

→ TRPO

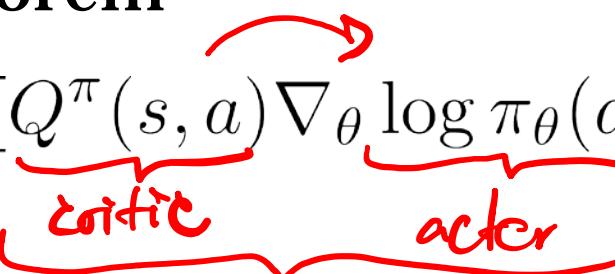
# Trust Region methods

Konpat Preechakul  
Chulalongkorn University  
November 2019

# Recap policy gradient

- Policy gradient theorem

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

✓ Variance problem 

- Advantage

$$\underline{A}(s, a) = \underline{Q}(s, a) - \underline{V}(s)$$

- Variance reduction
- A2C

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

$$\uparrow J(\theta) = \sum p(s_0) \pi^\pi(s_0)$$

# Recap policy gradient

- Off-policy gradient
  - Off-policy critic ✓  $Q \leftarrow$ ,  $A \leftarrow$
  - Off-policy actor ✓

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

*some*

↑  
TDH

The handwritten annotations include a checkmark next to 'Off-policy critic' and 'Off-policy actor'. A red checkmark is placed next to the equation. Brackets under 'd^b(s)' and 'Q^{\pi}(s, a)' are labeled 's' and 'a' respectively. A bracket under the entire term 'd^b(s) \sum\_a Q^{\pi}(s, a)' is labeled 'some'. To the right of the equation, there is a handwritten 'TDH' with an upward arrow and a red bracket labeled 'C' with a downward arrow.

- Deterministic policy gradient

- DDPG

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

The handwritten annotations include a circled 'a' under 'd^{\pi}(s)', a circled 'a' under 'Q\_{\phi}(s, a)|\_{a=\pi(s)}', and a circled 's' under 'nabla\_{\theta} \pi\_{\theta}(s)'. There is also a circled 'pi' with an upward arrow.

# **Today's topic**

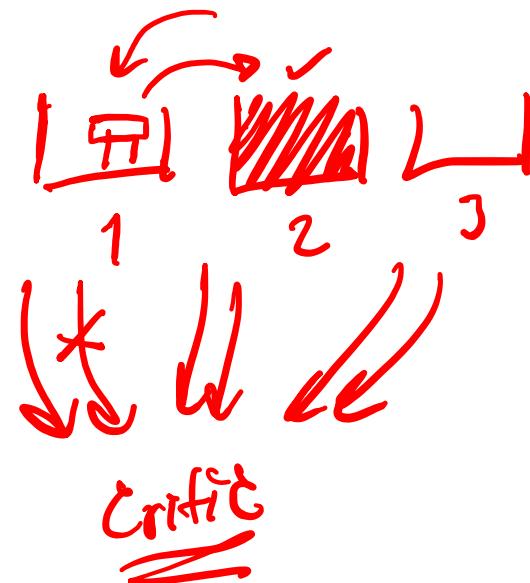
# Today topics

- ✓ Why policy gradient fails?
- More robust policy gradient
  - Trust region methods
  - 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

