabla_a Q_{\phi}(s, a)|_{a=\pi(s)} \nabla_{\theta} \pi_{\theta}(s)$$

Today's topic

Today topics

- Why policy gradient fails?
- More robust policy gradient
 - Trust region methods
 - Approximation
- Trust region on critic
- Natural gradients

Why policy gradient fails?

Bad critic

- Critic is not oracle, it has its flaws
- **What are some flaws?**

Bad critic

- Critic is not oracle, it has its flaws
- How to reduce the flaws?
- **Critic is prone to forgetting**
- When does critic forget?

Critic is prone to forgetting

- On-policy training “limits” kind of data the critic sees
- If the data is concentrated in “late game”, the critic forgets “early game” states
 - It is likely to give gibberish Q to the actor
- Underlines using replay, parallel actors

Bad critic

- Critic is not oracle, it has its flaws
- Critic is prone to forgetting
- **How to know when to trust?**

When to trust critic?

- Do critic have “confidence”?

When to trust critic?

- Do critic have “confidence”?
- How to get one?

When to trust critic?

- Do critic have “confidence”?
- How to get one?
- Confidence is challenging:
 - Critic outputs variance?

When to trust critic?

- Do critic have “confidence”?
- How to get one?
- Confidence is challenging:
 - Critic outputs variance?
 - Multiple critics?

When to trust critic?

- Do critic have “confidence”?
- How to get one?
- Confidence is challenging:
 - Critic outputs variance?
 - Multiple critics?
 - Critic dropout?

Bad critic

- Critic is not oracle, it has its flaws
- Critic is prone to forgetting
- How to know when to trust?
- **How to limit the trust?**
- Not updating too much ...

Not updating too much

Not updating too much

- How much is too much?
 - The right LR?
 - What is the right LR?
 - If the critic is “abruptly” large, no small LR is small enough
-
- We need to be serious about update!

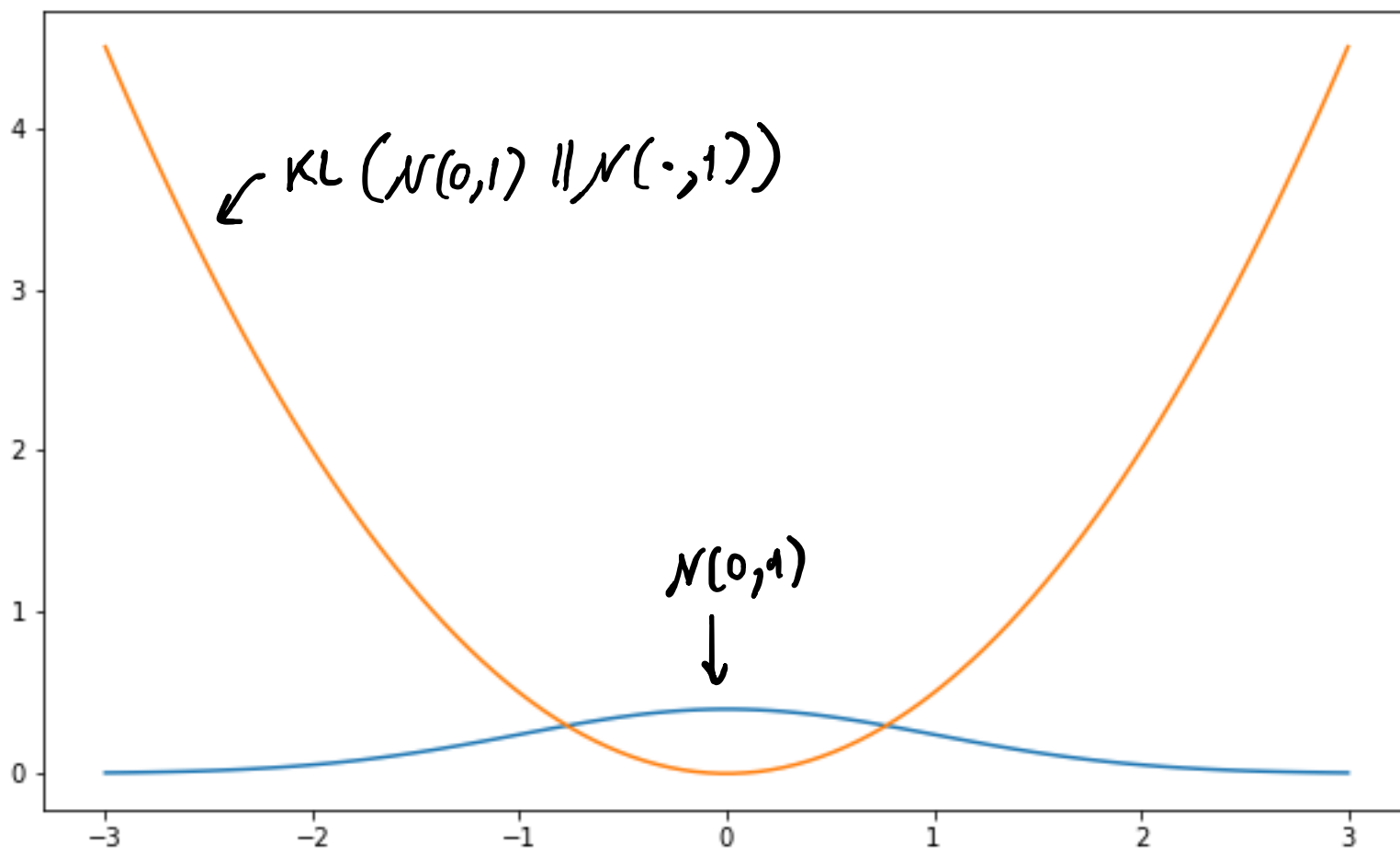
Quantifying “much”

