### Lecture 7: Advanced Policy Gradients

#### Zhi Wang & Chunlin Chen

Department of Control Science and Intelligence Engineering Nanjing University

### Natural policy gradient

- Policy gradient in distribution space
- Solve the constrained optimization problem
- Natural gradient
- 2 Trust region policy optimization (TRPO)
- 3 Proximal policy optimization (PPO)

### Review: Vanilla policy gradient (REINFORCE)

# REINFORCE algorithm: Loop: 1. sample $\{\tau^i\}$ from $\pi_{\theta}(a_t|s_t)$ (run the policy) 2. $\nabla_{\theta}J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t^i|s_t^i) \sum_{t'=t}^{T} \gamma^{t'-t} r(s_t^i, a_t^i)$ 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta}J(\theta)$



# Problems of vanilla policy gradient (REINFORCE)

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(a_{t}^{i} | s_{t}^{i}) Q^{\pi}(s_{t}^{i}, a_{t}^{i})$$
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

- Hard to select the step size  $\alpha$ 
  - Too big step: Bad policy → data collected under bad policy → we cannot recover (in Supervised Learning, data does not depend on neural network weights)
  - Too small step: Not efficient use of experience (in Supervised Learning, data can be trivially re-used)

# Problems of vanilla policy gradient (REINFORCE)



• Small changes in the policy parameters can unexpectedly lead to big **changes** in the policy

• The step size in gradient descent results from solving the following optimization problem, e.g., using line search

$$d^* = \operatorname*{arg\,max}_{||d|| \le \epsilon} J(\theta + d)$$

- Euclidean distance in parameter space
- Stochastic gradient descent (SGD)

$$\theta \leftarrow \theta + d^*$$

### Hard to pick the threshold $\epsilon$

• It is hard to predict the result on the parameterized distribution

• Especially for nonlinear function approximators, e.g., neural networks





• Gradient descent in parameter space

$$d^* = \operatorname*{arg\,max}_{||d|| \le \epsilon} J(\theta + d)$$

• Natural gradient descent: the step size in parameter space is determined by considering the KL divergence in the distributions before and after the update

$$d^* = \underset{d}{\arg\max} J(\theta + d), \quad s.t. \operatorname{D}_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta + d}) \leq \epsilon$$

- KL divergence in distribution space
- Easier to pick the distance threshold!!!

### Distance for probability distributions

• How to calculate the distance between two points in a 2D coordinate?

distance = 
$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

• Euclidean distance



 How to calculate the distance between two probability distributions, p(x) and q(x)? • A measure of how one probability distribution, p(x), is different from a second, reference probability distribution, q(x)

$$D_{\mathrm{KL}}(p(x)||q(x)) = \sum_{i} p(x_i) \log \frac{p(x_i)}{q(x_i)}$$
$$D_{\mathrm{KL}}(p(x)||q(x)) = \int_{x} p(x) \log \frac{p(x)}{q(x)} \,\mathrm{d}x$$

• A KL divergence of 0 indicates that the two distributions are identical



• Suppose two Gaussian distributions:

$$p(x) \sim \mathcal{N}(\mu_1, \sigma_1^2), \quad q(x) \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

• What is  $D_{KL}(p(x)||q(x))$ ?

$$\log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

$$p(x) \sim \mathcal{N}(\mu_1, \sigma_1^2), \quad \mathbb{E}_{p(x)}[x] = \mu_1, \quad var_{p(x)}[x] = \mathbb{E}[(x - \mu_1)^2] = \sigma_1^2$$

$$\begin{aligned} \mathbf{D}_{\mathrm{KL}}(p(x)||q(x)) &= \mathbb{E}_{p(x)}[\log p(x) - \log q(x)] \\ &= \mathbb{E}_{p(x)} \left[ -\log(\sqrt{2\pi\sigma_1}) - \frac{(x-\mu_1)^2}{2\sigma_1^2} + \log(\sqrt{2\pi\sigma_2}) + \frac{(x-\mu_2)^2}{2\sigma_2^2} \right] \\ &= \log \frac{\sigma_2}{\sigma_1} - \frac{\mathbb{E}_{p(x)}[(x-\mu_1)^2]}{2\sigma_1^2} + \frac{\mathbb{E}_{p(x)}[(x-\mu_1+\mu_1-\mu_2)^2]}{2\sigma_2^2} \\ &= \log \frac{\sigma_2}{\sigma_1} - \frac{\sigma_1^2}{2\sigma_1^2} + \frac{\mathbb{E}_{p(x)}[(x-\mu_1)^2 + 2(x-\mu_1)(\mu_1-\mu_2) + (\mu_1-\mu_2)^2]}{2\sigma_2^2} \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1-\mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \end{aligned}$$

