CS 8803 Deep Reinforcement Learning

Lec 3: Intro to RL Fall 2024

Animesh Garg

Summary: MDP Equations

• Value iteration equation:

$$V_{k+1}(s) \leftarrow \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V_k(s') \right]$$

• Policy evaluation equation:

$$V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$$

• Policy iteration equation:

٦

$$\pi_{i+1}(s) = \arg\max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$

The Bellman Equations



Convergence when Solving MDPs

- Redefine value update as general Bellman Utility update
 - Recursive update or utility (sum of discounted reward) $V_{k+1}^{\pi_i}(s) \leftarrow \sum_{s'} T(s, \pi_i(s), s') \left[R(s, \pi_i(s), s') + \gamma V_k^{\pi_i}(s') \right]$ \downarrow $U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum P(s'|s, a) U_i(s')$
- How does this converge?
 - Assume fixed policy $\pi_i(s)$.
 - *R(s)* is the short term reward of being in *s*

Convergence when Solving MDPs

- How does this update rule converge? $U_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum P(s'|s, a) U_i(s')$
- Re-write update: $U_{i+1} \leftarrow BU_i$
 - *B* is a linear operator (like a matrix)
 - *U* is a vector
- Interested in delta between Utilities:

$$||BU_{i+1} - BU_i|| \le \gamma ||U_{i+1} - U_i||$$

Bellman Recursion

 $||U_{i+1} - U_i||$

v' = Av

Detour: Convergence of Bellman backup

Let $||V - V'|| = \max_{s} |V(s) - V'(s)|$ be the infinity norm

$$\begin{split} \|BV_{k} - BV_{j}\| &= \left\| \max_{a} \left(R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V_{k}(s') \right) - \max_{a'} \left(R(s, a') + \gamma \sum_{s' \in S} P(s'|s, a') V_{j}(s') \right) \right\| \\ &\leq \max_{a} \left\| \left(R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V_{k}(s') - R(s, a) - \gamma \sum_{s' \in S} P(s'|s, a) V_{j}(s') \right) \right\| \\ &= \max_{a} \left\| \gamma \sum_{s' \in S} P(s'|s, a) (V_{k}(s') - V_{j}(s')) \right\| \\ &\leq \max_{a} \left\| \gamma \sum_{s' \in S} P(s'|s, a) \| V_{k} - V_{j} \| \right) \\ &= \max_{a} \left\| \gamma \| V_{k} - V_{j} \| \sum_{s' \in S} P(s'|s, a) \right\| \\ &= \gamma \| V_{k} - V_{j} \| \end{split}$$

Note: Even if all inequalities are equalities, this is still a contraction if $\gamma < 1$

Convergence when Solving MDPs

• How does this delta converge?

$$|BU_{i+1} - BU_i|| \le \gamma ||U_{i+1} - U_i||$$

- Utility error estimate reduced by γ each iteration:
- Total Utilities are bounded,

$$\sum_{i=0}^{\infty} R_{max} \gamma^i \qquad \pm rac{R_{max}}{(1-\gamma)}$$

- Consider minimum initial error:
 - (Max norm)
- Max error: reduce by discount each step.

$$\|U_0 - U\| \le \frac{2R_{max}}{(1-\gamma)}$$

Utility Error Bound

• Error at step 0: $\|U_0 - U\| \leq \frac{2R_{max}}{(1 - \gamma)}$

• Error at step N:
$$\|U_N - U\| = \gamma^N \cdot \frac{2R_{max}}{(1 - \gamma)} < \epsilon$$

• Steps for error below ϵ :

$$N = \frac{\log\left(\frac{2R_{max}}{\epsilon(1-\gamma)}\right)}{\log\left(\frac{1}{\gamma}\right)}$$

MDP Convergence Visualized

- Value iteration converges exponentially (with discount factor)
- Policy iteration will converge linearly to 0.



Summary: MDP Algorithms

- So you want to....
 - Compute optimal values: use value iteration or policy iteration
 - Compute values for a particular policy: use policy evaluation
 - Turn your values into a policy: use policy extraction (one-step lookahead)
- These all look the same!
 - They basically are they are all variations of Bellman updates
 - They all use one-step lookahead expectimax fragments
 - They differ only in whether we plug in a fixed policy or max over actions

Definitions

Terminology & notation

O₁



Imitation Learning





Reward functions



which action is better or worse?