- Policies are probability functions
- We find a “distance” measure on the probabilities
- KL divergence:

$$\text{KL}(P\|Q) = \sum_x P(x) [\log P(x) - \log Q(x)]$$

KL Divergence in picture

$$\text{KL}(P\|Q) = \sum_x P(x) [\log P(x) - \log Q(x)]$$



Quantify “much” in RL

- Context of RL, for a state S :

$$\text{KL}_s(\theta||\theta') = \sum_a \pi_\theta(a|s) [\log \pi_\theta(a|s) - \log \pi_{\theta'}(a|s)]$$

- For all states:

$$\text{KL}(\theta||\theta') = \sum_s d^\pi(s) \text{KL}_s(\theta||\theta')$$

**Do policy gradient while
making sure not going too
far in terms of KL**


New objective for policy update

- From the policy gradient objective

$$J(\theta) = \sum_s P_{s_0} V^\pi(s)$$

$$\theta \leftarrow \theta + d \quad d = \alpha \nabla_\theta J(\theta)$$

- A new update direction should be: *don't worry too much*

$$d^* = \operatorname{argmax}_d J(\theta + d) \quad \text{s.t.} \quad \text{KL}(\theta \| \theta + d) = c$$


New update direction

$$d^* = \operatorname{argmax}_d J(\theta + d) \quad \text{s.t.} \quad \text{KL}(\theta \| \theta + d) = c$$

- A constrained optimization
- We relax them using **Lagrangian**:

$$\mathcal{L}(d, \lambda) = J(\theta + d) + \lambda (\text{KL}(\theta \| \theta + d) - c)$$

- Optimal d is at the critical point

$$\nabla_{d, \lambda} \mathcal{L} = 0$$

Solving for the update direction

$$\mathcal{L}(d, \lambda) = J(\theta + d) + \lambda (\text{KL}(\theta \| \theta + d) - c)$$

