

# Lecture 5: Policy Gradient I

Emma Brunskill

CS234 Reinforcement Learning

Winter 2025

- With many slides from or derived from David Silver and John Schulman and Pieter Abbeel
- Additional reading: Sutton and Barto 2018 Chp. 13

# L5N1 Refresh Your Knowledge. Comparing Policy Performance

- Consider doing experience replay over a finite, but extremely large, set of  $(s,a,r,s')$  tuples). Q-learning is initialized to 0 everywhere and all rewards are positive. Select all that are true
  - ➊ Assume all tuples were gathered from a fixed, deterministic policy  $\pi$ . Then in the tabular setting, if each tuple is sampled at random and used to do a Q-learning update, and this is repeated an infinite number of times, then there exists a learning rate schedule so that the resulting estimate will converge to the true  $Q^\pi$ .
  - ➋ In situation (1) (the first option above) the resulting Q estimate will be identical to if one computed an estimated dynamics model and reward model using maximum likelihood evaluation from the tuples, and performed policy evaluation using the estimated dynamics and reward models.
  - ➌ If one uses DQN to populate the experience replay set of tuples, then doing experience replay with DQN is always guaranteed to converge to the optimal Q function.
  - ➍ Not sure

# L5N1 Refresh Your Knowledge. Comparing Policy Performance

- Consider doing experience replay over a finite, but extremely large, set of  $(s,a,r,s')$  tuples). Q-learning is initialized to 0 everywhere and all rewards are positive. Select all that are true
  - Assume all tuples were gathered from a fixed, deterministic policy  $\pi$ . Then in the tabular setting, if each tuple is sampled at random and used to do a Q-learning update, and this is repeated an infinite number of times, then there exists a learning rate schedule so that the resulting estimate will converge to the true  $Q^\pi(s, \pi(s))$ .  ~~$Q^\pi(s,a)$~~  false
  - In situation (1) (the first option above) the resulting Q estimate will be identical to if one computed an estimated dynamics model and reward model using maximum likelihood evaluation from the tuples, and performed policy evaluation using the estimated dynamics and reward models.
  - If one uses DQN to populate the experience replay set of tuples, then doing experience replay with DQN is always guaranteed to converge to the optimal Q function.
  - Not sure

$$Q^\pi(s,a)$$

Answer: 1 is true and 2 is true. 3 is false.

# Class Structure

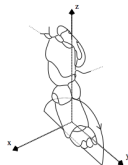
- Last time: Learning to Control in Tabular MDPs to Deep RL / Generalization to scale RL
- **This time: Policy Search**
- Next time: Policy Search Cont.

# RL Algorithms Involve

- Optimization
- Delayed consequences
- Exploration
- Generalization

# Do We Need "RL" at All? Can we Just do Online Optimization?

- Policy gradient methods have been **very** influential
- In NLP (Sequence Level Training with Recurrent Neural Networks built on REINFORCE)
- End-to-End Training of Deep Visuomotor Policies  
<https://arxiv.org/abs/1504.00702>
- In homework 2 you will be implementing Proximal Policy Optimization (PPO) which was used in training ChatGPT



**Figure:** Early example of policy gradient methods: training a AIBO to have a faster walk. Paper: Kohl and Stone, ICRA 2004

# Policy-Based Reinforcement Learning

- In the last lecture we approximated the value or action-value function using parameters  $w$ ,

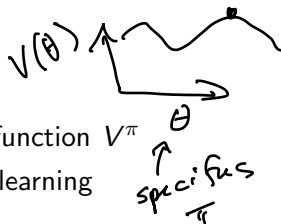
$$V_w(s) \approx V^\pi(s)$$

$$Q_w(s, a) \approx Q^\pi(s, a)$$

- A policy was generated directly from the value function
  - e.g. using  $\epsilon$ -greedy
- In this lecture we will directly parametrize the policy, and will typically use  $\theta$  to show parameterization:

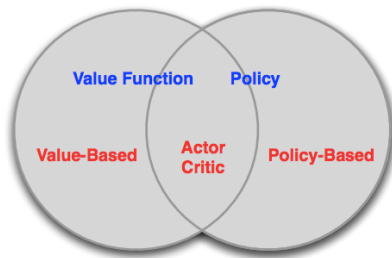
$$\pi_\theta(s, a) = \mathbb{P}[a|s; \theta]$$

- Goal is to find a policy  $\pi$  with the highest value function  $V^\pi$
- We will focus again on model-free reinforcement learning



# Value-Based and Policy-Based RL

- Value Based
  - learned Value Function
  - Implicit policy (e.g.  $\epsilon$ -greedy)
- Policy Based
  - No Value Function
  - Learned Policy
- Actor-Critic
  - Learned Value Function
  - Learned Policy





# Types of Policies to Search Over

- So far have focused on deterministic policies or  $\epsilon$ -greedy policies
- Now we are thinking about direct policy search in RL, will focus heavily on stochastic policies

# Example: Rock-Paper-Scissors



- Two-player game of rock-paper-scissors
  - Scissors beats paper
  - Rock beats scissors
  - Paper beats rock
- Let state be history of prior actions (rock, paper and scissors) and if won or lost
- Is deterministic policy optimal? Why or why not?

