Trust Region Policy Optimization

Desiderata for policy optimization method:

Stable, monotonic improvement. (How to choose stepsizes?)

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

Good sample efficiency

Step Sizes

Why are step sizes a big deal in RL?

- Supervised learning
 - Step too far \rightarrow next updates will fix it
- Reinforcement learning
 - Step too far \rightarrow bad policy
 - Next batch: collected under bad policy
 - Can't recover, collapse in performance!



Surrogate Objective

- Let $\eta(\pi)$ denote the expected return of π
- \blacktriangleright We collect data with $\pi_{old}.$ Want to optimize some objective to get a new policy π
- Define $L_{\pi_{\mathrm{old}}}(\pi)$ to be the "surrogate objective" ¹

$$egin{aligned} \mathcal{L}(\pi) &= \mathbb{E}_{\pi_{\mathrm{old}}}\left[rac{\pi(a \mid s)}{\pi_{\mathrm{old}}(a \mid s)}\mathcal{A}^{\pi_{\mathrm{old}}}(s, a)
ight] \ &
abla_{ heta}\mathcal{L}(\pi_{ heta})ig|_{ heta_{\mathrm{old}}} &=
abla_{ heta}\eta(\pi_{ heta})ig|_{ heta_{\mathrm{old}}} & (ext{policy gradient}) \end{aligned}$$

 \blacktriangleright Local approximation to the performance of the policy; does not depend on parameterization of π

¹S. Kakade and J. Langford. "Approximately optimal approximate reinforcement learning". In: ICML. vol. 2. (2002, opp 267–274) 👍 🗄 🗧 🕤 🔍 🗠

Improvement Theory

- Theory: bound the difference between L_{πold}(π) and η(π), the performance of the policy
- ► Result: $\eta(\pi) \ge L_{\pi_{\text{old}}}(\pi) C \cdot \max_{s} \text{KL}[\pi_{\text{old}}(\cdot | s), \pi(\cdot | s)]$, where $c = 2\epsilon\gamma/(1-\gamma)^2$
- Monotonic improvement guaranteed (MM algorithm)



Practical Algorithm: TRPO

Constrained optimization problem

$$\max_{\pi} L(\pi), \text{ subject to } \overline{\mathsf{KL}}[\pi_{\mathrm{old}}, \pi] \leq \delta$$

where $L(\pi) = \mathbb{E}_{\pi_{\mathrm{old}}} \left[\frac{\pi(a \mid s)}{\pi_{\mathrm{old}}(a \mid s)} A^{\pi_{\mathrm{old}}}(s, a) \right]$

Construct loss from empirical data

$$\hat{L}(\pi) = \sum_{n=1}^{N} rac{\pi(a_n \mid s_n)}{\pi_{ ext{old}}(a_n \mid s_n)} \hat{A}_n$$

Make quadratic approximation and solve with conjugate gradient algorithm

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 悪 = のへで

J. Schulman, S. Levine, P. Moritz, et al. "Trust Region Policy Optimization". In: ICML. 2015

Practical Algorithm: TRPO

for iteration= $1, 2, \ldots$ do

Run policy for T timesteps or N trajectories

Estimate advantage function at all timesteps

Compute policy gradient g

Use CG (with Hessian-vector products) to compute $F^{-1}g$

Do line search on surrogate loss and KL constraint

end for

Practical Algorithm: TRPO

Applied to

► Locomotion controllers in 2D



Atari games with pixel input

"Proximal" Policy Optimization

Use penalty instead of constraint

$$\underset{\theta}{\mathsf{minimize}} \sum_{n=1}^{N} \frac{\pi_{\theta}(a_n \mid s_n)}{\pi_{\theta_{\text{old}}}(a_n \mid s_n)} \hat{A}_n - \beta \overline{\mathsf{KL}}[\pi_{\theta_{\text{old}}}, \pi_{\theta}]$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