$$\nabla_{d, \lambda} \mathcal{L} = 0$$

Problematic terms

Cannot calculate easily

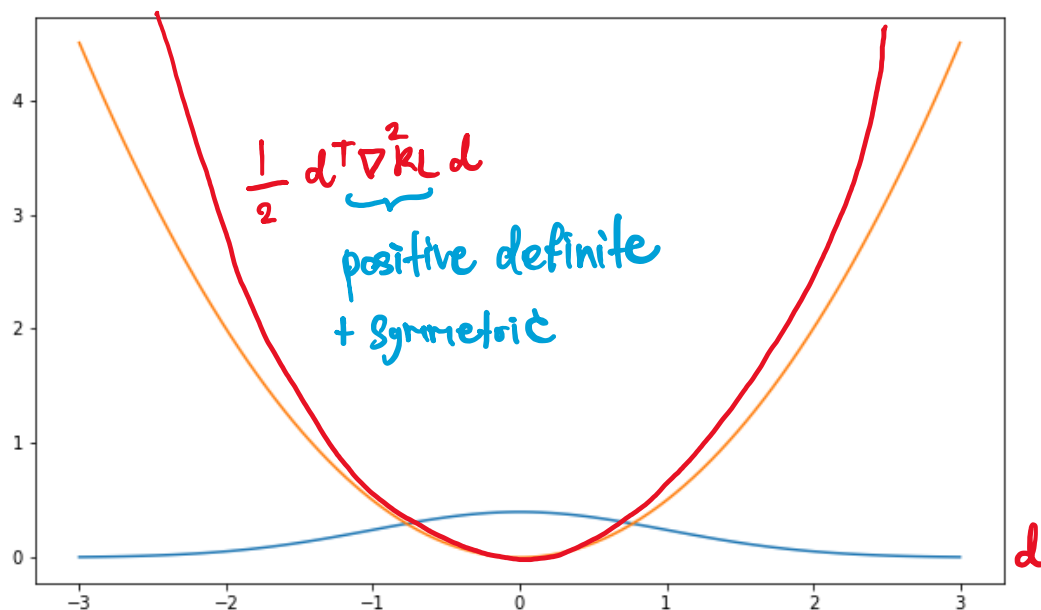
$$\nabla_d J(\theta + d)$$

$$\nabla_d \text{KL}(\theta \| \theta + d)$$

Problematic term

Taylor's to the second order:

$$\begin{aligned}\text{KL}(\theta \parallel \theta + d) &\approx \overline{\text{KL}}(\theta \parallel \theta + d) \\ &= \text{KL}(\theta \parallel \theta) + \nabla_{\theta'} \text{KL}(\theta \parallel \theta')|_{\theta'=\theta} d + \frac{1}{2} d^T \nabla_{\theta'}^2 \text{KL}(\theta \parallel \theta')|_{\theta'=\theta} d\end{aligned}$$



Solving for the update direction

New update direction

$$d = \frac{s}{\lambda}$$


What does it really mean



$$s = (\nabla_{\theta}^2 \text{KL})^{-1} \nabla_{\theta} J(\theta)$$

$$\lambda = \sqrt{\frac{s^T \nabla^2 \text{KL} s}{2c}}$$

previously
 $\nabla_{\theta} L(b)$



$$\theta \leftarrow \theta + d$$

Inverse of KL?

- Naïve inverse is not possible to calculate online

$$\nabla_{\theta}^2 \text{KL}^{-1} = \left(\mathbb{E}_s \left[\nabla_{\theta}^2 \text{KL}_s \right] \right)^{-1}$$

- We use:

$$\nabla_{\theta'}^2 \text{KL}(\pi_{\theta}(\cdot|s) \parallel \pi_{\theta'}(\cdot|s)) \Big|_{\theta'=\theta}$$

$$\nabla_{\theta}^2 \text{KL}^{-1} \approx \mathbb{E}_s \left[\left(\nabla_{\theta}^2 \text{KL}_s \right)^{-1} \right]$$

Finally, Limited trust PG is

$$\theta \leftarrow \theta + d \quad d = \frac{s}{\lambda} \quad \lambda = \sqrt{\frac{s^T \nabla^2 \text{KL} s}{2c}} \leftarrow \text{Trust region}$$

$$s = (\nabla_{\theta}^2 \text{KL})^{-1} \nabla_{\theta} J(\theta)$$

$$s \approx \mathbb{E}_s \left[(\nabla_{\theta}^2 \text{KL}_s)^{-1} \mathbb{E}_a [Q^{\pi}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)] \right]$$

Calculation of inverse is expensive

There are tricks to improve this even further

Bad critic

- Critic is not oracle, it has its flaws
- Critic is prone to forgetting
- How to know when to trust?
- **How to limit the trust?**

Bad critic

- Critic is not oracle, it has its flaws
 - Critic is prone to forgetting
 - How to know when to trust?
 - How to limit the trust?
-
- **Still doesn't guarantee policy improvement**

Policy improvement guarantee

Motivation

- We now have better update
- But, still need to know how large or small the “C” (trust parameter) is
- C could be “varying”
- Too large C could degrade policy
- Too small C is too conservative
- **We want to find C that is just right**

Forming the problem

- We want to “guarantee” policy improvement

$$J(\theta) = \mathbb{E}_{s_0 \sim P(s_0)} [V^\pi(s_0)]$$

$$J(\theta') \geq J(\theta)$$

- Objective becomes:

$$\operatorname{argmax}_{\theta'} J(\theta') - J(\theta)$$

- How to estimate?