 $r(\mathbf{s}, \mathbf{a})$: reward function tells us which states and actions are better **s**, **a**, $r(\mathbf{s}, \mathbf{a})$, and $p(\mathbf{s'}|\mathbf{s}, \mathbf{a})$ define Markov decision process



high reward



low reward

Algorithms

The anatomy of a reinforcement learning algorithm





Another example: RL by backprop





Value Functions

How do we deal with all these expectations?

$$E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

$$E_{\mathbf{s}_{1} \sim p(\mathbf{s}_{1})} \left[E_{\mathbf{a}_{1} \sim \pi(\mathbf{a}_{1} | \mathbf{s}_{1})} \left[r(\mathbf{s}_{1}, \mathbf{a}_{1}) + E_{\mathbf{s}_{2} \sim p(\mathbf{s}_{2} | \mathbf{s}_{1}, \mathbf{a}_{1})} \left[E_{\mathbf{a}_{2} \sim \pi(\mathbf{a}_{2} | \mathbf{s}_{2})} \left[r(\mathbf{s}_{2}, \mathbf{a}_{2}) + \dots | \mathbf{s}_{2} \right] | \mathbf{s}_{1}, \mathbf{a}_{1} \right] | \mathbf{s}_{1} \right] \right]$$
what if we knew this part?
$$Q(\mathbf{s}_{1}, \mathbf{a}_{1}) = r(\mathbf{s}_{1}, \mathbf{a}_{1}) + E_{\mathbf{s}_{2} \sim p(\mathbf{s}_{2} | \mathbf{s}_{1}, \mathbf{a}_{1})} \left[E_{\mathbf{a}_{2} \sim \pi(\mathbf{a}_{2} | \mathbf{s}_{2})} \left[r(\mathbf{s}_{2}, \mathbf{a}_{2}) + \dots | \mathbf{s}_{2} \right] | \mathbf{s}_{1}, \mathbf{a}_{1} \right]$$

$$E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] = E_{\mathbf{s}_{1} \sim p(\mathbf{s}_{1})} \left[E_{\mathbf{a}_{1} \sim \pi(\mathbf{a}_{1} | \mathbf{s}_{1})} \left[Q(\mathbf{s}_{1}, \mathbf{a}_{1}) | \mathbf{s}_{1} \right] \right]$$
easy to modify $\pi_{\theta}(\mathbf{a}_{1} | \mathbf{s}_{1})$ if $Q(\mathbf{s}_{1}, \mathbf{a}_{1})$ is known!
example: $\pi(\mathbf{a}_{1} | \mathbf{s}_{1}) = 1$ if $\mathbf{a}_{1} = \arg \max_{\mathbf{a}_{1}} Q(\mathbf{s}_{1}, \mathbf{a}_{1})$

Definition: Q-function

 $Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t \right]: \text{ total reward from taking } \mathbf{a}_t \text{ in } \mathbf{s}_t$

Definition: value function

 $V^{\pi}(\mathbf{s}_t) = \sum_{t'=t}^{T} E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$: total reward from \mathbf{s}_t

 $V^{\pi}(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)}[Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)]$

 $E_{\mathbf{s}_1 \sim p(\mathbf{s}_1)}[V^{\pi}(\mathbf{s}_1)]$ is the RL objective!

Using Q-functions and value functions

Idea 1: if we have policy π , and we know $Q^{\pi}(\mathbf{s}, \mathbf{a})$, then we can *improve* π :

set $\pi'(\mathbf{a}|\mathbf{s}) = 1$ if $\mathbf{a} = \arg \max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a})$

this policy is at least as good as π (and probably better)!

and it doesn't matter what π is

Idea 2: compute gradient to increase probability of good actions **a**:

if $Q^{\pi}(\mathbf{s}, \mathbf{a}) > V^{\pi}(\mathbf{s})$, then **a** is better than average (recall that $V^{\pi}(\mathbf{s}) = E[Q^{\pi}(\mathbf{s}, \mathbf{a})]$ under $\pi(\mathbf{a}|\mathbf{s})$) modify $\pi(\mathbf{a}|\mathbf{s})$ to increase probability of **a** if $Q^{\pi}(\mathbf{s}, \mathbf{a}) > V^{\pi}(\mathbf{s})$

These ideas are *very* important in RL; we'll revisit them again and again!

