# Deep Reinforcement Learning based on Policy Gradients

Ângelo G. Lovatto

`angelo.lovatto@usp.br`

June 27, 2019

**INSTITUTO DE MATEMÁTICA E ESTATÍSTICA**
UNIVERSIDADE DE SÃO PAULO

# Table of contents

# Policy Gradients Introduction

## Setup

Environment as a Markov Decision Process (MDP) [7]: $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \rho_0)$.

- $\mathcal{S}$ is the set of possible states;
- $\mathcal{A}$ is the set of possible actions;
- $\mathcal{P}(s' \mid s, a)$ is the next state distribution;
- $\mathcal{R}(s, a, s')$ is the reward function;
- $\rho_0(s)$ is the initial state distribution.

## Setup

A *trajectory* is a state-action sequence $\tau = (s_0, a_0, s_1 \ldots, a_{T-1}, s_T)$, whose probability under a stochastic policy $\pi$ is given by

$$P(\tau \mid \pi) \doteq \rho_0(s_0) \prod_{t=0}^{T-1} \pi(a_t \mid s_t) \mathcal{P}(s_{t+1} \mid s_t, a_t),$$

and its undiscounted *return* is given by

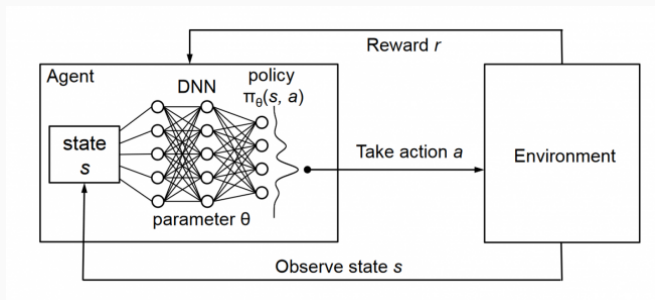$$R(\tau) \doteq \sum_{t=0}^{T-1} \mathcal{R}(s_t, a_t, s_{t+1})$$

The on-policy value functions are defined as usual:

$$V^\pi(s) \doteq \mathop{\mathbb{E}}_{\tau \sim \pi} [R(\tau) \mid s_0 = s], \quad Q^\pi(s, a) \doteq \mathop{\mathbb{E}}_{\tau \sim \pi} [R(\tau) \mid s_0 = s, a_0 = a].$$

3

We consider a smoothly parameterized policy class:

- fix a set of parametric stochastic policies $\{\pi_\theta; \theta \in \mathbb{R}^n\}$
- we assume $\nabla_\theta \pi_\theta(a \mid s), (s, a) \in \mathcal{S} \times \mathcal{A}$, exists and is finite

## Policy Gradient Methods - Problem

We'd like to find the parameters which maximize the expected return

$$\theta^* = \arg\max_\theta J(\theta), \text{ where } J(\theta) = \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} [R(\tau)].$$

Policy gradient methods optimize the policy using gradient ascent, e.g.,

$$\theta_{k+1} = \theta_k + \alpha \nabla J(\theta_k).$$

Problems:

- how to derive the gradient from the expectation?
- how to obtain sample estimates of the gradient?

## Score Function Estimator

Suppose that *x* is a random variable, *f* is a function, and we are interested in computing

$$\frac{\partial}{\partial \theta} \mathbb{E}_x [f(x)], \quad x \sim p(\cdot; \theta).$$

The *score function* (SF) estimator of the equation above is derived as follows:

$$\frac{\partial}{\partial \theta} \mathbb{E}_x [f(x)] = \frac{\partial}{\partial \theta} \int p(x; \theta) f(x) dx = \int \frac{\partial}{\partial \theta} p(x; \theta) f(x) dx$$

$$= \int p(x; \theta) \frac{\partial}{\partial \theta} \log p(x; \theta) f(x) dx = \mathbb{E}_x \left[ f(x) \frac{\partial}{\partial \theta} \log p(x; \theta) \right].$$

