

POLICY SEARCH METHODS 3

Herke van Hoof

STAYING CLOSE TO PREVIOUS POLICIES

- Small policy update steps tends to be more ‘safe’
- Estimated value function can be imprecise
(approximation or estimation errors)
- So we don’t want to fully trust the current best guess!

STAYING CLOSE TO PREVIOUS POLICIES

- Small policy update steps tends to be more ‘safe’
- How to measure the difference between policies?
- Policy gradient answer: l2 distance parameter vectors
- One way to derive this:

$$\begin{aligned}\underbrace{\theta^* - \theta_0}_{\text{s.t. } d\theta^T d\theta = c} &= \max_{d\theta} J(\theta_0 + d\theta) \\ &\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta \quad \text{s.t. } d\theta^T d\theta = c \\ &\propto \nabla_{\theta} J(\theta_0)\end{aligned}$$

STAYING CLOSE TO PREVIOUS POLICIES

- Small policy update steps tends to be more ‘safe’
- How to measure the difference between policies?
- Policy gradient answer: l2 distance parameter vectors
- One way to derive this:

$$\begin{aligned}\theta^* - \theta_0 &= \max_{d\theta} J(\theta_0 + d\theta) && \text{s.t. } d\theta^T d\theta = c \\ &\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta && \text{s.t. } d\theta^T d\theta = c\end{aligned}$$

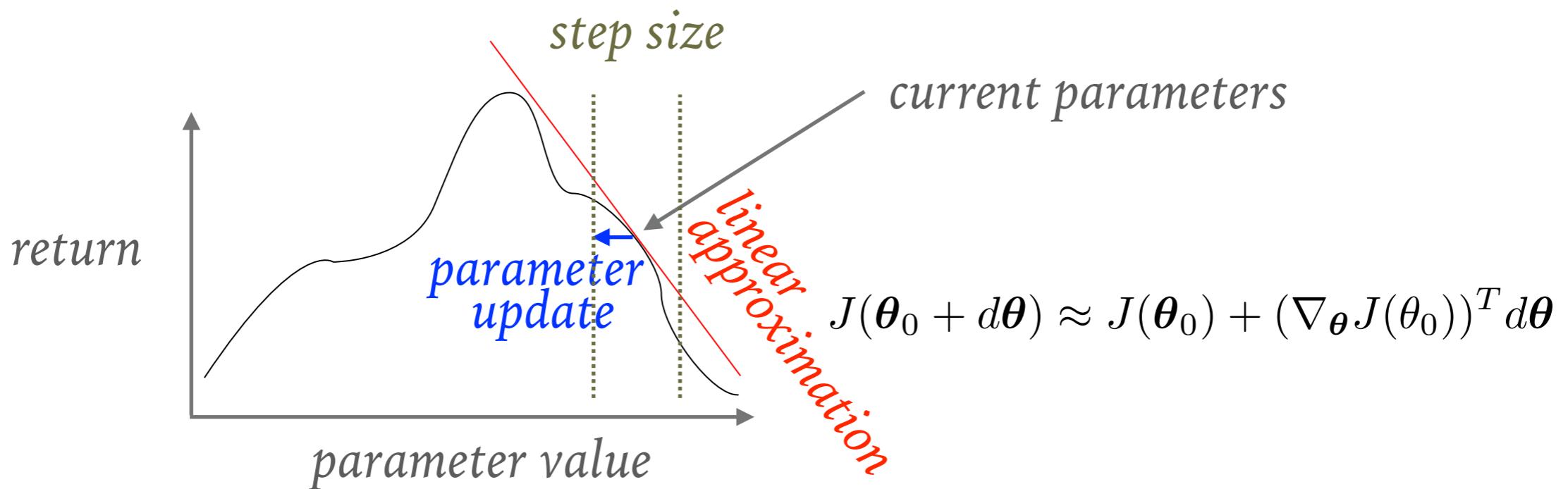
STAYING CLOSE TO PREVIOUS POLICIES

- Small policy update steps tends to be more ‘safe’
- How to measure the difference between policies?
- Policy gradient answer: l2 distance parameter vectors
- One way to derive this:

$$\begin{aligned}\theta^* - \theta_0 &= \max_{d\theta} J(\theta_0 + d\theta) && \text{s.t. } d\theta^T d\theta = c \\ &\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta && \text{s.t. } d\theta^T d\theta = c \\ &\propto \nabla_{\theta} J(\theta_0)\end{aligned}$$

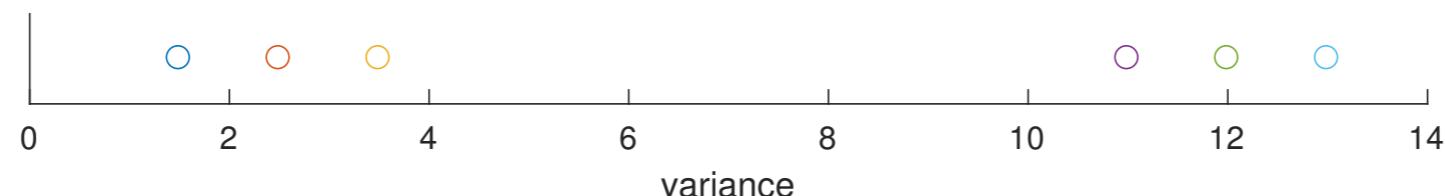
- Standard (‘vanilla’) policy gradients finds direction of most improvement per unit l2 distance in **parameters**

IN A PICTURE



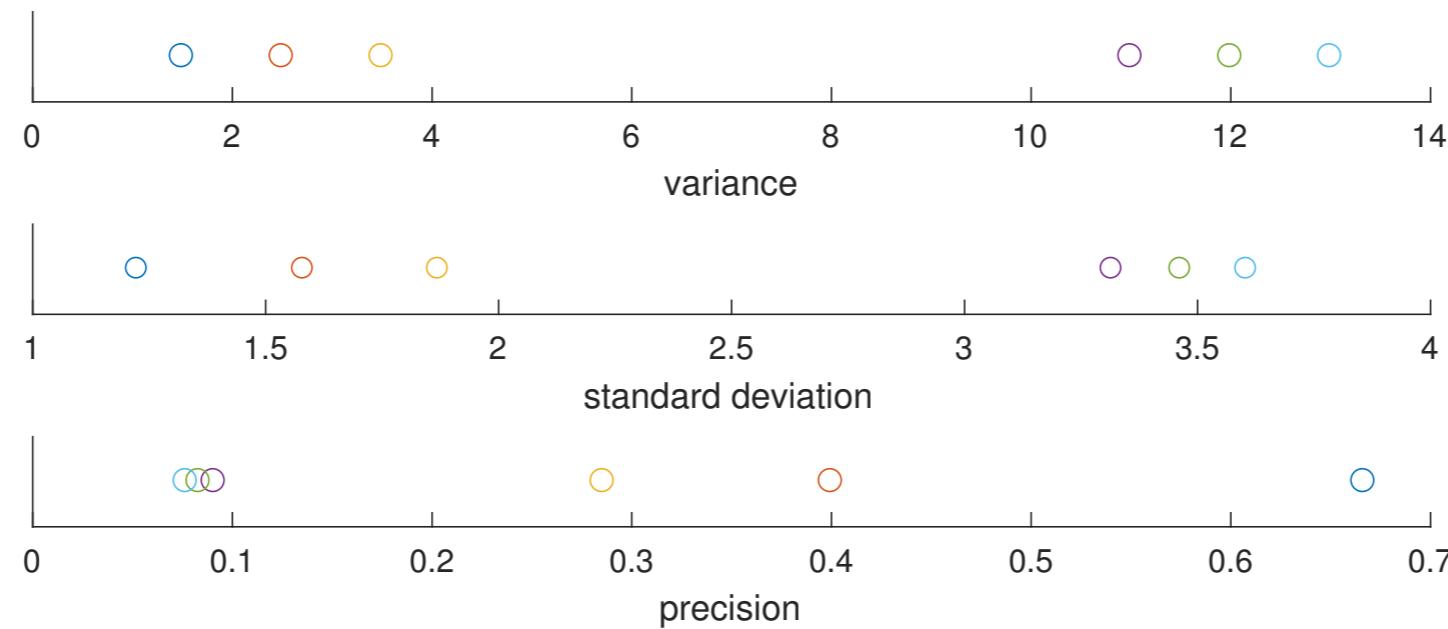
STAYING CLOSE TO PREVIOUS POLICIES

- Standard ('vanilla') policy gradients finds direction of most improvement per unit l2 distance in **parameters**
- This norm is sensitive to parametrisation:



STAYING CLOSE TO PREVIOUS POLICIES

- Standard ('vanilla') policy gradients finds direction of most improvement per unit l2 distance in **parameters**
- This norm is sensitive to parametrisation:



- How to express policy closeness covariantly?
- (covariant: independent of choice of parametrisation)

STAYING CLOSE TO PREVIOUS POLICIES

- Why do we want covariant norm (invariant to parametrisation)?
 - Don't waste time tuning parametrisation *thy parameters*
 - Parameters with different 'meaning': mean and precision
 - does a norm in this space make sense?
 - step size never right on all parameters if scale different
(have to take step small enough for most sensitive direction)
 - Correlations between parameters ignored
(feature modulated by more parameters easier to change)

STAYING CLOSE TO PREVIOUS POLICIES

- Why do we want covariant norm (invariant to parametrisation)?
 - Don't waste time tuning parametrisation
 - Parameters with different 'meaning': mean and precision
 - does a norm in this space make sense?
 - step size never right on all parameters if scale different
(have to take step small enough for most sensitive direction)
 - Correlations between parameters ignored
(feature modulated by more parameters easier to change)
- Conceptually, it's not the change in **parameters** we care about!
 - Limit change in **trajectories, states, and/or actions?**

STAYING CLOSE TO PREVIOUS POLICIES

- How to express policy closeness covariantly?
- Kullback-Leibler (KL) divergence is information-theoretic quantification of difference between distributions

$$D_{\text{KL}}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

- Minimal value of 0 when $p=q$
- KL is invariant under parameter transformations
- Idea: find direction of maximal improvement per unit of **KL** between policies

STAYING CLOSE TO PREVIOUS POLICIES

- Idea: find direction of maximal improvement per unit of KL between policies
- Several algorithms can be understood using this idea
 - Natural policy gradient
 - Trust region policy optimization (TRPO)

NATURAL POLICY GRADIENT

- Idea: make policy gradients covariant [Kakade 2002]
- this yields an algorithm that exploits structure of parameters
- Here, will look how it relates to KL [Bagnell 2003]

[Kakade 2002, Bagnell 2003]

NATURAL POLICY GRADIENT

- Recall vanilla policy gradients

$$\begin{aligned}\theta^* - \theta_0 &= \max_{d\theta} J(\theta_0 + d\theta) \\ &\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta \\ &\propto \underline{\nabla_{\theta} J(\theta_0)}\end{aligned}$$

s.t. $d\theta^T d\theta = c$

s.t. $d\theta^T d\theta = c$

- replace constraint by quadratic expansion of KL divergence

$$\begin{aligned}c &= \mathbb{E}_s [D_{\text{KL}}(\pi(a|s; \theta_0) \| \pi(a|s; \theta_0 + d\theta))] = \underline{\text{EKL}(d\theta)} \\ &\approx \cancel{\text{EKL}(\theta_0)} + d\theta^T \nabla_{d\theta} \text{EKL}(d\theta) + \frac{1}{2} d\theta^T (\nabla_{d\theta}^2 \text{EKL}(d\theta)) d\theta\end{aligned}$$

ge

(gradients and Hessians evaluated at $d\theta = 0$)