"Proximal" Policy Optimization

Use penalty instead of constraint

$$\operatorname{minimize}_{\theta} \sum_{n=1}^{N} \frac{\pi_{\theta}(a_n \mid s_n)}{\pi_{\theta_{\mathrm{old}}}(a_n \mid s_n)} \hat{A}_n - \beta \overline{\mathsf{KL}}[\pi_{\theta_{\mathrm{old}}}, \pi_{\theta}]$$

Pseudocode:

for iteration= $1, 2, \dots$ do

Run policy for T timesteps or N trajectories Estimate advantage function at all timesteps Do SGD on above objective for some number of epochs If KL too high, increase β . If KL too low, decrease β . end for

"Proximal" Policy Optimization

Use penalty instead of constraint

$$\operatorname{minimize}_{\theta} \sum_{n=1}^{N} \frac{\pi_{\theta}(a_n \mid s_n)}{\pi_{\theta_{\mathrm{old}}}(a_n \mid s_n)} \hat{A}_n - \beta \overline{\mathsf{KL}}[\pi_{\theta_{\mathrm{old}}}, \pi_{\theta}]$$

Pseudocode:

for iteration= $1, 2, \ldots$ do

Run policy for T timesteps or N trajectories Estimate advantage function at all timesteps Do SGD on above objective for some number of epochs If KL too high, increase β . If KL too low, decrease β . end for

ightarrow same performance as TRPO, but only first-order optimization

Variance Reduction Using Value Functions

Variance Reduction

Now, we have the following policy gradient formula:

$$abla_ heta \mathbb{E}_ au\left[R
ight] = \mathbb{E}_ au \left[\sum_{t=0}^{ au-1}
abla_ heta \log \pi(a_t \mid s_t, heta) \mathsf{A}^\pi(s_t, a_t)
ight]$$

- A^{π} is not known, but we can plug in \hat{A}_t , an advantage estimator
- Previously, we showed that taking

$$\hat{A}_t = r_t + r_{t+1} + r_{t+2} + \cdots - b(s_t)$$

for any function $b(s_t)$, gives an unbiased policy gradient estimator. $b(s_t) \approx V^{\pi}(s_t)$ gives variance reduction.

The Delayed Reward Problem

With policy gradient methods, we are confounding the effect of multiple actions:

$$\hat{A}_t = r_t + r_{t+1} + r_{t+2} + \cdots - b(s_t)$$

mixes effect of $a_t, a_{t+1}, a_{t+2}, \ldots$

- SNR of \hat{A}_t scales roughly as 1/T
 - ► Only a_t contributes to signal A^π(s_t, a_t), but a_{t+1}, a_{t+2},... contribute to noise.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへで

Variance Reduction with Discounts

- ► Discount factor *γ*, 0 < *γ* < 1, downweights the effect of rewars that are far in the future—ignore long term dependencies</p>
- We can form an advantage estimator using the *discounted return*:

$$\hat{A}_{t}^{\gamma} = \underbrace{r_{t} + \gamma r_{t+1} + \gamma^{2} r_{t+2} + \dots}_{\text{discounted return}} - b(s_{t})$$

reduces to our previous estimator when $\gamma = 1$.

So advantage has expectation zero, we should fit baseline to be *discounted* value function

$$\mathcal{W}^{\pi,\gamma}(s) = \mathbb{E}_{\tau} \left[\mathbf{r}_0 + \gamma \mathbf{r}_1 + \gamma^2 \mathbf{r}_2 + \dots \mid \mathbf{s}_0 = s
ight]$$

- Discount γ is similar to using a horizon of $1/(1-\gamma)$ timesteps
- \hat{A}_t^{γ} is a biased estimator of the advantage function

Value Functions in the Future

- Baseline accounts for and removes the effect of past actions
- Can also use the value function to estimate future rewards

 $r_t + \gamma V(s_{t+1})$ cut off at one timestep $r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2})$ cut off at two timesteps... $r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$ ∞ timesteps (no V)