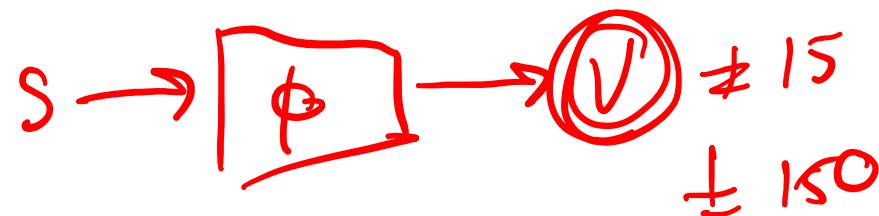
# Bad critic

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



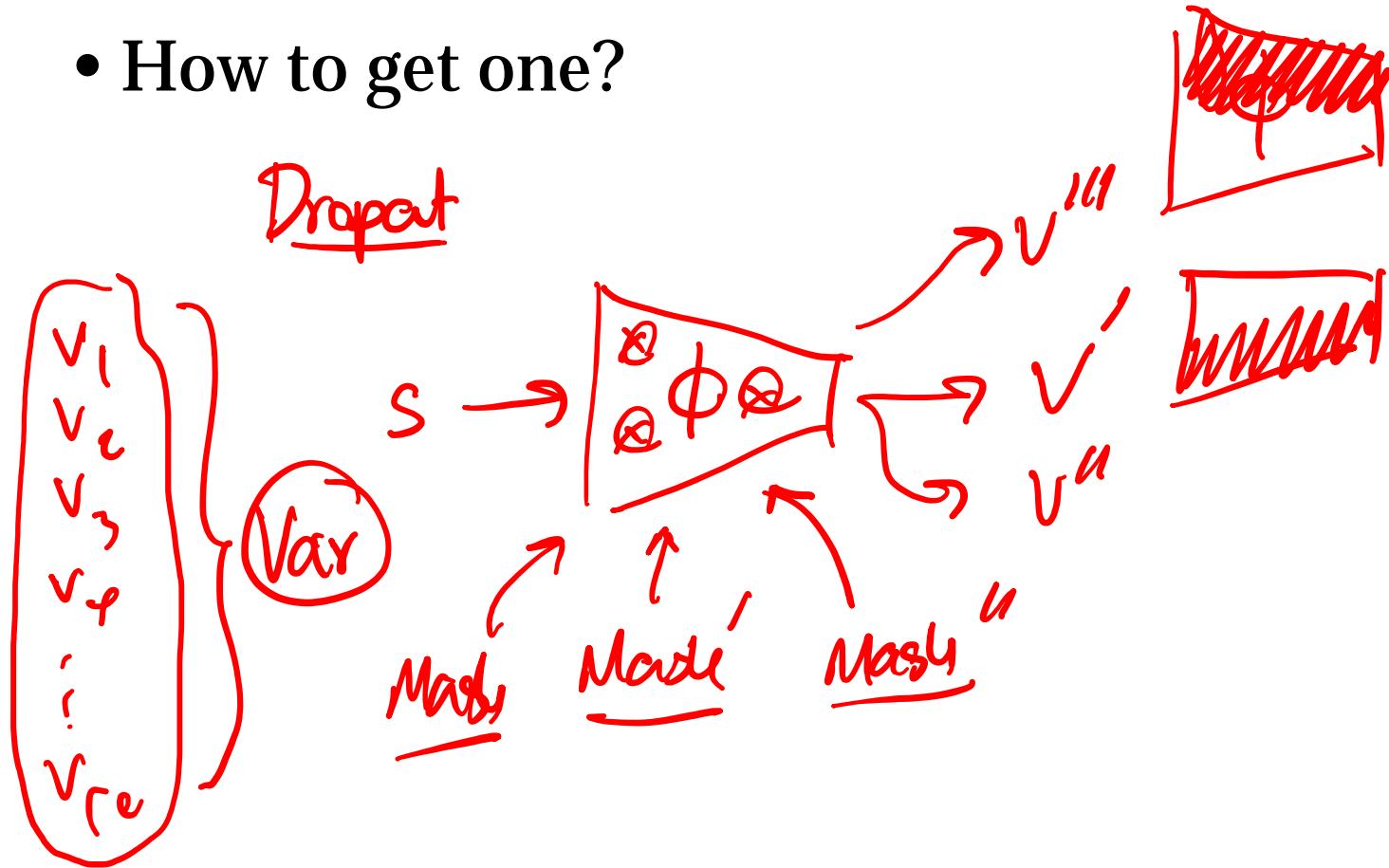
# When to trust critic?

- Do critic have “confidence”?



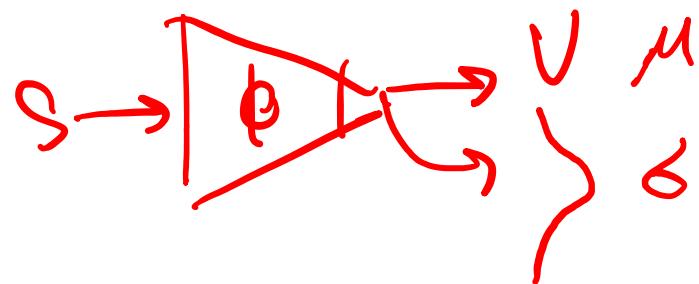
# When to trust critic?

- 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?

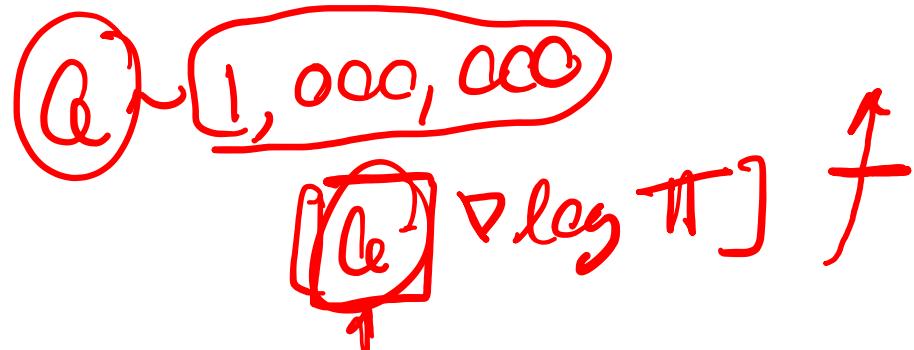
# When to trust critic?