### Natural policy gradient

- Policy gradient in distribution space
- Solve the constrained optimization problem
- Natural gradient
- 2 Trust region policy optimization (TRPO)
- 3 Proximal policy optimization (PPO)

• How to solve this constrained optimization problem?

$$d^* = \underset{d}{\arg\max} J(\theta + d), \quad s.t. \operatorname{D}_{\operatorname{KL}}(\pi_{\theta} || \pi_{\theta + d}) \le \epsilon$$

- What tool to use?
  - Turn the constrained optimization problem to an unconstrained one?

• How to solve this constrained optimization problem?

$$d^* = \operatorname*{arg\,max}_{d} J(\theta + d), \quad s.t. \, \mathrm{D}_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta + d}) \leq \epsilon$$

 Use the Lagrangian multiplier λ, turn to the unconstrained penalized objective

$$d^* = \operatorname*{arg\,max}_{d} J(\theta + d) - \lambda(\mathrm{D}_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta + d}) - \epsilon)$$

$$d^* = \underset{d}{\arg\max} J(\theta + d) - \lambda(\mathbf{D}_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta + d}) - \epsilon)$$

• First-order Taylor expansion for the loss

 $J(\theta + d) \approx J(\theta) + \nabla_{\theta'} J(\theta')|_{\theta' = \theta} \cdot d$ 

• Second-order Taylor expansion for the KL

$$\mathbf{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta+d}) \approx \frac{1}{2} d^T \cdot \nabla^2_{\theta'} \mathbf{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} \cdot d$$

 A representation of a function as an infinite sum of terms that are calculated from the values of the function's derivatives at a single point

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$
  
=  $f(a) + f'(a)(x-a) + \frac{f''(a)}{2} (x-a)^2 + \dots$ 

Examples

$$e^x = ?$$

$$\frac{1}{1-x} = ?$$

### Taylor series/expansion

 A representation of a function as an infinite sum of terms that are calculated from the values of the function's derivatives at a single point

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$
  
=  $f(a) + f'(a)(x-a) + \frac{f''(a)}{2} (x-a)^2 + \dots$ 

Examples

$$e^{x} = 1 + x + \frac{x^{2}}{2!} + \frac{x^{3}}{3!} + \dots$$
$$\frac{1}{1-x} = 1 + x + x^{2} + x^{3} + \dots$$

• Let  $\theta' = \theta + d$  is the independent variable

• That is, 
$$x = \theta'$$
,  $a = \theta$ ,  $x - a = d$ 

• What is the Taylor expansion of  $J(\theta + d)$ ?

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$
  
=  $f(a) + f'(a)(x-a) + \frac{f''(a)}{2}(x-a)^2 + \dots$ 

- Let  $\theta' = \theta + d$  is the independent variable
- That is,  $x = \theta'$ ,  $a = \theta$ , x a = d
- What is the Taylor expansion of  $J(\theta + d)$ ?
- First-order Taylor expansion for the loss:

 $J(\theta + d) \approx J(\theta) + \nabla_{\theta'} J(\theta')|_{\theta' = \theta} \cdot d$ 

• Let  $\theta' = \theta + d$  is the independent variable

• That is, 
$$x = \theta'$$
,  $a = \theta$ ,  $x - a = d$ 

• What is the Taylor expansion of  $D_{KL}(\pi_{\theta}||\pi_{\theta+d})$ ?

