

Policy Gradients methods

Ciprian Paduraru

University of Bucharest

Agenda

Part 1 (today):

- Recap & Motivation
- Policy objective and optimization
- Reinforce algorithm
- A3C and GAE
- Results

Part 2

- Other different ways to choose the policy network – gaussian, softmax, etc.
- Address the problem when the policy network is not differentiable
- A more formal look at Policy Theorem
- Examples from practice, real implementations

Part 3

- Importance of optimization method, step size tuning
- Use Importance sampling to have monotonic improvenets
- Lower bounds for local Approximation, Trust Regions, TRPO Algorithm
- PPO
- A look at Alpha-Go

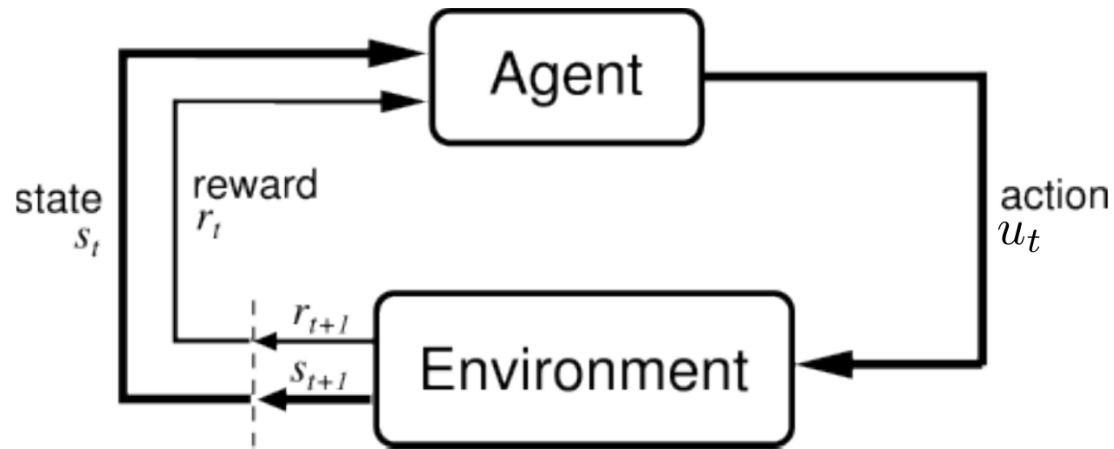
Lab 1: Experiment with base class algorithms on Cart-Pole

Lab 2: Experiment with various parameters on DeepMimic environment

1. Short recap

Last lectures:

- The general approach:



[Figure source: Sutton & Barto, 1998]

- Have looked at (exact or approximate) approaches for finding $V(s)$, $Q(s, a)$
- A policy was generated directly from the value function
$$\pi(s) = \arg \max_a Q(s, a)$$
- In this lecture, we'll **directly parameterize the policy** (Policy-based RL)
- Focus will still be on **model-free reinforcement learning**

- **Value-based**

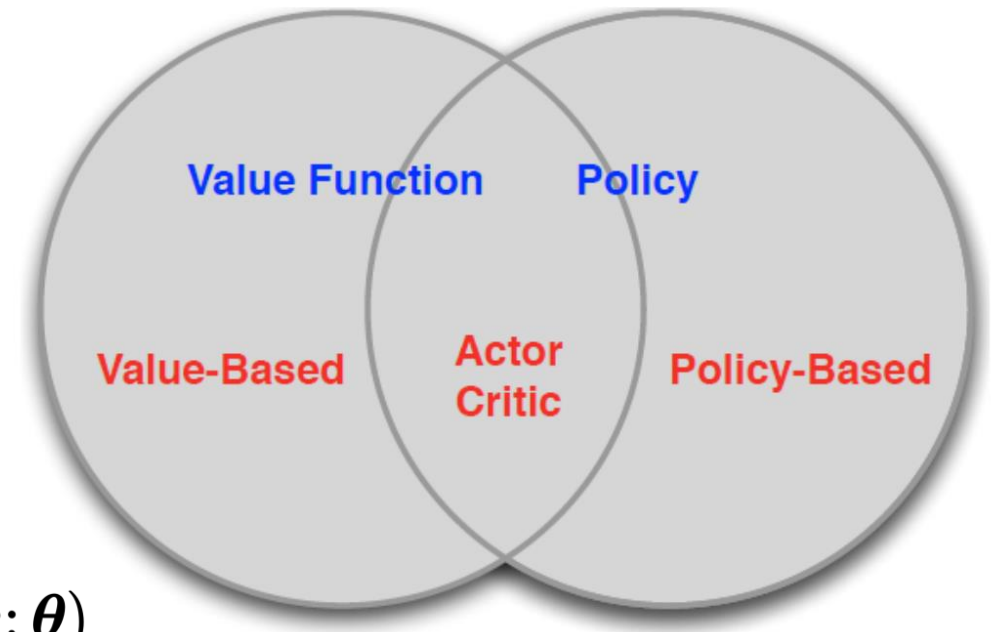
- Learn $V(s)$, $Q(s, a)$
- Extract the policy from $V(s)$, $Q(s, a)$

- **Policy-based**

- No intermediate learning of value functions
- Directly learn the policy $\pi_{\theta}(a | s) = P(a | s; \theta)$

- **Actor-Critic** = Combination of both worlds

- Learn $V(s)$, $Q(s, a)$ -> **CRITIC**
- Learn Policy $\pi_{\theta}(a | s)$ -> **ACTOR**
- Note: many times, better than pure Policy-based approaches !