Value Functions in the Future

Subtracting out baselines, we get advantage estimators

$$\hat{A}_{t}^{(1)} = r_{t} + \gamma V(s_{t+1}) - V(s_{t})$$
$$\hat{A}_{t}^{(2)} = r_{t} + r_{t+1} + \gamma^{2} V(s_{t+2}) - V(s_{t})$$
...
$$\hat{A}_{t}^{(\infty)} = r_{t} + \gamma r_{t+1} + \gamma^{2} r_{t+2} + \dots - V(s_{t})$$

- $\hat{A}_t^{(1)}$ has low variance but high bias, $\hat{A}_t^{(\infty)}$ has high variance but low bias.
- Using intermediate k (say, 20) gives an intermediate amount of bias and variance

Finite-Horizon Methods: Advantage Actor-Critic

A2C / A3C uses this fixed-horizon advantage estimator

V. Mnih, A. P. Badia, M. Mirza, et al. "Asynchronous Methods for Deep Reinforcement Learning". In: ICML (2016) + 4 = + 4 = + 6 - 2 - 9 0 0

Finite-Horizon Methods: Advantage Actor-Critic

- A2C / A3C uses this fixed-horizon advantage estimator
- Pseudocode

for iteration=1, 2, ... do Agent acts for T timesteps (e.g., T = 20), For each timestep t, compute

$$\hat{R}_t = r_t + \gamma r_{t+1} + \dots + \gamma^{T-t+1} r_{T-1} + \gamma^{T-t} V(s_t)$$
$$\hat{A}_t = \hat{R}_t - V(s_t)$$

 \hat{R}_t is target value function, in regression problem \hat{A}_t is estimated advantage function Compute loss gradient $g = \nabla_{\theta} \sum_{t=1}^{T} \left[-\log \pi_{\theta} (a_t \mid s_t) \hat{A}_t + c(V(s) - \hat{R}_t)^2 \right]$ g is plugged into a stochastic gradient descent variant, e.g., Adam. end for

A3C Video

◆□ > ◆□ > ◆三 > ◆三 > ○ ○ ○

A3C Results



$\mathsf{TD}(\lambda)$ Methods: Generalized Advantage Estimation

Recall, finite-horizon advantage estimators

$$\hat{A}_t^{(k)} = r_t + \gamma r_{t+1} + \dots + \gamma^{k-1} r_{t+k-1} + \gamma^k V(s_{t+k}) - V(s_t)$$

• Define the TD error
$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

By a telescoping sum,

$$\hat{A}_t^{(k)} = \delta_t + \gamma \delta_{t+1} + \dots + \gamma^{k-1} \delta_{t+k-1}$$

Take exponentially weighted average of finite-horizon estimators:

$$\hat{\mathcal{A}}^{\lambda} = \hat{\mathcal{A}}_t^{(1)} + \lambda \hat{\mathcal{A}}_t^{(2)} + \lambda^2 \hat{\mathcal{A}}_t^{(3)} + \dots$$

We obtain

$$\hat{A}_t^{\lambda} = \delta_t + (\gamma \lambda) \delta_{t+1} + (\gamma \lambda)^2 \delta_{t+2} + \dots$$

This scheme named generalized advantage estimation (GAE) in [1], though versions have appeared earlier, e.g., [2]. Related to TD(λ)

J. Schulman, P. Moritz, S. Levine, et al. "High-dimensional continuous control using generalized advantage estimation". In: ICML. 2015

H. Kimura and S. Kobayashi. "An Analysis of Actor/Critic Algorithms Using Eligibility Traces: Reinforcement Learning with Imperfect Value Function." In: ICML. 1998, pp. 278–286

Choosing parameters γ, λ

Performance as γ,λ are varied



◆□ > ◆□ > ◆臣 > ◆臣 > 善臣 - のへで

TRPO+GAE Video

◆□ > ◆□ > ◆三 > ◆三 > ○ ○ ○

Pathwise Derivative Policy Gradient Methods

Deriving the Policy Gradient, Reparameterized

• Episodic MDP:



▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● ● ● ● ●

Want to compute $\nabla_{\theta} \mathbb{E}[R_T]$. We'll use $\nabla_{\theta} \log \pi(a_t \mid s_t; \theta)$

Deriving the Policy Gradient, Reparameterized

• Episodic MDP:



Want to compute $\nabla_{\theta} \mathbb{E}[R_T]$. We'll use $\nabla_{\theta} \log \pi(a_t \mid s_t; \theta)$

• Reparameterize: $a_t = \pi(s_t, z_t; \theta)$. z_t is noise from fixed distribution.



