

6905: Human-Centric Machine Learning

Reinforcement Learning Overview

Prof. Scott Niekum

Reinforcement Learning:

Learning from experience in sequential decision problems

BC: (s, a^*) vs. RL: (s, a, s', r)

MDP: States: S

Actions: A

Transition probabilities: $T(s, a, s') = p(s' | s, a)$

Reward function $r(s, a)$

Start state distribution: d_0

Discount factor: γ

Markov Property: $p(s_t | s_{t-1}, a_{t-1}, \dots, s_0, a_0) = p(s_t | s_{t-1}, a_{t-1})$

Objective : find an optimal policy π^* that maximizes expected return

$$\text{Return} : R = \sum_{t=0}^{\infty} \gamma^t r_t$$

$$\text{Expected Return : } \underset{\text{on MDP}}{E_{\pi, s_0 \sim d_0}} [R]$$

Why discounting? Why geometric?

Side note : Sometimes average reward optimized for instead.

Why is RL hard?

Delayed effects / credit assignment

Curse of dimensionality / horizon

Exploration vs. Exploitation

Continuous states + actions

No fixed dataset

Generalization / distribution shift / nonstationarity

Designing state representation

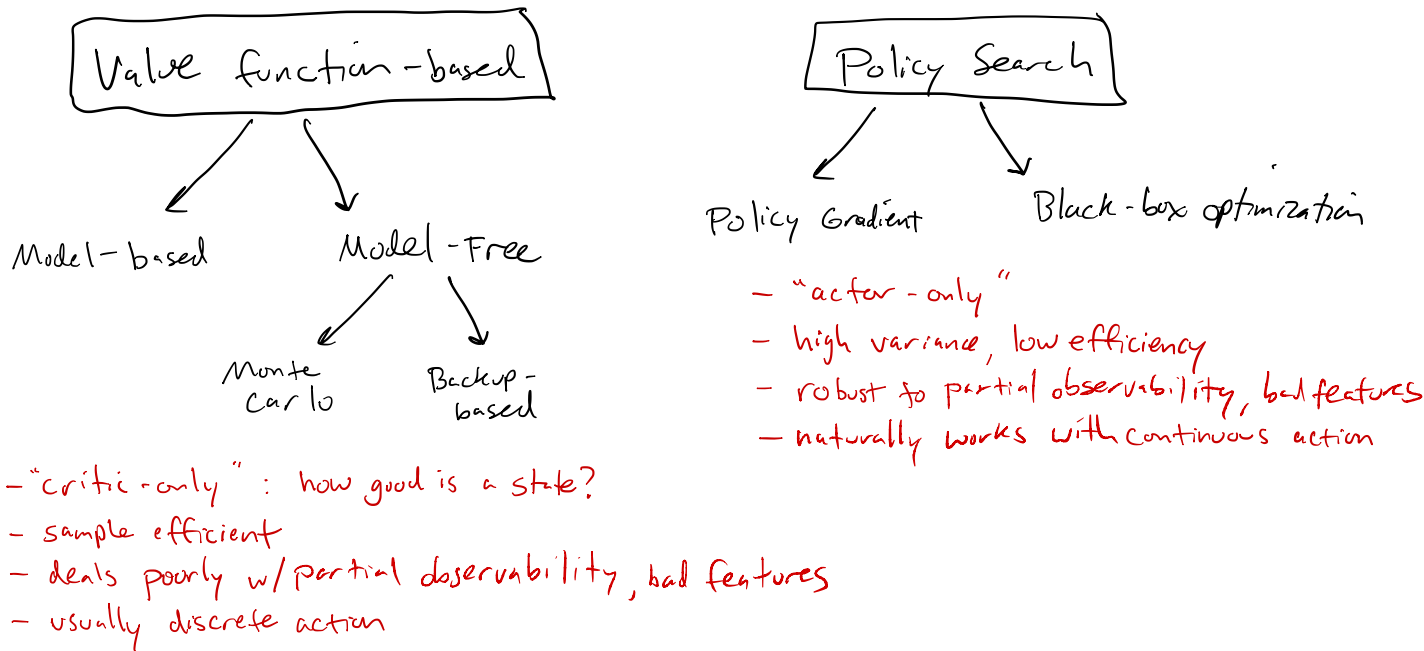
Defining reward

Hyperparameter tuning

Reproducibility

⋮

RL landscape



Actor-critic
combines VF + PS

- often best of both worlds
- efficient, robust
- continuous action

Value functions:

$$\begin{aligned} V_{\pi}(s) &= E_{\pi} [R_t | s_t = s] \text{ where } R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+i} \\ &= \sum_a \pi(s, a) \sum_{s'} T(s, a, s') [r(s, a) + \gamma V_{\pi}(s')] \end{aligned}$$

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi} [R_t | s_t = s, a_t = a] \\ &= \sum_{s'} T(s, a, s') [r(s, a) + \gamma V_{\pi}(s')] \end{aligned}$$

Bellman optimality equations:

$$V^* = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') [r(s, a) + \gamma V^*(s')]$$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [r(s, a) + \gamma V^*(s')]$$

Optimal policy:

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

Value function methods

Model-based : Dynamic programming to compute Q^*/V^*

Model-free :

Monte-Carlo : Average return samples

On the i^{th} visit to (s, a) , record return after: $R_{s,a}^i$

$$\hat{Q}(s, a) = \frac{1}{N} \sum_{i=0}^N R_{s,a}^i$$

$$Q\text{-learning: } \hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \alpha \left[\underbrace{r_t + \gamma \max_{a'} \hat{Q}(s', a')}_{\text{"observation"}} - \underbrace{\hat{Q}(s, a)}_{\text{prediction}} \right]$$

"Bootstrapping"

off-policy vs. on-policy

$$Q\text{-learning: } \hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha [r_t + \gamma \max_{a'} \hat{Q}(s',a') - \hat{Q}(s,a)]$$

Q-learning is **off-policy**. Actions can come from any policy, but \hat{Q} will still converge to Q^* (in the tabular case)

SARSA is **on-policy**. Learns and improves Q^π based on current policy π

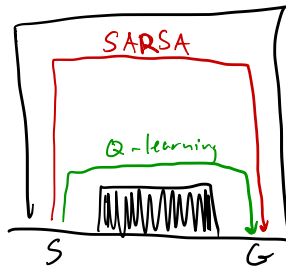
$$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha [r_t + \gamma \hat{Q}(s',\pi(s')) - \hat{Q}(s,a)]$$

e.g. where $\pi \sim$ epsilon-greedy ($\hat{Q}(s',\cdot)$)

π must become greedy over time for on-policy to converge to Q^* !
otherwise



Converged policies:
Cliff world with
E-greedy
exploration



Why would you use
on-policy algorithms?

→ Better convergence
guarantees under function
approximation!

Policy Search

Parameterized policy $f_{\theta}(s) \rightarrow a$ or $f_{\theta}(s, a) \rightarrow \mathbb{R}$

Example with linear function approx. w/ features ϕ :

$$f_{\theta}(s) = \theta_1 \phi_1(s) + \theta_2 \phi_2(s) + \dots + \theta_n \phi_n(s)$$

Methods:

① Policy gradient: approximate $\frac{dR}{d\theta}$ and use to update parameters

$$J(\theta) = E_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

$$\nabla_{\theta} J(\theta) = \sum_s \mu_{\pi}(s) \sum_a Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(s, a)$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_{\theta} J(\theta)$$

Need to approximate

REINFORCE does

Monte Carlo approx of Q^{π}

② Black-box optimization

- CMA-ES

- perturb + hill-climb

Actor - Critic

Like policy search, directly parameterizes policy

Like value function methods, also computes VF

Example:

- REINFORCE (policy search) uses MC estimation of $Q(s_t, a_t)$ when estimating $\nabla_{\theta} J(\theta)$

$$\text{e.g. } \hat{Q}(s_t, a_t) = r_t + \gamma r_{t+1} + \dots + \gamma^N r_{t+N}$$

- Critic can reduce variance by using:

$$\hat{Q}(s_t, a_t) = r_t + \gamma \hat{V}(s_{t+1})$$

✓ Continuous action

✓ Lower variance than policy search alone

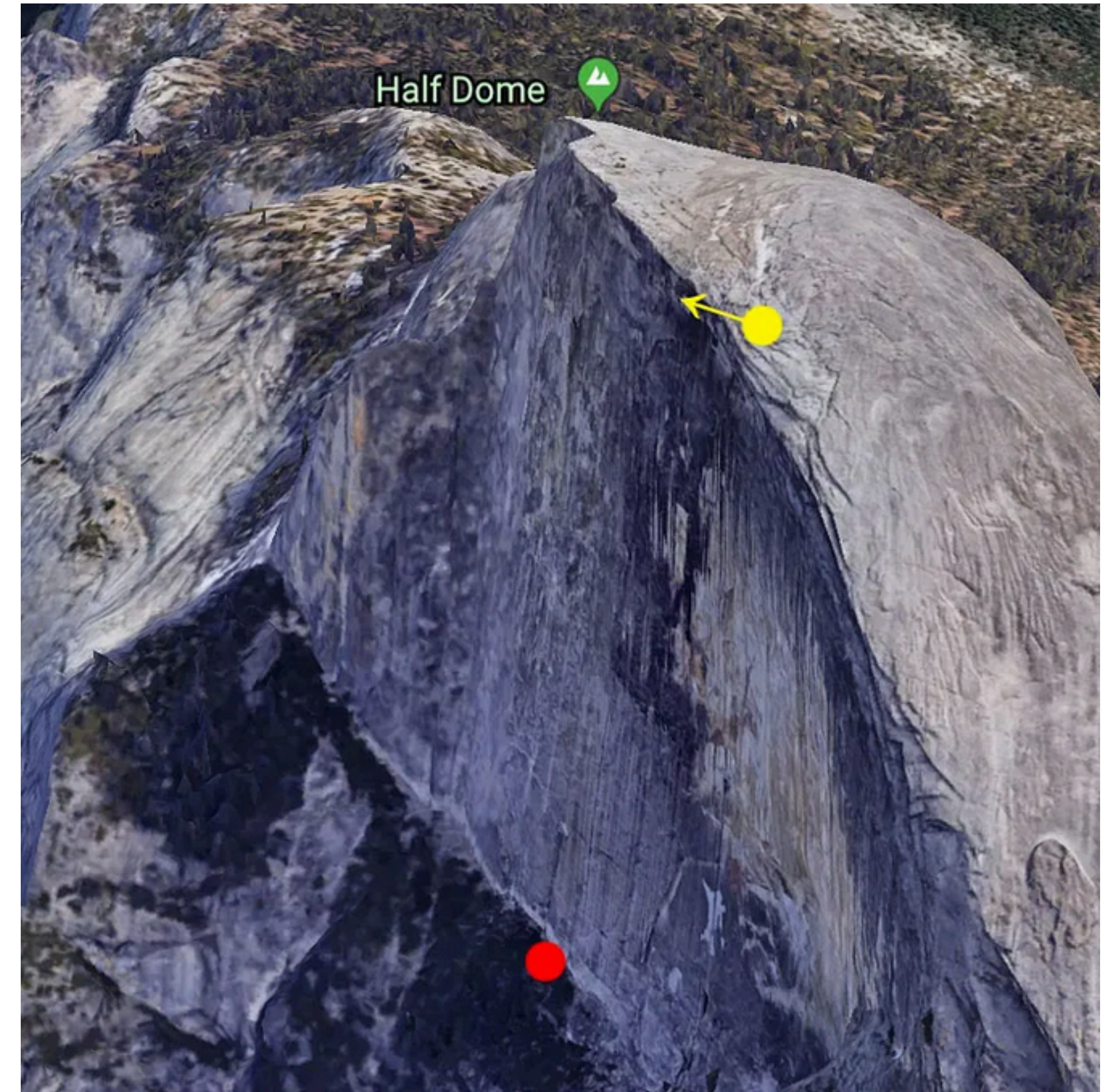
✗ More complex, sometimes annoying in practice