NATURAL POLICY GRADIENT

- Recall vanilla policy gradients

$$\theta^* - \theta_0 = \max_{d\theta} J(\theta_0 + d\theta)$$

$$\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta$$

$$\propto \nabla_{\theta} J(\theta_0)$$

s.t. $d\theta^T d\theta = c$

s.t. $d\theta^T d\theta = c$

- replace constraint by quadratic expansion of KL divergence

$$c = \mathbb{E}_s [D_{\text{KL}}(\pi(a|s; \theta_0) \| \pi(a|s; \theta_0 + d\theta))] = \text{EKL}(d\theta)$$

$$\approx \text{EKL}(\theta_0) + d\theta^T \nabla_{d\theta} \text{EKL}(d\theta) + \frac{1}{2} d\theta^T (\nabla_{d\theta}^2 \text{EKL}(d\theta)) d\theta$$

$C = 0$

$$+ \frac{1}{2} d\theta^T (\nabla_{d\theta}^2 \text{EKL}(d\theta)) d\theta$$

- since minimal value of 0 is reached if parameter doesn't change

NATURAL POLICY GRADIENT

$$c = \mathbb{E}_{\mathbf{s}} [D_{\text{KL}}(\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0) \| \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0 + d\boldsymbol{\theta})] = \text{EKL}(d\boldsymbol{\theta})$$

$$\approx \frac{1}{2} d\boldsymbol{\theta}^T (\nabla_{d\boldsymbol{\theta}}^2 \text{EKL}(d\boldsymbol{\theta})) d\boldsymbol{\theta}$$

- This is the squared norm with respect to matrix

$$F = \nabla_{d\boldsymbol{\theta}}^2 \text{EKL} = \mathbb{E}_{\mathbf{s}} \left[\underbrace{\nabla_{d\boldsymbol{\theta}}^2 D_{\text{KL}} (\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0) \| \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0 + d\boldsymbol{\theta}))}_{\text{FIM}} \right]$$

- F can be shown to be the expected Fisher information matrix of the policy
- F characterises information about parameters in action

[Kakade 2002, Bagnell 2003]

NATURAL POLICY GRADIENT

- Consider now the modified optimisation problem

$$\theta^* - \theta_0 = \max_{d\theta} J(\theta_0 + d\theta)$$

$$\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta$$

s.t. $d\theta^T F d\theta = c$

s.t. $d\theta^T F d\theta = c$

- solve constraint optimisation problem: Lagrangian

$$L(d\theta, \lambda) = \underbrace{J(\theta_0)}_{\text{blue}} + \underbrace{(\nabla_{\theta} J(\theta_0))^T d\theta}_{\text{red}} + \lambda \underbrace{(d\theta^T F d\theta - c)}_{\text{green}}$$

- At optimality, partial derivatives of L are 0

$$\frac{\partial L(d\theta, \lambda)}{\partial d\theta} = 0 \Rightarrow (\nabla_{\theta} J(\theta_0)) + \lambda \frac{1}{2} \cancel{F d\theta} \quad]$$

$$\frac{\partial L(d\theta, \lambda)}{\partial \lambda} = 0 \Rightarrow d\theta^T F d\theta - c$$

[Kakade 2002, Bagnell 2003]

NATURAL POLICY GRADIENT

- Consider now the modified optimisation problem

$$\theta^* - \theta_0 = \max_{d\theta} J(\theta_0 + d\theta) \quad \text{s.t. } d\theta^T \mathbf{F} d\theta = c$$

$$\approx \max_{d\theta} J(\theta_0) + (\nabla_{\theta} J(\theta_0))^T d\theta \quad \text{s.t. } d\theta^T \mathbf{F} d\theta = c$$

- solve constraint optimisation problem: Lagrangian

$$L(d\theta, \lambda) = J(\theta_0) + \nabla_{\theta} J(\theta_0)^T d\theta + \lambda(d\theta^T \mathbf{F} d\theta - c)$$

- At optimality, partial derivatives of L are 0

$$\frac{\partial L(d\theta, \lambda)}{\partial d\theta} = 0 \quad \rightarrow \quad (\nabla_{\theta} J(\theta_0)) + \lambda F d\theta = 0$$
$$\frac{\partial L(d\theta, \lambda)}{\partial \lambda} = 0 \quad \quad \quad d\theta^T F d\theta = c$$

[Kakade 2002, Bagnell 2003]

NATURAL POLICY GRADIENT

- So optimality conditions are

$$\begin{aligned} (\underbrace{\nabla_{\theta} J(\theta)}_{\text{blue}}) + \lambda F \boxed{d\theta} &= 0 \\ d\theta^T F d\theta &= c \end{aligned}$$

- From the first line, update direction

$$d\theta \underset{\text{blue}}{\propto} \underbrace{F^{-1} \nabla_{\theta} J(\theta)}_{\text{blue}}$$

- This is the **natural gradient**

(natural gradients in ML used at least since [Amari, 1998],
used in RL since [Kakade 2002])

NATURAL POLICY GRADIENT

- The policy is adapted using the **natural gradient** update rule:

$$\theta_{t+1} = \theta_t + \alpha F^{-1} \nabla_{\theta_t} J(\theta_t)$$

*Reinforce, --
'regular'*

- We can use any known approach for the **vanilla gradient**
- Will this always improve J ?
- For small enough step size, objective improves if

$$(\underline{\theta^* - \theta_0}) \underline{\nabla J(\theta_0)} \geq 0$$
$$\nabla_{\theta} J(\theta_0)^T F^{-1} \nabla J(\theta_0) \geq 0$$

NATURAL POLICY GRADIENT

- The policy is adapted using the **natural gradient** update rule:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \boxed{F^{-1}} \nabla_{\boldsymbol{\theta}_t} J(\boldsymbol{\theta}_t)$$

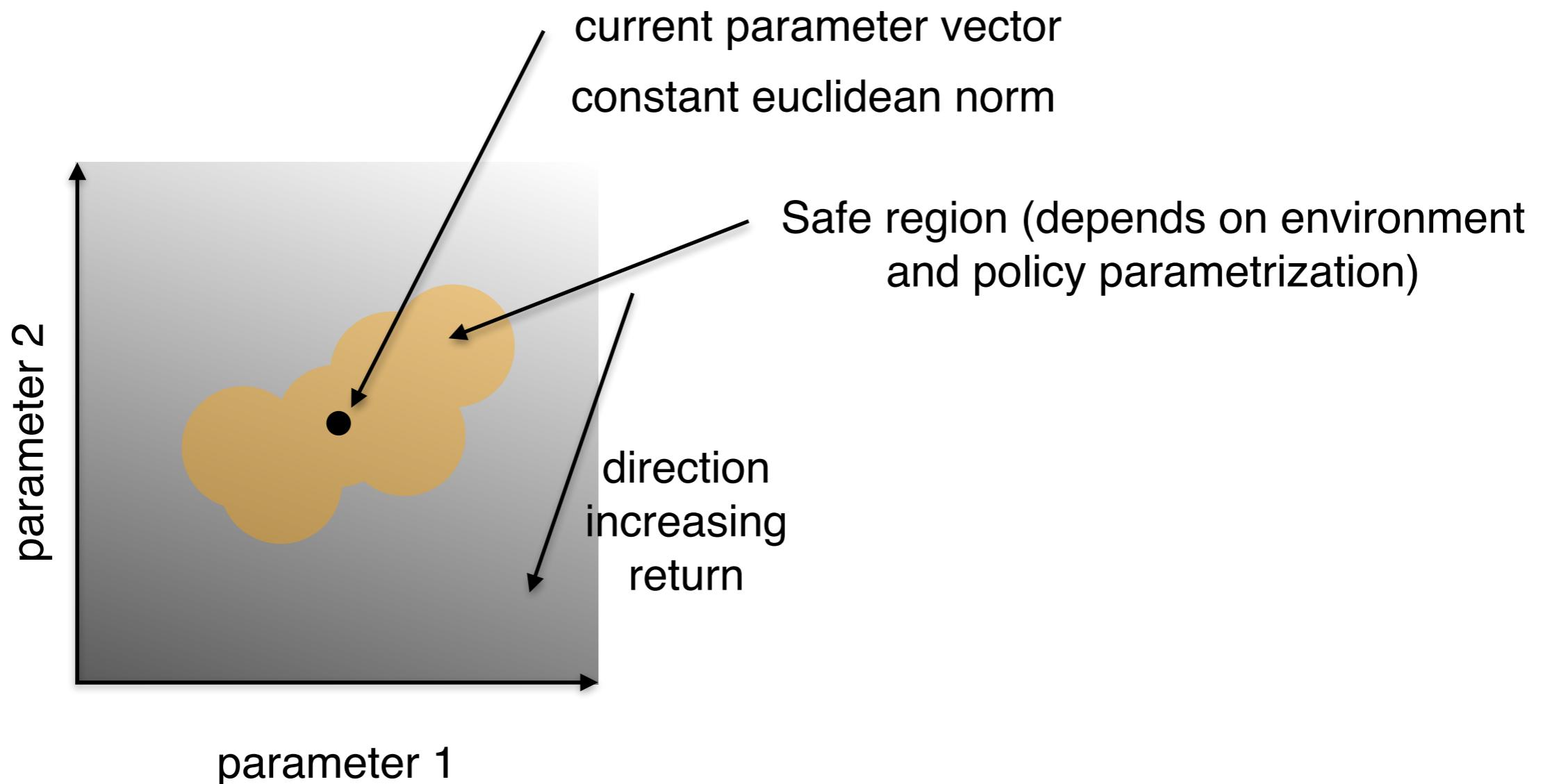
- We can use any known approach for the vanilla gradient
- Will this always improve J ?
- For small enough step size, objective improves if

$$(\boldsymbol{\theta}^* - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) > 0$$