- 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”

$$\pi \xrightarrow{\text{C}} \pi'$$

★ Policies are probability functions

- We find a “distance” measure on the probabilities

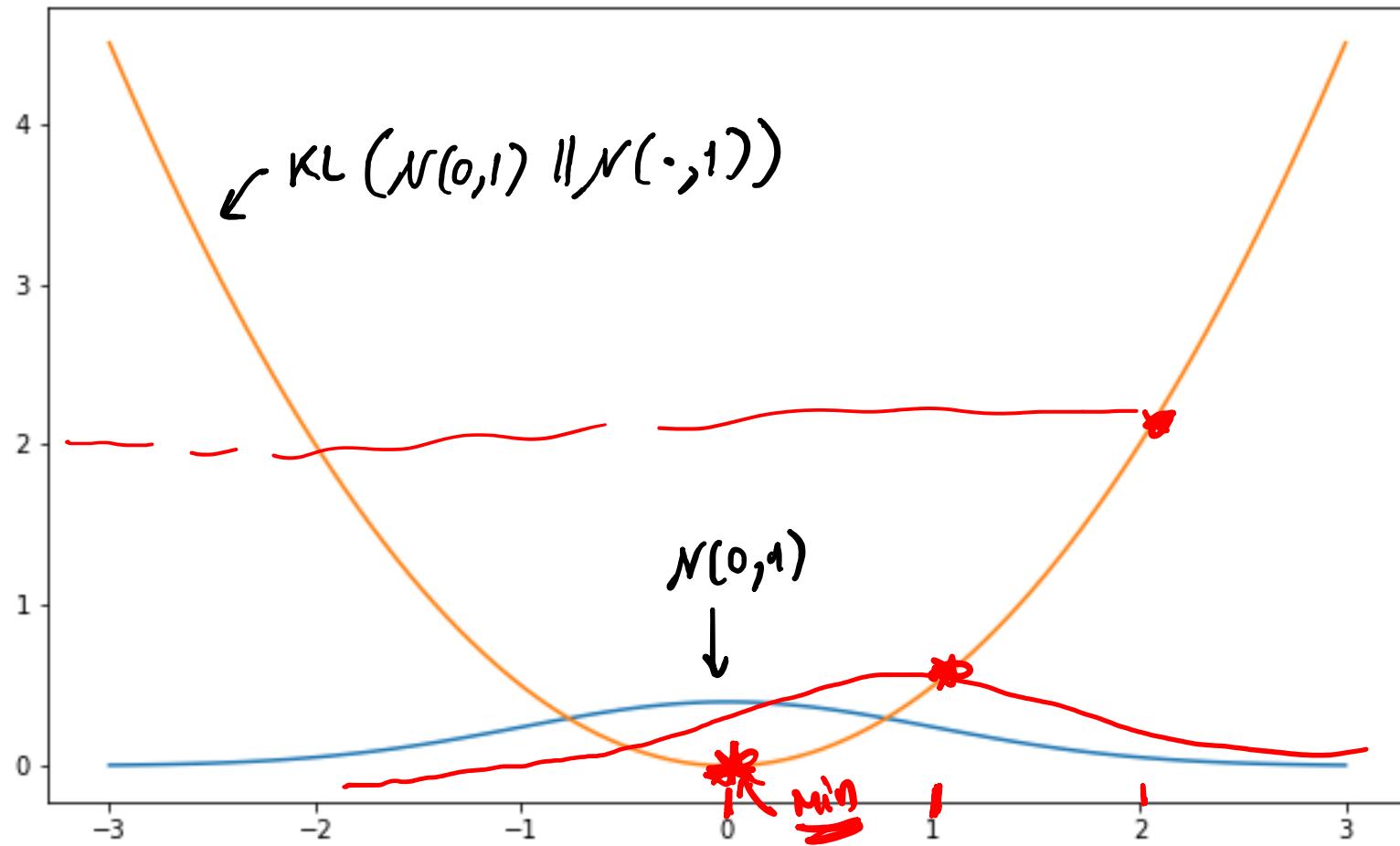
★ KL divergence:  $\Theta \rightsquigarrow \Delta$