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$
  
=  $f(a) + f'(a)(x-a) + \frac{f''(a)}{2}(x-a)^2 + \dots$ 

- Let  $\theta' = \theta + d$  is the independent variable
- That is,  $x = \theta'$ ,  $a = \theta$ , x a = d
- What is the Taylor expansion of  $D_{KL}(\pi_{\theta}||\pi_{\theta+d})$ ?
- Second-order Taylor expansion for  $D_{KL}(\pi_{\theta}||\pi_{\theta'})$ :

 $D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'}) \approx D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta}) + d^{T} \nabla_{\theta'} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} + \frac{1}{2} d^{T} \nabla_{\theta'}^{2} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} d$ 

### Taylor expansion of KL

 $D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'}) \approx D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta}) + d^{T} \nabla_{\theta'} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} + \frac{1}{2} d^{T} \nabla_{\theta'}^{2} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} d^{T} \nabla_{\theta'}^{2} D_{\mathrm{KL}}(\pi_{\theta}||\pi$ 

$$D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'}) = \int \pi_{\theta}(x) \log \frac{\pi_{\theta}(x)}{\pi_{\theta'}(x)} \, \mathrm{d}x = \underbrace{\int \pi_{\theta}(x) \log \pi_{\theta}(x) \, \mathrm{d}x}_{\text{independent of } \theta'} - \int \pi_{\theta}(x) \log \pi_{\theta'}(x) \, \mathrm{d}x$$

$$\nabla_{\theta'} \operatorname{D}_{\operatorname{KL}}(\pi_{\theta} || \pi_{\theta'}) |_{\theta'=\theta} = -\nabla_{\theta'} \int \pi_{\theta}(x) \log \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}$$
$$= -\int \pi_{\theta}(x) \nabla_{\theta'} \log \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}$$
$$= -\int \frac{\pi_{\theta}(x)}{\pi_{\theta'}(x)} \nabla_{\theta'} \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}$$
$$= -\nabla_{\theta'} \int \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}$$
$$= 0$$

$$D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'}) \approx D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta}) + d^{T} \nabla_{\theta'} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} + \frac{1}{2} d^{T} \nabla_{\theta'}^{2} D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta} d^{T} \nabla_{\theta'}^{2} D_{\mathrm{KL}}(\pi_{\theta}||\pi$$

$$\nabla_{\theta'}^{2} \operatorname{D}_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta'}) |_{\theta'=\theta} = -\int \pi_{\theta}(x) \nabla_{\theta'}^{2} \log \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}$$

$$= -\int \pi_{\theta}(x) \frac{\pi_{\theta'}(x) \nabla_{\theta'}^{2} \pi_{\theta'}(x) - \nabla_{\theta'} \pi_{\theta'}(x) \nabla_{\theta'} \pi_{\theta'}(x)^{T}}{\pi_{\theta'}(x)^{2}} \, \mathrm{d}x |_{\theta'=\theta}$$

$$= \underbrace{-\nabla_{\theta'}^{2} \int \pi_{\theta'}(x) \, \mathrm{d}x |_{\theta'=\theta}}_{0} + \int \pi_{\theta}(x) \nabla_{\theta'} \log \pi_{\theta'}(x) \nabla_{\theta'} \log \pi_{\theta'}(x)^{T} \mathrm{d}x |_{\theta'=\theta}$$

$$= \mathbb{E}_{x \sim \pi_{\theta}} [\nabla_{\theta'} \log \pi_{\theta'}(x) \nabla_{\theta'} \log \pi_{\theta'}(x)^{T} |_{\theta'=\theta}]$$

Z Wang & C Chen (NJU)

## Hessian of KL = Fisher information matrix (FIM)

• **Hessian**: A square matrix of second-order partial derivatives of a scalar-valued function, which describes the local curvature of a function of many variables

$$\boldsymbol{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

• Fisher information: a way of measuring the amount of information that an observable random variable X carries about an unknown parameter  $\theta$  upon which the probability of X depends

$$\boldsymbol{F}(\theta) = \mathbb{E}_{x \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(x) \nabla_{\theta} \log \pi_{\theta}(x)^{T}]$$

### Hessian of KL = Fisher information matrix (FIM)

#### • The FIM is exactly the Hessian matrix of KL divergence

$$\underbrace{\nabla_{\theta'}^2 \operatorname{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta}}_{\mathsf{Hessian of KL}} = \underbrace{\mathbb{E}_{x \sim \pi_{\theta}} \left[\nabla_{\theta'} \log \pi_{\theta'}(x) \nabla_{\theta'} \log \pi_{\theta'}(x)^T|_{\theta'=\theta}\right]}_{\mathsf{FIM}}$$

$$\operatorname{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'}) \approx \underbrace{\operatorname{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta})}_{0} + d^T \underbrace{\nabla_{\theta'} \operatorname{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta}}_{0} + \frac{1}{2} d^T \underbrace{\nabla_{\theta'}^2 \operatorname{D}_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta'})|_{\theta'=\theta}}_{F(\theta)} d$$

$$= \frac{1}{2} d^T F(\theta) d$$

$$= \frac{1}{2} (\theta' - \theta)^T F(\theta) (\theta' - \theta)$$

$$D_{\mathrm{KL}}(\pi_{\theta}||\pi_{\theta+d}) \approx \frac{1}{2} d^T F(\theta) d$$