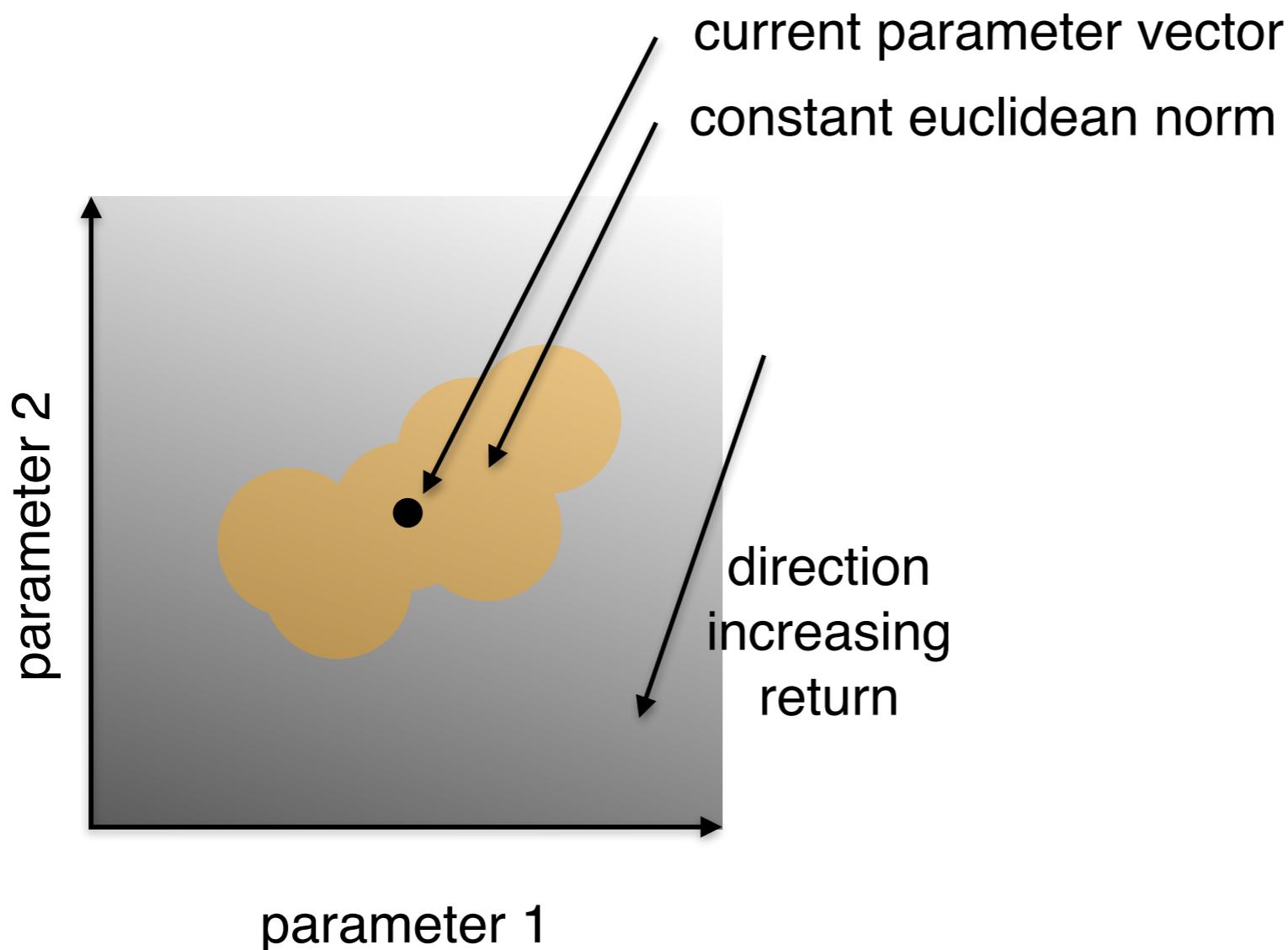
$$(\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0))^T F^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) \stackrel{?}{>} 0$$

- Since Fisher information is positive definite, answer is **yes**

NATURAL POLICY GRADIENT: SOME INTUITION

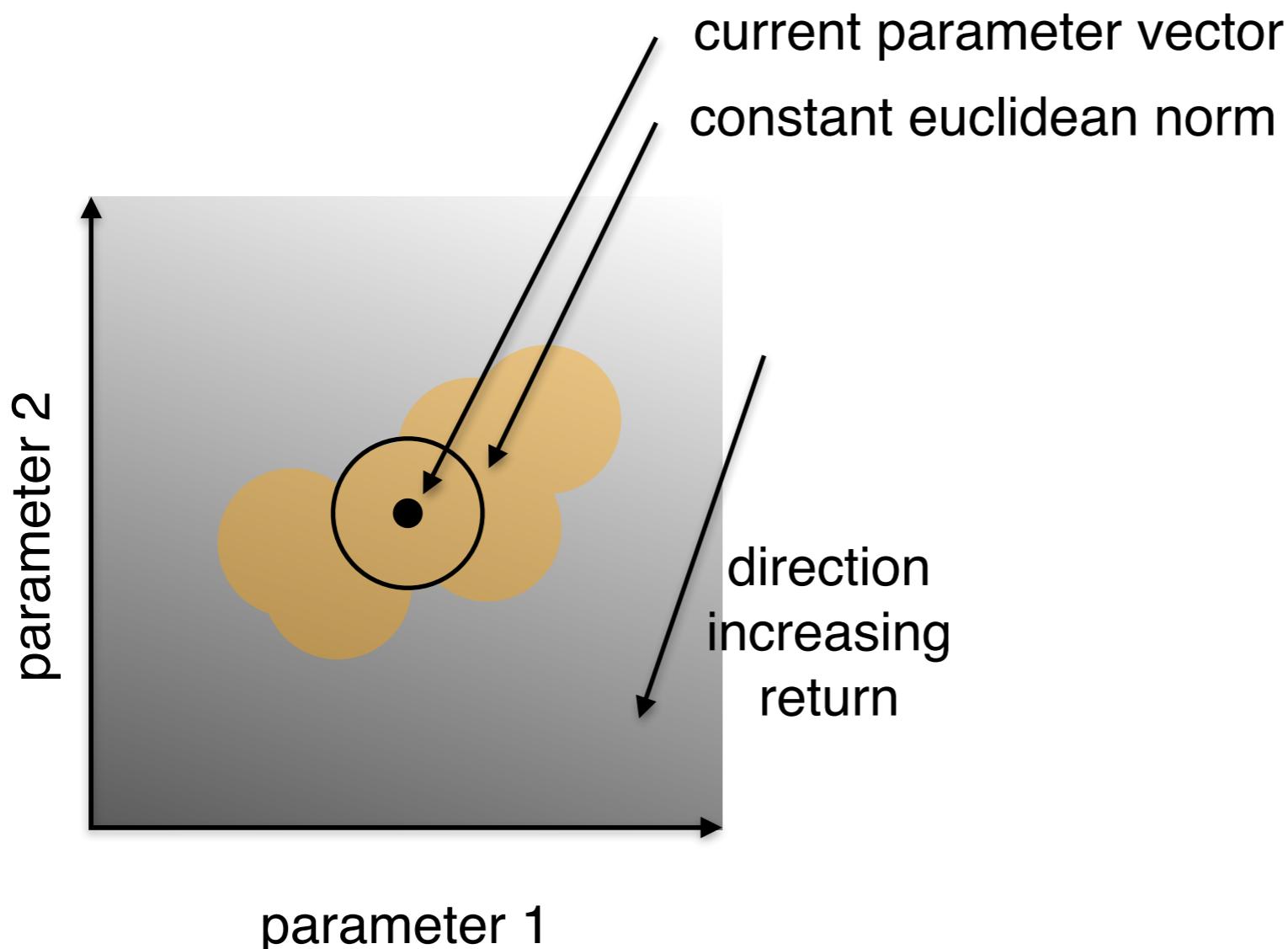


NATURAL POLICY GRADIENT: SOME INTUITION



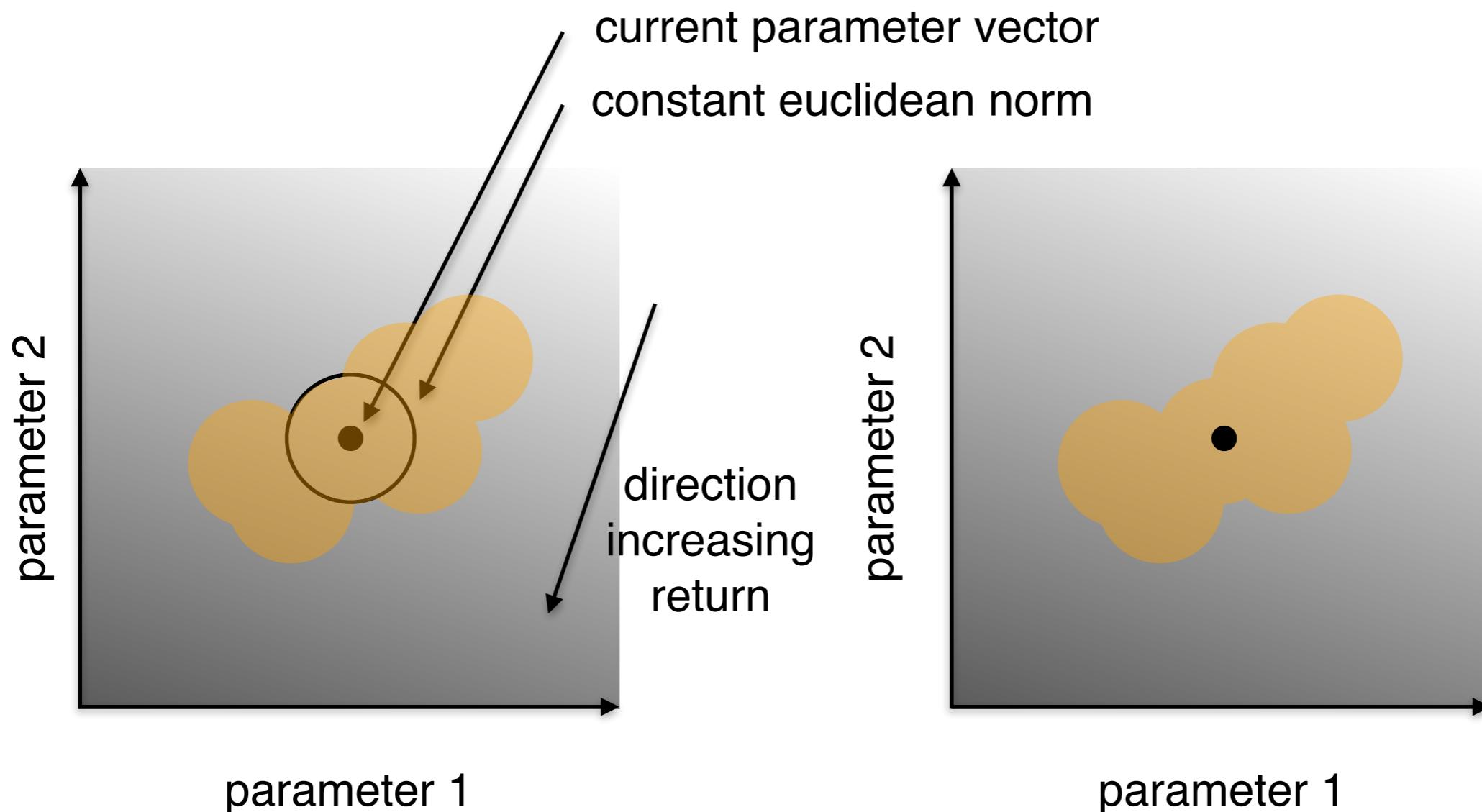
Regular gradient: increase step size till we notice we are outside of safe region