The problem of estimation

- We want $J(\theta')$
- We need:
 - Create a new policy
 - Evaluate the policy
- Aim:
 - Estimate the new policy from **what we have**

Write it in another form

Write it in another form

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P'(s_t)} \left[\mathbb{E}_{a_t \sim \pi} \left[\frac{\pi'(a_t|s_t)}{\pi(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right] \right]$$

Can we **lower bound** it while using **only what we have**?

$$J(\theta') - J(\theta) \geq \sum_t \mathbb{E}_{s_t \sim P(s_t)} [\mathbb{E}_{a_t \sim \pi} [\dots]]$$

We need to bound: $|P(s_t) - P'(s_t)|$

Bound the state probability

Intuition: $\pi' \approx \pi \rightarrow P' \approx P$

We will show only a “glimpse”

Assume π is deterministic $a_t = \pi(s_t)$

Let:

$$\pi'(a_t \neq \pi(s_t) | s_t) \leq \epsilon \quad \leftarrow \text{closeness}$$

Then:

$$P'(s_t) = \underbrace{(1 - \epsilon)^t P(s_t)}_{\text{taking action like } \pi} + \underbrace{(1 - (1 - \epsilon)^t) P_{\text{oth}}(s_t)}_{\text{taking at least 1 wrong action}}$$

Bound the probability

$$|P(s_t) - P'(s_t)|$$

General policy case

- It is also possible to show (not here):

$$|\pi'(a|s) - \pi(a|s)| \leq \epsilon$$

$$|P(s_t) - P'(s_t)| \leq \epsilon t$$

- In optimization, we don't have π'
- We need to estimate it, if so:
- Using KL instead would ease estimation

Using KL

- It is possible to show (not here):

$$|\pi'(a|s) - \pi(a|s)| \leq \sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \parallel \pi')} = \epsilon'$$

$$|P(s_t) - P'(s_t)| \leq \epsilon' t$$

Bound some function

$$|P(s_t) - P'(s_t)| \leq \epsilon t$$

$$\mathbb{E}_{s_t \sim P'(s_t)} [f(s_t)]$$

Returning to our objective

$$\mathcal{A}_\pi(\pi') = \mathbb{E}_{a_t \sim \pi} \left[\frac{\pi'(a_t|s_t)}{\pi(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right]$$

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P'(s_t)} [\mathcal{A}_\pi(\pi')]$$

Lower bound:

make ϵ small, ignore this!

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P(s_t)} [\mathcal{A}_\pi(\pi')] - \sum_t \epsilon t \mathcal{O}\left(\frac{r_{\max}}{1-\gamma}\right)$$

$$\text{s.t. } \sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \parallel \pi')} = \epsilon$$

New objective

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P(s_t)} [\mathcal{A}_\pi(\pi')] - \sum_t \epsilon t \mathcal{O}\left(\frac{r_{\max}}{1 - \gamma}\right)$$

s.t. $\sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \parallel \pi')} = \epsilon$

- Bound is very loose
- We should interpret as:
Keep KL small, policy improves!

Approaching optimization

$$J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P(s_t)} [\mathcal{A}_\pi(\pi')] - \sum_t \epsilon t \mathcal{O}\left(\frac{r_{\max}}{1 - \gamma}\right)$$
$$\text{s.t. } \sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \parallel \pi')} = \epsilon$$

Becomes:

$$\operatorname{argmax}_{\theta'} \mathcal{J}(\theta') \quad \text{s.t. } \text{KL}^{\max}(\pi \parallel \pi_{\theta'}) = \epsilon$$

Approaching optimization

$$\operatorname{argmax}_{\theta'} \mathcal{J}(\theta') \quad \text{s.t.} \quad \text{KL}^{\max}(\pi \parallel \pi_{\theta'}) = \epsilon$$

$\text{KL}^{\max}(\pi \parallel \pi_{\theta'})$ is impractical to estimate

We need to go for all S to get the max

Approximate as:

$$\operatorname{argmax}_{\theta'} \mathcal{J}(\theta') \quad \text{s.t.} \quad \mathbb{E}_s [\text{KL}_s(\pi \parallel \pi_{\theta'})] = \epsilon$$

Seem familiar?