The anatomy of a reinforcement learning algorithm



Types of Algorithms

Types of RL algorithms

$$\theta^{\star} = \arg\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

- Policy gradients: directly differentiate the above objective
- Value-based: estimate value function or Q-function of the optimal policy (no explicit policy)
- Actor-critic: estimate value function or Q-function of the current policy, use it to improve policy
- Model-based RL: estimate the transition model, and then...
 - Use it for planning (no explicit policy)
 - Use it to improve a policy
 - Something else

Model-based RL algorithms



Model-based RL algorithms

improve the policy

a few options

- 1. Just use the model to plan (no policy)
 - Trajectory optimization/optimal control (primarily in continuous spaces) essentially backpropagation to optimize over actions
 - Discrete planning in discrete action spaces e.g., Monte Carlo tree search
- 1. Backpropagate gradients into the policy
 - Requires some tricks to make it work
- 2. Use the model to learn a value function
 - Dynamic programming
 - Generate simulated experience for model-free learner

Value function based algorithms



Direct policy gradients



Actor-critic: value functions + policy gradients



Tradeoffs Between Algorithms

Why so many RL algorithms?

- Different tradeoffs
 - Sample efficiency
 - Stability & ease of use
- Different assumptions
 - Stochastic or deterministic?
 - Continuous or discrete?
 - Episodic or infinite horizon?
- Different things are easy or hard in different settings
 - Easier to represent the policy?
 - Easier to represent the model?



Comparison: sample efficiency

- Sample efficiency = how many samples do we need to get a good policy?
- Most important question: is the algorithm *off policy*?
 - Off policy: able to improve the policy without generating new samples from that policy
 - On policy: each time the policy is changed, even a little bit, we need to generate new samples



Comparison: sample efficiency



Why would we use a *less* efficient algorithm?

Wall clock time is not the same as efficiency!

Comparison: stability and ease of use

- Does it converge?
- And if it converges, to what?
- And does it converge every time?

Why is any of this even a question???

- Supervised learning: almost *always* gradient descent
- Reinforcement learning: often not gradient descent
 - Q-learning: fixed point iteration
 - Model-based RL: model is not optimized for expected reward
 - Policy gradient: *is* gradient descent, but also often the least efficient!
Comparison: stability and ease of use

- Value function fitting
 - At best, minimizes error of fit ("Bellman error")
 - Not the same as expected reward
 - At worst, doesn't optimize anything
 - Many popular deep RL value fitting algorithms are not guaranteed to converge to *anything* in the nonlinear case
- Model-based RL
 - Model minimizes error of fit
 - This will converge
 - No guarantee that better model = better policy
- Policy gradient
 - The only one that actually performs gradient descent (ascent) on the true objective

Comparison: assumptions

- Common assumption #1: full observability
 - Generally assumed by value function fitting methods
 - Can be mitigated by adding recurrence
- Common assumption #2: episodic learning
 - Often assumed by pure policy gradient methods
 - Assumed by some model-based RL methods
- Common assumption #3: continuity or smoothness
 - Assumed by some continuous value function learning methods
 - Often assumed by some model-based RL methods







Examples of Algorithms

Examples of specific algorithms

- Value function fitting methods
 - Q-learning, DQN
 - Temporal difference learning
 - Fitted value iteration
- Policy gradient methods
 - REINFORCE
 - Natural policy gradient
 - Trust region policy optimization
- Actor-critic algorithms
 - Asynchronous advantage actor-critic (A3C)
 - Soft actor-critic (SAC)
- Model-based RL algorithms
 - Dyna
 - Guided policy search

We'll learn about most of these in the next few weeks!

The goal of reinforcement learning



$$\underbrace{p_{\theta}(\mathbf{s}_{1}, \mathbf{a}_{1}, \dots, \mathbf{s}_{T}, \mathbf{a}_{T})}_{p_{\theta}(\tau)} = p(\mathbf{s}_{1}) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\theta^{\star} = \arg\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

The goal of reinforcement learning

$$\theta^{\star} = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

Evaluating the objective

$$\theta^{\star} = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$
$$J(\theta)$$



$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \approx \frac{1}{N} \sum_{i} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

Direct policy differentiation

$$\theta^{\star} = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$
$$J(\theta)$$

a convenient identity
$$p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) = p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} = \underline{\nabla_{\theta} p_{\theta}(\tau)}$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)] = \int p_{\theta}(\tau)r(\tau)d\tau$$
$$\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\nabla_{\theta} J(\theta) = \int \underline{\nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau} = \int \underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) r(\tau) d\tau} = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