NATURAL POLICY GRADIENT: SOME INTUITION



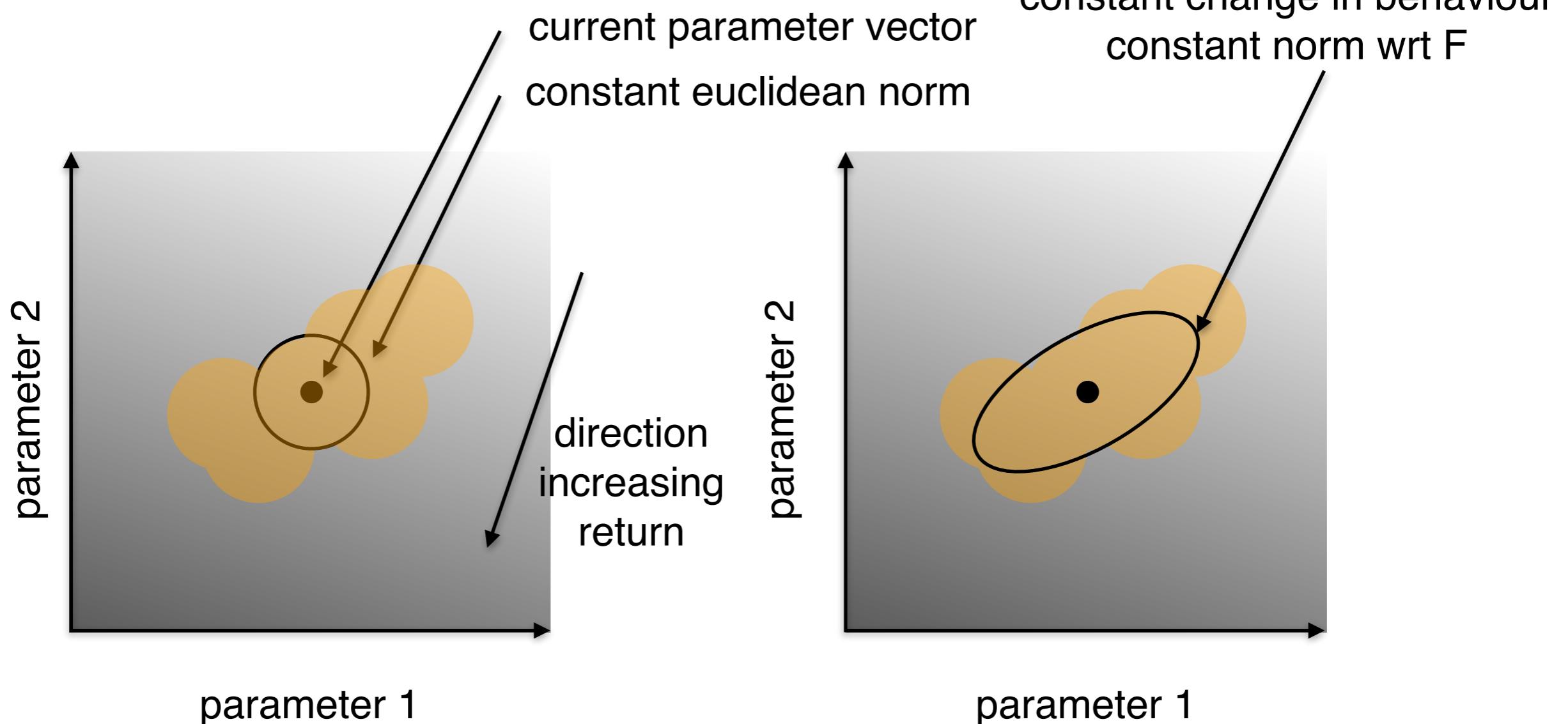
Regular gradient: increase step size till we notice we are outside of safe region

NATURAL POLICY GRADIENT: SOME INTUITION



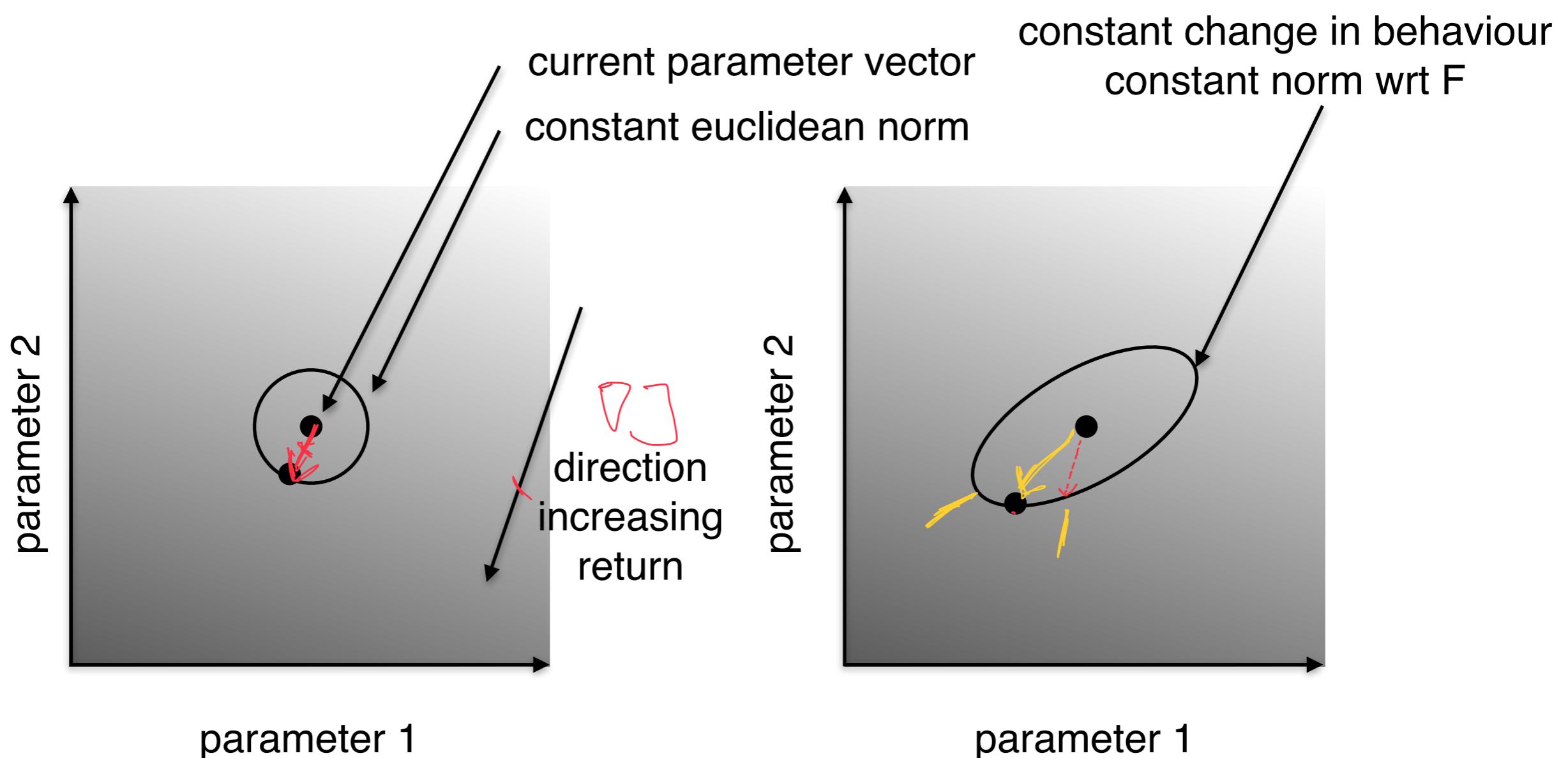
Do the same for natural gradient...

NATURAL POLICY GRADIENT: SOME INTUITION



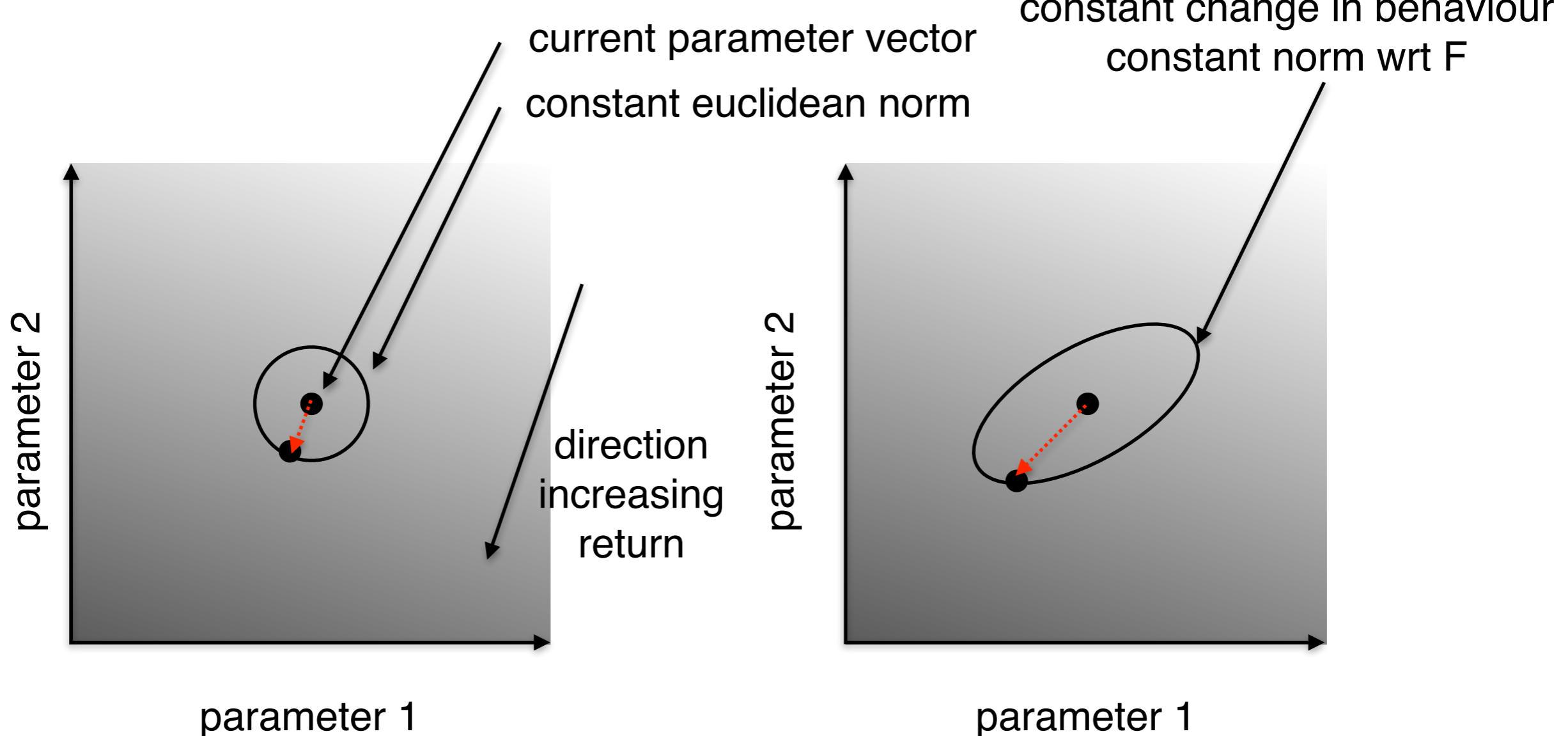
Do the same for natural gradient...