## Policy Gradient Derivation

The score function estimator for $J(\theta)$ is then

$$\nabla J(\theta) = \nabla_\theta \operatorname*{\mathbb{E}}_{\tau \sim \pi_\theta} [R(\tau)]$$

$$= \operatorname*{\mathbb{E}}_{\tau \sim \pi_\theta} [R(\tau) \nabla_\theta \log P(\tau \mid \theta)] ,$$

where

$$\nabla_\theta \log P(\tau \mid \theta) = \nabla_\theta \left( \rho_0(s_0) + \sum_{t=0}^{T-1} \big( \log \pi_\theta(a_t \mid s_t) + \log \mathcal{P}(s_{t+1} \mid s_t, a_t) \big) \right)$$

$$= \sum_{t=0}^{T-1} \underbrace{\nabla_\theta \log \pi_\theta(a_t \mid s_t)}_{\text{no dynamics model needed!}} .$$

Thus,

$$\nabla J(\theta) = \operatorname*{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} R(\tau) \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]$$

We can then estimate the gradient by averaging across sampled trajectories:

$$\nabla J(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^{m} \sum_{t=0}^{T-1} R(\tau_i) \nabla_\theta \log \pi_\theta(a_{i,t} \mid s_{i,t}).$$

Intuitively, taking a step in the gradient direction

- increases the probability of trajectories with positive $R(\tau)$,
- decreases the probability of trajectories with negative $R(\tau)$.

We can remove terms that don't depend on current action:

$$
\mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} R(\tau) \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]
$$

$$
= \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t \mid s_t) \sum_{t'=0}^{T-1} \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \right]
$$

$$
= \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t \mid s_t) \left( \sum_{i=0}^{t-1} \mathcal{R}(s_{i+1} \mid s_i, a_i) + \sum_{j=t}^{T-1} \mathcal{R}(s_{j+1} \mid s_j, a_j) \right) \right]
$$

$$
= \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t \mid s_t) \underbrace{\sum_{t'=t}^{T-1} \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'})}_{\text{reward-to-go}} \right]
$$

## Reducing Variance: Adding a Baseline

Furthermore, a state-dependent function can be subtracted

$$g \approx \frac{1}{m} \sum_{i=1}^{m} \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_{i,t} \,|\, s_{i,t}) \left( \sum_{t'=t}^{T-1} \mathcal{R}(s_{i,t'+1} \,|\, s_{i,t'}, a_{i,t'}) - b(s_{i,t}) \right).$$

which is also known as a *control variate*. It can be shown that

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t \,|\, s_t) b(s_t) \right] = 0,$$

thus, the estimate is still unbiased. This form of is also known as the REINFORCE estimator.

# Vanilla Policy Gradient and GAE

## General Policy Gradient

In general, the policy gradient can be expressed as

$$g = \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \Phi_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]$$

where $\Phi_t$ may be one of the following.

- $R(\tau)$
- $\sum_{k=t}^{T-1} \mathcal{R}(s_{k+1} \mid s_k, a_k)$
- $\sum_{k=t}^{T-1} \mathcal{R}(s_{k+1} \mid s_k, a_k) - b(s_t)$

- $Q^{\pi_\theta}(s_t, a_t)$
- $A^{\pi_\theta}(s_t, a_t)$
- $r_{t+1} + V^{\pi_\theta}(s_{t+1}) - V^{\pi_\theta}(s_t)$

The advantage function, $A^{\pi_\theta}(s_t, a_t) = Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)$, yields very low variance in practice.

$$g = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} A^{\pi_\theta}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]$$

Intuitively, taking a step in the direction of this gradient:

- increases the probability of better-than-average actions;
- decreases the probability of worse-than-average actions.

We don't have access to the true value functions and sampled rewards have high variance. We can address this in two ways:

- using discounted returns
- using a learned value function (critic)

Consider discounted variants of the value functions

$$V^{\pi,\gamma}(s) \doteq \mathop{\mathbb{E}}_{\tau \sim \pi} \left[ R^\gamma(\tau) \,|\, s_0 = s \right],$$

$$Q^{\pi,\gamma}(s,a) \doteq \mathop{\mathbb{E}}_{\tau \sim \pi} \left[ R^\gamma(\tau) \,|\, s_0 = s, a_0 = a \right],$$

$$A^{\pi,\gamma}(s,a) = Q^{\pi,\gamma}(s,a) - V^{\pi,\gamma}(s),$$

where $R^\gamma(\tau) \doteq \sum_{t=0}^{T-1} \gamma^t \mathcal{R}(s_t, a_t, s_{t+1})$.