- KL divergence is roughly analogous to a distance measure between distributions
- Fisher information serves as a local distance metric between distributions: how much you change the distribution if you move the parameters a little bit in a given direction

### Back to solving the KL constrained problem

$$d^* = \arg\max_{d} J(\theta + d) - \lambda (D_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta + d}) - \epsilon)$$
  

$$\approx \arg\max_{d} J(\theta) + \nabla_{\theta'} J(\theta')|_{\theta' = \theta} \cdot d - \lambda (\frac{1}{2} d^T \nabla_{\theta'}^2 D_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta'})|_{\theta' = \theta} d - \epsilon)$$
  

$$= \arg\max_{d} \nabla_{\theta'} J(\theta')|_{\theta' = \theta} \cdot d - \frac{1}{2} \lambda d^T F(\theta) d$$

• Set the gradient to 0:

$$0 = \frac{\partial}{\partial d} \left( \nabla_{\theta'} J(\theta') |_{\theta'=\theta} \cdot d - \frac{1}{2} \lambda d^T \mathbf{F}(\theta) d \right)$$
$$= \nabla_{\theta'} J(\theta') |_{\theta'=\theta} - \lambda \mathbf{F}(\theta) d$$

$$d^* = \frac{1}{\lambda} \boldsymbol{F}^{-1}(\theta) \nabla_{\theta'} J(\theta')|_{\theta'=\theta} = \frac{1}{\lambda} \boldsymbol{F}^{-1}(\theta) \nabla_{\theta} J(\theta)$$

Z Wang & C Chen (NJU)

### Natural policy gradient

- Policy gradient in distribution space
- Solve the constrained optimization problem
- Natural gradient

### 2 Trust region policy optimization (TRPO)

### **3** Proximal policy optimization (PPO)

### Natural gradient descent

• The natural gradient:

$$\widetilde{\nabla}_{\theta} J(\theta) = \boldsymbol{F}^{-1}(\theta) \underbrace{\nabla_{\theta} J(\theta)}_{\hat{g}}$$

• Natural gradient ascent:

$$\theta' = \theta + \alpha \cdot \boldsymbol{F}^{-1}(\theta)\hat{g}$$

• How to determine the learning rate  $\alpha$ :

$$D_{\mathrm{KL}}(\pi_{\theta} || \pi_{\theta} + d) \approx \frac{1}{2} (\theta' - \theta)^{T} \boldsymbol{F}(\theta) (\theta' - \theta) \leq \epsilon$$
$$\frac{1}{2} (\alpha \hat{g})^{T} \boldsymbol{F}(\alpha \hat{g}) = \epsilon$$
$$\alpha = \sqrt{\frac{2\epsilon}{\hat{g}^{T} \boldsymbol{F} \hat{g}}}$$

### Geometric interpretation of natural policy gradient

#### • Find the steepest direction for parameter updating



### Natural gradient descent $\rightarrow$ Natural policy gradient (NPG)

#### Algorithm 1 Natural Policy Gradient

Input: initial policy parameters  $\theta_0$ for k = 0, 1, 2, ... do Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$ Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm Form sample estimates for

• policy gradient  $\hat{g}_k$  (using advantage estimates)

• and KL-divergence Hessian / Fisher Information Matrix  $\hat{H}_k$ 

Compute Natural Policy Gradient update:

$$heta_{k+1} = heta_k + \sqrt{rac{2\,m{\epsilon}}{\hat{m{g}}_k^{\, T}\hat{m{H}}_k \ \hat{m{g}}_k}} \hat{m{H}}_k^{-1} \hat{m{g}}_k$$

end for

- Originated from natural gradient descent in supervised learning
- Very expensive to compute the inverse of Hessian matrix for a large number of parameters