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

$$d(x,y) \neq d(y,x)$$

# KL Divergence in picture

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



# Quantify “much” in RL

- Context of RL, for a state  $S$ :

$$\pi \rightarrow \pi'$$
$$KL(\pi \parallel \pi') = c$$

$$KL_s(\theta \parallel \theta') = \sum_a \pi_\theta(a|s) [\log \pi_\theta(a|s) - \log \pi_{\theta'}(a|s)]$$
$$KL(\pi(\cdot|s) \parallel \pi'(\cdot|s))$$

- For all states:

$$\underline{KL}(\theta \parallel \theta') = \sum_s d^\pi(s) KL_s(\theta \parallel \theta')$$

$\uparrow$

$E_{s \sim d^\pi}$

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

$$KL(\pi || \pi') = \zeta$$

# New objective for policy update

\* From the policy gradient objective

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

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

$$E_{s,a} [Q(s,a) \nabla \log \pi(\cdot)]$$

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

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

# New update direction

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

- A constrained optimization

$$f(d) \rightarrow \nabla_d f(d) = 0$$

- We relax it 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

$$\textcircled{1} \quad \mathcal{L}(d, \lambda) = J(\theta + d) + \lambda (\text{KL}(\theta \| \theta + d) - c)$$

$$\textcircled{2} \quad \nabla_{\underline{d}, \lambda} \mathcal{L} = 0$$

$$\nabla_d \mathcal{L}(d, \lambda) = \nabla_d J(\theta + d) + \lambda \nabla_d \text{KL}(\theta \| \theta + d) = 0$$

$$\nabla_d \mathcal{L}(d, \lambda) = \cancel{\text{KL}(\theta \| \theta + d)} - c = 0$$

$$\boxed{\nabla_\theta J(\theta)}$$

# Problematic terms

Cannot calculate easily

①  $\nabla_d J(\theta + d)$   $\nabla J_\theta$  ✓

$$J \approx \bar{J}(\theta + d) = J(\theta) + \nabla_\theta J(\theta) \cdot d$$

$$\nabla_d \bar{J}(\theta + d) = \nabla_\theta J(\theta)$$

②  $\nabla_d \text{KL}(\theta || \theta + d)$

$$\begin{aligned} &\approx \\ &\text{KL}(\theta || \theta + d) = \cancel{\text{KL}(\theta || \theta)} = 0 \end{aligned}$$

$$+ \cancel{\nabla_\theta \text{KL}(\theta || \theta) d} = 0$$

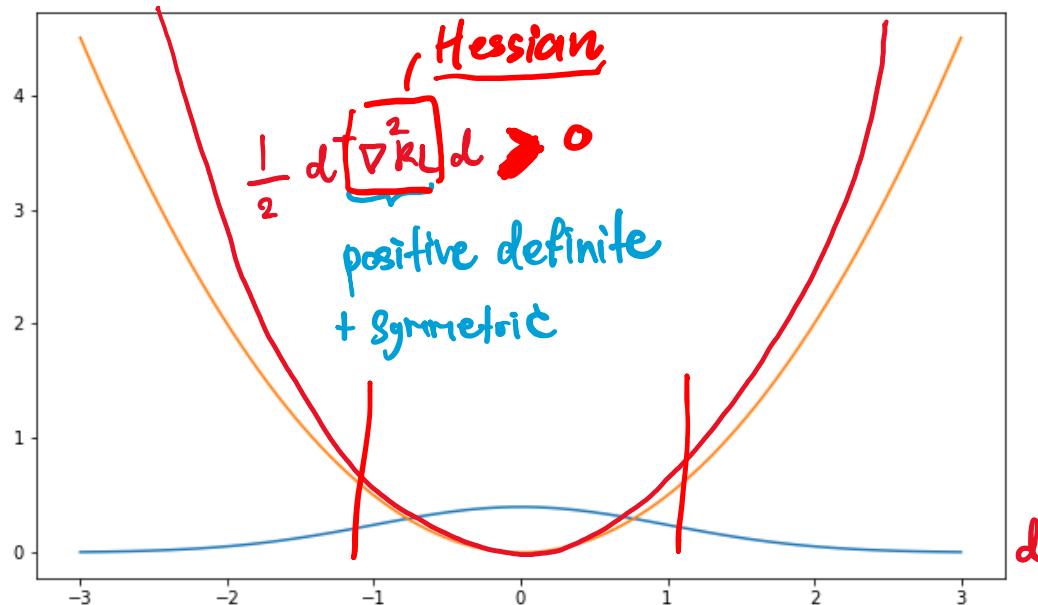
$$+ \frac{1}{2} d^\top \nabla_\theta^2 \text{KL} d$$