Deriving the Policy Gradient, Reparameterized

► Episodic MDP:



Want to compute $\nabla_{\theta} \mathbb{E}[R_T]$. We'll use $\nabla_{\theta} \log \pi(a_t \mid s_t; \theta)$

• Reparameterize: $a_t = \pi(s_t, z_t; \theta)$. z_t is noise from fixed distribution.



• Only works if $P(s_2 | s_1, a_1)$ is known $\ddot{\frown}$

Using a Q-function



$$\frac{\mathrm{d}}{\mathrm{d}\theta} \mathbb{E}\left[R_{T}\right] = \mathbb{E}\left[\sum_{t=1}^{T} \frac{\mathrm{d}R_{T}}{\mathrm{d}a_{t}} \frac{\mathrm{d}a_{t}}{\mathrm{d}\theta}\right] = \mathbb{E}\left[\sum_{t=1}^{T} \frac{\mathrm{d}}{\mathrm{d}a_{t}} \mathbb{E}\left[R_{T} \mid a_{t}\right] \frac{\mathrm{d}a_{t}}{\mathrm{d}\theta}\right]$$
$$= \mathbb{E}\left[\sum_{t=1}^{T} \frac{\mathrm{d}Q(s_{t}, a_{t})}{\mathrm{d}a_{t}} \frac{\mathrm{d}a_{t}}{\mathrm{d}\theta}\right] = \mathbb{E}\left[\sum_{t=1}^{T} \frac{\mathrm{d}}{\mathrm{d}\theta}Q(s_{t}, \pi(s_{t}, z_{t}; \theta))\right]$$

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ ▲目 ● ● ●

SVG(0) Algorithm

• Learn Q_{ϕ} to approximate $Q^{\pi,\gamma}$, and use it to compute gradient estimates.

SVG(0) Algorithm

- Learn Q_{ϕ} to approximate $Q^{\pi,\gamma}$, and use it to compute gradient estimates.
- Pseudocode:

for iteration=1,2,... do Execute policy π_{θ} to collect T timesteps of data Update π_{θ} using $g \propto \nabla_{\theta} \sum_{t=1}^{T} Q(s_t, \pi(s_t, z_t; \theta))$ Update Q_{ϕ} using $g \propto \nabla_{\phi} \sum_{t=1}^{T} (Q_{\phi}(s_t, a_t) - \hat{Q}_t)^2$, e.g. with $\mathsf{TD}(\lambda)$ end for

N. Heess, G. Wayne, D. Silver, et al. "Learning continuous control policies by stochastic value gradients". dn: NIPS 🖓 015 (🚊 + (🛓 - () 🙊) () ()

SVG(1) Algorithm



- Instead of learning Q, we learn
 - State-value function $V \approx V^{\pi,\gamma}$
 - Dynamics model f, approximating $s_{t+1} = f(s_t, a_t) + \zeta_t$
- Given transition (s_t, a_t, s_{t+1}) , infer $\zeta_t = s_{t+1} f(s_t, a_t)$

$$\blacktriangleright Q(s_t, a_t) = \mathbb{E}\left[r_t + \gamma V(s_{t+1})\right] = \mathbb{E}\left[r_t + \gamma V(f(s_t, a_t) + \zeta_t)\right], \text{ and } a_t = \pi(s_t, \theta, \zeta_t)$$

$SVG(\infty)$ Algorithm



- ► Just learn dynamics model f
- Given whole trajectory, infer all noise variables
- Freeze all policy and dynamics noise, differentiate through entire deterministic computation graph