2. Properties of Policy-based RL, motivation

Advantages

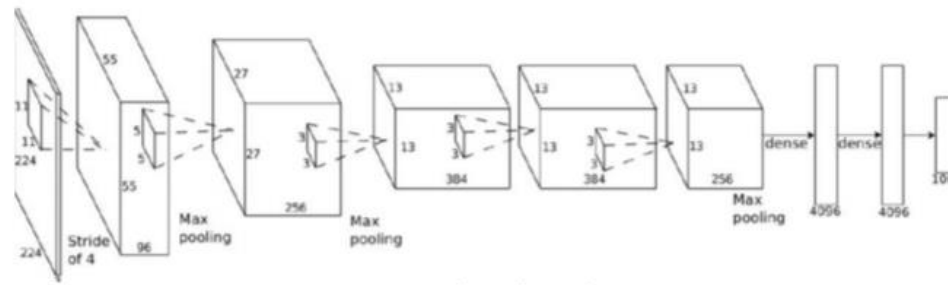
- Better convergence properties – more details in The Sutton and Barto reference, Ch 13.3
- Effectiveness in high-dimensional or continuous action spaces

Examples:

- Robotics: Isn't it harder to compute $V(s)$, $Q(s,a)$ for every state / action of a robot rather than a policy ?
- In IRL (Inverse Reinforcement Learning) settings,
isn't it simpler to think in terms of a policy directly instead of computing the value for each state ?
Think about the self-driving cars field !



\mathbf{o}_t



$\pi_{\theta}(\mathbf{a}_t | \mathbf{o}_t)$



\mathbf{a}_t

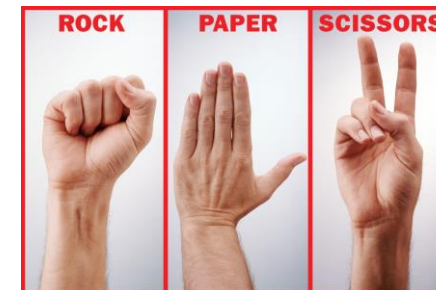
- Can learn **stochastic policies** (remember previously with MDPs we had deterministic policies).

Why is a stochastic policy needed ?

- State representation is not Markov
- Adversarial / non-stationary domain problems

Example: think about *rock-paper-scissors* game. Can someone exploit a deterministic strategy ? 😊.

<https://www.timeforkids.com/k1/rock-paper-scissors/>



Disadvantages

- Typically converges to a local instead of global optimum (however also true for value-based RL with function approximation)
- Evaluating policy (expected reward) is typically inefficient and high variance

3. Policy optimization

- Denote **trajectory** $\tau = (s_0, u_0, r_1, s_2, u_1, r_2, \dots, u_{H-1}, r_H, s_H, \dots)$
- Actions along τ are taken using policy $\pi_\theta(a|s)$, based on parameters θ used by a function π_θ
- Methods to evaluate the policy:

➤ In episodic environments, can use the *start value* of the policy:

$$J(\theta) = V^{\pi_\theta}(s_1) = E_{\pi_\theta}[v_1]$$

➤ In continuing environments, stationary distribution of π_θ over states can be used: $d^{\pi_\theta}(s)$

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

- Note: We'll discuss the episodic case, but all results can be reused in the non-episodic ones
- Also, considering finite horizons of length H.

3. Policy optimization

- **Trajectory** $\tau = (s_0, u_0, r_1, s_2, u_1, r_2, \dots, u_{H-1}, r_H, s_H, \dots)$
- **Reward** of a trajectory: $R(\tau) = \sum_{t=0}^H R(s_t, u_t)$
- **Objective** function (expected sum of rewards along trajectories samples from π_θ):

$$U(\theta) = \mathbb{E} \left[\sum_{t=0}^H R(s_t, u_t); \pi_\theta \right] = \sum_{\tau} P(\tau; \theta) R(\tau)$$

- Our **goal** find the optimal θ :

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

3.1 Non-differentiable optimization methods

- Having an objective function and a parameters space, we can treat the policy-based RL problem as an optimization one.
- There are some **gradient free methods** such as :
 - Hill climbing
 - Simplex / amoeba / Nelder Mead
 - Genetic algorithms
 - Cross-Entropy method (CEM)
 - Covariance Matrix Adaptation (CMA)
 - Evolution strategies
- **Their advantage:** These allows the policy parametrization to be non-differentiable and are often easy to parallelize


Disadvantage: these methods ignore the temporal structure of rewards

- Updates consider only the total episodes' rewards, do not break up the reward for each state in the trajectory !

3.2 Policy gradient

- We assume policy π_θ is differentiable whenever it is non-zero

Softmax Policy class

- Weight actions by using l.c. of features and parameters $\phi(s, a)^T \theta$
- Probability of actions , $\pi_\theta(s, a) = \frac{e^{\phi(s, a)^T \theta}}{\left(\sum_a e^{\phi(s, a)^T \theta}\right)}$  $\nabla_\theta \log \pi_\theta(s, a) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, \cdot)]$

Full proof:

Using the log identity $\log(x/y) = \log(x) - \log(y)$ we can write

$$\log(\pi_\theta(s, a)) = \log(e^{\phi(s, a)^T \theta}) - \log\left(\sum_{k=1}^N e^{\phi(s, a_k)^T \theta}\right)$$

“Likelihood Ratio” Policy Gradient

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

- We take the gradient with respect to θ

$$\begin{aligned} \nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \end{aligned}$$

“Likelihood Ratio” Policy Gradient

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

- We take the gradient with respect to θ

$$\begin{aligned} \nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \end{aligned}$$

“Likelihood Ratio” Policy Gradient

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

- We take the gradient with respect to θ

$$\begin{aligned} \nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} R(\tau) \end{aligned}$$

“Likelihood Ratio” Policy Gradient

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

- We take the gradient with respect to θ

$$\begin{aligned} \nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} R(\tau) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau) \end{aligned}$$

“Likelihood Ratio” Policy Gradient

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

- We take the gradient with respect to θ

$$\begin{aligned} \nabla_{\theta} U(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} R(\tau) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau) \\ &= \mathbb{E}_{\tau} [\nabla_{\theta} \log P(\tau; \theta) R(\tau)] \end{aligned}$$