# Problematic term

Taylor's to the second order:

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

$$\rightarrow \underbrace{\text{KL}(\theta\|\theta)}_{\text{Hessian}} + \underbrace{\nabla_{\theta'} \text{KL}(\theta\|\theta')|_{\theta'=\theta} d}_{\text{positive definite}} + \frac{1}{2} \underbrace{d^T \nabla_{\theta'}^2 \text{KL}(\theta\|\theta')|_{\theta'=\theta} d}_{\text{+ Symmetric}}$$



$$\pi \xrightarrow{c} \pi'$$

# Solving for the update direction

$$0 = \nabla_d J(\theta + d) + \lambda \nabla_d KL(\theta \| \theta + d) \quad \left( \frac{\cancel{d^T \nabla^2_{\theta\theta} d}}{2} \right)$$

$$0 \approx \nabla_\theta J(\theta) + \lambda \nabla_\theta^2 KL(\theta \| \theta) \cdot d \quad \left( \frac{\cancel{\lambda \nabla^2_{\theta\theta} d^2}}{2} \right)$$

$$-\nabla_\theta J(\theta) = \lambda \nabla_\theta^2 KL \cdot d$$

$$-\underbrace{(\nabla_\theta^2 KL)^{-1} \nabla_\theta J}_{S} = d$$

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

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

$$\underline{KL(\theta \| \theta + d)} - c = 0$$

$$\frac{1}{2} d^T \nabla^2_{\theta\theta} d - c = 0$$

$$\frac{1}{2} \left( \frac{S}{\lambda} \right)^T \nabla^2_{\theta\theta} \left( \frac{S}{\lambda} \right) = \frac{2c}{\lambda} =$$

$$\frac{s^T \nabla^2_{\theta\theta} s}{2c} = \lambda$$

$$\sqrt{\frac{s^T \nabla^2_{\theta\theta} s}{2c}} = \lambda$$

# 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

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

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

# Inverse of KL?

- Naïve inverse is not possible to calculate online

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

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

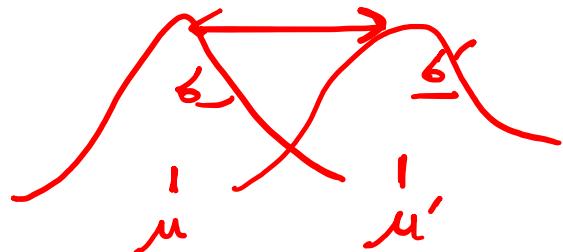
- We use:

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

$$\nabla^2 L(\theta) \leftarrow$$

# Finally, Limited trust PG is

$$\theta \leftarrow \theta + d \quad d = \frac{s}{\lambda}$$



$$\lambda = \sqrt{\frac{s^T \nabla^2 \text{KL} s}{2C}} \xrightarrow{\text{Trust region}}$$

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

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

**PG**

\* Calculation of inverse is expensive

There are tricks to improve this even further

TRPO

# 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

$$\pi \xrightarrow{C} \pi'$$

→ 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

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

$$\underline{J(\theta')} \geq \underline{J(\theta)}$$

- Objective becomes:

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

- How to estimate?  $J(\theta')$

# The problem of estimation

- We want  $J(\theta')$
  - We need:
    - ① Create a new policy
    - ② Evaluate the policy
  - Aim:
    - Estimate the new policy from **what we have**
- $\checkmark a, s \sim \pi$   
 $\times a, s \sim \pi$
-

# Write it in another form

$$J(\theta') - J(\theta) = \underbrace{E_{s_0} [V^{\pi'}(s_0)]}_{\text{---}} - \underbrace{E_{s_0} [V^{\pi}(s_0)]}_{\text{---}}$$

$$\textcircled{1} E_{s_0} [V^{\pi}(s_0)] = E_{s_0} \left[ E_{T \sim \pi} \left[ \sum_{t=0}^{\infty} r^t v(s_t) - \sum_{t=1}^{\infty} \gamma^t v(s_t) \right] \right]$$

$$\downarrow E_{T \sim \pi'} \left[ \sum_{t=0}^{\infty} r^t v(s_t) - \gamma^{t+1} v(s_{t+1}) \right]$$

$$\textcircled{2} E_{s_0} [V^{\pi'}(s_0)] = E_{s_0} \left[ E_{T \sim \pi'} \left[ \sum_{t=0}^{\infty} r^t r_t \right] \right]$$

$$a \approx r + v' \quad q - v = r + v' - v$$

$$v = v \quad A =$$

$$(J(\theta') - J(\theta)) = E_{T \sim \pi'} \left[ \sum_{t=0}^{\infty} r^t r_t - \gamma^t v(s_t) + \gamma^{t+1} v(s_{t+1}) \right]$$