SVG Results

Applied to 2D robotics tasks



Overall: different gradient estimators behave similarly



N. Heess, G. Wayne, D. Silver, et al. "Learning continuous control policies by stochastic value gradients". dm NIPS 🔁 015 (🚊) (🛓)

Deterministic Policy Gradient

- For Gaussian actions, variance of score function policy gradient estimator goes to infinity as variance goes to zero
- But SVG(0) gradient is fine when $\sigma \rightarrow 0$

$$\nabla_{\theta} \sum_{t} Q(s_t, \pi(s_t, \theta, \zeta_t))$$

- Problem: there's no exploration.
- Solution: add noise to the policy, but estimate Q with TD(0), so it's valid off-policy
- Policy gradient is a little biased (even with Q = Q^π), but only because state distribution is off—it gets the right gradient at every state

D. Silver, G. Lever, N. Heess, et al. "Deterministic policy gradient algorithms". In: ICML. 2014

Deep Deterministic Policy Gradient

 Incorporate replay buffer and target network ideas from DQN for increased stability

Deep Deterministic Policy Gradient

- Incorporate replay buffer and target network ideas from DQN for increased stability
- Use lagged (Polyak-averaging) version of Q_{ϕ} and π_{θ} for fitting Q_{ϕ} (towards $Q^{\pi,\gamma}$) with TD(0)

$$\hat{Q}_t = r_t + \gamma Q_{\phi'}(s_{t+1}, \pi(s_{t+1}; \theta'))$$

Deep Deterministic Policy Gradient

- Incorporate replay buffer and target network ideas from DQN for increased stability
- Use lagged (Polyak-averaging) version of Q_{ϕ} and π_{θ} for fitting Q_{ϕ} (towards $Q^{\pi,\gamma}$) with TD(0)

$$\hat{Q}_t = r_t + \gamma Q_{\phi'}(s_{t+1}, \pi(s_{t+1}; \theta'))$$

Pseudocode:

for iteration=1, 2, ... do Act for several timesteps, add data to replay buffer Sample minibatch Update π_{θ} using $g \propto \nabla_{\theta} \sum_{t=1}^{T} Q(s_t, \pi(s_t, z_t; \theta))$ Update Q_{ϕ} using $g \propto \nabla_{\phi} \sum_{t=1}^{T} (Q_{\phi}(s_t, a_t) - \hat{Q}_t)^2$, end for



Applied to 2D and 3D robotics tasks and driving with pixel input



T. P. Lillicrap, J. J. Hunt, A. Pritzel, et al. "Continuous control with deep reinforcement learning". In: ICLR:(2015) 🖶 + 4 🗄 + 4 👼 + 4 \emptyset

Policy Gradient Methods: Comparison

Two kinds of policy gradient estimator

- ▶ REINFORCE / score function estimator: $\nabla \log \pi(a \mid s) \hat{A}$.
 - Learn Q or V for variance reduction, to estimate \hat{A}
- Pathwise derivative estimators (differentiate wrt action)
 - SVG(0) / DPG: $\frac{d}{da}Q(s,a)$ (learn Q)
 - SVG(1): $\frac{d}{da}(r + \gamma V(s'))$ (learn f, V)
 - SVG(∞): $\frac{\mathrm{d}}{\mathrm{d}a_t}(r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots)$ (learn f)
- Pathwise derivative methods more sample-efficient when they work (maybe), but work less generally due to high bias