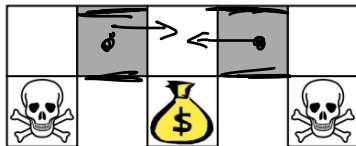
# Example: Rock-Paper-Scissors, Vote



- Two-player game of rock-paper-scissors
  - Scissors beats paper
  - Rock beats scissors
  - Paper beats rock
- Let state be history of prior actions (rock, paper and scissors) and if won or lost

Deterministic policy is easily exploited by an adversary. System is not Markov. A uniform random policy is optimal (Nash equilibrium).

# Example: Aliased Gridworld (1)



- The agent cannot differentiate the grey states
- Consider features of the following form (for all N, E, S, W)

$$\phi(s, a) = \mathbb{1}(\text{wall to } N, a = \text{move } E)$$

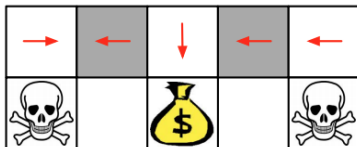
- Compare value-based RL, using an approximate value function

$$Q_{\theta}(s, a) = f(\phi(s, a); \theta)$$

- To policy-based RL, using a parametrized policy

$$\pi_{\theta}(s, a) = g(\phi(s, a); \theta)$$

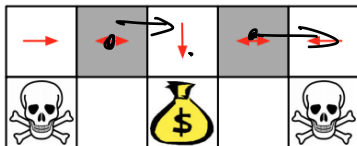
## Example: Aliased Gridworld (2)



- Under aliasing, an optimal **deterministic** policy will either
  - move W in both grey states (shown by red arrows)
  - move E in both grey states
- Either way, it can get stuck and never reach the money
- Value-based RL learns a near-deterministic policy
  - e.g. greedy or  $\epsilon$ -greedy
- So it will traverse the corridor for a long time

## Example: Aliased Gridworld (3)

$$p(s'|s, a)$$



- An optimal **stochastic** policy will randomly move E or W in grey states

$$\pi_{\theta}(\text{wall to N and S, move E}) = 0.5$$

$$\pi_{\theta}(\text{wall to N and S, move W}) = 0.5$$

- It will reach the goal state in a few steps with high probability
- Policy-based RL can learn the optimal stochastic policy

# Policy Objective Functions

before  $|A|^{15}$

- Goal: given a policy  $\pi_\theta(s, a)$  with parameters  $\theta$ , find best  $\theta$
- But how do we measure the quality for a policy  $\pi_\theta$ ?
- In episodic environments can use policy value at start state  $V(s_0, \theta)$
- For simplicity, today will mostly discuss the episodic case, but can easily extend to the continuing / infinite horizon case

# Policy optimization

- Policy based reinforcement learning is an **optimization** problem
- Find policy parameters  $\theta$  that maximize  $V(s_0, \theta)$



# Policy optimization

- Policy based reinforcement learning is an **optimization** problem
- Find policy parameters  $\theta$  that maximize  $V(s_0, \theta)$
- Can use gradient free optimization
  - Hill climbing
  - Simplex / amoeba / Nelder Mead
  - Genetic algorithms
  - Cross-Entropy method (CEM)
  - Covariance Matrix Adaptation (CMA)

# Human-in-the-Loop Exoskeleton Optimization (Zhang et al. Science 2017)

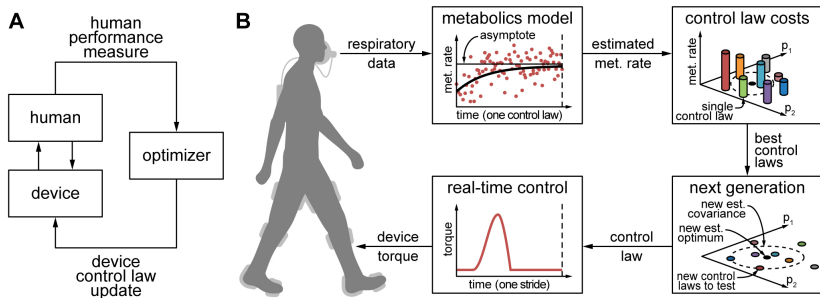


Figure: Zhang et al. Science 2017

- Optimization was done using CMA-ES, variation of covariance matrix evaluation

# Gradient Free Policy Optimization

- Can often work embarrassingly well: "discovered that evolution strategies (ES), an optimization technique that's been known for decades, rivals the performance of standard reinforcement learning (RL) techniques on modern RL benchmarks (e.g. Atari/MuJoCo)" (<https://blog.openai.com/evolution-strategies/>)