Direct policy differentiation

$$\theta^{\star} = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)]$$

$$\log \text{ of both}$$

$$\frac{p_{\theta}(\mathbf{s}_{1}, \mathbf{a}_{1}, \dots, \mathbf{s}_{T}, \mathbf{a}_{T})}{p_{\theta}(\tau)} = p(\mathbf{s}_{1}) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\log p_{\theta}(\tau) = \log p(\mathbf{s}_{1}) + \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\nabla_{\theta} \left[\log p(\mathbf{s}_{1}) + \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right]$$

recall: $J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left| \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right| \approx \frac{1}{N} \sum_{t} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$ $\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left| \left(\sum_{t=1}^{I} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left(\sum_{t=1}^{I} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right|$ fit a model to estimate return $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{i=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$ generate samples (i.e. run the policy) $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ improve the

policy

REINFORCE algorithm:

1. sample
$$\{\tau^i\}$$
 from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Evaluating the policy gradient

Understanding Policy Gradients

Evaluating the policy gradient

recall:
$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \approx \frac{1}{N} \sum_{i} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$
what is this?







 \mathbf{a}_t

Comparison to maximum likelihood

policy gradient:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

maximum likelihood:
$$\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$



Example: Gaussian policies

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

example: $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = \mathcal{N}(f_{\text{neural network}}(\mathbf{s}_t); \Sigma)$ $\log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = -\frac{1}{2} ||f(\mathbf{s}_t) - \mathbf{a}_t||_{\Sigma}^2 + \text{const}$ $\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = -\frac{1}{2} \Sigma^{-1} (f(\mathbf{s}_t) - \mathbf{a}_t) \frac{df}{d\theta}$

REINFORCE algorithm:

1. sample
$$\{\tau^i\}$$
 from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

What did we just do?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau_{i}) r(\tau_{i})}_{T}$$
$$\sum_{t=1}^{T} \nabla_{\theta} \log_{\theta} \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$
good stuff is made more likely

bad stuff is made less likely

simply formalizes the notion of "trial and error"!

REINFORCE algorithm:

1. sample
$$\{\tau^i\}$$
 from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

maximum likelihood:
$$\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log \pi_{\theta}(\tau_i)$$



Partial observability



$$\mathbf{a}_t$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{o}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Markov property is not actually used!

Can use policy gradient in partially observed MDPs without modification

What is wrong with the policy gradient?



Review

- Evaluating the RL objective
 - Generate samples
- Evaluating the policy gradient
 - Log-gradient trick
 - Generate samples
- Understanding the policy gradient
 - Formalization of trial-and-error
- Partial observability
 - Works just fine
- What is wrong with policy gradient?



Reducing Variance

Reducing variance

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Causality: policy at time t' cannot affect reward at time t when t < t'

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{\substack{t' \in \mathbf{t} \\ \mathbf{t}' \in \mathbf{t}}}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$

"reward to go"
 $\hat{Q}_{i,t}$

Baselines

a convenient identity $p_{ heta}(\tau) \nabla_{ heta} \log p_{ heta}(\tau) = \nabla_{ heta} p_{ heta}(\tau)$



$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b]$$
$$b = \frac{1}{N} \sum_{i=1}^{N} r(\tau) \qquad \text{but... are we allowed to do that??}$$

$$E[\nabla_{\theta} \log p_{\theta}(\tau)b] = \int p_{\theta}(\tau)\nabla_{\theta} \log p_{\theta}(\tau)b \, d\tau = \int \nabla_{\theta} p_{\theta}(\tau)b \, d\tau = b\nabla_{\theta} \int p_{\theta}(\tau)d\tau = b\nabla_{\theta} 1 = 0$$

subtracting a baseline is *unbiased* in expectation!

average reward is *not* the best baseline, but it's pretty good!

Analyzing variance

can we write down the variance?

 $\operatorname{Var}[x] = E[x^2] - E[x]^2$ $\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]$ $\operatorname{Var} = E_{\tau \sim p_{\theta}(\tau)} [(\nabla_{\theta} \log p_{\theta}(\tau)(r(\tau) - b))^{2}] - E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau)(r(\tau) - b)]^{2}$ this bit is just $E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$ (baselines are unbiased in expectation) $\frac{d\operatorname{Var}}{db} = \frac{d}{db}E[g(\tau)^2(r(\tau)-b)^2] = \frac{d}{db}\left(E[g(\tau)^2r(\tau)^2] - 2E[g(\tau)^2r(\tau)b] + b^2E[g(\tau)^2]\right)$ $= -2E[q(\tau)^{2}r(\tau)] + 2bE[q(\tau)^{2}] = 0$ This is just expected reward, but weighted $b = \frac{E[g(\tau)^2 r(\tau)]}{E[g(\tau)^2]} \quad \longleftarrow \quad$ by gradient magnitudes!