Policy improvement guarantee:

$$\operatorname{argmax}_{\theta'} \mathcal{J}(\theta') \quad \text{s.t.} \quad \mathbb{E}_s [\text{KL}_s(\pi \| \pi_{\theta'})] = \epsilon$$

 How to calculate $\nabla_{\theta'} \mathcal{J}(\theta')$

Limited trust PG:

$$\operatorname{argmax}_d J(\theta + d) \quad \text{s.t.} \quad \text{KL}(\theta \| \theta + d) = c$$

Estimating the gradient

$$\mathcal{J}(\theta') = \sum_t \mathbb{E}_{s_t, a_t \sim \pi} \left[\frac{\pi_{\theta'}(a_t | s_t)}{\pi(a_t | s_t)} \gamma^t A^\pi(s_t, a_t) \right]$$

Same old Taylor trick:

$$\nabla_{\theta'} \mathcal{J}(\theta') \approx \nabla_{\theta'} \mathcal{J}(\theta')|_{\theta'=\theta}$$

Policy improvement guarantee

$$s \approx \mathbb{E}_s \left[\left(\nabla_{\theta}^2 \text{KL}_s \right)^{-1} \mathbb{E}_a \left[Q^{\pi}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s) \right] \right]$$

- Turns out to be the same as “Limited trust PG”
- Limit trust \Rightarrow Policy improvement with high chance
- We still cannot “guarantee” we don’t know epsilon
- There is another interpretation of $\nabla_{\theta}^2 \text{KL}_s$

Natural gradients

Motivation

- Gradient descent is not always “steepest”
- Under a more realistic assumption, a better gradient could be derived
- **Leading to faster convergence**
- But higher computation

What is steepest descent?

- Given a fixed budget, what is the update that “reduces” the loss the most
- Euclidean space:

$$\operatorname{argmin}_d L(\theta + d) \quad \text{s.t.} \quad \|d\| = c$$

- If the update “d” is of norm “c”, what should be its direction?

Gradient descent is steepest in Euclidean

- d is steepest if it reduces the loss fastest

$$d^* = \operatorname{argmin}_d L(\theta + d) \quad \text{s.t.} \quad \|d\| = c$$

- Lagrangian

$$\mathcal{L}(d, \lambda) = L(\theta + d) + \lambda(d^T d - c^2)$$

- Solve for critical point:

$$\nabla_d \mathcal{L} = \nabla_d L(\theta + d) + \lambda d = 0$$

$$\nabla_\lambda \mathcal{L} = d^T d - c^2 = 0$$

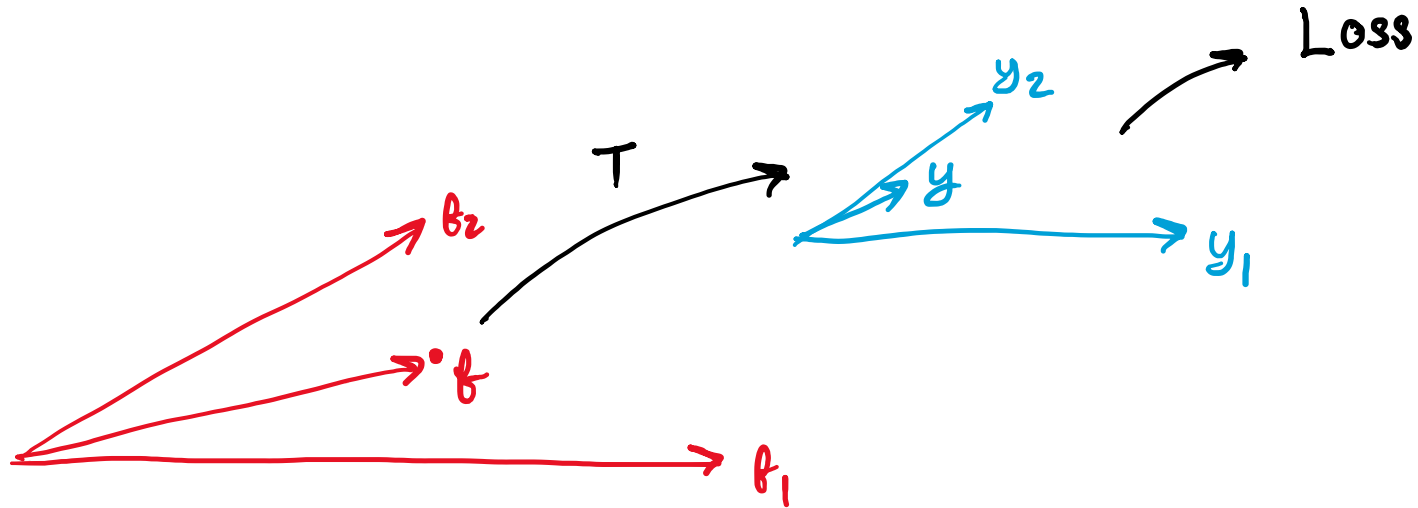
Gradient descent is steepest in Euclidean