# Gradient Free Policy Optimization

- Often a great simple baseline to try
- Benefits
  - Can work with any policy parameterizations, including non-differentiable
  - Frequently very easy to parallelize
- Limitations
  - Often less sample efficient because it ignores temporal structure

# Policy optimization

- Policy based reinforcement learning is an **optimization** problem
- Find policy parameters  $\theta$  that maximize  $V(s_0, \theta)$
- Can use gradient free optimization:
- Greater efficiency often possible using gradient
  - Gradient descent
  - Conjugate gradient
  - Quasi-newton
- We focus on gradient descent, many extensions possible
- And on methods that exploit sequential structure

$$V^{\pi(\theta)}(s)$$

- Define  $V(\theta) = V(s_0, \theta)$  to make explicit the dependence of the value on the policy parameters [but don't confuse with value function approximation, where parameterized value function]
- Assume episodic MDPs (easy to extend to related objectives, like average reward)

# Policy Gradient

- Define  $V^{\pi_{\theta}} = V(s_0, \theta)$  to make explicit the dependence of the value on the policy parameters
- Assume episodic MDPs
- Policy gradient algorithms search for a *local* maximum in  $V(s_0, \theta)$  by ascending the gradient of the policy, w.r.t parameters  $\theta$

*best in the policy class*

$$\Delta\theta = \alpha \nabla_{\theta} V(s_0, \theta)$$

- Where  $\nabla_{\theta} V(s_0, \theta)$  is the **policy gradient**

$$\nabla_{\theta} V(s_0, \theta) = \begin{pmatrix} \frac{\partial V(s_0, \theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial V(s_0, \theta)}{\partial \theta_n} \end{pmatrix}$$

- and  $\alpha$  is a step-size parameter

# Value of a Parameterized Policy

$$\gamma = (s_0, \pi(s_0), r \dots \dots s_H, \pi, r(s_H, \pi(s_H)))$$

- Now assume policy  $\pi_\theta$  is differentiable whenever it is non-zero and we know the gradient  $\nabla_\theta \pi_\theta(s, a)$  *we could add  $g^t$*
- Recall policy value is  $V(s_0, \theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T R(s_t, a_t); \pi_\theta, s_0 \right]$  where the expectation is taken over the states & actions visited by  $\pi_\theta$
- We can re-express this in multiple ways

$$V(s_0, \theta) = \sum_a \pi_\theta(a|s_0) Q(s_0, a, \theta)$$

$$= \sum_\gamma p(\gamma|\theta) R(\gamma)$$

*also used  $G(\gamma)$  to denote  $\sum$  of rewards*

*approx with sampling*

*freq  $\gamma$  sampled using  $\pi_\theta$*



# Value of a Parameterized Policy

- Assume policy  $\pi_\theta$  is differentiable whenever it is non-zero and we can compute the gradient  $\nabla_\theta \pi_\theta(s, a)$
- Recall policy value is  $V(s_0, \theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T R(s_t, a_t); \pi_\theta, s_0 \right]$  where the expectation is taken over the states & actions visited by  $\pi_\theta$
- We can re-express this in multiple ways
  - $V(s_0, \theta) = \sum_a \pi_\theta(a|s_0) Q(s_0, a, \theta)$
  - $V(s_0, \theta) = \sum_\tau P(\tau; \theta) R(\tau)$ 
    - where  $\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$  is a state-action trajectory,
    - $P(\tau; \theta)$  is used to denote the probability over trajectories when executing policy  $\pi(\theta)$  starting in state  $s_0$ , and
    - $R(\tau) = \sum_{t=0}^T R(s_t, a_t)$  the sum of rewards for a trajectory  $\tau$
- To start will focus on this latter definition. See Chp 13.1-13.3 of SB for a nice discussion starting with the other definition

# Likelihood Ratio Policies

- Denote a state-action trajectory as  $\tau = (s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$
- Use  $R(\tau) = \sum_{t=0}^T R(s_t, a_t)$  to be the sum of rewards for a trajectory  $\tau$
- Policy value is

$$V(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^T R(s_t, a_t); \pi_\theta \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

- where  $P(\tau; \theta)$  is used to denote the probability over trajectories when executing policy  $\pi(\theta)$
- In this new notation, our goal is to find the policy parameters  $\theta$ :

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

# Likelihood Ratio Policy Gradient

- Goal is to find the policy parameters  $\theta$ :