Consider discounted variants of the value functions

$$V^{\pi,\gamma}(s) \doteq \mathop{\mathbb{E}}_{\tau\sim\pi}\left[R^\gamma(\tau)\,|\,s_0 = s\right],$$

$$Q^{\pi,\gamma}(s,a) \doteq \mathop{\mathbb{E}}_{\tau\sim\pi}\left[R^\gamma(\tau)\,|\,s_0 = s, a_0 = a\right],$$

$$A^{\pi,\gamma}(s,a) = Q^{\pi,\gamma}(s,a) - V^{\pi,\gamma}(s),$$

where $R^\gamma(\tau) \doteq \sum_{t=0}^{T-1}\gamma^t\mathcal{R}(s_t,a_t,s_{t+1})$.

Note that $A^{\pi,\gamma}(s,a)$ is equivalent to the *TD residual,*

$$A^{\pi,\gamma}(s,a) = \mathbb{E}\left[r_{t+1} + \gamma V^{\pi,\gamma}(s_{t+1}) - V^{\pi,\gamma}(s_t)\,|\,s_t = s, a_t = a\right].$$

In fact, any n-step TD residual is valid.

$$
\begin{aligned}
A^{\pi,\gamma}&(s, a) \\
&= \mathbb{E}\left[r_{t+1} + \gamma V^{\pi,\gamma}(s_{t+1}) - V^{\pi,\gamma}(s_t) \,|\, s_t = s, a_t = a\right] \\
&= \mathbb{E}\left[r_{t+1} + \gamma r_{t+2} + \gamma^2 V^{\pi,\gamma}(s_{t+2}) - V^{\pi,\gamma}(s_t) \,|\, s_t = s, a_t = a\right] \\
&= \mathbb{E}\left[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V^{\pi,\gamma}(s_{t+3}) - V^{\pi,\gamma}(s_t) \,|\, s_t = s, a_t = a\right] \\
&\;\;\vdots \\
&= \mathbb{E}\left[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots + \gamma^{T-1} r_T - V^{\pi,\gamma}(s_t) \,|\, s_t = s, a_t = a\right]
\end{aligned}
$$

In practice, we use sample estimates and learned critics

$$\hat{A}_t^{(1)} = r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$$
$$\hat{A}_t^{(2)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 V_\phi(s_{t+2}) - V_\phi(s_t)$$
$$\hat{A}_t^{(3)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V_\phi(s_{t+3}) - V_\phi(s_t)$$
$$\vdots$$
$$\hat{A}_t^{(T)} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots + \gamma^{T-1} r_T - V_\phi(s_t)$$

Analogous to the n-step return $G_t^{(n)}$ used in *TD* methods.

The generalized advantage estimator is defined as follows

$$
\begin{aligned}
\hat{A}_t^{\text{GAE}(\gamma,\lambda)} \doteq\ & (r_{t+1} + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)) && (1-\lambda) \\
& + (r_{t+1} + \gamma r_{t+2} + \gamma^2 V_\phi(s_{t+2}) - V_\phi(s_t)) && (1-\lambda)\lambda \\
& + (r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V_\phi(s_{t+3}) - V_\phi(s_t)) && (1-\lambda)\lambda^2 \\
& \ \vdots \\
& + (r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots + \gamma^{T-1} r_T - V_\phi(s_t)) && \lambda^{T-t-1}
\end{aligned}
$$

Analogous to the lambda return $G_t^\lambda$ used in $TD(\lambda)$ (forward-view).