NATURAL POLICY GRADIENT: SOME INTUITION



As natural gradient can (hopefully!) better fit safe region, bigger steps possible

NATURAL POLICY GRADIENT: SOME INTUITION

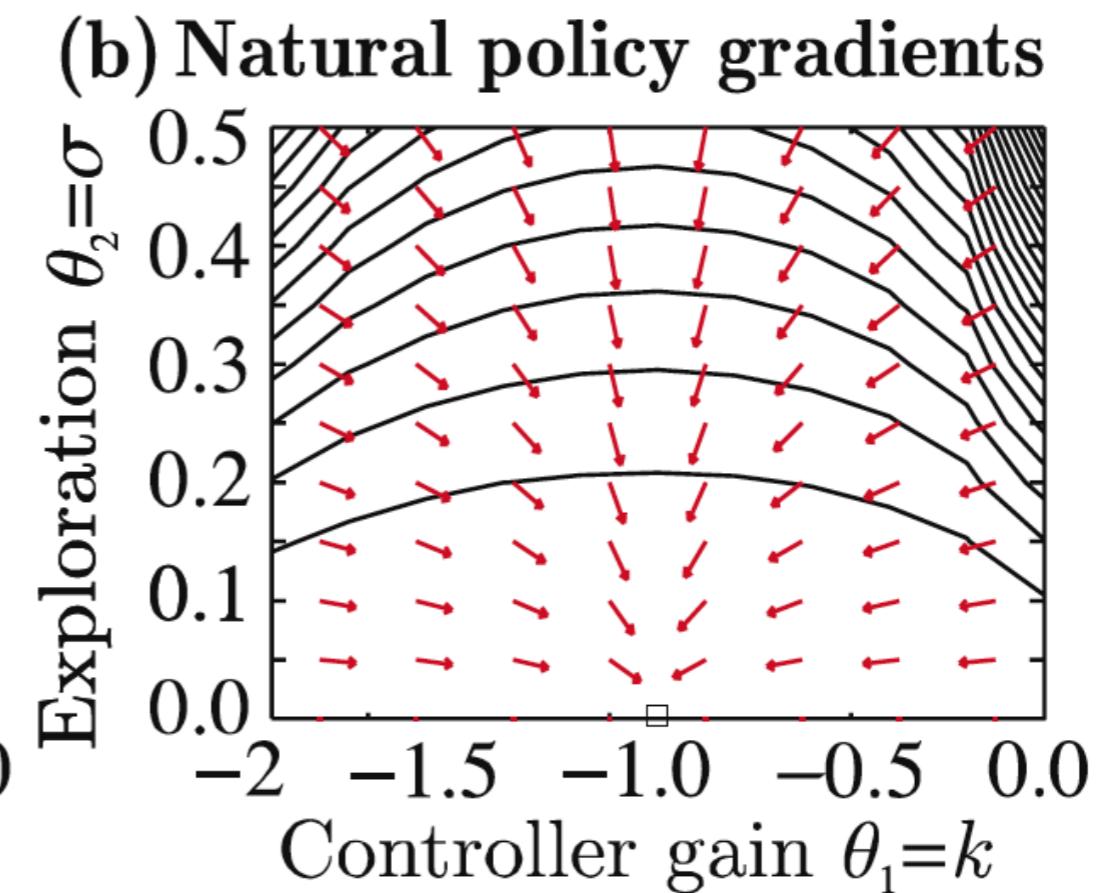
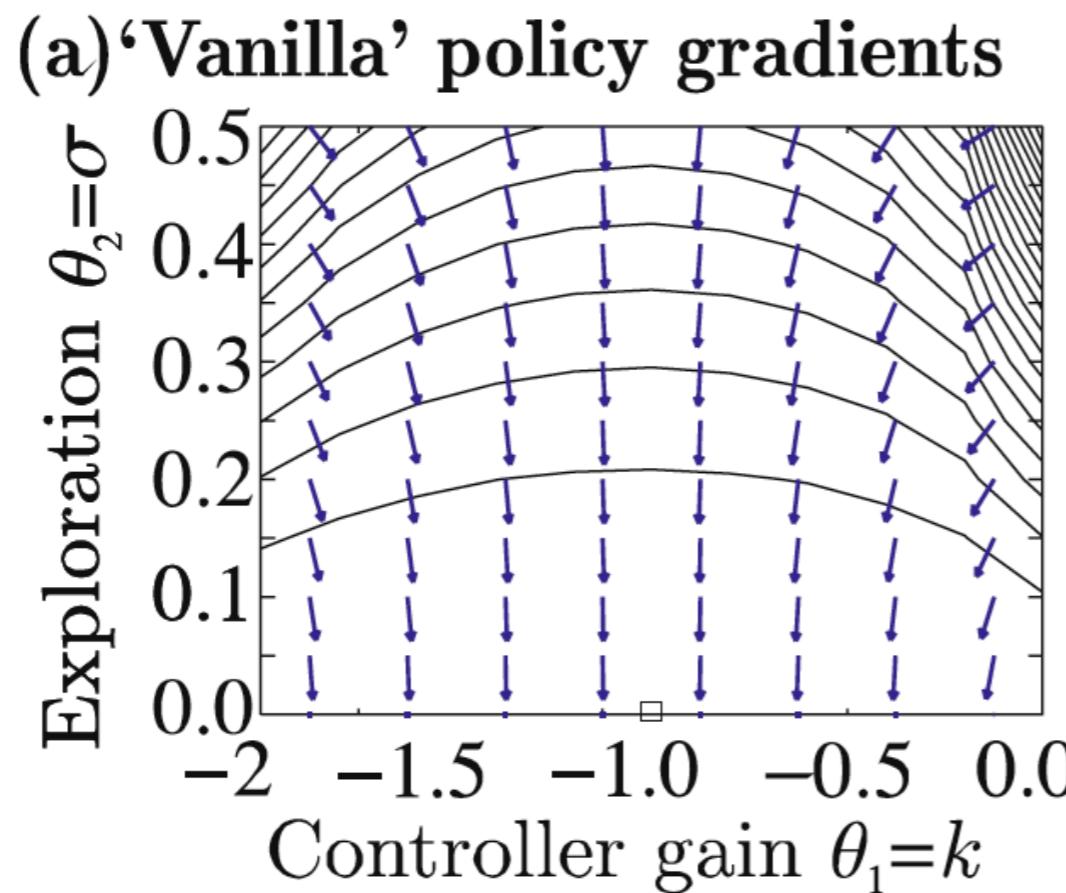


Regular gradient: direction of steepest increase expected return

Natural gradient: within 90° of that direction: will improve objective
‘steepest return per unit policy change’

NATURAL ACTOR CRITIC

- Natural gradients can help where the likelihood is almost flat



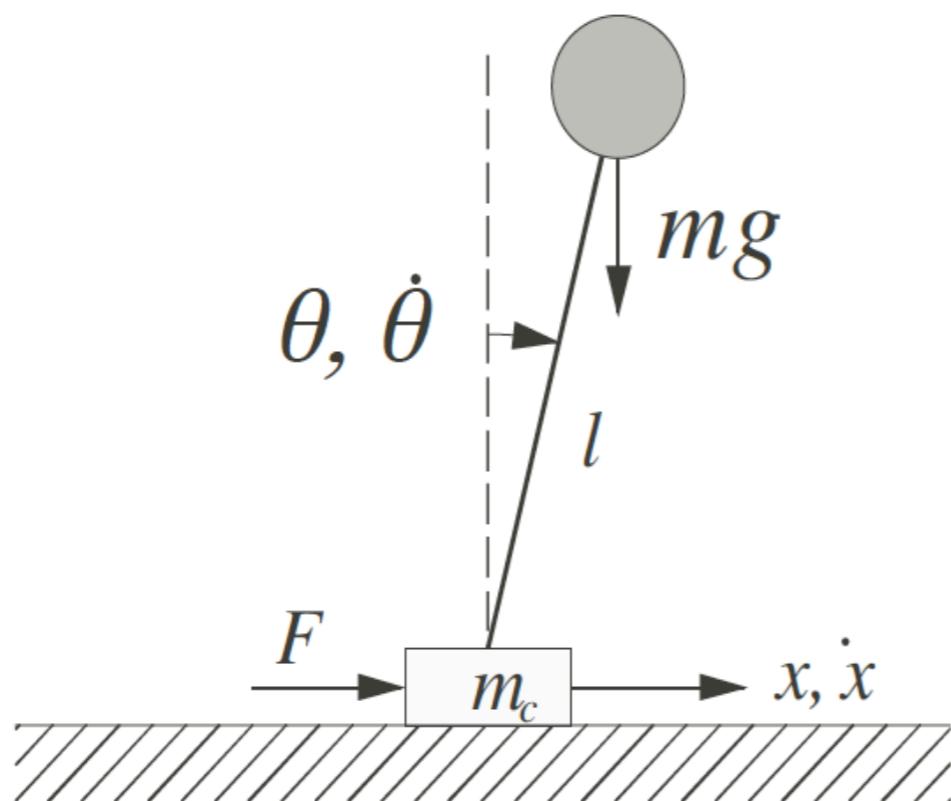
[Peters 2008]

NATURAL POLICY GRADIENT

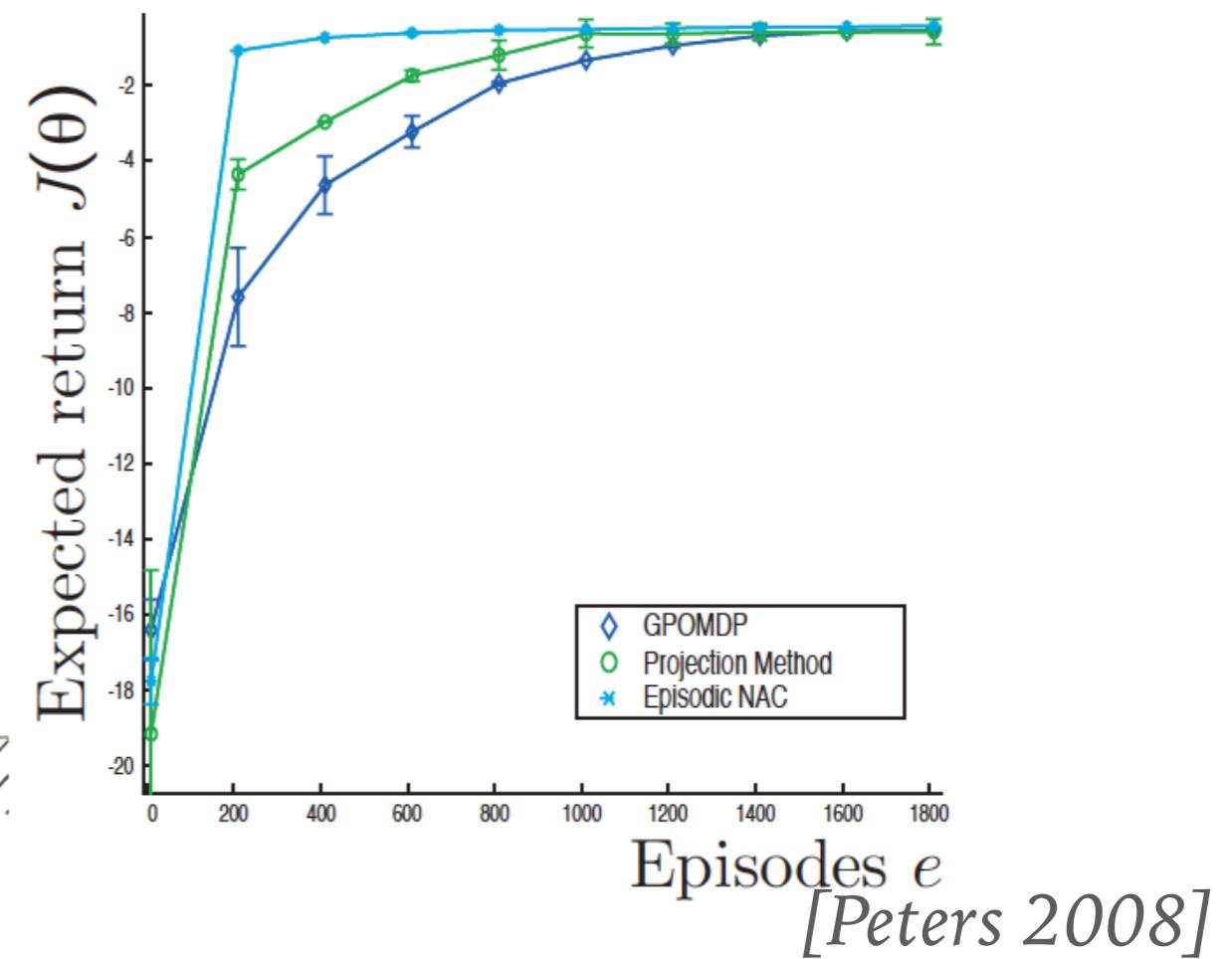
- Natural policy gradients can be used in actor-critic set-up

(1) Cart-Pole Comparison

(a) Physical system



(b) Performance



NATURAL POLICY GRADIENT

- Advantages
 - Usually needs less training than regular policy gradients
 - Inherits advantageous properties from vanilla gradients
- Limitations
 - Need Fisher information matrix
 - Known for some standard distributions, e.g. Gaussian
 → *Computational*.
 - Inherits disadvantages from PG (e.g., high variance)