- The gradient
  - Constrain parameter update in parameter space (using Euclidean distance)
- The natural gradient
  - Constrain parameter update in distribution space (using KL divergence)
  - The meaning of "natural": the distance metric is **invariant** to function parameterization
- Fisher information matrix (FIM)
  - Second-order information: a local distance metric between distributions
  - The FIM is exactly the Hessian matrix of KL divergence
  - Expensive to compute for a large number of parameters

### Natural policy gradient

- Policy gradient in distribution space
- Solve the constrained optimization problem
- Natural gradient

### 2 Trust region policy optimization (TRPO)

### 3 Proximal policy optimization (PPO)

### Trust region policy optimization (TRPO)

- John Schulman, Sergey Levine, Philipp Moritz, Michael Jordan, and Pieter Abbeel, **Trust Region Policy Optimization**, ICML, 2015.
- The family of statistical learning
  - $\bullet~$  John Schulman  $\rightarrow~$  Pieter Abbeel  $\rightarrow~$  Andrew Ng  $\rightarrow~$  Michael Jordan

#### John Schulman's Homepage

I'm a research scientist at OpenAI. I co-lead the reinforcement learning (RL) team, where we work on (1) designing better RL algorithms that enable agents to learn much faster in novel situations; (2) designing better training environments that teach agents transferrable skills. We mostly use games as a testbed.

Previously, I received my PhD in Computer Science from UC Berkeley, where I had the good fortune of being advised by Pieter Abbeel. Prior to my recent work in RL, I spent some time working on robotics, enabling robots to tie knots and stitches and plan movement using trajectory optimization.

- Publications
- Presentations
- Code
- Awards

Email: joschu@openai.com.



### Trust region policy optimization (TRPO)

	Michael I. Jordan		M FOLLOW	Cited by		VIEW ALL
	Professor of EECS and Professor of Statistics, <u>University of California, Berkeley</u> Verified email at cs.berkeley.edu - Homepage				All	Since 2014
	machine learning statistics computational biology a	artificial intelligence optimization		Citations h-index i10-index	165762 160 540	84682 114 425
TITLE		CITED BY	YEAR			17000
Latent dirichlet allocation DM Blei, AY Ng, MI Jordan Journal of machine Learning research 3 (Jan), 993-1022		29247	2003			12750
On spectral clustering: Analysis and an algorithm AY Ng, Mi Jordan, Y Weiss Advances in neural information processing systems, 849-856		7927	2002			4250
Adaptive mixtures RA Jacobs, MI Jordan, Neural computation 3 (	of local experts. SJ Nowlan, GE Hinton 1), 79-87	4089	1991	2012 2013 2014 2	015 2016 2017 201	18 2019





### TRPO - The KL constrained problem

• The objective function:

$$\begin{array}{ll} \underset{\theta}{\operatorname{maximize}} & \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right] \\ \\ \text{subject to} & \hat{\mathbb{E}}_t \left[ \mathrm{D}_{\mathrm{KL}}[\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)] \right] \leq \delta \end{array}$$

• Also worth considering using a penalty instead of a constraint:

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] - \beta \hat{\mathbb{E}}_t \left[ \text{D}_{\text{KL}}[\pi_{\theta_{old}}(\cdot|s_t), \pi_{\theta}(\cdot|s_t)] \right]$$

• Again the KL penalized problem!

#### Algorithm 3 Trust Region Policy Optimization

Input: initial policy parameters 
$$\theta_0$$

for k = 0, 1, 2, ... do

Collect set of trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$ 

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm Form sample estimates for

- policy gradient  $\hat{g}_k$  (using advantage estimates)
- and KL-divergence Hessian-vector product function  $f(v) = \hat{H}_k v$

Use CG with  $n_{cg}$  iterations to obtain  $x_k \approx \hat{H}_k^{-1} \hat{g}_k$ Estimate proposed step  $\Delta_k \approx \sqrt{\frac{2\delta}{x_k^T \hat{H}_k x_k}} x_k$ 