It turns out that $\hat{A}_t^{GAE(\gamma,\lambda)}$ can be computed as follows

$$\hat{A}_t^{GAE(\gamma,\lambda)} = \sum_{l=0}^{T-1} (\gamma\lambda)^l \delta_{t+l}^V, \text{ where } \delta_t^V = r_{t+1} + \gamma V(s_{t+1}) - V(s_t).$$

We can recover the one-step TD estimator and the monte carlo by setting $\lambda$ to 0 and 1 respectively:

$$\hat{A}_t^{GAE(\gamma,0)} = \delta_t^V \qquad = r_{t+1} + \gamma V(s_{t+1}) - V(s_t),$$
$$\hat{A}_t^{GAE(\gamma,1)} = \sum_{l=0}^{T-1} \gamma^l \delta_{t+l}^V = R_t^\gamma - V(s_t).$$

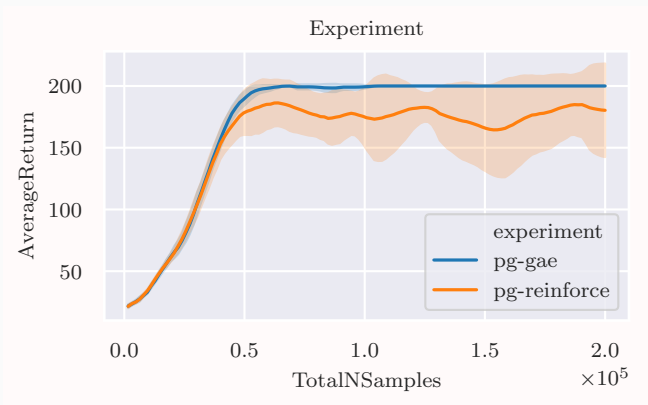As in $TD(\lambda)$, $\lambda$ makes a compromise between bias and variance.

**Algorithm 1** Vanilla Policy Gradient + GAE

1: Initialize policy parameter $\theta_0$ and value function parameter $\phi_0$
2: **for** k=0, 1, 2, . . . **do**
3:      Collect $N$ transitions $(s_t, a_t, r_{t+1})$ following $\pi_{\theta_k}$
4:      Compute $\delta_t^V$ at all timesteps $t \in \{1, 2, \ldots, N\}$
5:      Compute $\hat{A}_t = \sum_{l=0}^{T-1} (\gamma\lambda)^l \delta_{t+l}^V$ for every timestep
6:      Fit baseline by minimizing $\|V_\phi(s_t) - R_t^\gamma\|^2$
7:      Compute $\theta_{k+1}$ with VPG update
8: **end for**

**Figure 1:** Learning curves in the CartPole-v0 environment, where REINFORCE uses a baseline of zero and VPG uses GAE(0.99,0.97)

# Natural Policy Gradient

# Problems with the Vanilla Gradient

$$g = \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} A^{\pi_\theta}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]$$

- The data distribution depends on the parameters $\theta$
- Most optimizers work well in settings with i.i.d. data
- Hard to choose step size $\alpha$
    - too big $\rightarrow$ bad policy $\rightarrow$ data collected under bad policy $\rightarrow$ collapse in performance
    - too small: inefficient use of data collected

Consider a family of policies with parametrization:

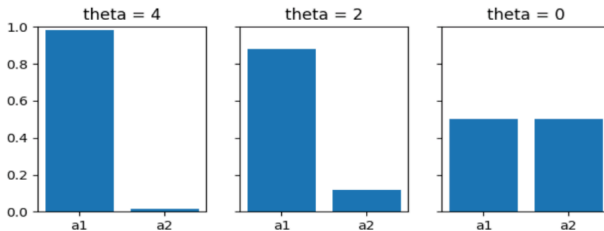$$\pi_\theta(a) = \begin{cases} \sigma(\theta) & a = 1 \\ 1 - \sigma(\theta) & a = 2 \end{cases}$$



Figure: Small changes in the policy parameters can unexpectedly lead to **big** changes in the policy.

Figure source: `http://www.andrew.cmu.edu/course/10-703/`

The step size in gradient ascent comes from solving the following optimization problem, e.g., using line search

$$d^* = \arg\max_d J(\theta + d)$$

$$\text{s.t.} \quad \|d\| \leq \epsilon$$

However, the distance in parameters, as we've seen, gives no insight into the resulting policy change.

The step size in natural gradient ascent comes from solving the following optimization problem,

$$d^* = \arg \max_d J(\theta + d)$$

$$\text{s.t.} \quad D(\pi_\theta \| \pi_{\theta+d}) \leq \epsilon$$

It is be better to define the constraint in terms of a "distance" between successive policies.