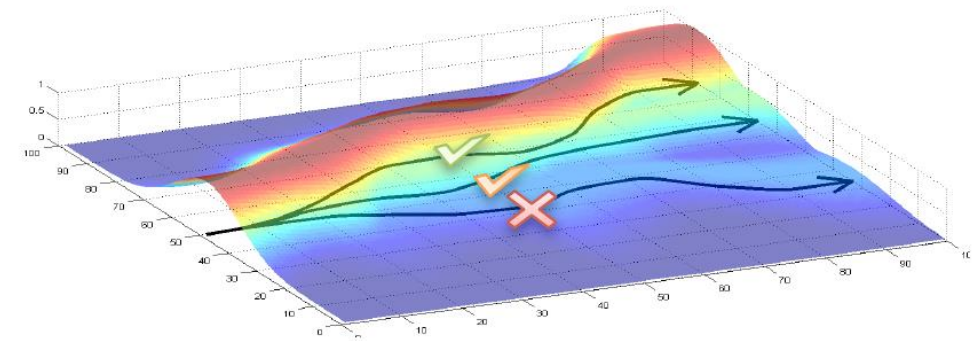
$$\nabla_{\theta} U(\theta) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau) = \mathbb{E}_{\tau} [\nabla_{\theta} \log P(\tau; \theta) R(\tau)]$$

- We can now approximate with the empirical estimate for m sample paths under policy – Monte Carlo (MC) sampling method from the previous courses !

- Collect m episodes then compute the average

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

- **Intuition** - The gradient tries to:
 - Increase the probability of paths with positive reward
 - Decrease the probability of paths with negative rewards



- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \underbrace{\nabla_{\theta} \log P(\tau^{(i)}; \theta)}_{\text{Let's take a step closer at this}} R(\tau^{(i)})$$

Let's take a step closer at this

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \underbrace{\nabla_{\theta} \log P(\tau^{(i)}; \theta)}_{\text{Let's take a step closer at this}} R(\tau^{(i)})$$

Let's take a step closer at this

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$



Initial state distribution



Model's dynamic



Policy samples

- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

(using log property of products to sum)

$$= \nabla_{\theta} \left[\mu(s_0) + \sum_{t=0}^{H-1} \log P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

$$\text{(using log property of products to sum)} \quad = \nabla_{\theta} \left[\mu(s_0) + \sum_{t=0}^{H-1} \log P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

$$\text{(first two don't depend on } \theta) \quad = \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)})$$

- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

$$\text{(using log property of products to sum)} \quad = \nabla_{\theta} \left[\mu(s_0) + \sum_{t=0}^{H-1} \log P(s_{t+1}^{(i)} \mid s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)}) \right]$$

$$\text{(first two don't depend on } \theta \text{)} \quad = \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)})$$

$$= \sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} \mid s_t^{(i)})$$

- Let's go further: Decompose paths into states and actions along trajectory $\tau = (s_0, u_0, r_1, s_1, u_1, r_2, \dots, u_{H-1}, r_H, s_H)$

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\mu(s_0) \cdot \prod_{t=0}^{H-1} P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)}) \cdot \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right]$$

(using log property of products to sum)

$$= \nabla_{\theta} \left[\mu(s_0) + \sum_{t=0}^{H-1} \log P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right]$$

(first two don't depend on θ)

$$= \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})$$

$$= \sum_{t=0}^H \underbrace{\nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{Also named } \textit{score function}}$$

Note: No model dynamics is no longer required !

No initial state distribution !

Let's recap !

- We have an **unbiased** estimate of the gradient
- **No need to access model dynamics**

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P\left(\tau^{(i)}; \theta\right) R\left(\tau^{(i)}\right)$$

- **Unbiased but very noisy** (remember the MC properties from previous courses)
- There are a couple of ways to address this and reduce variance in practice:
 - Exploit the temporal structure
 - Baseline method
 - Next course: KL-Divergence trust region / natural gradient

Optimization 1: Exploit the temporal structure

- Current estimation: $\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P \left(\tau^{(i)}; \theta \right) R \left(\tau^{(i)} \right)$

Optimization 1: Exploit the temporal structure

- Current estimation:
$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P \left(\tau^{(i)}; \theta \right) R \left(\tau^{(i)} \right)$$

(expanding the trajectory's steps as before)
$$= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \right) \left(\sum_{t=0}^{H-1} R \left(s_t^{(i)}, u_t^{(i)} \right) \right)$$

Optimization 1: Exploit the temporal structure

- Current estimation:
$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P \left(\tau^{(i)}; \theta \right) R \left(\tau^{(i)} \right)$$

(expanding the trajectory's steps as before)
$$= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \right) \left(\sum_{t=0}^{H-1} R \left(s_t^{(i)}, u_t^{(i)} \right) \right)$$

(splitting rewards before/after t)
$$= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left[\left(\sum_{k=0}^{t-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) + \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right] \right)$$

Any observation

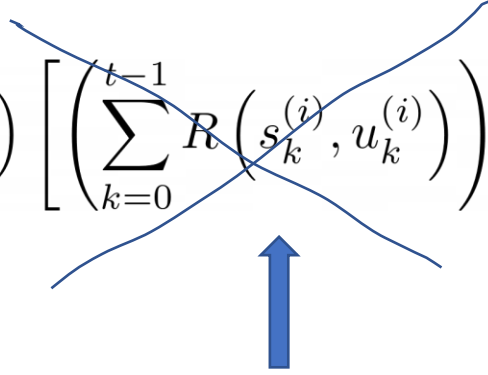
?

Optimization 1: Exploit the temporal structure

- Current estimation: $\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P \left(\tau^{(i)}; \theta \right) R \left(\tau^{(i)} \right)$

(expanding the trajectory's steps as before) $= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \right) \left(\sum_{t=0}^{H-1} R \left(s_t^{(i)}, u_t^{(i)} \right) \right)$

(splitting rewards before/after t) $= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left[\left(\sum_{k=0}^{t-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) + \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right] \right)$



This term (previous rewards) do not depend on u_t !
 Whatever probability we choose now, this remains
 fixed ! (More formal definition in seminar)

- Improved version, less variance $\hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$

REINFORCE algorithm

REINFORCE (Williams, 1992).

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$$

REINFORCE algorithm

REINFORCE (Williams, 1992).

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \underbrace{\left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right)}_{G_t \text{ notation from MC course}} \right)$$

- At each trajectory, and step we could update the parameters of the policy by:

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

REINFORCE algorithm

REINFORCE (Williams, 1992).