Review

- The high variance of policy gradient
- Exploiting causality
 - Future doesn't affect the past
- Baselines
 - Unbiased!
- Analyzing variance
 - Can derive optimal baselines



Off-Policy Policy Gradients

Policy gradient is on-policy

 $\theta^{\star} = \arg \max_{\theta} J(\theta)$

 $J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)]$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$
this is trouble...

- Neural networks change only a little bit with each gradient step
- On-policy learning can be extremely inefficient!

REINFORCE algorithm: 1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot) 2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$ 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Off-policy learning & importance sampling

 $\theta^{\star} = \arg \max_{\theta} J(\theta)$

 $J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)]$

what if we don't have samples from $p_{\theta}(\tau)$? (we have samples from some $\bar{p}(\tau)$ instead)

$$J(\theta) = E_{\tau \sim \bar{p}(\tau)} \left[\frac{p_{\theta}(\tau)}{\bar{p}(\tau)} r(\tau) \right]$$
$$p_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\frac{p_{\theta}(\tau)}{\bar{p}(\tau)} = \frac{p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}{p(\mathbf{s}_1) \prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} = \frac{\prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}$$

| importance sampling |
|---|
| $E_{x \sim p(x)}[f(x)] = \int p(x)f(x)dx$ |
| $= \int \frac{q(x)}{q(x)} p(x) f(x) dx$ |
| $= \int q(x) \frac{p(x)}{q(x)} f(x) dx$ |
| $= E_{x \sim q(x)} \left[\frac{p(x)}{q(x)} f(x) \right]$ |

Deriving the policy gradient with IS

 $\theta^{\star} = \arg\max_{\theta} J(\theta)$

 $J(\theta) = E_{\tau \sim p_{\theta}(\tau)}[r(\tau)]$

can we estimate the value of some new parameters θ' ?

 $J(\theta') = E_{\tau \sim p_{\theta}(\tau)} \begin{bmatrix} p_{\theta'}(\tau) \\ p_{\theta}(\tau) \end{bmatrix}$ the only bit that depends on θ'

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim p_{\theta}(\tau)} \left[\frac{\nabla_{\theta'} p_{\theta'}(\tau)}{p_{\theta}(\tau)} r(\tau) \right] = E_{\tau \sim p_{\theta}(\tau)} \left[\frac{p_{\theta'}(\tau)}{p_{\theta}(\tau)} \nabla_{\theta'} \log p_{\theta'}(\tau) r(\tau) \right]$$

now estimate locally, at $\theta = \theta'$: $\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$

a convenient identity

$$p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} p_{\theta}(\tau)$$

The off-policy policy gradient

A first-order approximation for IS (preview)

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^{T} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \left(\prod_{t'=1}^{t} \frac{\pi_{\theta'}(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'})} \right) \left(\sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right]$$
exponential in *T*...

let's write the objective a bit differently...

on-policy policy gradient:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$$

off-policy policy gradient: $\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\pi_{\theta'}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$
We'll see why this is
reasonable
later in the course! ignore this part

Implementing Policy Gradients

Policy gradient with automatic differentiation

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}}{\mathbf{1}}$$
 pretty inefficient to compute these explicitly!

How can we compute policy gradients with automatic differentiation?

We need a graph such that its gradient is the policy gradient!

maximum likelihood:
$$\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \qquad J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$

Just implement "pseudo-loss" as a weighted maximum likelihood:

$$\tilde{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$$
cross entropy (discrete) or squared error (Gaussian)

Policy gradient with automatic differentiation

Pseudocode example (with discrete actions):

Maximum likelihood:

```
# Given:
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
loss = tf.reduce_mean(negative_likelihoods)
gradients = loss.gradients(loss, variables)
```

Policy gradient with automatic differentiation

Pseudocode example (with discrete actions):

```
Policy gradient:
```

```
# Given:
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# q_values - (N*T) x 1 tensor of estimated state-action values
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
weighted_negative_likelihoods = tf.multiply(negative_likelihoods, q_values)
loss = tf.reduce_mean(weighted_negative_likelihoods)
gradients = loss.gradients(loss, variables)
```



Policy gradient in practice

- Remember that the gradient has high variance
 - This isn't the same as supervised learning!
 - Gradients will be really noisy!
- Consider using much larger batches
- Tweaking learning rates is very hard
 - Adaptive step size rules like ADAM can be OK-ish
 - We'll learn about policy gradient-specific learning rate adjustment methods later!