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Take the gradient with respect to  $\theta$ :

$$\begin{aligned} \nabla_{\theta} V(\theta) &= \nabla_{\theta} \sum P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} R(\tau) \nabla_{\theta} P(\tau; \theta) \quad R(\tau) \text{ indep of } \theta \\ &= \sum_{\tau} R(\tau) \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) \quad \nabla_{\theta} \log P(\tau; \theta) = \frac{1}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) \\ &= \sum_{\tau} R(\tau) P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) \end{aligned}$$

# Likelihood Ratio Policy Gradient

- Goal is to find the policy parameters  $\theta$ :

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Take the gradient with respect to  $\theta$ :

$$\begin{aligned} \nabla_{\theta} V(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \underbrace{\frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)}}_{\text{likelihood ratio}} \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta) \end{aligned}$$

# Likelihood Ratio Policy Gradient

- Goal is to find the policy parameters  $\theta$ :

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Take the gradient with respect to  $\theta$ :

$$\nabla_{\theta} V(\theta) = \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta)$$

- Approximate using  $m$  sample trajectories under policy  $\pi_{\theta}$ :

$$\nabla_{\theta} V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta)$$

# Decomposing the Trajectories Into States and Actions

- Approximate using  $m$  sample paths under policy  $\pi_\theta$ :  $\mu(s_0) = \text{prob}_{\text{env}} s_0$

$$\nabla_\theta V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_\theta \log P(\tau^{(i)})$$

$$\begin{aligned} \nabla_\theta \log P(\tau^{(i)}; \theta) &= \nabla_\theta \log \left[ \mu(s_0) \prod_{t=0}^{\tau-1} \pi_\theta(a_t | s_t) p(s_{t+1} | a_t, s_{0:t}, a_{0:t}) \right] \\ &= \nabla_\theta \log \mu(s_0) + \sum_{t=0}^{\tau-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \\ &\quad + \sum_{t=0}^{\tau-1} \nabla_\theta \log p(s_{t+1} | a_t, s_{0:t}, a_{0:t}) \\ &= \sum_{t=0}^{\tau-1} \nabla_\theta \log \pi_\theta(a_t | s_t) \end{aligned}$$

# Decomposing the Trajectories Into States and Actions

- Approximate using  $m$  sample paths under policy  $\pi_\theta$ :

$$\nabla_\theta V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_\theta \log P(\tau^{(i)})$$

$$\begin{aligned} \nabla_\theta \log P(\tau^{(i)}; \theta) &= \nabla_\theta \log \left[ \underbrace{\mu(s_0)}_{\text{Initial state distrib.}} \prod_{t=0}^{T-1} \underbrace{\pi_\theta(a_t|s_t)}_{\text{policy}} \underbrace{P(s_{t+1}|s_t, a_t)}_{\text{dynamics model}} \right] \\ &= \nabla_\theta \left[ \log \mu(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t|s_t) + \log P(s_{t+1}|s_t, a_t) \right] \\ &= \sum_{t=0}^{T-1} \underbrace{\nabla_\theta \log \pi_\theta(a_t|s_t)}_{\text{no dynamics model required!}} \end{aligned}$$

# Decomposing the Trajectories Into States and Actions

- Approximate using  $m$  sample paths under policy  $\pi_\theta$ :

$$\nabla_\theta V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_\theta \log P(\tau^{(i)})$$

$$\begin{aligned} \nabla_\theta \log P(\tau^{(i)}; \theta) &= \nabla_\theta \log \left[ \underbrace{\mu(s_0)}_{\text{Initial state distrib.}} \prod_{t=0}^{T-1} \underbrace{\pi_\theta(a_t|s_t)}_{\text{policy}} \underbrace{P(s_{t+1}|s_t, a_t)}_{\text{dynamics model}} \right] \\ &= \nabla_\theta \left[ \log \mu(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t|s_t) + \log P(s_{t+1}|s_t, a_t) \right] \\ &= \sum_{t=0}^{T-1} \underbrace{\nabla_\theta \log \pi_\theta(a_t|s_t)}_{\text{score function}} \end{aligned}$$



# Score Function

- A score function is the derivative of the log of a parameterized probability / likelihood
- Example: let  $\pi(s; \theta)$  be the probability of state  $s$  under parameter  $\theta$
- Then the score function would be

$$\nabla_{\theta} \log \pi(s; \theta) \tag{1}$$

- For many policy classes, it is not hard to compute the score function

# Softmax Policy

- Weight actions using linear combination of features  $\phi(s, a)^T \theta$
- Probability of action is proportional to exponentiated weight

$$\pi_{\theta}(s, a) = e^{\phi(s, a)^T \theta} / \left( \sum_a e^{\phi(s, a)^T \theta} \right)$$