✓^{data - 1 step}

STAYING CLOSE TO PREVIOUS POLICIES

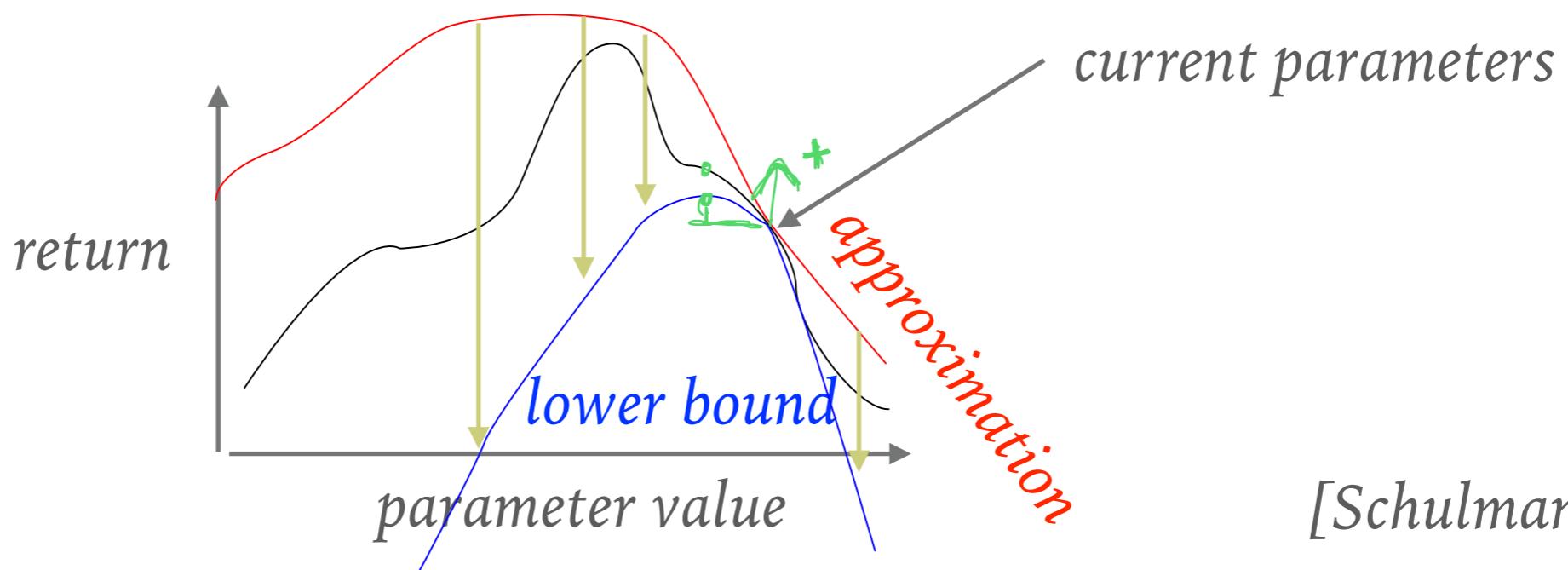
- Idea: use KL to specify how the policy can change in one step
 - Natural policy gradient
 - Trust region policy optimization (TRPO)

TRUST REGION POLICY OPTIMISATION

- Trust region: region where approximation is valid
- Schulman’s “Trust region policy optimisation” defines a trust region based RL algorithm inspired by a theoretical lower bound
- Theoretical starting point quite different than policy gradients, but practical implementation is quite similar

TRPO: THEORY

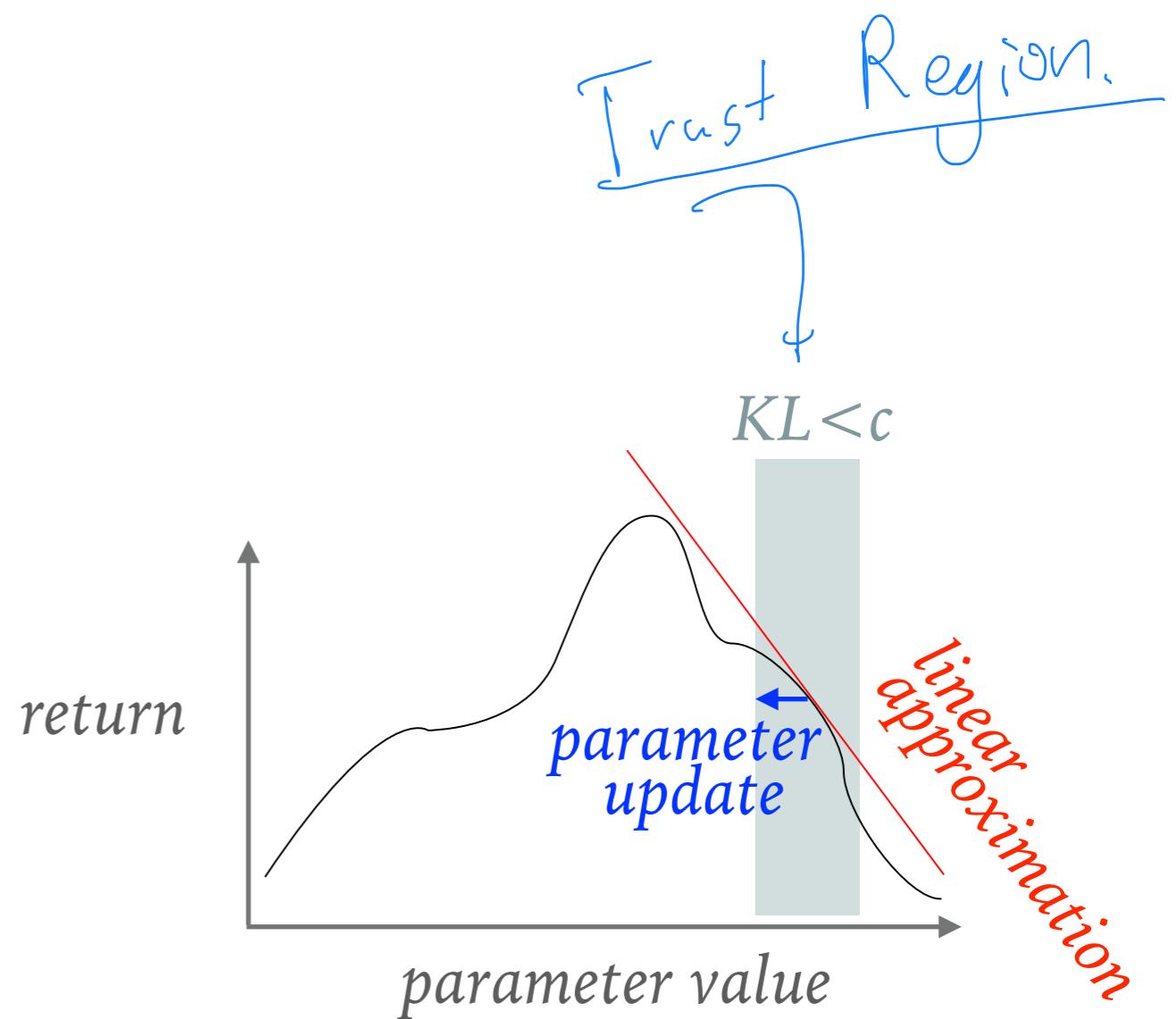
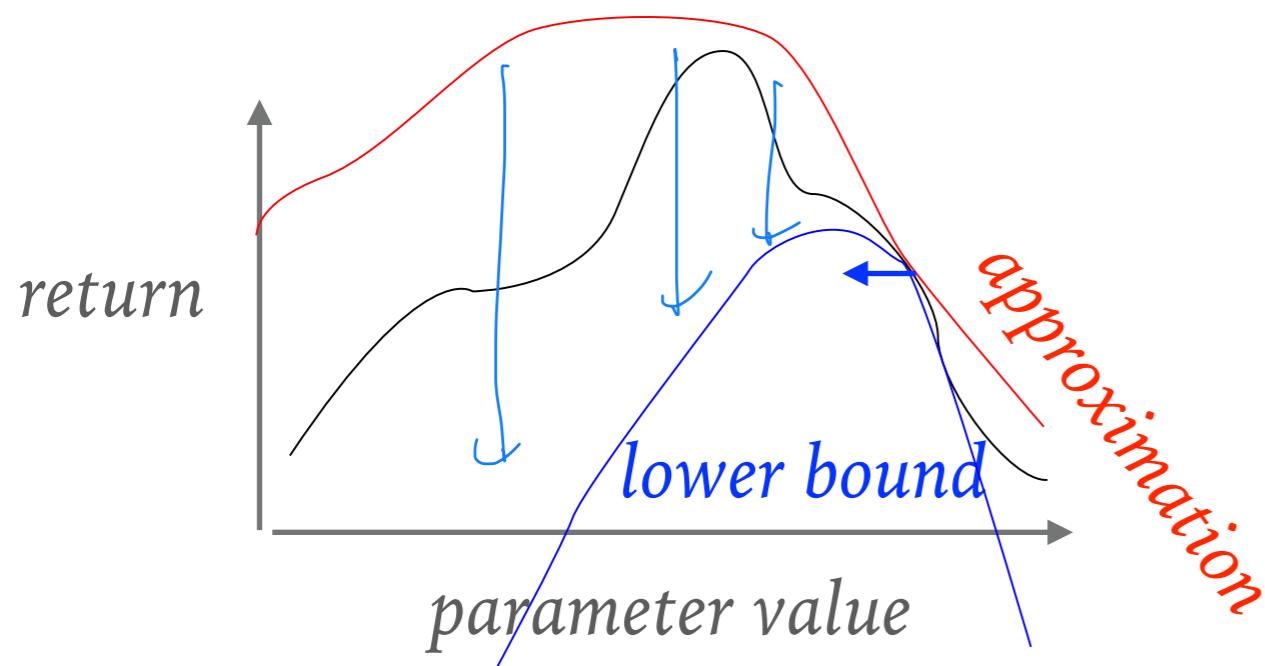
- Idea: take big steps while guaranteeing improvement
1. approximate the return function
(ignoring the difference between μ_b and μ_π)
 2. apply KL-based penalty term to yield lower bound
 3. maximize this lower bound (no need to specify step size!)
- state distributions of
behavior &
new policy.



[Schulman 2016]

TRPO: PRACTICE

- Approximate theoretical update by maximising a linearised approximation under a KL constraint inspired by the KL-based *penalty term*



CONNECTION TO NATURAL GRADIENTS

- Maximising linearised performance under KL constraint is very reminiscent of our perspective on natural gradients!
- In natural gradients, we found the update direction with maximal improvement per unit KL
$$\mathbf{d}\theta \propto \tilde{\nabla} J = F^{-1} \nabla_{d\theta} J(\boldsymbol{\theta}_0 + d\boldsymbol{\theta})$$
taking a fixed-size step β .
- The optimal step size might not be same everywhere in parameter space
- Instead: derive step-size from KL constraint.