Policy Gradient Methods: Comparison

Task	Random	REINFORCE	TNPG	RWR	REPS	TRPO	CEM	CMA-ES	DDPG
Cart-Pole Balancing Inverted Pendulum ⁴ Mountain Car Acrobot Double Inverted Pendulum ⁴	$\begin{array}{r} 77.1\pm0.0\\-153.4\pm0.2\\-415.4\pm0.0\\-1904.5\pm1.0\\149.7\pm0.1\end{array}$	$\begin{array}{rrrr} 4693.7\pm&14.0\\ 13.4\pm&18.0\\ -67.1\pm&1.0\\ -508.1\pm&91.0\\ 4116.5\pm&65.2 \end{array}$	$\begin{array}{rrrr} 3986.4 & \pm \ 748.9 \\ 209.7 & \pm \ 55.5 \\ \textbf{-66.5} & \pm \ \textbf{4.5} \\ -395.8 \pm 121.2 \\ \textbf{4455.4} & \pm \ \textbf{37.6} \end{array}$	$\begin{array}{r} \textbf{4861.5} \pm \textbf{12.3} \\ 84.7 \pm 13.8 \\ -79.4 \pm 1.1 \\ -352.7 \pm 35.9 \\ 3614.8 \pm 368.1 \end{array}$	$\begin{array}{c} 565.6\pm137.6\\-113.3\pm&4.6\\-275.6\pm166.3\\-1001.5\pm&10.8\\446.7\pm114.8\end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 4815.4\pm&4.8\\38.2\pm&25.7\\-66.0\pm&2.4\\-436.8\pm&14.7\\2566.2\pm178.9\end{array}$	$\begin{array}{rrrr} 2440.4\pm 568.3\\ -40.1\pm & 5.7\\ -85.0\pm & 7.7\\ -785.6\pm & 13.1\\ 1576.1\pm & 51.3 \end{array}$	$\begin{array}{r} 4634.4 \ \pm \ 87.8 \\ 40.0 \ \pm \ 244.6 \\ -288.4 \ \pm \ 170.3 \\ \textbf{-223.6} \ \pm \ \textbf{5.8} \\ 2863.4 \ \pm \ 154.0 \end{array}$
Swimmer* Hopper 2D Walker Half-Cheetah Ant* Simple Humanoid Full Humanoid	$\begin{array}{c} -1.7\pm0.1\\ 8.4\pm0.0\\ -1.7\pm0.0\\ -90.8\pm0.3\\ 13.4\pm0.7\\ 41.5\pm0.2\\ 13.2\pm0.1\end{array}$	$\begin{array}{rrrr} 92.3 \pm & 0.1 \\ 714.0 \pm & 29.3 \\ 506.5 \pm & 78.8 \\ 1183.1 \pm & 69.2 \\ 548.3 \pm & 55.5 \\ 128.1 \pm & 34.0 \\ 262.2 \pm & 10.5 \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{cccc} 60.7\pm & 5.5\\ 553.2\pm & 71.0\\ 136.0\pm & 15.9\\ 376.1\pm & 28.2\\ 37.6\pm & 3.1\\ 93.3\pm & 17.4\\ 46.7\pm & 5.6\end{array}$	$\begin{array}{rrrr} 3.8 \pm & 3.3 \\ 86.7 \pm & 17.6 \\ -37.0 \pm & 38.1 \\ 34.5 \pm & 38.0 \\ 39.0 \pm & 9.8 \\ 28.3 \pm & 4.7 \\ 41.7 \pm & 6.1 \end{array}$	$\begin{array}{rrrrr} 96.0 & \pm & 0.2 \\ 1183.3 & \pm & 150.0 \\ 1353.8 & \pm & 820.1 \\ 1914.0 & \pm & 120.1 \\ 730.2 & \pm & 61.3 \\ 269.7 & \pm & 40.3 \\ 287.0 & \pm & 23.4 \end{array}$	$\begin{array}{c} 68.8 \pm & 2.4 \\ 63.1 \pm & 7.8 \\ 84.5 \pm & 19.2 \\ 330.4 \pm 274.8 \\ 49.2 \pm & 5.9 \\ 60.6 \pm & 12.9 \\ 36.9 \pm & 2.9 \end{array}$	$\begin{array}{cccc} 64.9 \pm & 1.4 \\ 20.3 \pm & 14.3 \\ 77.1 \pm & 24.3 \\ 441.3 \pm & 107.6 \\ 17.8 \pm & 15.5 \\ 28.7 \pm & 3.9 \\ \text{N/A} \pm & \text{N/A} \end{array}$	$\begin{array}{rrrr} 85.8 \pm & 1.8 \\ 267.1 \pm & 43.5 \\ 318.4 \pm 181.6 \\ 2148.6 \pm 702.7 \\ 326.2 \pm & 20.8 \\ 99.4 \pm & 28.1 \\ 119.0 \pm & 31.2 \end{array}$
Cart-Pole Balancing (LS)* Inverted Pendulum (LS) Mountain Car (LS) Acrobot (LS)*	$\begin{array}{c} 77.1 \pm 0.0 \\ -122.1 \pm 0.1 \\ -83.0 \pm 0.0 \\ -393.2 \pm 0.0 \end{array}$	$\begin{array}{r} 420.9\pm265.5\\-13.4\pm&3.2\\-81.2\pm&0.6\\-128.9\pm&11.6\end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 68.9 \pm & 1.5 \\ -107.4 \pm & 0.2 \\ -81.7 \pm & 0.1 \\ -235.9 \pm & 5.3 \end{array}$	$\begin{array}{rrrr} 898.1\pm&22.1\\-87.2\pm&8.0\\-82.6\pm&0.4\\-379.5\pm&1.4\end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{r} 227.0\pm223.0\\-81.2\pm&33.2\\\textbf{-68.9}\pm&\textbf{1.3}\\-149.5\pm&15.3\end{array}$	$\begin{array}{rrrr} 68.0 \pm & 1.6 \\ -62.4 \pm & 3.4 \\ \textbf{.73.2} \ \pm \ \textbf{0.6} \\ -159.9 \pm & 7.5 \end{array}$	