• The score function is  $\nabla_{\theta} \log \pi_{\theta}(s, a) =$

$$\begin{aligned} & \nabla_{\theta} \left[ \log e^{\phi(s, a)^T \theta} / \sum_a e^{\phi(s, a)^T \theta} \right] \\ &= \nabla_{\theta} \left[ \phi(s, a)^T \theta - \log \sum_a e^{\phi(s, a)^T \theta} \right] \\ &= \phi(s, a) - \frac{1}{\sum_a e^{\phi(s, a)^T \theta}} \sum_a \nabla_{\theta} e^{\phi(s, a)^T \theta} \\ &= \phi(s, a) - \frac{1}{\sum_a e^{\phi(s, a)^T \theta}} \sum_a \phi(s, a) e^{\phi(s, a)^T \theta} \end{aligned}$$

# Softmax Policy

- Weight actions using linear combination of features  $\phi(s, a)^T \theta$
- Probability of action is proportional to exponentiated weight

$$\pi_{\theta}(s, a) = e^{\phi(s, a)^T \theta} / \left( \sum_a e^{\phi(s, a)^T \theta} \right)$$

- The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \nabla_{\theta} \left[ \log \left[ \frac{e^{\phi(s, a)^T \theta}}{\sum_a e^{\phi(s, a)^T \theta}} \right] \right] \quad (2)$$

$$= \nabla_{\theta} \left[ \phi(s, a)^T \theta - \log \left[ \sum_a e^{\phi(s, a)^T \theta} \right] \right] \quad (3)$$

$$= \phi(s, a) - \frac{\sum_a \nabla_{\theta} e^{\phi(s, a)^T \theta}}{\sum_a e^{\phi(s, a)^T \theta}} \quad (4)$$

$$= \phi(s, a) - \frac{\sum_a \phi(s, a) e^{\phi(s, a)^T \theta}}{\sum_a e^{\phi(s, a)^T \theta}} \quad (5)$$

$$= \phi(s, a) - \mathbb{E}_{\pi_{\theta}}[\phi(s, \cdot)] \quad (6)$$

- Weight actions using linear combination of features  $\phi(s, a)^T \theta$
- Probability of action is proportional to exponentiated weight

$$\pi_{\theta}(s, a) = e^{\phi(s, a)^T \theta} / (\sum_a e^{\phi(s, a)^T \theta})$$

- The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \phi(s, a) - \mathbb{E}_{\pi_{\theta}}[\phi(s, \cdot)]$$

# Gaussian Policy

- In continuous action spaces, a Gaussian policy is natural
- Mean is a linear combination of state features  $\mu(s) = \phi(s)^T \theta$
- Variance may be fixed  $\sigma^2$ , or can also be parametrised
- Policy is Gaussian  $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$

- Deep neural networks (and other models where can compute the gradient) can also be used to represent the policy

# Likelihood Ratio / Score Function Policy Gradient

- Putting this together
- Goal is to find the policy parameters  $\theta$ :

$$\arg \max_{\theta} V(\theta) = \arg \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Approximate with empirical estimate for  $m$  sample paths under policy  $\pi_{\theta}$  using score function:

$$\begin{aligned} \nabla_{\theta} V(\theta) &\approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta) \\ &= (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \end{aligned}$$

- Do not need to know dynamics model

# L5N2 Check Your Understanding L5: Score functions

$$\nabla_{\theta} V(\theta) = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

The likelihood ratio / score function policy gradient (select one):

- (a) requires reward functions that are differentiable
- (b) can only be used with Markov decision processes
- (c) Is useful mostly for infinite horizon tasks
- (a) and (b)
- a,b and c
- None of the above
- Not sure

# L5N2 Check Your Understanding L5: Score functions Solution

$$\nabla_{\theta} V(\theta) = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

The likelihood ratio / score function policy gradient (select one):

- (a) requires reward functions that are differentiable
- (b) can only be used with Markov decision processes
- (c) Is useful mostly for infinite horizon tasks
- (a) and (b)
- a,b and c
- None of the above
- Not sure

None of the above

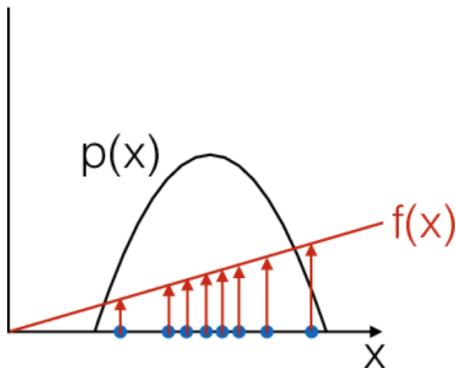


# Score Function Gradient Estimator: Intuition

- Consider generic form of  $R(\tau^{(i)})\nabla_{\theta} \log P(\tau^{(i)}; \theta)$ :  
 $\hat{g}_i = f(x_i)\nabla_{\theta} \log p(x_i|\theta)$
- $f(x)$  measures how good the sample  $x$  is.
- Moving in the direction  $\hat{g}_i$  pushes up the logprob of the sample, in proportion to how good it is
- *Valid even if  $f(x)$  is discontinuous, and unknown, or sample space (containing  $x$ ) is a discrete set*