Solve for critical point:

$$\nabla_d \mathcal{L} = \nabla_d L(\theta + d) + \lambda d = 0$$

$$\nabla_\lambda \mathcal{L} = d^T d - c^2 = 0$$

If we have another space above



- What is now the steepest descent wrt. the space above?

Norm of the new space

$$\|b\| = \langle b, b \rangle = b^T b$$

$$T: b \rightarrow y \quad y = T b \quad \text{Linear map}$$

$$\|y\| = \langle y, y \rangle = y^T y$$

$$= (Tb)^T (Tb)$$

$$= b^T T^T T b \neq b^T b$$

Norm of the new space

Steepest is “subjective” because “norm” is subjective

$$d^* = \underset{d}{\operatorname{argmin}} L(\theta + d) \quad \text{s.t.} \quad \|d\| = c$$



different meaning

$$d \neq -\frac{\nabla L}{\lambda}$$

Is there a more “natural” space than parameter space?

- Policy is a probability function
- It is more natural to think in “**space of probability functions**”
- What is the steepest descent in the probability function space?

Steepest in prob. fn. space

- Define the “distance” in the function space
- KL Divergence comes into mind:

$$\text{KL}(P\|Q) = \sum_x P(x) [\log P(x) - \log Q(x)]$$

- Steepest descent can get from solving:

$$d^* = \underset{d}{\operatorname{argmax}} J(\theta + d) \quad \text{s.t.} \quad \text{KL}(\theta\|\theta + d) = c$$

Connection

- Limited trust policy gradient
- Policy improvement guarantee
- Steepest descent on probability function space
(Natural gradient)

They are doing the same thing

Related works

- Natural policy gradient
- Trust region policy optimization (TRPO)
 - We present here a “mini” version
- Proximal policy optimization (PPO)
 - An approximation of TRPO
 - Works well and easy to implement

More on policy gradient

Action dependent baseline

- **PG** has high variance because it uses “indirect gradient”

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

- **DPG** has lower variance because it can “backprop”

$$\nabla J(\theta) = \sum_s d^\pi(s) \nabla_a Q_\phi(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$

Action dependent baseline

- **Can we combine the two?**
- **Q-Prop**

$$\begin{aligned}\nabla J(\theta) = & \mathbb{E}_{s,a} [(Q^\pi(s, a) - \bar{Q}_\theta(s, a)) \nabla_\theta \log \pi_\theta(a|s)] \\ & + \mathbb{E}_s [\nabla_a Q_\phi(s, a)|_{a=u_\theta(s)} \nabla_\theta u_\theta(s)]\end{aligned}$$

$$u_\theta(s) = \sum_a \pi_\theta(a|s) a$$

- **Taylor expansion (first order)**

$$\bar{Q}_\phi(s, a) = Q_\phi(s, u_\theta(s)) + \nabla_a Q_\phi(s, a)|_{a=u_\theta(s)} (a - u_\theta(s))$$

Policy gradient from minimizing KL

- If we look at Q as “unnormalized” policy
 - A little bit sharper of Q is $\exp(Q)$
 - This is our target policy
- We could use a KL:

$$\pi = \operatorname{argmin}_{\pi \in \Pi} D_{KL} \left(\pi(\cdot|s) \left\| \frac{\exp(Q(\cdot, s))}{Z} \right. \right)$$

- Minimizing KL is an optimization task

Policy gradient from minimizing KL

- KL policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_s \left[\nabla_{\theta} D_{KL} \left(\pi(\cdot|s) \parallel \frac{\exp(Q(\cdot, s))}{Z} \right) \right]$$

- Z is a constant, ignored
- Policy improve to Q
- Policy eval: Q gets even sharper
- Repeat