Perform backtracking line search with exponential decay to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^j \Delta_k$$

end for

#### Algorithm 2 Line Search for TRPO

Compute proposed policy step  $\Delta_k = \sqrt{\frac{2\delta}{\hat{g}_k^T \hat{h}_k^{-1} \hat{g}_k}} \hat{H}_k^{-1} \hat{g}_k$ for j = 0, 1, 2, ..., L do Compute proposed update  $\theta = \theta_k + \alpha^j \Delta_k$ if  $\mathcal{L}_{\theta_k}(\theta) \ge 0$  and  $\bar{D}_{KL}(\theta || \theta_k) \le \delta$  then accept the update and set  $\theta_{k+1} = \theta_k + \alpha^j \Delta_k$ break end if end for

• Still very **expensive** to compute the **inverse of Hessian matrix** for a large number of parameters

### Natural policy gradient

- Policy gradient in distribution space
- Solve the constrained optimization problem
- Natural gradient

2 Trust region policy optimization (TRPO)

### 3 Proximal policy optimization (PPO)

### Proximal policy optimization (PPO): Clipped objective

• The surrogate objective function:

$$\mathcal{L}^{\mathrm{IS}}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

• Form a lower bound via clipped importance ratios

$$\mathcal{L}^{\mathsf{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \mathsf{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]$$

- Prevent large changes of policies, constrain the policy update
- Achieve similar performance to TRPO without second-order information (no Fisher matrix!)



## Proximal policy optimization (PPO): Adaptive KL penalty

Input: initial policy parameters  $\theta_0$ , initial KL penalty  $\beta_0$ , target KL-divergence  $\delta$ for  $k = 0, 1, 2, \dots$  do Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$ Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm Compute policy update  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$ by taking K steps of minibatch SGD (via Adam) if  $\overline{D}_{KL}(\theta_{k+1}||\theta_k) > 1.5\delta$  then  $\beta_{k+1} = 2\beta_k$ else if  $\overline{D}_{KL}(\theta_{k+1}||\theta_k) \leq \delta/1.5$  then  $\beta_{k+1} = \beta_k/2$ Don't use second order approximation for KI which is end if expensive, use standard gradient descent end for

- Penalty coefficient  $\beta$  changes between iterations to approximately enforce KL-divergence constraint
- Achieve similar performance to TRPO without second-order information (no Fisher matrix!)

- TRPO: again the KL penalty problem
  - Natural policy gradient + Monotonic policy improvement + Line search
  - Still need to compute the natural gradient with Hessian matrix

PPO

- Achieve TRPO-like performance without second-order computation
- Clipped objective, adaptive KL penalty

$$\begin{array}{ll} \underset{\theta}{\operatorname{maximize}} & \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} \hat{A}_t \right] \\ \\ \text{subject to} & \hat{\mathbb{E}}_t \left[ \mathrm{D}_{\mathrm{KL}} [\pi_{\theta_{old}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)] \right] \leq \delta \end{array}$$

- You should be able to...
  - Know how to derive the natural policy gradient
  - Be aware of several advanced algorithms, e.g., TRPO, PPO
  - Enhance your mathematical skills

- Lecture 9 of CS285 at UC Berkeley, Deep Reinforcement Learning, Decision Making, and Control
  - http://rail.eecs.berkeley.edu/deeprlcourse/static/slides/lec-9.pdf
- Classic papers
  - Peters & Schaal (2008). Reinforcement learning of motor skills with policy gradients: very accessible overview of optimal baselines and natural gradient.
- DRL policy gradient papers
  - Schulman, L., Moritz, Jordan, Abbeel (2015). **Trust region policy optimization**: deep RL with natural policy gradient and adaptive step size.
  - Schulman, Wolski , Dhariwal , Radford, Klimov (2017). Proximal policy optimization algorithms: deep RL with importance sampled policy gradient.
  - Y. Duan, et al., Benchmarking Deep Reinforcement Learning for Continuous Control, *ICML*, 2016.

### Homework 1

- Study the policy gradient algorithm in detail
- ullet Implement the series of policy gradient algorithms on problems 1 & 2
  - Problem 1: the point maze navigation, continuous state-action space  $(s, a \in \mathbb{R}^2, s \in [-0.5, 0.5]^2, a \in [-0.1, 0.1]^2)$
  - Problem 2: the MuJoCo HalfCheetah, make the robot run forward
  - Must use vanilla policy gradient and natural policy gradient, encourage to use TRPO and PPO
- Write a report introducing the algorithms and your experimentation
  - Explanations, steps, evaluation results, visualizations...
  - Submit the code and the report to *huzican0419@gmail.com*





# THE END