# Score Function Gradient Estimator: Intuition

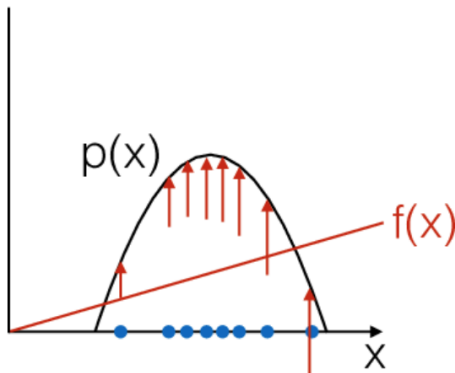
$$\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i | \theta)$$



# Score Function Gradient Estimator: Intuition

*for vs. reward?*

$$\hat{g}_i = f(x_i) \nabla_{\theta} \log p(x_i | \theta)$$



# Policy Gradient Theorem

- The policy gradient theorem generalizes the likelihood ratio approach

## Theorem

*For any differentiable policy  $\pi_\theta(s, a)$ ,  
for any of the policy objective function  $J = J_1$ , (episodic reward),  $J_{avR}$   
(average reward per time step), or  $\frac{1}{1-\gamma} J_{avV}$  (average value),  
the policy gradient is*

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)]$$

- Chapter 13.2 in SB has a nice derivation of the policy gradient theorem for episodic tasks and discrete states

# Table of Contents

## 4 Policy Gradient Algorithms and Reducing Variance

- Temporal Structure
- Baseline
- Alternatives to MC Returns

# Likelihood Ratio / Score Function Policy Gradient

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Unbiased but very noisy
- Fixes that can make it practical
  - Temporal structure
  - Baseline

# Policy Gradient: Use Temporal Structure

- Previously:

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

- We can repeat the same argument to derive the gradient estimator for a single reward term  $r_{t'}$ .

viewed at time step  $t'$

$$\nabla_{\theta} \mathbb{E}[r_{t'}] = \mathbb{E} \left[ \underbrace{r_{t'}}_{\text{viewed at time step } t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

$r_2 \mathbb{P}(a_2 | s_0)$

- To see this, recall  $V(s_0, \theta) = \mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^T R(s_t, a_t); \pi_{\theta}, s_0 \right]$  where the expectation is taken over the states & actions visited by  $\pi_{\theta}$
- $P(\gamma) R(\gamma)$ 
 $s_0, a_1, r_1, s_2, a_2, r_2, \dots$

# Policy Gradient: Use Temporal Structure

- Previously:

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

- We can repeat the same argument to derive the gradient estimator for a single reward term  $r_{t'}$ .

$$\nabla_{\theta} \mathbb{E}[r_{t'}] = \mathbb{E} \left[ r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- Summing this formula over t, we obtain

$$V(\theta) = \nabla_{\theta} \mathbb{E}[R] = \mathbb{E} \left[ \sum_{t'=0}^{T-1} r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

eqn



# Policy Gradient: Use Temporal Structure

- Previously:  $V(\theta) = \mathbb{E}_{\tau} [R(\tau) | \theta]$

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

- We can repeat the same argument to derive the gradient estimator for a single reward term  $r_{t'}$ .

$$\nabla_{\theta} \mathbb{E}[r_{t'}] = \mathbb{E} \left[ r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

- Summing this formula over  $t$ , we obtain

$$\begin{aligned} \nabla_{\theta} V(\theta) &= \nabla_{\theta} \mathbb{E}[R] = \mathbb{E} \left[ \sum_{t'=0}^{T-1} r_{t'} \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right] \\ &= \mathbb{E} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t, s_t) \sum_{t'=t}^{T-1} r_{t'} \right] \end{aligned}$$

later decisions don't influence past rewards

reduce  $\sum_{t'=t}^{T-1} r_{t'} *$

# Policy Gradient: Use Temporal Structure

- Recall for a particular trajectory  $\tau^{(i)}$ ,  $\sum_{t'=t}^{T-1} r_{t'}^{(i)}$  is the return  $G_t^{(i)}$

$$\nabla_{\theta} \mathbb{E}[R] \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t, s_t) G_t^{(i)}$$

# Monte-Carlo Policy Gradient (REINFORCE)

- Leverages likelihood ratio / score function and temporal structure

$$\Delta\theta_t = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$$

## **REINFORCE:**

Initialize policy parameters  $\theta$  arbitrarily

**for** each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_{\theta}$  **do**

**for**  $t = 1$  to  $T - 1$  **do**

$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) G_t$

**endfor**

**endfor**

**return**  $\theta$

# Likelihood Ratio / Score Function Policy Gradient

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m R(\tau^{(i)}) \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Unbiased but very noisy
- Fixes that can make it practical
  - Temporal structure
  - **Baseline**
  - Alternatives to using Monte Carlo returns  $R(\tau^{(i)})$  as targets