CS 690: Human-Centric Machine Learning

Prof. Scott Niekum

Trust region methods

First-order optimization can be dangerous



Trust regions



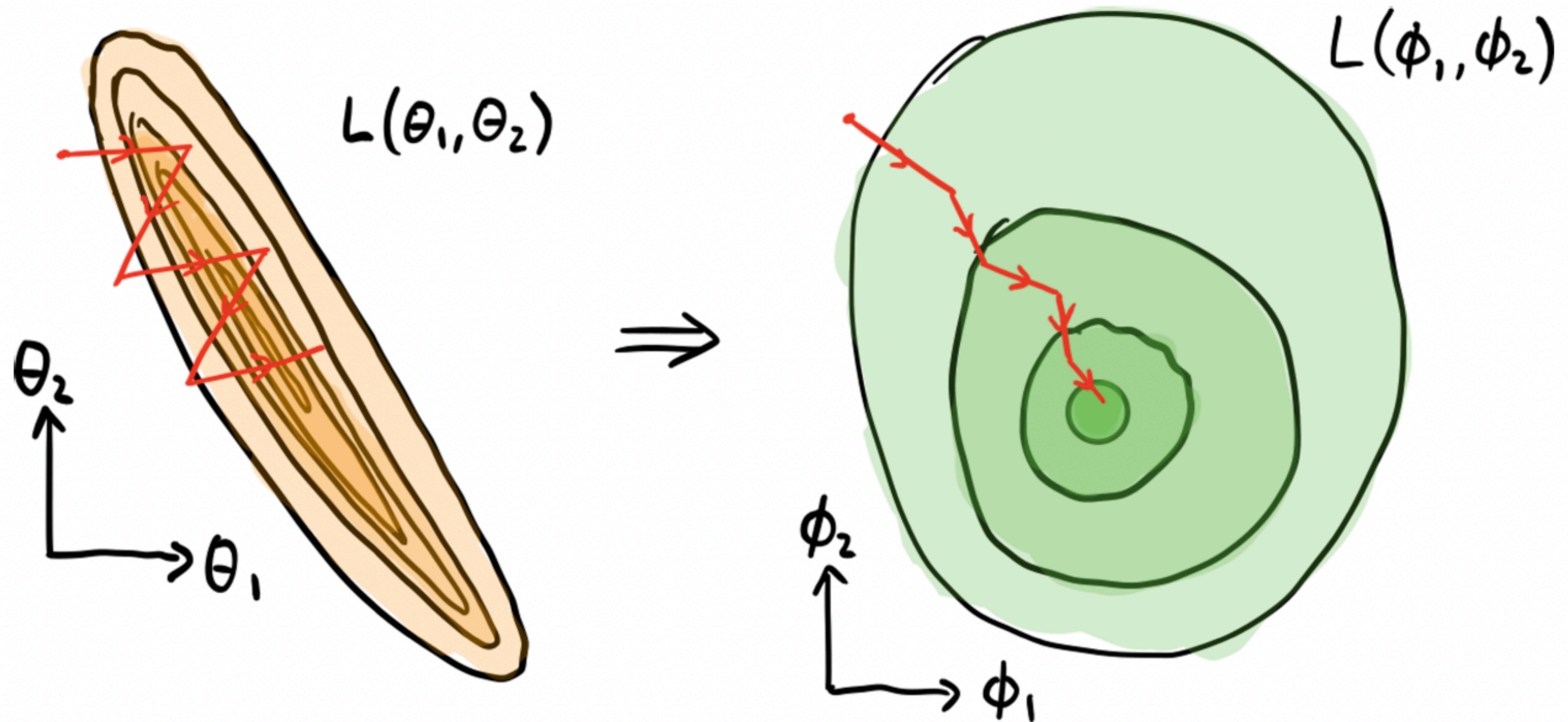
Line search
(like gradient ascent)