The Kullback-Leibler (KL) divergence provides a way to measure dissimilarities between distributions.

$$D_{\text{KL}}\left(P \,\|\, Q\right) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

$$D_{\text{KL}}\left(p \,\|\, q\right) = \int_{-\infty}^{+\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

We define the KL divergence between policies as follows.

$$\overline{D}_{\text{KL}}\left(\theta_k \,\|\, \theta\right) \doteq \mathop{\mathbb{E}}_{s \sim \pi_{\theta_k}} \left[D_{\text{KL}}\left(\pi_{\theta_k}(\cdot \,|\, s) \,\|\, \pi_\theta(\cdot \,|\, s)\right)\right].$$

The objective then becomes

$$\theta_{k+1} = \arg\max_{\theta} J(\theta_k)$$
$$\text{s.t. } \overline{D}_{KL}\left(\theta_k \,\|\, \theta\right) \leq \delta,$$

- Easier to pick the distance threshold
- Invariant to the parameterization of the policy
- Hard to calculate in practice

Using first- and second-order Taylor expansions for the objective and constraint respectively yields

$$\theta_{k+1} = \arg\max_\theta J(\theta_k) + g^\mathsf{T}(\theta - \theta_k)$$
$$\text{s.t.} \quad \tfrac{1}{2}(\theta - \theta_k)^\mathsf{T} F(\theta_k)(\theta - \theta_k) \leq \delta.$$

where $F(\theta_k)$ denotes the Hessian of the KL divergence, also known as the Fisher Information Matrix,

$$F(\theta_k) = \nabla_\theta^2 \overline{D}_{\mathsf{KL}}\left(\theta_k \,\|\, \theta\right)|_{\theta=\theta_k}.$$

A Lagrange multiplier argument gives the solution to this constrained optimization problem:

$$\theta_{k+1} = \theta_k + \underbrace{\sqrt{\frac{2\delta}{g^\mathsf{T} F(\theta_k)^{-1} g}}}_{\text{scaling term}} \underbrace{F(\theta_k)^{-1} g}_{\text{npg}},$$

Computing $F(\theta_k)^{-1}g$ can be expensive:

- DNN policy can have thousands of parameters
- Storing and inverting the Hessian becomes impractical

Computing $F(\theta_k)^{-1}g$ can be expensive:

- DNN policy can have thousands of parameters
- Storing and inverting the Hessian becomes impractical
- Can instead solve for $F(\theta_k)x = g$ using conjugate gradient
  - Does not form the whole matrix
  - Requires only a function $f(v) = F(\theta_k)v$

Computing $F(\theta_k)^{-1}g$ can be expensive:

- DNN policy can have thousands of parameters
- Storing and inverting the Hessian becomes impractical
- Can instead solve for $F(\theta_k)x = g$ using conjugate gradient
  - Does not form the whole matrix
  - Requires only a function $f(v) = F(\theta_k)v$
- We can evaluate $F(\theta_k)v$ using automatic differentiation.

$$F(\theta_k)v = \nabla_\theta \left( \left( \nabla_\theta \overline{D}_{\mathrm{KL}} \left( \theta_k \parallel \theta \right) \right)^{\mathsf{T}} v \right) \Big|_{\theta = \theta_k}.$$

# NPG Algorithm

---

**Algorithm 2** Natural Policy Gradient

---

1: Initialize policy parameter $\theta_0$
2: **for** k=0, 1, 2, . . . **do**
3:     Collect a set of trajectories following $\pi_{\theta_k}$
4:     Estimate $\hat{A}_t$'s using any advantage estimation algorithm
5:     Compute the policy gradient $g$
6:     Use CG and Fisher vector products to obtain $F(\theta_k)^{-1}g$
7:     Update the policy parameters

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g^\mathsf{T} F(\theta_k)^{-1} g}} F(\theta_k)^{-1} g,$$

8: **end for**

---

# Trust Region Policy Optimization

Natural policy gradient gives us a good step size, however

- How do we ensure improvement in the objective?
- Can we estimate the performance of a new policy?