Review

- Policy gradient is on-policy
- Can derive off-policy variant
 - Use importance sampling
 - Exponential scaling in T
 - Can ignore state portion (approximation)
- Can implement with automatic differentiation – need to know what to backpropagate
- Practical considerations: batch size, learning rates, optimizers



Advanced Policy Gradients
What *else* is wrong with the policy gradient?



$$r(\mathbf{s}_t, \mathbf{a}_t) = -\mathbf{s}_t^2 - \mathbf{a}_t^2$$
$$\log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = -\frac{1}{2\sigma^2} (k\mathbf{s}_t - \mathbf{a}_t)^2 + \text{const} \qquad \theta = (k, \sigma)$$



(image from Peters & Schaal 2008)





Covariant/natural policy gradient

 $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta) \qquad \qquad \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$

some parameters change probabilities a lot more than others!

$$\theta' \leftarrow \arg \max_{\theta'} (\theta' - \theta)^T \nabla_{\theta} J(\theta) \text{ s.t. } \frac{\|\theta' - \theta\|^2 \le \epsilon}{\|\theta\|^2 \le \epsilon}$$

controls how far we go

can we *rescale* the gradient so this doesn't happen?

$$\theta' \leftarrow \arg \max_{\theta'} (\theta' - \theta)^T \nabla_{\theta} J(\theta) \text{ s.t. } D(\pi_{\theta'}, \pi_{\theta}) \leq \epsilon$$

parameterization-independent divergence measure usually KL-divergence: $D_{\text{KL}}(\pi_{\theta'} || \pi_{\theta}) = E_{\pi_{\theta'}}[\log \pi_{\theta} - \log \pi_{\theta'}]$

 $D_{\mathrm{KL}}(\pi_{\theta'} \| \pi_{\theta}) \approx (\theta' - \theta)^T \mathbf{F}(\theta' - \theta) \qquad \mathbf{F} = E_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s})^T]$ Fisher-information matrix can estimate with samples



Covariant/natural policy gradient

 $D_{\mathrm{KL}}(\pi_{\theta'} \| \theta_{\pi}) \approx (\theta' - \theta)^T \mathbf{F}(\theta' - \theta) \qquad \mathbf{F} = E_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a} | \mathbf{s})^T]$

 $\theta' \leftarrow \arg \max_{\theta'} (\theta' - \theta)^T \nabla_{\theta} J(\theta) \text{ s.t. } D(\pi_{\theta'}, \pi_{\theta}) \leq \epsilon$

 $\theta \leftarrow \theta + \alpha \mathbf{F}^{-1} \nabla_{\theta} J(\theta)$

natural gradient: pick α

trust region policy optimization: pick ϵ

can solve for optimal α while solving $\mathbf{F}^{-1} \nabla_{\theta} J(\theta)$

conjugate gradient works well for this

see Schulman, L., Moritz, Jordan, Abbeel (2015) Trust region policy optimization



Advanced policy gradient topics

- What more is there?
- Next time: introduce value functions and Q-functions
- Later in the class: more on natural gradient and automatic step size adjustment

Policy gradients suggested readings

- Classic papers
 - Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning: introduces REINFORCE algorithm
 - Baxter & Bartlett (2001). Infinite-horizon policy-gradient estimation: temporally decomposed policy gradient (not the first paper on this! see actor-critic section later)
 - Peters & Schaal (2008). Reinforcement learning of motor skills with policy gradients: very accessible overview of optimal baselines and natural gradient
- Deep reinforcement learning policy gradient papers
 - Levine & Koltun (2013). Guided policy search: deep RL with importance sampled policy gradient (unrelated to later discussion of guided policy search)
 - Schulman, L., Moritz, Jordan, Abbeel (2015). Trust region policy optimization: deep RL with natural policy gradient and adaptive step size
 - Schulman, Wolski, Dhariwal, Radford, Klimov (2017). Proximal policy optimization algorithms: deep RL with importance sampled policy gradient

Acknowledgements

Slides adapted from

CS 188 UC Berkeley Pieter Abbeel, Dan Klein et al.

CS 285 UC Berkeley Sergey Levine

CSC 498 Univ of Toronto Animesh Garg