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \underbrace{\left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right)}_{G_t} \right)$$

G_t notation from MC course

Pseudocode

Initialize policy parameters θ 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 θ

3.3 Differentiable policy classes

- Assume policy π_θ is differentiable whenever it is non-zero
- Many choices for differentiable policy classes. Popular ones:
 - Softmax
 - Gaussian
 - Neural networks

Softmax Policy class – good for *discrete action spaces*

- Weight actions by using a linear comb of features and parameters $\phi(s, a)^T \theta$
- Probability of actions , $\pi_\theta(s, a) = \frac{e^{\phi(s, a)^T \theta}}{\left(\sum_a e^{\phi(s, a)^T \theta}\right)}$  $\nabla_\theta \log \pi_\theta(s, a) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, \cdot)]$

Full proof:

Using the log identity $\log(x/y) = \log(x) - \log(y)$ we can write

$$\log(\pi_\theta(s, a)) = \log(e^{\phi(s, a)^T \theta}) - \log\left(\sum_{k=1}^N e^{\phi(s, a_k)^T \theta}\right)$$

$$\nabla_{\theta} \log(\pi_{\theta}(s, a)) = \nabla_{\theta} \log(e^{\phi(s, a)^{\top} \theta}) - \nabla_{\theta} \log\left(\sum_{k=1}^N e^{\phi(s, a_k)^{\top} \theta}\right)$$

$$left = \nabla_{\theta} \log(e^{\phi(s, a)^{\top} \theta}) = \nabla_{\theta} \phi(s, a)^{\top} \theta = \phi(s, a)$$

The right term simplifies as follows:

Using the chain rule:

$$\nabla_x \log(f(x)) = \frac{\nabla_x f(x)}{f(x)}$$

We can write:

$$right = \nabla_{\theta} \log\left(\sum_{k=1}^N e^{\phi(s, a_k)^{\top} \theta}\right) = \frac{\nabla_{\theta} \sum_{k=1}^N e^{\phi(s, a_k)^{\top} \theta}}{\sum_{k=1}^N e^{\phi(s, a_k)^{\top} \theta}}$$

Taking the gradient of the numerator we get:

$$right = \frac{\sum_{k=1}^N \phi(s, a_k) e^{\phi(s, a_k)^{\top} \theta}}{\sum_{k=1}^N e^{\phi(s, a_k)^{\top} \theta}}$$

Substituting the definition of $\pi_{\theta}(s, a)$ we can simplify to:

$$right = \sum_{k=1}^N \phi(s, a_k) \pi_{\theta}(s, a_k)$$

Given the definition of Expected Value:

$$E[X] = X \cdot P = x_1 p_1 + x_2 p_2 + \dots + x_n p_n$$

Which in English is just the sum of each feature times its probability.

$$X = \text{features} = \phi(s, a)$$

$$P = \text{probabilities} = \pi_{\theta}(s, a)$$

So now we can write the expected value of the features:

$$right = E_{\pi_{\theta}}[\phi(s, \cdot)]$$

where \cdot means all possible actions.

Putting it all together:

$$\nabla_{\theta} \log(\pi_{\theta}(s, a)) = left - right = \phi(s, a) - E_{\pi_{\theta}}[\phi(s, \cdot)]$$

Interpretation: difference between the features for certain states and each individual action minus mean of features for the same state and all possible actions.

Gaussian Policy - *continuous action spaces* (e.g. robotics)

- Mean could be a linear combination of state features $\mu(s) = \phi(s)^T \theta$
- Variance σ^2 could be parametrized or fixed
- Policy is then a Gaussian such that actions are sampled:

$$a \sim \mathcal{N}(\mu(s), \sigma^2)$$

- Derivative (i.e., score function) becomes:

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

- Note: Nice interpretation again ! An action is more probable if close to the mean estimated by the parameters of the state.
 - It works for continuous cases !
- Neural networks – most common

Optimization 2: Baseline trick

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$$

Optimization 2: Baseline trick

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$$



Idea: Extract a **baseline value** to improve the variance further

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$

Optimization 2: Baseline trick

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$$



Idea: Extract a **baseline value** to improve the variance further

$$\begin{aligned} \nabla U(\theta) \approx \hat{g} &= \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right) \\ &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t) \left(\sum_{k=t}^{H-1} R(s_k, u_k) - b(s_t) \right) \right] \text{ (Just coming back a bit to } E \text{ form)} \end{aligned}$$

Optimization 2: Baseline trick

- We are here:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) \right) \right)$$



Idea: Extract a **baseline value** to improve the variance further

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$

- We are here:

(Just coming back a bit to E form)

- For any choice of baseline $b(s_t)$, gradient estimator is unbiased.
NOTE: it should only depend on previous states, not future, i.e., those affected by future decisions

- For any choice of baseline $b(s_t)$, gradient estimator is unbiased.

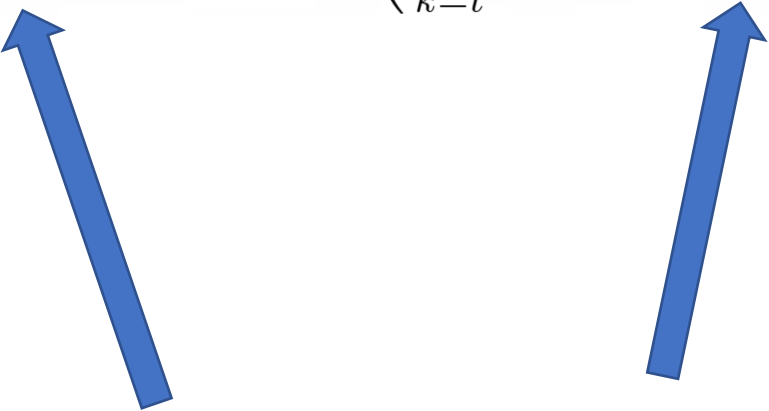
NOTE: it should only depend on previous states, not future, i.e., those affected by future decisions