Cart-Pole Balancing (NO)* Inverted Pendulum (NO) Mountain Car (NO) Acrobot (NO)*	$\begin{array}{c} 101.4\pm0.1\\-122.2\pm0.1\\-83.0\pm0.0\\-393.5\pm0.0\end{array}$	$\begin{array}{c} 616.0\pm210.8\\ 6.5\pm&1.1\\ -74.7\pm&7.8\\ \textbf{\cdot186.7}\pm&\textbf{31.3} \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 93.8\pm&1.2\\-110.0\pm&1.4\\-81.7\pm&0.1\\-233.1\pm&0.4\end{array}$	$\begin{array}{rrrr} 99.6\pm&7.2\\-119.3\pm&4.2\\-82.9\pm&0.1\\-258.5\pm&14.0\end{array}$	$\begin{array}{c} 606.2 \pm 122.2 \\ 10.4 \ \pm \ 2.2 \\ \textbf{-60.2} \ \pm \ 2.0 \\ \textbf{-149.6} \ \pm \ \textbf{8.6} \end{array}$	$\begin{array}{rrrr} 181.4\pm&32.1\\-55.6\pm&16.7\\-67.4\pm&1.4\\-213.4\pm&6.3\end{array}$	$\begin{array}{rrrr} 104.4\pm&16.0\\-80.3\pm&2.8\\-73.5\pm&0.5\\-236.6\pm&6.2\end{array}$	
Cart-Pole Balancing (SI)* Inverted Pendulum (SI) Mountain Car (SI) Acrobot (SI)*	$\begin{array}{r} 76.3 \pm 0.1 \\ -121.8 \pm 0.2 \\ -82.7 \pm 0.0 \\ -387.8 \pm 1.0 \end{array}$	$\begin{array}{r} 431.7\pm274.1\\-5.3\pm&5.6\\-63.9\pm&0.2\\\textbf{\cdot169.1}\pm&\textbf{32.3}\end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 69.0 \pm & 2.8 \\ -108.7 \pm & 4.7 \\ -81.4 \pm & 0.1 \\ -233.2 \pm & 2.6 \end{array}$	$\begin{array}{r} 702.4 \pm 196.4 \\ -92.8 \pm \ 23.9 \\ -80.7 \pm \ 2.3 \\ -216.1 \pm \ 7.7 \end{array}$	$\begin{array}{rrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrrr$	$\begin{array}{rrrr} 746.6\pm&93.2\\-51.8\pm&10.6\\-63.9\pm&1.0\\-250.2\pm&13.7\end{array}$	$\begin{array}{rrrr} 71.6\pm&2.9\\-63.1\pm&4.8\\-66.9\pm&0.6\\-245.0\pm&5.5\end{array}$	
Swimmer + Gathering Ant + Gathering Swimmer + Maze Ant + Maze	${}^{0.0\pm0.0}_{-5.8\pm5.0}_{0.0\pm0.0}_{0.0\pm0.0}_{0.0\pm0.0}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -0.1\pm & 0.1\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0\end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -0.4\pm & 0.1\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0\end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -5.5\pm & 0.5\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0\end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -6.7\pm & 0.7\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0 \end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -0.4\pm & 0.0\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0\end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -4.7\pm & 0.7\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0 \end{array}$	$\begin{array}{ccc} 0.0\pm & 0.0 \\ N/A\pm & N/A \\ 0.0\pm & 0.0 \\ N/A\pm & N/A \end{array}$	$\begin{array}{cccc} 0.0\pm & 0.0\\ -0.3\pm & 0.3\\ 0.0\pm & 0.0\\ 0.0\pm & 0.0 \end{array}$