$$- E_{T \sim \pi} \left[ \sum_{t=0}^{\infty} r^t A^{\pi}(s_t, a_t) \right]$$

$\rightarrow \pi$

$$J(\pi') - J(\pi) = E_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} r^t A^\pi(s_t, a_t) \right] \quad s_0 \ s_1 \ s_2 \dots s_\infty$$

$$= \sum_{t=0}^{\infty} E_{\substack{s_t \sim \pi' \\ \uparrow}} \left[ E_{a_t \sim \pi'} \left[ r^t A^\pi(s_t, a_t) \right] \right]$$

$$= \sum_{t=0}^{\infty} \underbrace{E_{\substack{s_t \sim \pi' \\ \uparrow}} \left[ E_{a_t \sim \pi'} \left[ \underbrace{\frac{\pi'(a_t | s_t)}{\pi(a_t | s_t)} \cdot r^t A^\pi(s_t, a_t)}_{\text{II}} \right] \right]}_{\text{I}}$$

# 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:  $|P(s_t) - P'(s_t)| \leq \boxed{\phantom{00}}$

# Bound the state probability

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

We will show only a “glimpse”

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

Let:

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

Then:  $\tau: s_0 \xrightarrow{\pi'} s_1 \xrightarrow{\pi'} s_2 \dots \xrightarrow{\pi'} s_\infty$

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

# Bound the probability

$$\begin{aligned}
 & |P(s_t) - P'(s_t)| \\
 &= \left| p(s_t) - (1-\varepsilon)^t p(s_t) - (1-(1-\varepsilon)^t) p_{\text{oth}}(s_t) \right| \\
 &= (1-(1-\varepsilon)^t) \left| p(s_t) - p_{\text{oth}}(s_t) \right| \\
 &\leq (1-(1-\varepsilon)^t) 1 \\
 &\leq \underline{\varepsilon t}
 \end{aligned}$$

$\varepsilon = 0.1$   
 $(1-\varepsilon)^t = 0.9^t$   
 $t$

$$\pi(s) = a \longrightarrow \pi(a|s)$$

# General policy case

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

$$\rightarrow \textcircled{A} \underbrace{|\pi'(a|s) - \pi(a|s)|}_{\text{KL}} \leq \textcircled{\epsilon}$$

$$\rightarrow \checkmark \underbrace{|P(s_t) - P'(s_t)|}_{\text{KL}} \leq \textcircled{\epsilon t}$$

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

# Using KL

- It is possible to show (not here; Pollard, 2000):

$$\underbrace{|\pi'(a|s) - \pi(a|s)|}_{\text{red underline}} \leq \sqrt{\frac{1}{2} \text{KL}_{\overbrace{s}^{\text{red arrow}}}^{\max}(\pi \| \pi')} = \epsilon'$$

$$\underbrace{|P(s_t) - P'(s_t)|}_{\text{red underline}} \leq \epsilon' t$$

# Bound some function

$$\begin{aligned} & \mathbb{E}_{\underbrace{s_t \sim P'(s_t)}_{\text{1}}} [f(s_t)] \geq \mathbb{E}_{\underbrace{s \sim \pi}_{\text{2}}} [f(s)] \\ & = \sum_{s_t} p(s_t) [f(s_t)] \\ & \geq \sum_{s_t} p(s_t) f(s_t) - |P(s_t) - \underline{P'(s_t)}| \cdot \max_s f(s) \\ & \geq \sum_{s_t} \underline{P(s_t)} f(s_t) - \underline{\epsilon t \cdot \max_s f(s)} \end{aligned}$$

$$\frac{r_{\max} + \gamma r_{\max} + \gamma^2 r_{\max} + \dots + \gamma^t r_{\max}}{1-\gamma} = \frac{r_{\max}}{1-\gamma}$$

# Returning to our objective

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

①  $J(\theta') - J(\theta) = \sum_t \mathbb{E}_{s_t \sim P'(\cdot|s_t)} [\mathcal{A}_\pi(\pi')]$

Lower bound:

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

make  $\epsilon$  small, ignore this!

small  $\epsilon \rightarrow$  policy improvement s.t.  $\sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \| \pi')} = \epsilon$

# 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.  $\sqrt{\frac{1}{2} \text{KL}^{\max}(\pi \| \pi')} = \epsilon$

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

# Approaching optimization

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

$\Downarrow$

$$\mathcal{J}(\theta')$$

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

Becomes:

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

# Approaching optimization

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

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

We need to go for all S to get the max

Approximate as:

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

↑ *Limited trust PG*

# Seem familiar?

① Policy improvement guarantee:

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

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

② Limited trust PG:

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

# Estimating the gradient

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

$\nabla_{\theta'} \log \pi$

Same old Taylor trick:

$$\begin{aligned} \nabla_{\theta'} \mathcal{J}(\theta') &\approx \nabla_{\theta'} \mathcal{J}(\theta')|_{\theta'=\theta} \\ \nabla_{\theta'} \mathcal{J}(\theta') &\approx \sum_t \mathbb{E}_\pi \left[ \frac{\nabla_{\theta'} \pi_\theta}{\pi_\theta} \gamma^t A^\pi \right] \Big|_{\theta'=\theta} \\ &= \sum_t \mathbb{E}_\pi \left[ \frac{\nabla_{\theta'} \pi_\theta}{\pi_\theta} \gamma^t A^\pi \right] = \mathbb{E}_\pi \left[ \frac{\gamma^t A^\pi}{\pi_\theta} \nabla_{\theta'} \log \pi_\theta \right] \end{aligned}$$

$\nabla \log \pi$   
 $\frac{F_{a,s|\pi}}{Q^\pi}$

① Limited trust PG

## ② Policy improvement guarantee

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

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

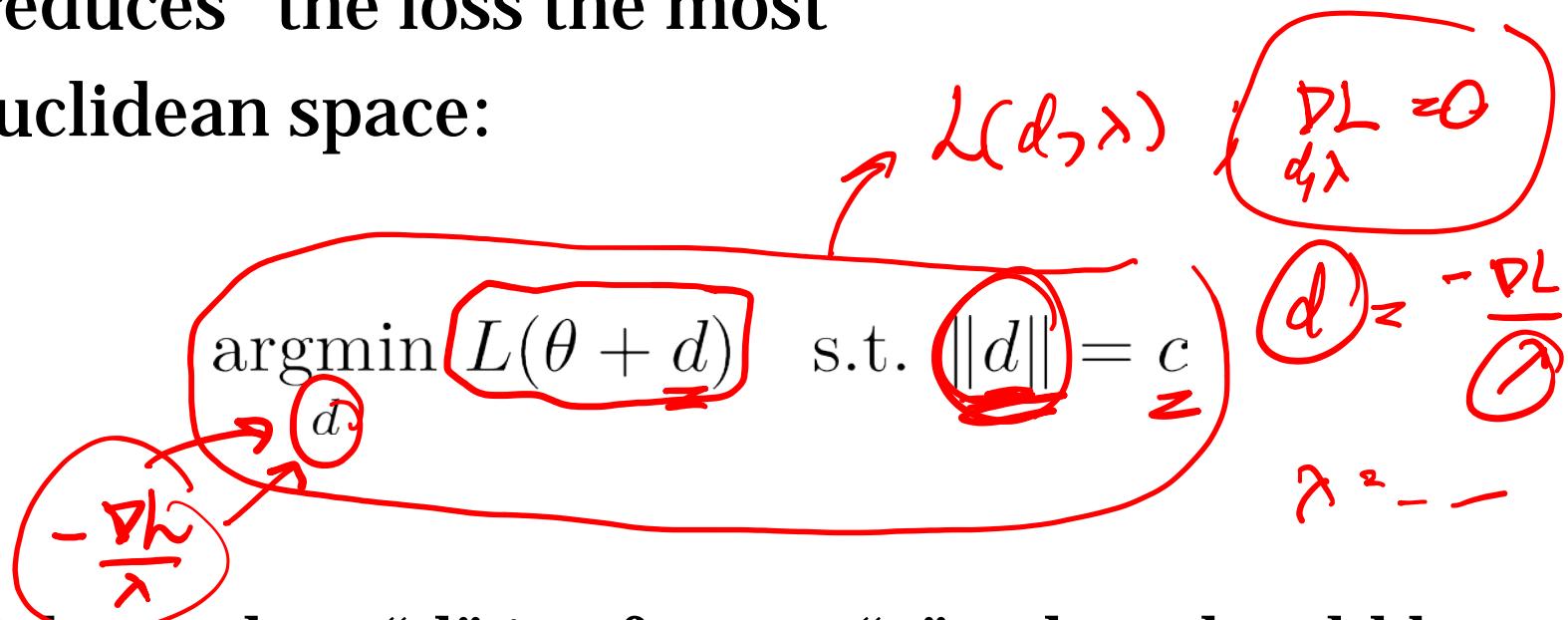
# **Natural gradients**

# 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:



- 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

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

- Lagrangian

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

- Solve for critical point:

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

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

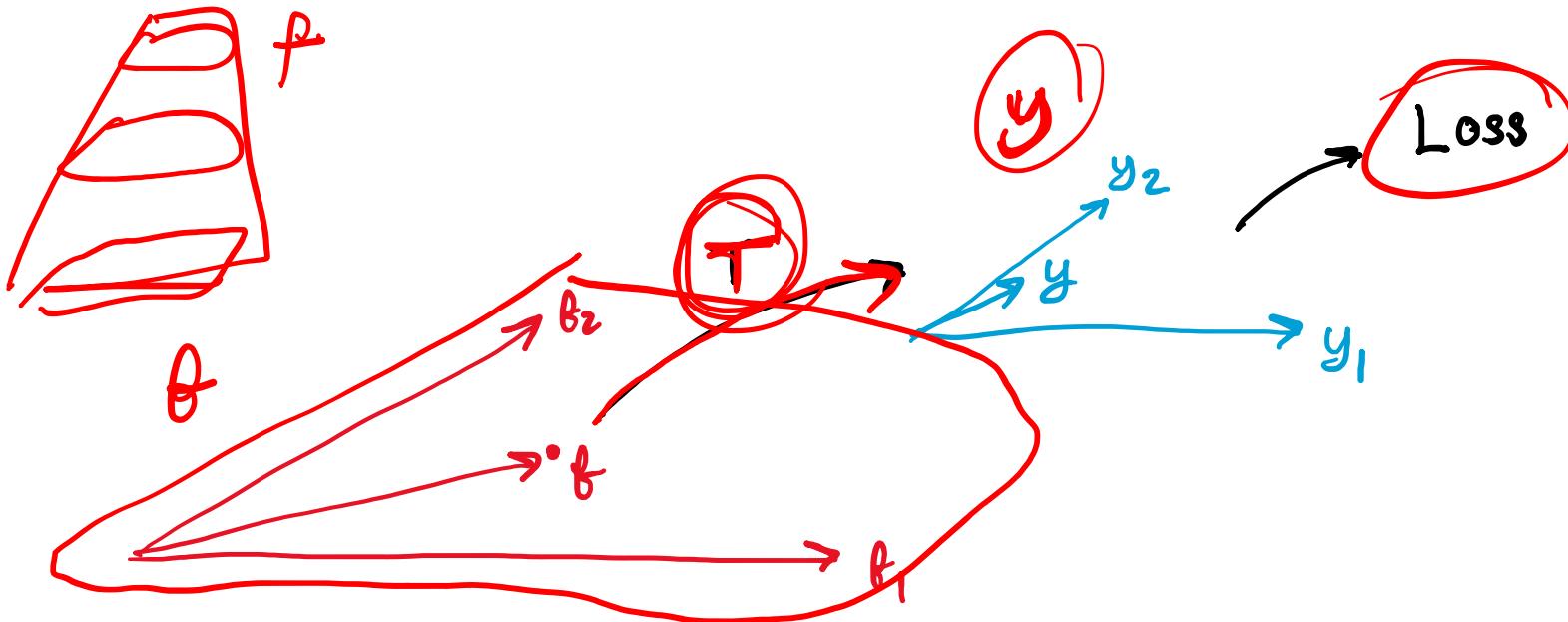
# Gradient descent is steepest in Euclidean

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?

$$\underset{d}{\operatorname{argmin}} \ L(\underline{f} + \underline{d}) \quad \text{s.t. } \underline{\|\underline{d}\|_2 = c}$$

# Norm of the new space

$$\|\underline{f}\| = \langle \underline{f}, \underline{f} \rangle = \underline{f}^T \underline{f}$$



$$T: \underline{f} \rightarrow \underline{y} \quad \boxed{\underline{y} = T \underline{f}} \quad \text{Linear map}$$

$$\begin{aligned} \|\underline{y}\| &= \langle \underline{y}, \underline{y} \rangle = \underline{y}^T \underline{y} \\ &= (\underline{T\underline{f}})^T (\underline{T\underline{f}}) \\ &= \underline{\underline{\underline{f}^T T^T T f}} \neq \underline{f^T f} \\ \|\underline{f}\|_y &\neq \|\underline{f}\|_f \end{aligned}$$

# Norm of the new space

Steepest is “subjective” because “norm” is subjective

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



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?

$$\left( \nabla_{\theta}^2 \mathcal{L} \right)^{-1}$$

# 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:

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

# Connection

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

**They are doing the same thing**

# Related works

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

$$[\nabla_{KL}]^\top$$

(2017)

# **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

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

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

- 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:

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

- Minimizing KL is an optimization task

# Policy gradient from minimizing KL

- KL policy gradient:

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

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