We must prove that:
$$\mathbb{E}_{\tau} \left[\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t) b(s_t) \right] = 0 \quad \text{or even simpler, that} \quad \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t) b(s_t)] = 0$$

Proof that by introducing baseline gradient is still unbiased:

$$\begin{aligned} & \mathbb{E}_{\tau} [\nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t) b(s_t)] \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} [\mathbb{E}_{s_{(t+1):T}, a_{t:(T-1)}} [\nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t) b(s_t)]] \text{ (split expectation before and after } t \text{)} \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} [b(s_t) \mathbb{E}_{s_{(t+1):H}, a_{t:(H-1)}} [\nabla_{\theta} \log \pi_{\theta} (u_t \mid s_t)]] \text{ (baseline doesn't affect inner E)} \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} [b(s_t) \mathbb{E}_{u_t} [\nabla_{\theta} \log \pi (u_t \mid s_t)]] \text{ (remove irrelevant } s \text{ iteration)} \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} \left[b(s_t) \sum_u \pi_{\theta} (u_t \mid s_t) \frac{\nabla_{\theta} \pi_{\theta} (u_t \mid s_t)}{\pi_{\theta} (u_t \mid s_t)} \right] \text{ (iterate over actions)} \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} \left[b(s_t) \sum_u \nabla_{\theta} \pi_{\theta} (u_t \mid s_t) \right] \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} \left[b(s_t) \nabla_{\theta} \sum_u \pi_{\theta} (u_t \mid s_t) \right] \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} [b(s_t) \nabla_{\theta} 1] \\ &= \mathbb{E}_{s_0:t, a_0:(t-1)} [b(s_t) \cdot 0] = 0 \end{aligned}$$

- Current state:

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$


- **Intuition at this point:** Increase the probability of actions proportionally to how much its returns are better than the estimated return under the current policy !
- Also called **the Advantage** of that action over baseline

How to choose the baseline ?

How to choose the baseline ?

- Constant baseline: $b = \mathbb{E} [R(\tau)] \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$

How to choose the baseline ?

- Constant baseline: $b = \mathbb{E} [R(\tau)] \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$
- Optimal constant baseline (according to variance): $b = \frac{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2 R(\tau^{(i)})}{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2}$

[**Read:** Greensmith, Bartlett, Baxter, JMLR 2004 for variance reducing techniques.]

How to choose the baseline ?

- Constant baseline: $b = \mathbb{E} [R(\tau)] \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$
- Optimal constant baseline (according to variance): $b = \frac{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2 R(\tau^{(i)})}{\sum_i (\nabla_{\theta} \log P(\tau^{(i)}; \theta))^2}$
- Time-based average: $b_t = \frac{1}{m} \sum_{i=1}^m \sum_{k=t}^{H-1} R(s_k^{(i)}, u_k^{(i)})$
- State-based expected return: $b(s_t) = \mathbb{E} [r_t + r_{t+1} + r_{t+2} + \dots + r_{H-1}] = \mathbb{E} [G_t]$

[**Read:** Greensmith, Bartlett, Baxter, JMLR 2004 for variance reducing techniques.]

“Vanilla” Policy Gradient pseudocode

(the **template** of all policy gradients algorithms)

Initialize policy parameter θ , baseline b

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

 Collect a set of trajectories τ by executing the current policy

 for each trajectory τ_i

 At each timestep t in trajectory τ^i

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

 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 $\nabla_{\theta} \log \pi_{\theta}(u_t \mid s_t) \hat{A}_t$.

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

How can we improve the baseline choice ?

- Current state: We estimate the baseline target as returns, G_t , from MC samples

$$\sum_i \sum_t \|b(s_t) - G_t^i\|^2$$

- This is unbiased, but high variance, collected from individual roll-outs.
- Ideas from previous courses:
 - Reduce variance by adding bias using **bootstrapping** and **discounting**
 - Using **function approximation**

How can we improve the baseline choice ?

- Now we use : $b(s_t) = \mathbb{E} [r_t + r_{t+1} + r_{t+2} + \dots + r_{H-1}]$

- Recall **State-value function** – can serve as a **baseline estimator**

$$V^{\pi, \gamma}(s) = \mathbb{E}_{\pi} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s]$$

- Recall **Q-function (state-action-value)**:

$$Q^{\pi, \gamma}(s, u) = \mathbb{E} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, u_0 = u]$$

How can we improve the baseline choice ?

- Now we use : $b(s_t) = \mathbb{E}[r_t + r_{t+1} + r_{t+2} + \dots + r_{H-1}]$

- Recall **State-value function** – can serve as a **baseline estimator**

$$V^{\pi, \gamma}(s) = \mathbb{E}_{\pi}[r_0 + \gamma r_1 + \gamma^2 r_2 \dots \mid s_0 = s]$$

- Recall **Q-function (state-action-value)**:

$$\begin{aligned} Q^{\pi, \gamma}(s, u) &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma V^{\pi}(s_1) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 V^{\pi}(s_2) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 V^{\pi}(s_3) \mid s_0 = s, u_0 = u] \\ &= \dots \end{aligned}$$

How can we improve the baseline choice ?

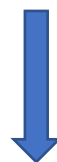
- Current state:

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$

How can we improve the baseline choice ?

- Current state:

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$



Replacing with bootstrapping and action value estimation from the previous slide

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) (Q^{\pi, \gamma}(s_t, u_t) - V^{\pi, \gamma}(s_t)) \right)$$

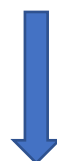
Estimated value of
taking action u_t in
state s_t

Estimated value of
state s_t

How can we improve the baseline choice ?

- Current state:

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) \left(\sum_{k=t}^{H-1} R \left(s_k^{(i)}, u_k^{(i)} \right) - b(s_t) \right) \right)$$



Replacing with bootstrapping and action value estimation from the previous slide

$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) (Q^{\pi, \gamma}(s_t, u_t) - V^{\pi, \gamma}(s_t)) \right)$$

Estimated value of
taking action u_t in
state s_t

Estimated value of
state s_t

Replacing with the advantage function:

what is the advantage over average if agent takes action u_t in state s_t

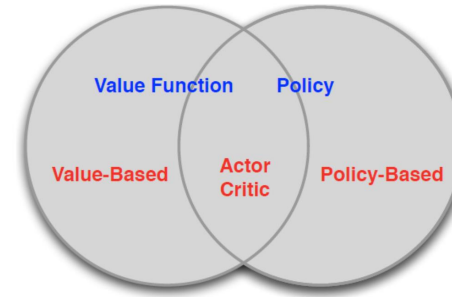


$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) A^{\pi, \gamma}(s_t, u_t) \right)$$

How can we improve the baseline choice ?

