Actor-Critic Algorithms

CS 285: Deep Reinforcement Learning, Decision Making, and Control Sergey Levine

Class Notes

- 1. HW2 is out!
 - Due 9/30
 - Start early! This one is harder
- 2. Remember to start forming final project groups
 - Proposal due 9/25

Today's Lecture

- 1. Improving the policy gradient with a critic
- 2. The policy evaluation problem
- 3. Discount factors
- 4. The actor-critic algorithm
- Goals:
 - Understand how policy evaluation fits into policy gradients
 - Understand how actor-critic algorithms work

Recap: policy gradients



"reward to go"

Improving the policy gradient

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

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

 $\hat{Q}_{i,t}$: estimate of expected reward if we take action $\mathbf{a}_{i,t}$ in state $\mathbf{s}_{i,t}$ can we get a better estimate?

$$Q(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]: \text{ true } expected \text{ reward-to-go}$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$



What about the baseline?

 $Q(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]: \text{ true } expected \text{ reward-to-go}$ $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) (\mathcal{O}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})) - \mathcal{O}(\mathbf{s}_{i,t}))$ $b_{t} = a_{v}^{1} \exp_{i} \left[\exp(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right] \text{ average what}?$

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



State & state-action value functions

 $Q(\mathbf{x}_{i},\mathbf{x}_{i},\mathbf{x}_{i}) = \sum_{t'} \sum_{t' t' t' t'} E_{\mathbf{x}_{i}} \left[r[(\mathbf{x}_{i},\mathbf{x}_{i'},\mathbf$

 $V^{\pi}(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}[Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)]$: total reward from \mathbf{s}_t

 $A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi}(\mathbf{s}_t)$: how much better \mathbf{a}_t is

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



the better this estimate, the lower the variance

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

unbiased, but high variance single-sample estimate

Value function fitting

$$Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]$$
$$V^{\pi}(\mathbf{s}_{t}) = E_{\mathbf{a}_{t} \sim \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t})} [Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t})]$$
$$A^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) - V^{\pi}(\mathbf{s}_{t})$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

fit what to what?

 $Q^{\pi}, V^{\pi}, A^{\pi}?$

 $Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \approx \underbrace{\sum_{t=t}^{T} \mathbf{a}_{t}}_{T} \underbrace{\mathbb{E}}_{\pi_{t}} \underbrace{\mathbb{E}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{E}}}_{\mathbf{s}_{t}} \underbrace{\mathbb{$

let's just fit $V^{\pi}(\mathbf{s})!$





fit V^{π}

Policy evaluation

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

 $J(\theta) = E_{\mathbf{s}_1 \sim p(\mathbf{s}_1)}[V^{\pi}(\mathbf{s}_1)]$

how can we perform policy evaluation?

Monte Carlo policy evaluation (this is what policy gradient does)

 $V^{\pi}(\mathbf{s}_t) \approx \sum r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$

 $V^{\pi}(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \qquad (\text{requires us to reset the simulator})$



Monte Carlo evaluation with function approximation

$$V^{\pi}(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

not as good as this: $V^{\pi}(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$

but still pretty good!

training data:
$$\left\{ \left(\mathbf{s}_{i,t}, \sum_{t'=t}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) \right\}$$

supervised regression:
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$





Can we do better?

ideal target: $y_{i,t} = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t} \right] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \sum_{t'=t}^{T} \sum_{t'=t}^{T} \left[r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t'}) + \sum_{t'=t}^{T} \sum_{t'=t}^{T} \left[r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t'})$ Monte Carlo target: $y_{i,t} = \sum_{t'=t}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$ directly use previous fitted value function!

uncerty use previous inteed value function

training data:
$$\left\{ \left(\mathbf{s}_{i,t}, r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1}) \right) \right\}$$

supervised regression:
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$

sometimes referred to as a "bootstrapped" estimate

Policy evaluation examples

TD-Gammon, Gerald Tesauro 1992



Figure 2. An illustration of the normal opening position in backgammon. TD-Gammon has sparked a near-universal conversion in the way experts play certain opening rolls. For example, with an opening roll of 4-1, most players have now switched from the traditional move of 13-9, 6-5, to TD-Gammon's preference, 13-9, 24-23. TD-Gammon's analysis is given in Table 2.



Figure 1. An illustration of the multilayer perception architecture used in TD-Gammon's neural network. This architecture is also used in the popular backpropagation learning procedure. Figure reproduced from [9].

AlphaGo, Silver et al. 2016



reward: game outcome

value function $\hat{V}^{\pi}_{\phi}(\mathbf{s}_t)$: expected outcome given board state reward: game outcome

value function $\hat{V}^{\pi}_{\phi}(\mathbf{s}_t)$: expected outcome given board state

An actor-critic algorithm

batch actor-critic algorithm:

$$V^{\pi}(\mathbf{s}y_{i,t}) \approx \sum_{i=\pm t}^{T} \mathcal{H}(\mathbf{s}_{i}) [\mathcal{H}(\hat{\mathbf{s}}_{i,t}) | \mathbf{s}_{i,t})]$$
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$



Aside: discount factors

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}^{\pi}_{\phi}(\mathbf{s}_{i,t+1})$$
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}^{\pi}_{\phi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$

what if T (episode length) is ∞ ? \hat{V}^{π}_{ϕ} can get infinitely large in many cases lox real time autonomous execution

episodic tasks

Iteration 2000



continuous/cyclical tasks

simple trick: better to get rewards sooner than later

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1})$$

discount factor $\gamma \in [0, 1]$ (0.99 works well)