CONNECTION TO NATURAL GRADIENTS

- Instead: derive step-size from KL constraint. Write $\mathbf{d}\theta = \beta \tilde{\nabla} J$ and look again at the Taylor expansion of the expected KL

$$c = \mathbb{E}_{\mathbf{s}} [D_{\text{KL}}(\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0) \| \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0 + \mathbf{d}\boldsymbol{\theta})] = \text{EKL}(\mathbf{d}\boldsymbol{\theta})$$

$$\approx \text{EKL}(\mathbf{0}) + \mathbf{d}\boldsymbol{\theta}^T \nabla_{\mathbf{d}\boldsymbol{\theta}} \text{EKL}(\mathbf{d}\boldsymbol{\theta}) + \frac{1}{2} \mathbf{d}\boldsymbol{\theta}^T (\nabla_{\mathbf{d}\boldsymbol{\theta}}^2 \text{EKL}(\mathbf{d}\boldsymbol{\theta})) \mathbf{d}\boldsymbol{\theta}$$

$$C = \frac{1}{2} \cdot \beta (\tilde{\nabla} J)^T (\nabla_{\mathbf{d}\boldsymbol{\theta}} \text{EKL}(\mathbf{d}\boldsymbol{\theta})) \beta (\tilde{\nabla} J)$$

$$\beta = \sqrt{\frac{1}{C} \cdot \frac{1}{2} \cdot (\tilde{\nabla} J)^T (\nabla_{\mathbf{d}\boldsymbol{\theta}} \text{EKL}(\mathbf{d}\boldsymbol{\theta})) (\tilde{\nabla} J)}$$

CONNECTION TO NATURAL GRADIENTS

$$F^{-1} \nabla \tilde{J}$$

- Instead: derive step-size from KL constraint. Write $\mathbf{d}\theta = \beta \tilde{\nabla} J$ and look again at the Taylor expansion of the expected KL

$$c = \mathbb{E}_{\mathbf{s}} [D_{\text{KL}}(\pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0) \| \pi(\mathbf{a}|\mathbf{s}; \boldsymbol{\theta}_0 + \mathbf{d}\boldsymbol{\theta})] = \text{EKL}(\mathbf{d}\boldsymbol{\theta})$$

$$\begin{aligned} & \approx \text{EKL}(\mathbf{0}) + \mathbf{d}\boldsymbol{\theta}^T \nabla_{\mathbf{d}\boldsymbol{\theta}} \text{EKL}(\mathbf{d}\boldsymbol{\theta}) + \frac{1}{2} \mathbf{d}\boldsymbol{\theta}^T (\nabla_{\mathbf{d}\boldsymbol{\theta}}^2 \text{EKL}(\mathbf{d}\boldsymbol{\theta})) \mathbf{d}\boldsymbol{\theta} \\ &= 0 \quad + \frac{1}{2} \beta^2 (\tilde{\nabla} J)^T \underbrace{(\nabla_{\mathbf{d}\boldsymbol{\theta}}^2 \text{EKL}(\mathbf{d}\boldsymbol{\theta})) (\tilde{\nabla} J)}_F \end{aligned}$$

- $c \approx \beta^2 (\tilde{\nabla} J)^T F (\tilde{\nabla} J) / 2$ can now be solved to get β

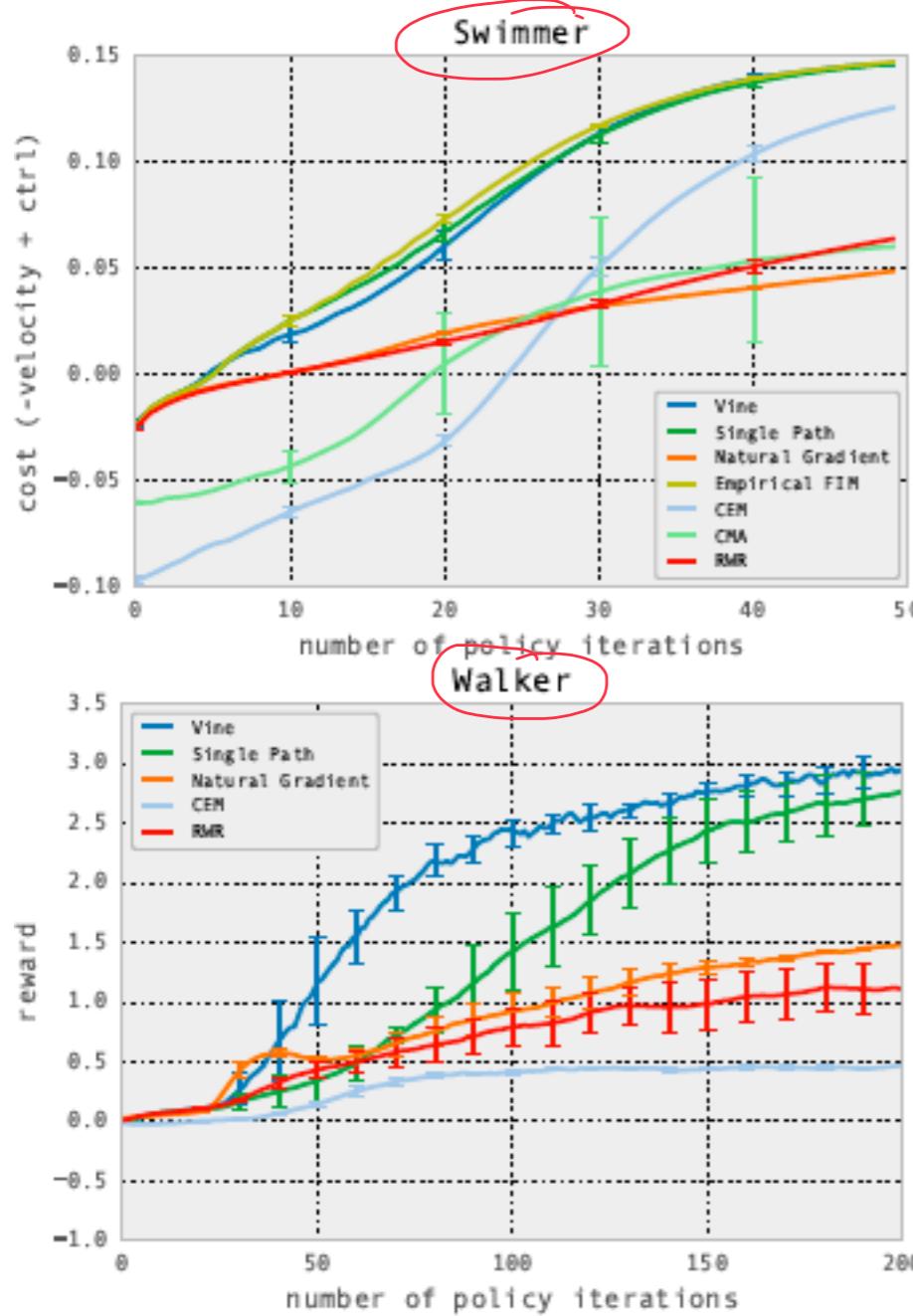
- Based on approximation of objective and KL

(*approximate*)

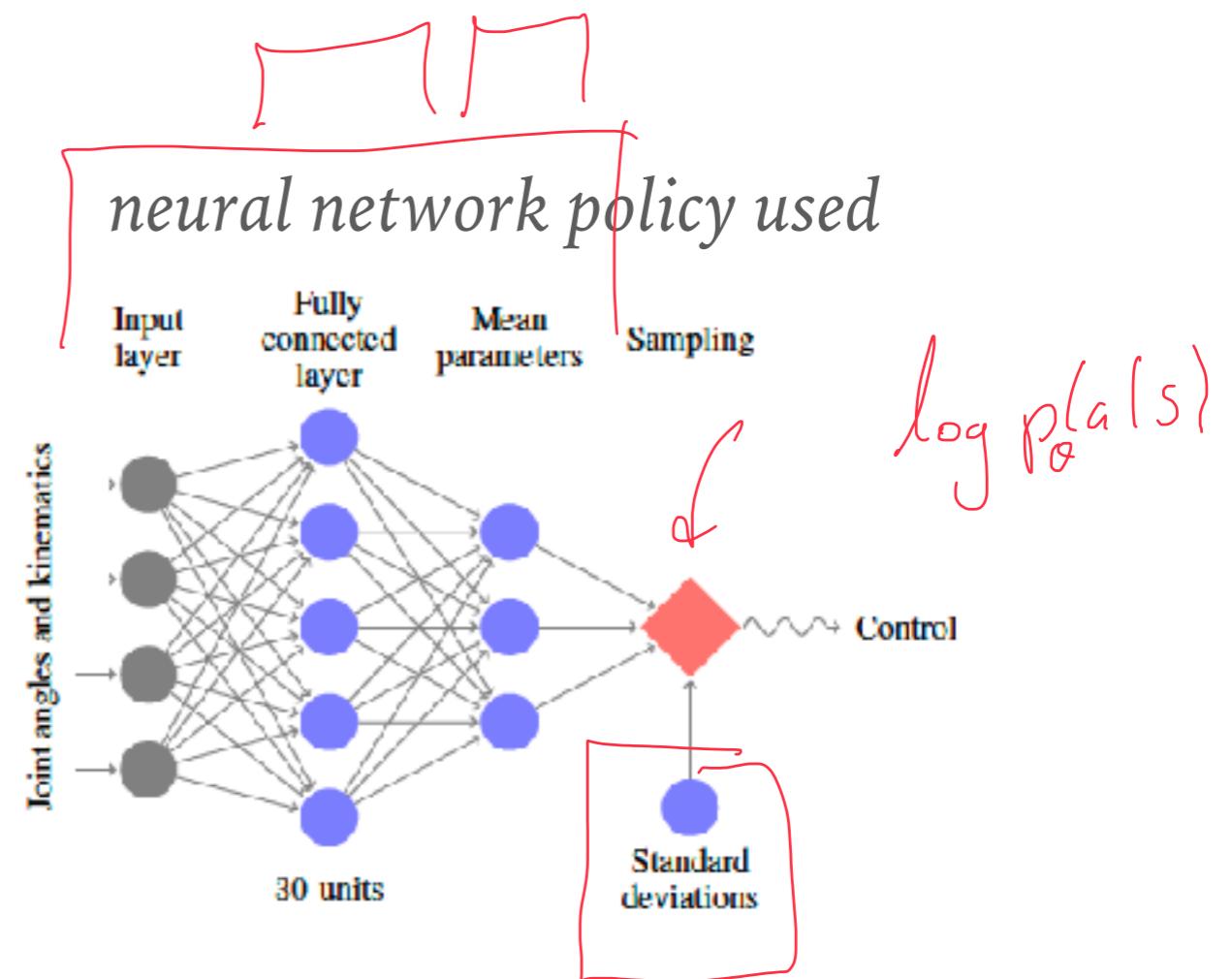
- Make sure constraint is met using analytic objective & KL

- Additional tricks to work with large number of parameters (NN)

TRPO EVALUATION



*different TRPO variants
direct policy search (meta heuristic)
natural gradients,
reward-weighted regression*



[Schulman, 2015]