- Actor-Critic method:

- Estimation of V and Q done by a **critic**
- Policy decisions are taken by an **actor**



- **A3C** (Mnih et al. ICML 2016) commonly used actor-critic method

- **Critic** can select any blend between TD and MC estimators for estimating the state-action value function Q .

$$\begin{aligned} Q^{\pi, \gamma}(s, u) &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma V^{\pi}(s_1) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 V^{\pi}(s_2) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 V^{\pi}(s_3) \mid s_0 = s, u_0 = u] \\ &= \dots \end{aligned}$$

- In A3C the look-ahead number of steps is a hyperparameter. E.g., $k=5$

How can we improve the baseline choice ?

Let's work on multiple look-ahead steps

- Recall:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) A^{\pi, \gamma}(s_t, u_t) \right)$$
$$A^{\pi, \gamma}(s_t, u_t) = Q^{\pi, \gamma}(s_t, u_t) - V^{\pi, \gamma}(s_t)$$

How can we improve the baseline choice ?

Let's work on multiple look-ahead steps

- Recall:
$$\nabla U(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta} \left(u_t^{(i)} \mid s_t^{(i)} \right) A^{\pi, \gamma}(s_t, u_t) \right)$$

$$A^{\pi, \gamma}(s_t, u_t) = Q^{\pi, \gamma}(s_t, u_t) - V^{\pi, \gamma}(s_t)$$

Look-ahead index



Simplifying notations (getting out π, γ)
and explicating multiple steps inside Q

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

$$A_t^{(2)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) - V(s_t)$$

.....

$$A_t^{(inf)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) + \dots - V(s_t)$$

Check

Look-ahead index



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

$$A_t^{(2)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) - V(s_t)$$

.....


$$A_t^{(inf)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) + \dots - V(s_t)$$

Select which ones are true:

- ☐ $A_t^{(1)}$ has low variance and low bias
- ☐ $A_t^{(1)}$ has high variance and low bias
- ☐ $A_t^{(inf)}$ has low variance and high bias
- ☐ $A_t^{(inf)}$ has high variance and low bias

Check

Look-ahead index


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

$$A_t^{(2)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) - V(s_t)$$

.....

$$A_t^{(inf)} = r_t + \gamma V(s_{t+1}) + \gamma^2 V(s_{t+2}) + \dots - V(s_t)$$

Select which ones are true:

☐ $A_t^{(1)}$ has low variance and low bias

☐ $A_t^{(1)}$ has high variance and low bias

☐ $A_t^{(inf)}$ has low variance and high bias

☒ $A_t^{(inf)}$ has high variance and low bias

Notes:

➤ $A_t^{(1)}$ has low variance but high bias !

➤ As we perform multiple look-ahead steps we increase the variance but lower the bias !

How can we improve the baseline choice ?

Let's improve further !

- **Generalized Advantage Estimation (GAE)** [Schulman et al, ICLR 2016]
- \sim TD(lambda) / eligibility traces [Sutton and Barto, 1990]

$$\begin{aligned} Q^{\pi, \gamma}(s, u) &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma V^{\pi}(s_1) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 V^{\pi}(s_2) \mid s_0 = s, u_0 = u] \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 V^{\pi}(s_3) \mid s_0 = s, u_0 = u] \\ &= \dots \end{aligned}$$

How can we improve the baseline choice ?

Let's improve further ! **Key Idea**: why not averaging all look-ahead steps in estimating Q ?

- **Generalized Advantage Estimation (GAE)** [Schulman et al, ICLR 2016]
- \sim TD(lambda) / eligibility traces [Sutton and Barto, 1990]

$$\begin{aligned} Q^{\pi, \gamma}(s, u) &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, u_0 = u] && (1 - \lambda) \\ &= \mathbb{E}[r_0 + \gamma V^{\pi}(s_1) \mid s_0 = s, u_0 = u] && (1 - \lambda)\lambda \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 V^{\pi}(s_2) \mid s_0 = s, u_0 = u] && (1 - \lambda)\lambda^2 \\ &= \mathbb{E}[r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 V^{\pi}(s_3) \mid s_0 = s, u_0 = u] && (1 - \lambda)\lambda^3 \\ &= \dots \end{aligned}$$



\hat{Q}

Averaged weighted of all steps

Actor-Critic with A3C + GAE

- Use two networks: one for policy, one for value estimation
- **Note:** Can update both independently, e.g., can use k-steps for V, full roll-out for π

Init $\pi_{\theta_0}, V_{\phi_0}^\pi$

Run episode = 1,2,....

Collect roll-outs $\{s, u, s', r\}$

Estimate $\hat{Q}_i(s, u)$ using *GAE*

Update the two networks:

$$\phi_{i+1} \leftarrow \min_{\phi} \sum_{(s,u,s',r)} \left\| \hat{Q}_i(s, u) - V_{\phi}^{\pi}(s) \right\|_2^2 + \kappa \|\phi - \phi_i\|_2^2$$

$$\theta_{i+1} \leftarrow \theta_i + \alpha \frac{1}{m} \sum_{k=1}^m \sum_{t=0}^{H-1} \nabla_{\theta} \log \pi_{\theta_i} \left(u_t^{(k)} \mid s_t^{(k)} \right) \underbrace{\left(\hat{Q}_i \left(s_t^{(k)}, u_t^{(k)} \right) - V_{\phi_i}^{\pi} \left(s_t^{(k)} \right) \right)}_{\hat{A}_t}$$