# Desired Properties of a Policy Gradient RL Algorithm

- Goal: Converge as quickly as possible to a local optima
  - Incurring reward / cost as execute policy, so want to minimize number of iterations / time steps until reach a good policy

# Table of Contents

- Temporal Structure

## 5 Policy Gradient Algorithms and Reducing Variance

- Baseline
- Alternatives to MC Returns

# Policy Gradient: Introduce Baseline

- Reduce variance by introducing a *baseline*  $b(s)$

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left( \sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right]$$

- For any choice of  $b$ , gradient estimator is unbiased.
  - Near optimal choice is the expected return,
- (if  $b$  is only a function of  $s$ )

$$b(s_t) \approx \mathbb{E}[r_t + r_{t+1} + \dots + r_{T-1}]$$

- Interpretation: increase logprob of action  $a_t$  proportionally to how much returns  $\sum_{t'=t}^{T-1} r_{t'}$  are better than expected

# Baseline $b(s)$ Does Not Introduce Bias—Derivation

$$\mathbb{E}_\tau [\nabla_\theta \log \pi(a_t | s_t; \theta) b(s_t)]$$

$$= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [\mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_\theta \log \pi(a_t | s_t; \theta) b(s_t)]]$$

$$= \mathbb{E}_{s_{0:t}, a_{0:t-1}} [b(s_t) \mathbb{E}_{s_{(t+1):T}, a_{t:T}} [\nabla_\theta \log \pi(a_t | s_t, \theta)]]$$

$$= \quad " \quad [b(s_t) \mathbb{E}_{a_t} \nabla_\theta \log \pi(a_t | s_t, \theta)]$$

$$= \quad " \quad [b(s_t) \mathbb{E}_{a_t} \frac{1}{\pi(a_t | s_t, \theta)} \nabla_\theta \pi(a_t | s_t, \theta)]$$

$$= \quad " \quad [b(s_t) \sum_{a_t} \pi(a_t | s_t, \theta) \frac{1}{\pi(a_t | s_t, \theta)} \nabla_\theta \pi(a_t | s_t, \theta)]$$

$$= \quad " \quad [b(s_t) \sum_{a_t} \nabla_\theta \pi(a_t | s_t, \theta)]$$

$$= \quad " \quad [b(s_t) \nabla_\theta \underbrace{\sum_{a_t} \pi(a_t | s_t, \theta)}_{=1}]$$



# Baseline $b(s)$ Does Not Introduce Bias–Derivation

$$\begin{aligned} & \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi(a_t | s_t; \theta) b(s_t)] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ \mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_{\theta} \log \pi(a_t | s_t; \theta) b(s_t)] \right] \text{ (break up expectation)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_{\theta} \log \pi(a_t | s_t; \theta)] \right] \text{ (pull baseline term out)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \mathbb{E}_{a_t} [\nabla_{\theta} \log \pi(a_t | s_t; \theta)] \right] \text{ (remove irrelevant variables)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \sum_a \pi_{\theta}(a_t | s_t) \frac{\nabla_{\theta} \pi(a_t | s_t; \theta)}{\pi_{\theta}(a_t | s_t)} \right] \text{ (likelihood ratio)} \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \sum_a \nabla_{\theta} \pi(a_t | s_t; \theta) \right] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} \left[ b(s_t) \nabla_{\theta} \sum_a \pi(a_t | s_t; \theta) \right] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [b(s_t) \nabla_{\theta} 1] \\ &= \mathbb{E}_{s_{0:t}, a_{0:(t-1)}} [b(s_t) \cdot 0] = 0 \end{aligned}$$

# "Vanilla" Policy Gradient Algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

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

Collect a set of trajectories by executing the current policy

At each timestep  $t$  in each trajectory  $\tau^i$ , compute

Return  $G_t^i = \sum_{t'=t}^{T-1} r_{t'}^i$ , and

Advantage estimate  $\hat{A}_t^i = G_t^i - b(s_t)$ .

Re-fit the baseline, by minimizing  $\sum_i \sum_t \|b(s_t) - G_t^i\|^2$ ,

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

Which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$ .

(Plug  $\hat{g}$  into SGD or ADAM)

**endfor**

# Other Choices for Baseline?

Initialize policy parameter  $\theta$ , baseline  $b$

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

Collect a set of trajectories by executing the current policy

At each timestep  $t$  in each trajectory  $\tau^i$ , compute

Return  $G_t^i = \sum_{t'=t}^{T-1} r_{t'}^i$ , and

Advantage estimate  $\hat{A}_t^i = G_t^i - b(s_t)$ .