$$
\begin{aligned}
J(\tilde{\pi}) - J(\pi) &= J(\tilde{\pi}) - \mathbb{E}_{s_0 \sim \rho_0} \left[ V^\pi(s_0) \right] \\
&= J(\tilde{\pi}) + \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{i=1}^{T-1} V^\pi(s_i) - \sum_{j=0}^{T-1} V^\pi(s_j) \right] \\
&= J(\tilde{\pi}) + \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{T-1} V^\pi(s_{t+1}) - V^\pi(s_t) \right] \\
&= \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{T-1} \mathcal{R}(s_t, a_t) + V^\pi(s_{t+1}) - V^\pi(s_t) \right] \\
&= \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{T-1} A^\pi(s_t, a_t) \right]
\end{aligned}
$$

TRPO makes the following approximation to the result of $J(\tilde{\pi}) - J(\pi)$:

$$\mathop{\mathbb{E}}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{T-1} A^\pi(s_t, a_t) \right] \approx \mathop{\mathbb{E}}_{\tau \sim \pi} \left[ \sum_{t=0}^{T-1} \frac{\tilde{\pi}(a_t \mid s_t)}{\pi(a_t \mid s_t)} A^\pi(s_t, a_t) \right].$$

That is, we approximate performance of a new policy relative to the old one using data gathered by the latter. In the case of parametric policies, the surrogate objective we seek to optimize is

$$\mathcal{L}(\theta_k, \theta) \doteq \mathop{\mathbb{E}}_{\tau \sim \pi_{\theta_k}} \left[ \sum_{t=0}^{T-1} \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_k}(a_t \mid s_t)} A^{\pi_{\theta_k}}(s_t, a_t) \right].$$

## TRPO: Problem

We could optimize $\mathcal{L}(\theta_k, \theta)$ freely, however

- we're substituting the trajectory distribution of the old policy for the new policy
- the approximation might not be accurate if the policies differ too much

It turns out that we can bound this approximation error

$$|J(\theta) - (\underbrace{J(\theta_k) + \mathcal{L}(\theta_k, \theta)}_{\approx J(\theta)})| \le C\sqrt{\overline{D}_{\mathsf{KL}}\left(\theta_k \| \theta\right)}$$

## TRPO: Problem

TRPO iteratively optimizes the surrogate objective function around a local neighborhood of the current policy, as measured by KL divergence:

$$\theta_{k+1} = \arg\max_{\theta} \ \mathbb{E}_{\tau \sim \pi_{\theta_k}} \left[ \sum_{t=0}^{T-1} \frac{\pi_\theta(a_t \mid s_t)}{\pi_{\theta_k}(a_t \mid s_t)} A^{\pi_{\theta_k}}(s_t, a_t) \right]$$

$$\text{s.t.} \ \ \overline{D}_{KL} \left( \theta_k \, \| \, \theta \right) \leq \delta \ .$$

The constraint $\delta$ defines a *trust region* around our local approximation to the new policy's relative performance.

## TRPO: Solution

It can be shown that $\nabla_\theta \mathcal{L}(\theta_k, \theta)$ evaluated at $\theta = \theta_k$ is equivalent to the policy gradient. Thus, the appropriate Taylor expansions give

$$\theta_{k+1} = \arg\max_\theta g^\mathsf{T}(\theta - \theta_k)$$
$$\text{s.t. } \tfrac{1}{2}(\theta - \theta_k)^\mathsf{T} F(\theta_k)(\theta - \theta_k) \leq \delta,$$

yielding the same solution as in NPG. TRPO, however, tries to ensure

- improvement in the objective $J(\theta)$,
- satisfaction of the KL constraint.

Final update:

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{g^\mathsf{T} F(\theta_k)^{-1} g}} F(\theta_k)^{-1} g,$$

where $\alpha \in (0, 1)$ is a backtracking coefficient.