Trust region

Image credit: https://medium.com/@jonathan_hui/rl-trust-region-policy-optimization-trpo-explained-a6ee04e04e9

Sensitivity to parameterization

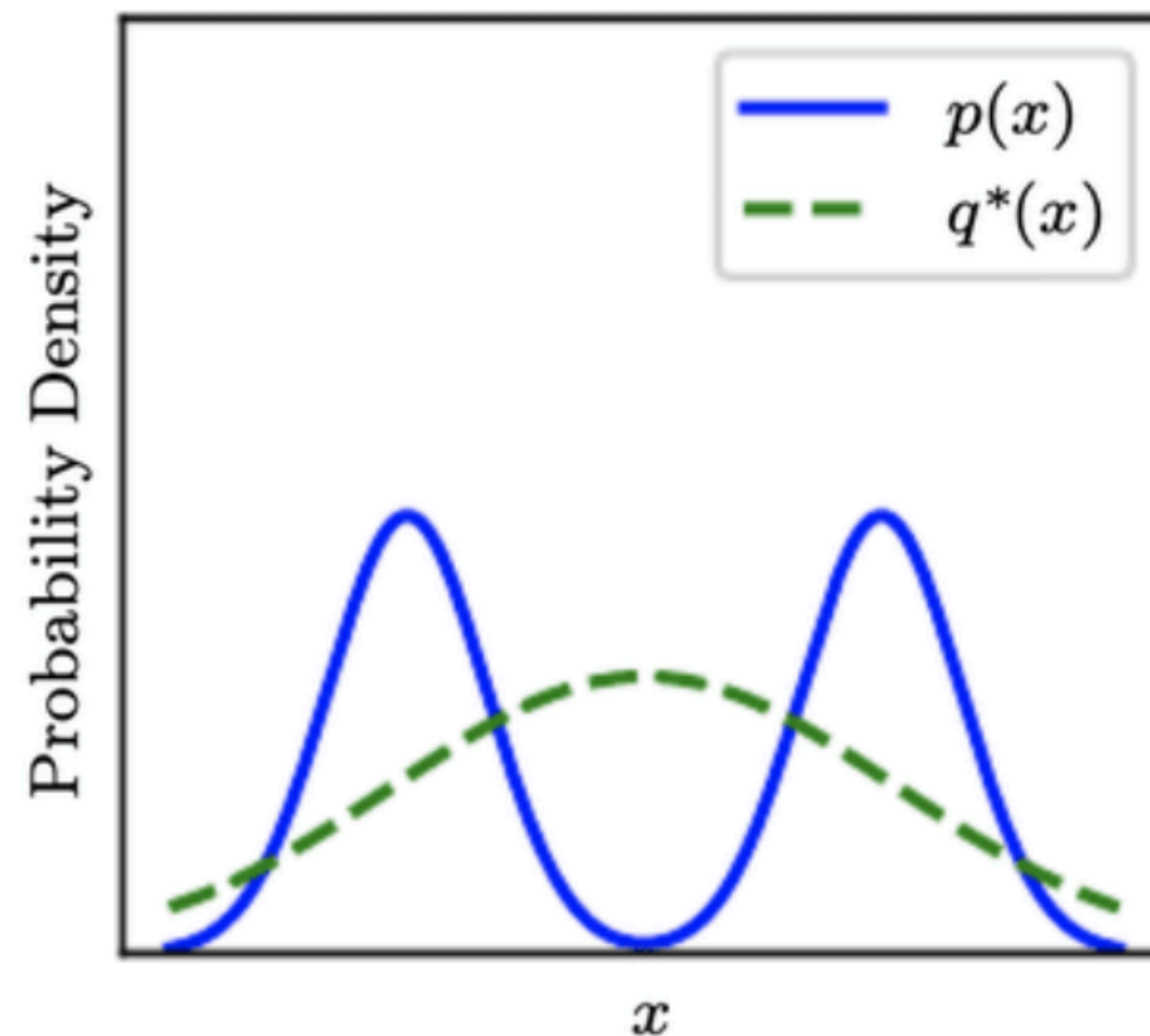


KL Divergence

$$D_{KL}(P \parallel Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(X)}{Q(X)} \right]$$