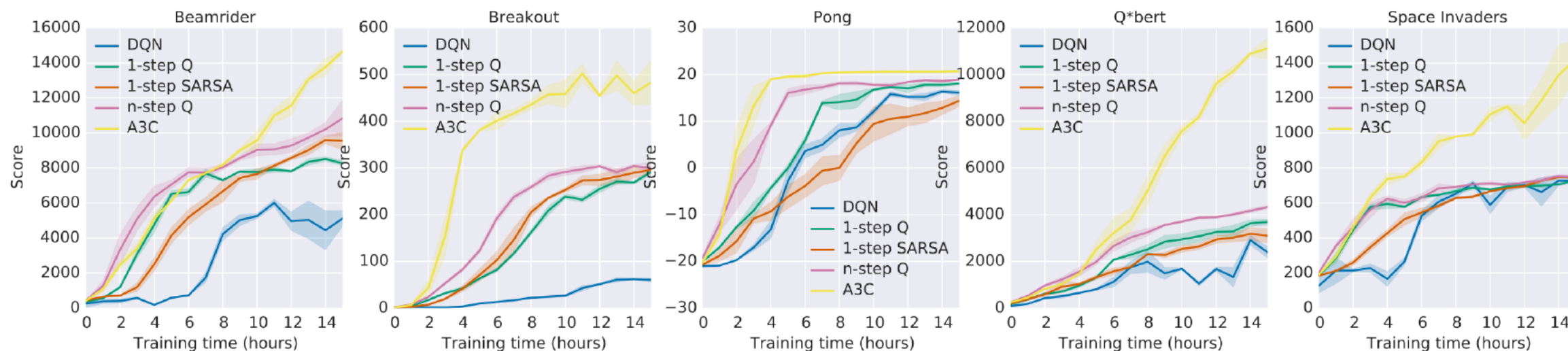
(advantage at time t)

Some results on ATARI games:

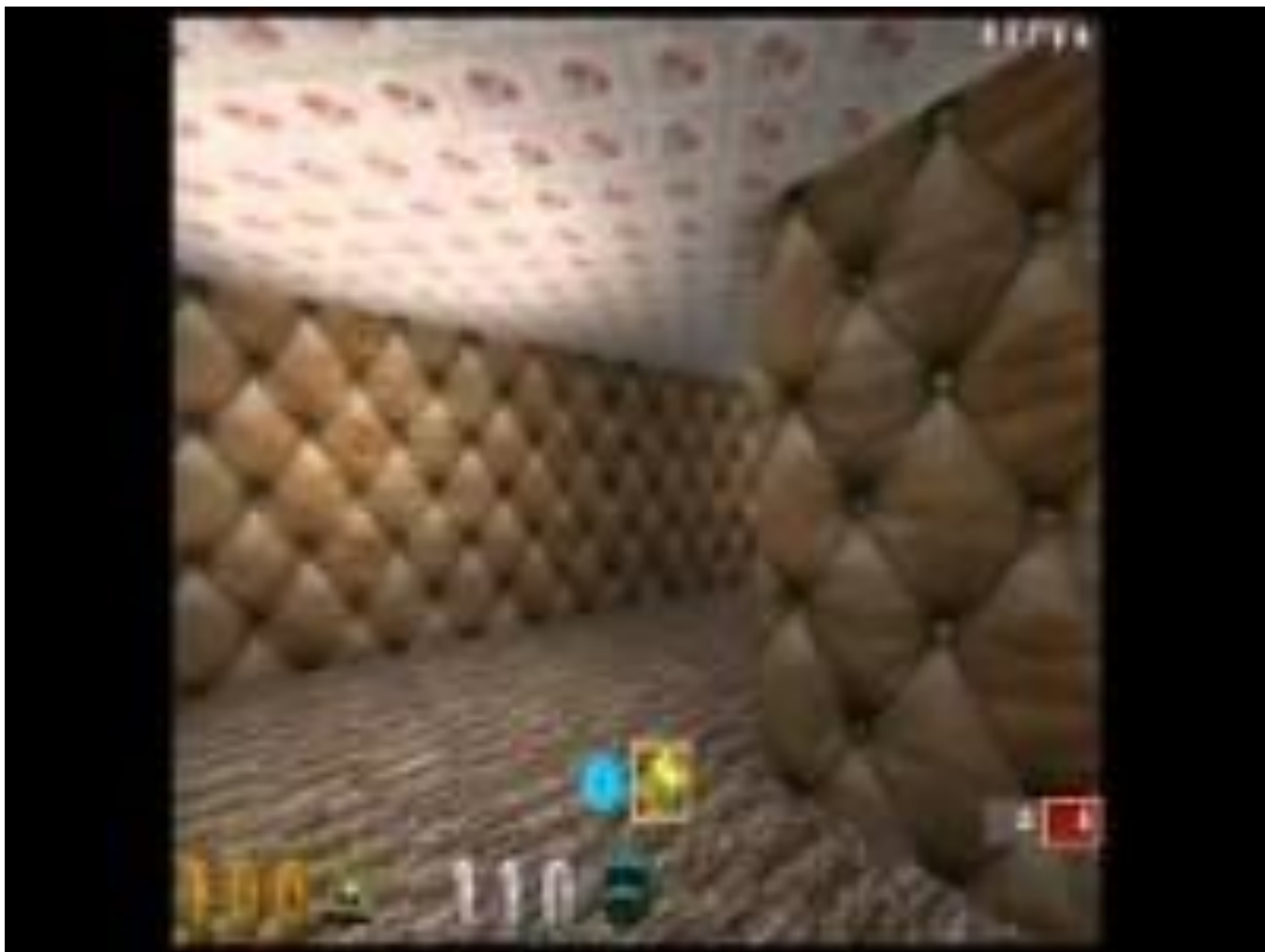
[Mnih et al, ICML 2016]