TRPO EXAMPLE



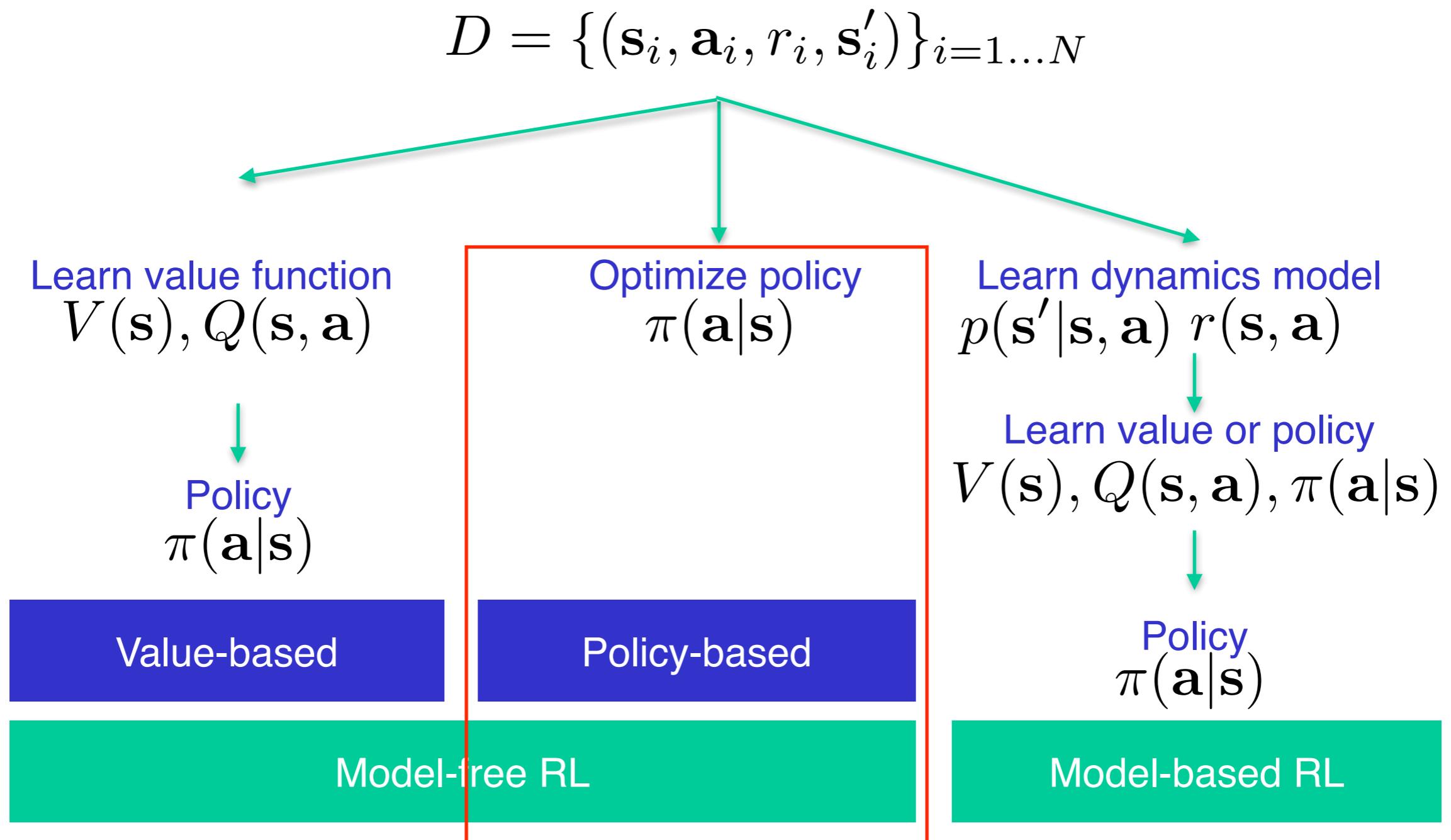
TRPO with generalised advantage estimate, [Schulman 2016]

TRPO EXAMPLE



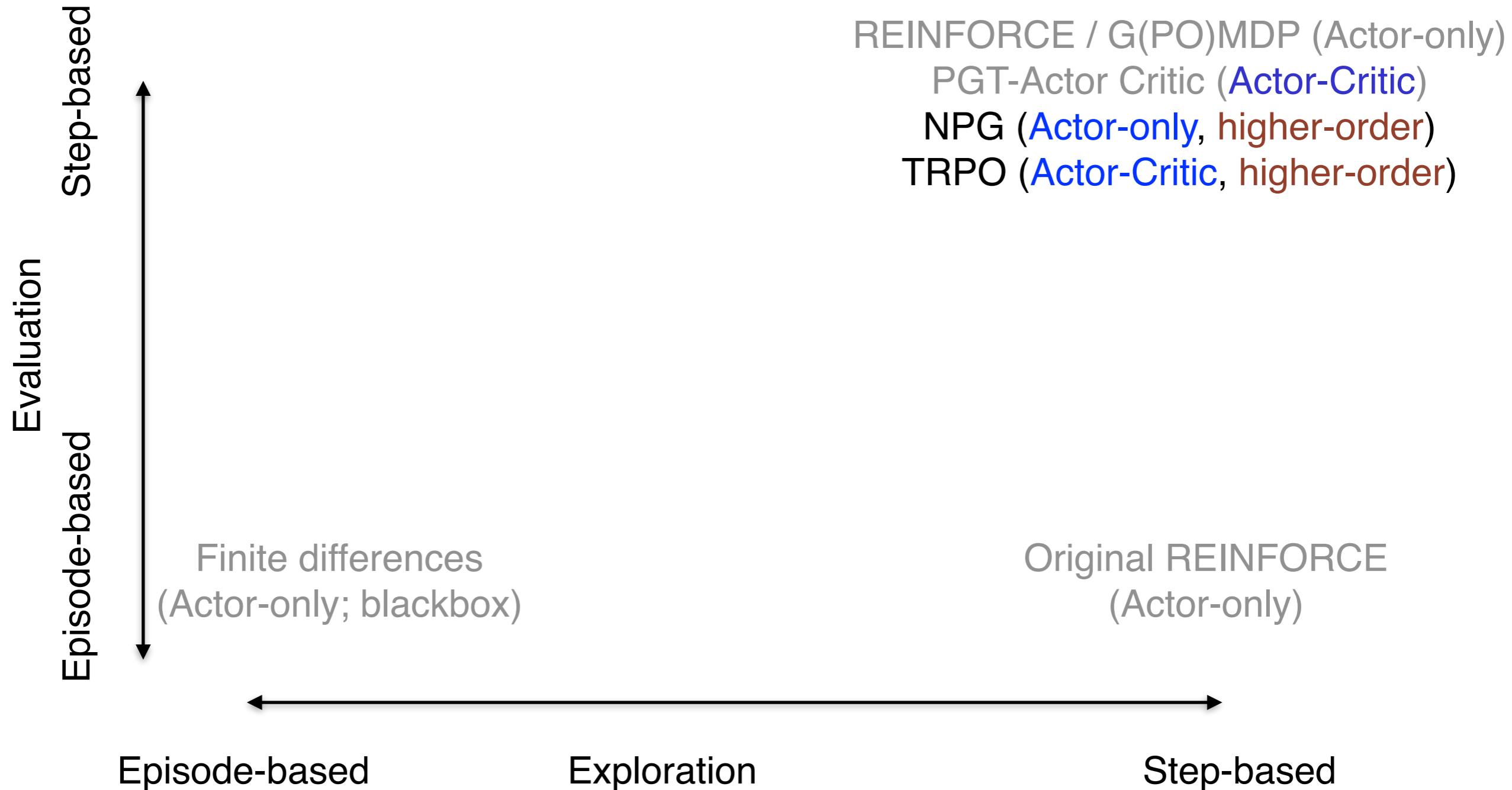
TRPO with generalised advantage estimate, [Schulman 2016]

Big picture: How to learn policies

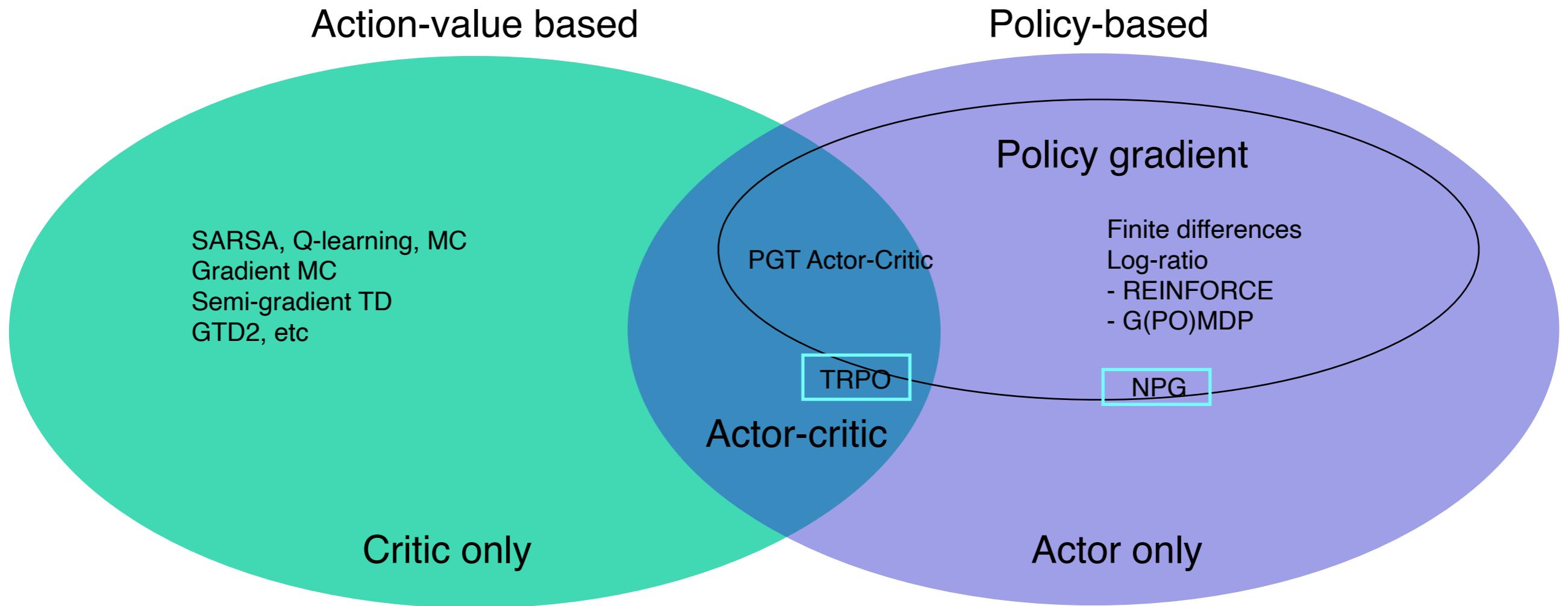


Thanks to Jan Peters

Comparison



Policies and action-values



CONCLUSIONS

- NPG, TRPO: all use improved metric for policy updates, exploit structure of policy
- Generally allows taking larger steps in policy space than ‘vanilla’ methods
→ Both have non-trivial computational cost compared to P-G.
- NPG: easier to implement, still manual step size
- TRPO: even larger steps (faster), step size from KL constraint

WHAT YOU SHOULD REMEMBER

- Advantage of covariant representation of distances?
- Advantage of specifying constraint instead of stepsize?
- Why do we need a constraint / penalty / stepsize?

Thanks for your attention!

Feedback?

h.c.vanhoof@uva.nl

REFERENCES

- [Bagnell 2003] Bagnell, J.A. and Schneider, J. Covariant policy search. IJCAI.
- [Kakade 2002] Kakade, S., 2002. A natural policy gradient. In: NeurIPS
- [Peters 2008] Peters, J. and Schaal, S., 2008. Natural actor-critic. Neurocomputing, 71(7), pp.1180-1190
- [Sutton 2000] Sutton, R.S., McAllester, D., Singh, S. & Mansour, Y. Policy Gradient Methods for Reinforcement Learning with Function Approximation. In: NeurIPS.
- [Schulman 2015] Schulman, J., Levine, S., Abbeel, P., Jordan, M.I. and Moritz, P. Trust Region Policy Optimization. In: ICML
- [Schulman 2016] Schulman, J., Moritz, P., Levine, S., Jordan, M.I. and Abbeel, P. High-dimensional continuous control using generalised advantage estimates. In ICLR.