Y. Duan, X. Chen, R. Houthooft, et al. "Benchmarking Deep Reinforcement Learning for Continuous Control" + Inx HML (2016) < 🚊 - 🔗 < 🖓

Stochastic Computation Graphs

Gradients of Expectations

Want to compute $\nabla_{\theta} \mathbb{E}[F]$. Where's θ ?

- ▶ In distribution, e.g., $\mathbb{E}_{x \sim p(\cdot \mid \theta)}[F(x)]$
 - $\triangleright \nabla_{\theta} \mathbb{E}_{x} [f(x)] = \mathbb{E}_{x} [f(x) \nabla_{\theta} \log p_{x}(x; \theta)].$
 - Score function estimator
 - ► Example: REINFORCE policy gradients, where *x* is the trajectory
- Outside distribution: $\mathbb{E}_{z \sim \mathcal{N}(0,1)} [F(\theta, z)]$

$$\nabla_{\theta} \mathbb{E}_{z} \left[f(x(z,\theta)) \right] = \mathbb{E}_{z} \left[\nabla_{\theta} f(x(z,\theta)) \right].$$

- Pathwise derivative estimator
- Example: SVG policy gradient
- > Often, we can reparametrize, to change from one form to another
- What if F depends on θ in complicated way, affecting distribution and F?

M. C. Fu. "Gradient estimation". In: Handbooks in operations research and management science 13 (2006) cpp. 57 📴 616 4 🚊 🕨 4 🚊 🖉 4 🖓 4 🖓

Stochastic Computation Graphs

- Stochastic computation graph is a DAG, each node corresponds to a deterministic or stochastic operation
- Can automatically derive unbiased gradient estimators, with variance reduction



Stochastic Computation Graphs



J. Schulman, N. Heess, T. Weber, et al. "Gradient Estimation Using Stochastic Computation Graphs". In: MIPS. 2015 + 4 🗄 + 4 🗟 + 4 🗟 - 🖓 🤄

Worked Example



- L = c + e. Want to compute $\frac{d}{d\theta} \mathbb{E}[L]$ and $\frac{d}{d\phi} \mathbb{E}[L]$.
- Treat stochastic nodes (b, d) as constants, and introduce losses logprob * (futurecost) at each stochastic node
- Obtain unbiased gradient estimate by differentiating surrogate:

Surrogate
$$(\theta, \psi) = \underbrace{c + e}_{(1)} + \underbrace{\log p(\hat{b} \mid a, d)\hat{c}}_{(2)}$$

(1): how parameters influence cost through deterministic dependencies

(2): how parameters affect distribution over random variables.