---

**Algorithm 3** Line Search for TRPO

---

1: Compute proposed policy step $\Delta\theta = \sqrt{\frac{2\delta}{g^\mathsf{T} F(\theta_k)^{-1} g}} F(\theta_k)^{-1} g$
2: **for** $i = 0, 1, 2, \ldots, L$ **do**
3:      Compute proposed update $\theta = \theta_k + \alpha^i \Delta\theta$
4:      **if** $\mathcal{L}(\theta_k, \theta) \geq 0$ and $\overline{D}_{\mathsf{KL}}(\theta_k \,\|\, \theta) \leq \delta$ **then**
5:          accept the update and set $\theta_{k+1} = \theta_k + \alpha^i \Delta\theta$
6:          break
7:      **end if**
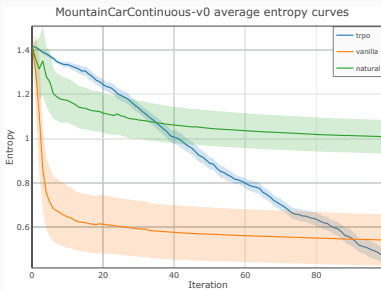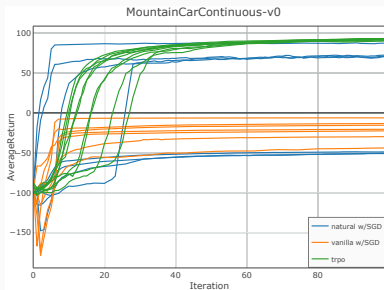8: **end for**

---

## TRPO: Algorithm

**Algorithm 4** TRPO

1: Initialize policy parameter $\theta_0$
2: **for** k=0, 1, 2, . . . **do**
3:     Collect a set of trajectories following $\pi_{\theta_k}$
4:     Estimate $\hat{A}_t$'s using any advantage estimation algorithm
5:     Compute the policy gradient $g$
6:     Use CG and Fisher vector products to obtain $F(\theta_k)^{-1}g$
7:     Compute the proposed update $\Delta\theta = \sqrt{\frac{2\delta}{g^{\intercal}F(\theta_k)^{-1}g}}F(\theta_k)^{-1}g$,
8:     Perform backtracking line search to obtain final update

$$\theta_{k+1} = \theta_k + \alpha^i \Delta\theta$$

9: **end for**

# Comparisons

Questions?

Recall that

$$
\mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} R(\tau) \nabla_\theta \log \pi_\theta(a_t \mid s_t) \right]
$$

$$
= \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t \mid s_t) \sum_{t'=0}^{T-1} \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \right]
$$

$$
= \sum_{t=0}^{T-1} \sum_{t'=0}^{T-1} \mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t \mid s_t) \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \right].
$$

We'll show that, for $t' < t$ (when the reward comes before the action),

$$
\mathop{\mathbb{E}}_{\tau \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t \mid s_t) \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \right] = 0.
$$

Let $f(t, t') = \nabla_\theta \log \pi_\theta(a_t \mid s_t) \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'})$. Note that this term only depends on timesteps $t, t'$ and $t' + 1$. Thus, we can use the *marginal likelihood*

$$
\begin{aligned}
\mathbb{E}_{\tau \sim \pi_\theta} [f(t, t')] &= \int_\tau P(\tau \mid \theta) f(t, t') \\
&= \int_{s_t, a_t, s_{t'}, a_{t'}, s_{t'+1}} P(s_t, a_t, s_{t'}, a_{t'}, s_{t'+1} \mid \theta) f(t, t') \\
&= \mathbb{E}_{s_t, a_t, s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta} [f(t, t')]
\end{aligned}
$$

We can break the expectation using the chain rule of probability:

$$\underset{s_t, a_t, s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta}{\mathbb{E}} [f(t, t')] = \underset{s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta}{\mathbb{E}} \left[ \underset{s_t, a_t \sim \pi_\theta}{\mathbb{E}} [f(t, t') \mid s_{t'}, a_{t'}, s_{t'+1}] \right].$$