Likelihood Ratio Policy Gradient

n-step Advantage Estimation



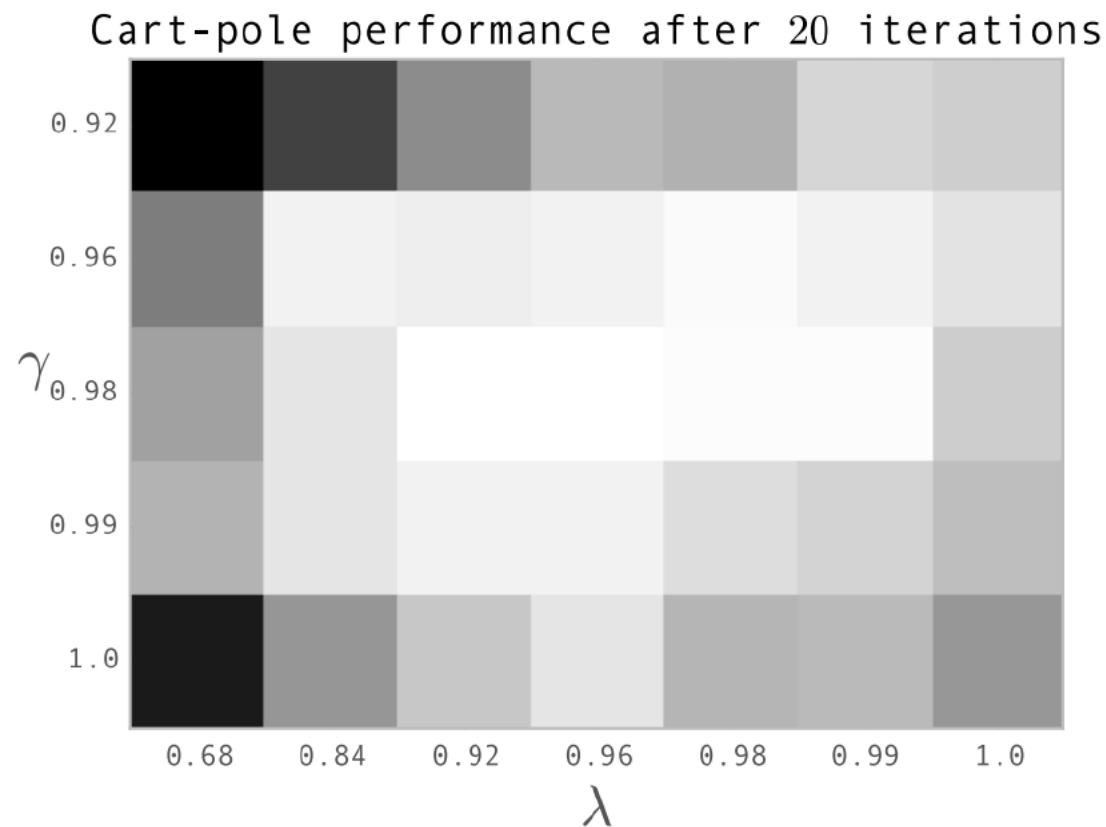
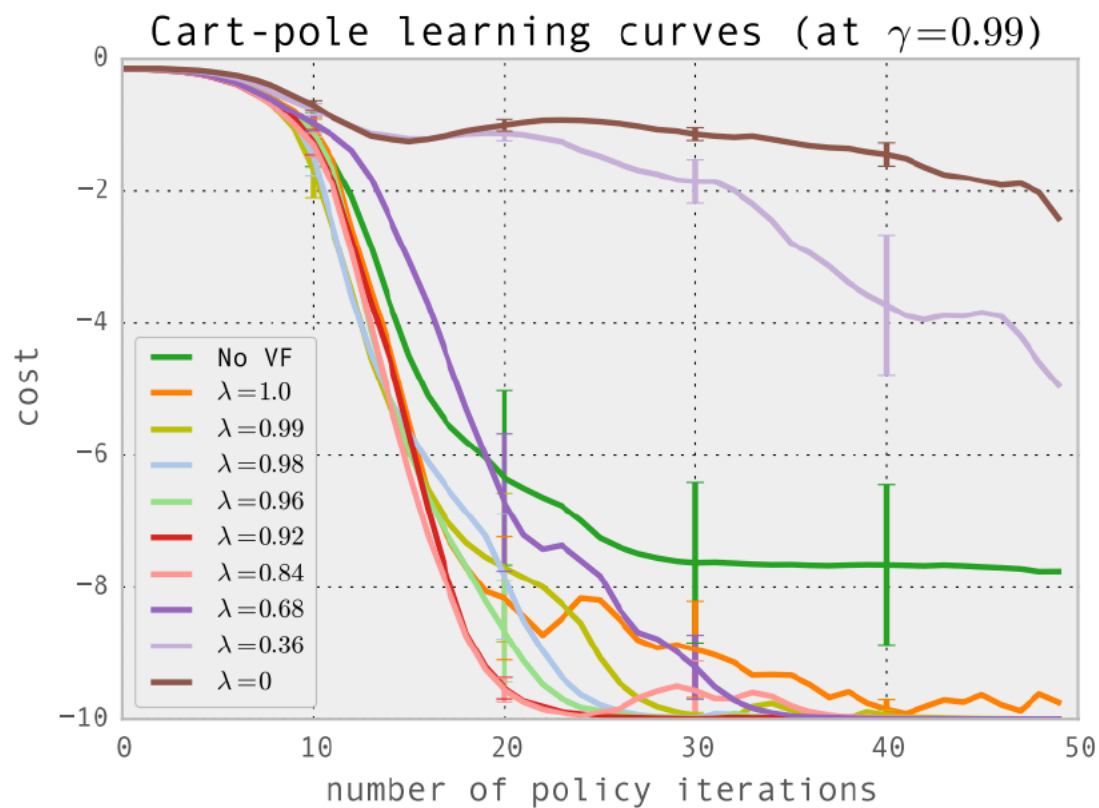
DeepMine A3C Labyrinth



The video shows an agent collecting rewards in previously unseen mazes using only raw pixels as input. The agent was trained using the Asynchronous Advantage Actor-Critic (A3C) algorithm and was only rewarded for picking up apples and orange portals during training.

Paper link -
<http://arxiv.org/pdf/1602.01783.pdf>

Cart-Pole, GAE: Effect of gamma and lambda



[Schulman et al, 2016 -- GAE]

Thank you !
Questions