Re-fit the baseline, by minimizing  $\sum_i \sum_t \|b(s_t) - G_t^i\|^2$ ,

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

Which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$ .

(Plug  $\hat{g}$  into SGD or ADAM)

**endfor**

# Choosing the Baseline: Value Functions

- Recall Q-function / state-action-value function:

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi} [r_0 + \gamma r_1 + \gamma^2 r_2 \cdots | s_0 = s, a_0 = a]$$

- State-value function can serve as a great baseline

$$\begin{aligned} V^{\pi}(s) &= \mathbb{E}_{\pi} [r_0 + \gamma r_1 + \gamma^2 r_2 \cdots | s_0 = s] \\ &= \mathbb{E}_{a \sim \pi} [Q^{\pi}(s, a)] \end{aligned}$$

# Table of Contents

- Temporal Structure
- Baseline

## 6 Policy Gradient Algorithms and Reducing Variance

- Alternatives to MC Returns

- Policy gradient:

$$\nabla_{\theta} \mathbb{E}[R] \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t, s_t) (G_t^{(i)} - b(s_t))$$

- Fixes that improve simplest estimator
  - Temporal structure (shown in above equation)
  - Baseline (shown in above equation)
  - **Alternatives to using Monte Carlo returns  $G_t^i$  as estimate of expected discounted sum of returns for the policy parameterized by  $\theta$ ?**

# Choosing the Target

- $G_t^i$  is an estimation of the value function at  $s_t$  from a single roll out
- Unbiased but high variance
- Reduce variance by introducing bias using bootstrapping and function approximation
  - Just like we saw for TD vs MC, and value function approximation

- Estimate of  $V/Q$  is done by a **critic**
- **Actor-critic** methods maintain an explicit representation of policy and the value function, and update both
- A3C (Mnih et al. ICML 2016) is a very popular actor-critic method



# Policy Gradient Formulas with Value Functions

- Recall:

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left( \sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right]$$

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] \approx \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) (Q(s_t, a_t; \mathbf{w}) - b(s_t)) \right]$$

- Letting the baseline be an estimate of the value  $V$ , we can represent the gradient in terms of the state-action advantage function

$$\nabla_{\theta} \mathbb{E}_{\tau}[R] \approx \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \hat{A}^{\pi}(s_t, a_t) \right]$$

- where the advantage function  $A^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$

# Choosing the Target: N-step estimators

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Note that critic can select any blend between TD and MC estimators for the target to substitute for the true state-action value function.

# Choosing the Target: N-step estimators

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- Note that critic can select any blend between TD and MC estimators for the target to substitute for the true state-action value function.

$$\hat{R}_t^{(1)} = r_t + \gamma V(s_{t+1})$$

$$\hat{R}_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 V(s_{t+2}) \quad \dots$$

$$\hat{R}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

- If subtract baselines from the above, get advantage estimators

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

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots - V(s_t)$$

# L5N3 Check Your Understanding: Blended Advantage Estimators

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- If subtract baselines from the above, get advantage estimators

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

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} + \dots - V(s_t)$$

- Select all that are true
- $\hat{A}_t^{(1)}$  has low variance & low bias.
- $\hat{A}_t^{(1)}$  has high variance & low bias.
- $\hat{A}_t^{(\infty)}$  low variance and high bias.
- $\hat{A}_t^{(\infty)}$  high variance and low bias.
- Not sure

# LN3 Check Your Understanding: Blended Advantage Estimators Answers

$$\nabla_{\theta} V(\theta) \approx (1/m) \sum_{i=1}^m \sum_{t=0}^{T-1} R_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)})$$

- If subtract baselines from the above, get advantage estimators

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

$$\hat{A}_t^{(\text{inf})} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+1} + \dots - V(s_t)$$

Solution:  $\hat{A}_t^{(1)}$  has low variance & high bias.  $\hat{A}_t^{(\infty)}$  high variance but low bias.

# "Vanilla" Policy Gradient Algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

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

Collect a set of trajectories by executing the current policy

At each timestep  $t$  in each trajectory  $\tau^i$ , compute

*Advantage estimate*  $\hat{A}_{it}^n$

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

Which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_{it}^n$ .

**(Plug  $\hat{g}$  into SGD or ADAM)**

**endfor**

- Note, we can choose which blended estimator  $\hat{A}^n$  to use

# Current Summary of Benefits of Policy-Based RL

## Advantages:

- Better convergence properties
- Effective in high-dimensional or continuous action spaces
- Can learn stochastic policies

## Disadvantages:

- Typically converge to a local rather than global optimum
- Evaluating a policy can be inefficient and high variance (though baseline and temporal structure helps)

# Class Structure

- Last time: Deep Model-free Value Based RL
- **This time: Policy Search**
- Next time: Policy Search Cont.