Substituting $f(t, t')$ in the equation above and noting that the reward term is constant w.r.t. the inner expectation:

$$\underset{s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta}{\mathbb{E}} \left[ \underset{s_t, a_t \sim \pi_\theta}{\mathbb{E}} [f(t, t') \mid s_{t'}, a_{t'}, s_{t'+1}] \right]$$

$$= \underset{s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta}{\mathbb{E}} \left[ \underset{s_t, a_t \sim \pi_\theta}{\mathbb{E}} [\nabla_\theta \log \pi_\theta(a_t \mid s_t) \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \mid s_{t'}, a_{t'}, s_{t'+1}] \right]$$

$$= \underset{s_{t'}, a_{t'}, s_{t'+1} \sim \pi_\theta}{\mathbb{E}} \left[ \mathcal{R}(s_{t'+1} \mid s_{t'}, a_{t'}) \underset{s_t, a_t \sim \pi_\theta}{\mathbb{E}} [\nabla_\theta \log \pi_\theta(a_t \mid s_t) \mid s_{t'}, a_{t'}, s_{t'+1}] \right].$$

Finally, since $t' < t$ we can break the inner expectation further using the chain rule of probability

$$\mathbb{E}_{s_t, a_t \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t \,|\, s_t) \,|\, s_{t'}, a_{t'}, s_{t'+1} \right]$$

$$= \mathbb{E}_{s_t \sim \pi_\theta} \left[ \mathbb{E}_{a_t \sim \pi_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t \,|\, s_t) \,|\, s_t \right] \Big| s_{t'}, a_{t'}, s_{t'+1} \right]$$

$$= 0$$

Which follows since for any* probability distribution,

$$\int_x P_\theta(x) \nabla_\theta \log P_\theta(x) = \int_x \nabla_\theta P_\theta(x) = \nabla_\theta \underbrace{\int_x P_\theta(x)}_{=1} = 0.$$

# Second-Order Taylor Expansion of KL

$$D_{\text{KL}}\left(\theta' \parallel \theta\right) \approx D_{\text{KL}}\left(\theta' \parallel \theta'\right) + \left(\nabla_\theta D_{\text{KL}}\left(\theta' \parallel \theta\right)|_{\theta=\theta'}\right)^{\mathsf{T}}\left(\theta - \theta'\right)$$
$$+ \tfrac{1}{2}(\theta - \theta')^{\mathsf{T}} H(\theta - \theta'),$$

where

$$\nabla_\theta D_{\text{KL}}\left(\theta' \parallel \theta\right)|_{\theta=\theta'} = -\int_X p_{\theta'}(x)\frac{\nabla_\theta p_\theta(x)|_{\theta=\theta'}}{p_{\theta'}(x)}$$

$$= -\int_X p_{\theta'}(x)\nabla_{\theta'}\log p_{\theta'}(x)dx = -\int_X \nabla_{\theta'}p_{\theta'}(x)dx = 0$$

and

$$H = \nabla_\theta^2 D_{\text{KL}}\left(\theta' \parallel \theta\right)|_{\theta=\theta'}$$

$$= -\int_X p_{\theta'}(x)\nabla_{\theta'}^2 \log p_{\theta'}(x)dx$$

$$= -\mathbb{E}_{p_{\theta'}}[\nabla_\theta^2 \log p_{\theta'}(x)],$$

J. Achiam, D. Held, A. Tamar, and P. Abbeel.
**Constrained policy optimization.**
In *ICML*, volume 70 of *Proceedings of Machine Learning Research*, pages 22–31. PMLR, 2017.

S. Kakade and J. Langford.
**Approximately optimal approximate reinforcement learning.**
In *ICML*, pages 267–274. Morgan Kaufmann, 2002.

S. M. Kakade.
**A natural policy gradient.**
In *Advances in neural information processing systems*, pages 1531–1538, 2002.

H. Mao, M. Alizadeh, I. Menache, and S. Kandula.
**Resource management with deep reinforcement learning.**
In *Proceedings of the 15th ACM Workshop on Hot Topics in Networks*, pages 50–56. ACM, 2016.

J. Schulman, N. Heess, T. Weber, and P. Abbeel.
**Gradient estimation using stochastic computation graphs.**
In *NIPS*, pages 3528–3536, 2015.

J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel.
**High-dimensional continuous control using generalized advantage estimation.**
In *ICLR*, 2016.

R. S. Sutton and A. G. Barto.
*Reinforcement learning: An introduction.*
MIT press, 2018.