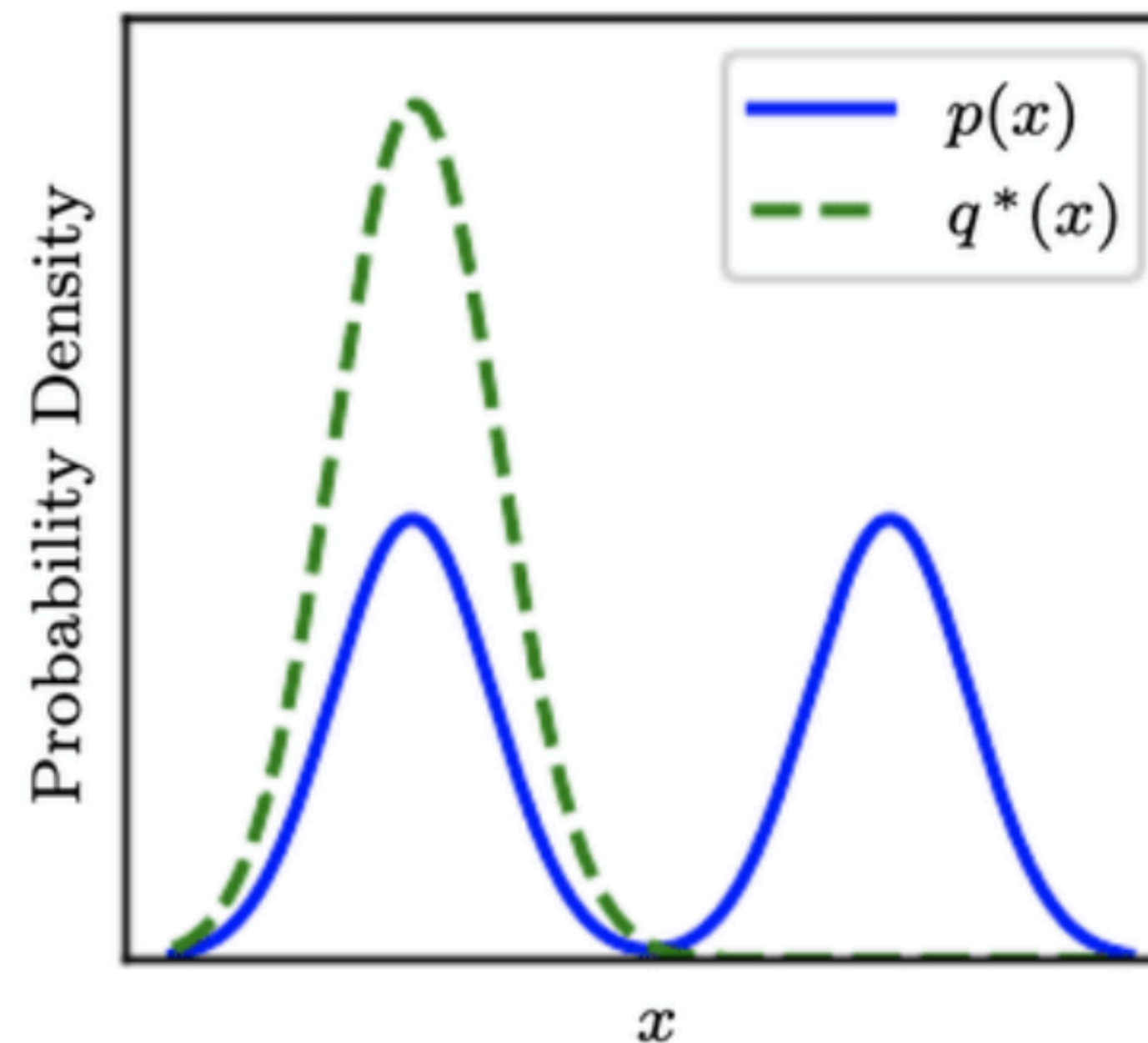
Forward KL
Mean-seeking

$$q^* = \operatorname{argmin}_q D_{KL}(p \parallel q)$$



Reverse KL
Mode-seeking

$$q^* = \operatorname{argmin}_q D_{KL}(q \parallel p)$$



TRPO basic idea

Algorithm 1 Policy iteration algorithm guaranteeing non-decreasing expected return η

Initialize π_0 .

for $i = 0, 1, 2, \dots$ until convergence **do**

 Compute all advantage values $A_{\pi_i}(s, a)$.

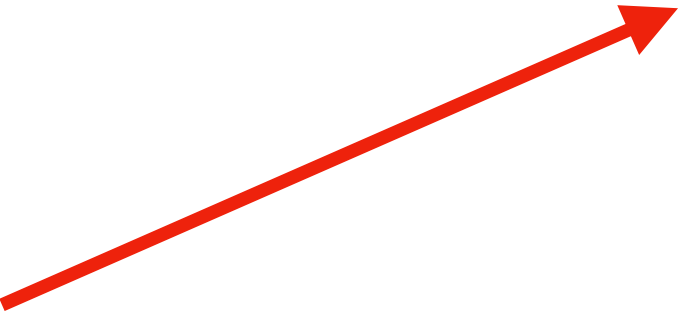
 Solve the constrained optimization problem


$$\pi_{i+1} = \arg \max_{\pi} [L_{\pi_i}(\pi) - C D_{\text{KL}}^{\max}(\pi_i, \pi)]$$


$$\text{where } C = 4\epsilon\gamma/(1 - \gamma)^2$$

$$\text{and } L_{\pi_i}(\pi) = \eta(\pi_i) + \sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a)$$

end for


$$\begin{aligned} & \underset{\theta}{\text{maximize}} \quad L_{\theta_{\text{old}}}(\theta) \\ & \text{subject to} \quad D_{\text{KL}}^{\max}(\theta_{\text{old}}, \theta) \leq \delta. \end{aligned}$$


$$\begin{aligned} & \underset{\theta}{\text{maximize}} \quad L_{\theta_{\text{old}}}(\theta) \\ & \text{subject to} \quad \boxed{\bar{D}_{\text{KL}}^{\rho_{\theta_{\text{old}}}}}(\theta_{\text{old}}, \theta) \leq \delta. \end{aligned}$$



Quadratic approximation:
Fisher information matrix

PPO

On-policy actor-critic

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{s, a \sim \pi_{\theta_k}} [L(s, a, \theta_k, \theta)],$$

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), g(\epsilon, A^{\pi_{\theta_k}}(s, a)) \right),$$

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0. \end{cases}$$

Soft actor-critic

Off-policy actor-critic

- PPO can't learn from ***offline data*** or reuse data, since it is on policy
- PPO policies tend to get more deterministic over time, leading exploration to collapse and learning to stagnate
- SAC can reuse past experience or offline data since it is off policy
- SAC uses entropy maximization to learn the most random policy possible that performs well, leading to natural exploration and robustness to estimation errors