 γ changes the MDP:



Aside: discount factors for policy gradients

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1})$$

$$\hat{A}^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$\hat{\mathcal{L}}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$
with critic:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t}) \right)$$

what about (Monte Carlo) policy gradients?

option 1:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
not the same!
option 2:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} \gamma^{t-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \gamma^{t-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
(later steps matter less)

Which version is the right one?

option 1:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
option 2:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \gamma^{t-1} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$
this is what we actually use...
why?
iteration 2000
$$(\mathbf{s}_{1} + \mathbf{s}_{2} + \mathbf{s}_{3} + \mathbf{s}_{4} + \mathbf{s}_{4$$

Further reading: Philip Thomas, Bias in natural actor-critic algorithms. ICML 2014

Actor-critic algorithms (with discount)

batch actor-critic algorithm:

online actor-critic algorithm:

1. take action
$$\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{s})$$
, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
2. update \hat{V}^{π}_{ϕ} using target $r + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}')$
3. evaluate $\hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}') - \hat{V}^{\pi}_{\phi}(\mathbf{s})$
4. $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}|\mathbf{s}) \hat{A}^{\pi}(\mathbf{s}, \mathbf{a})$
5. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

Break

Architecture design

online actor-critic algorithm:

1. take action
$$\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{s})$$
, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
2. update \hat{V}^{π}_{ϕ} using target $r + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}')$
3. evaluate $\hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}') - \hat{V}^{\pi}_{\phi}(\mathbf{s})$
4. $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}|\mathbf{s}) \hat{A}^{\pi}(\mathbf{s}, \mathbf{a})$
5. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



+ simple & stable

- no shared features between actor & critic



Online actor-critic in practice

online actor-critic algorithm:

1. take action
$$\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{s})$$
, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
2. update \hat{V}^{π}_{ϕ} using target $r + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}')$ works best with a batch (e.g., parallel workers)
3. evaluate $\hat{A}^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}') - \hat{V}^{\pi}_{\phi}(\mathbf{s})$
4. $\nabla_{\theta} J(\theta) \approx \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}|\mathbf{s}) \hat{A}^{\pi}(\mathbf{s}, \mathbf{a})$
5. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

synchronized parallel actor-critic



asynchronous parallel actor-critic



Critics as state-dependent baselines

Actor-critic:

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

+ lower variance (due to critic)- not unbiased (if the critic is not perfect)

Policy gradient:
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - b \right)$$

+ no bias

- higher variance (because single-sample estimate)

can we use \hat{V}^{π}_{ϕ} and still keep the estimator unbiased?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t}) \right)$$

+ no bias

+ lower variance (baseline is closer to rewards)

Control variates: action-dependent baselines

$$\begin{split} 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] \\ V^{\pi}(\mathbf{s}_{t}) &= E_{\mathbf{a}_{t} \sim \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t})} \left[Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \\ A^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) - V^{\pi}(\mathbf{s}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - V_{\phi}^{\pi}(\mathbf{s}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - Q_{\phi}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{s}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{s}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t}, \mathbf{s}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t'}, \mathbf{s}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t'}, \mathbf{s}_{t'}) &= \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t'}, \mathbf{s}_{t'}) \\ \hat{A}^{\pi}(\mathbf{s}_{t'}$$

use a critic *without* the bias (still unbiased), provided second term can be evaluated Gu et al. 2016 (Q-Prop)

Eligibility traces & n-step returns

$$\hat{A}_{\mathrm{C}}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+1}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$$

 $\hat{A}_{\mathrm{MC}}^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_t)$

- + lower variance
- higher bias if value is wrong (it always is)

+ no bias

- higher variance (because single-sample estimate)



 $\hat{A}_{n}^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{t+n} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t}) + \gamma^{n} \hat{V}_{\phi}^{\pi}(\mathbf{s}_{t+n})$

choosing n > 1 often works better!

Generalized advantage estimation



Review

- Actor-critic algorithms:
 - Actor: the policy
 - Critic: value function
 - Reduce variance of policy gradient
- Policy evaluation
 - Fitting value function to policy
- Discount factors
 - Carpe diem Mr. Robot 🐶
 - ...but also a variance reduction trick
- Actor-critic algorithm design
 - One network (with two heads) or two networks
 - Batch-mode, or online (+ parallel)
- State-dependent baselines
 - Another way to use the critic
 - Can combine: n-step returns or GAE



Actor-critic examples

- High dimensional continuous control with generalized advantage estimation (Schulman, Moritz, L., Jordan, Abbeel '16)
- Batch-mode actor-critic
- Blends Monte Carlo and function approximator estimators (GAE)

Iteration 0



Actor-critic examples

- Asynchronous methods for deep reinforcement learning (Mnih, Badia, Mirza, Graves, Lillicrap, Harley, Silver, Kavukcuoglu '16)
- Online actor-critic, parallelized batch
- N-step returns with N = 4
- Single network for actor and critic



Actor-critic suggested readings

- Classic papers
 - Sutton, McAllester, Singh, Mansour (1999). Policy gradient methods for reinforcement learning with function approximation: actor-critic algorithms with value function approximation
- Deep reinforcement learning actor-critic papers
 - Mnih, Badia, Mirza, Graves, Lillicrap, Harley, Silver, Kavukcuoglu (2016). Asynchronous methods for deep reinforcement learning: A3C -- parallel online actor-critic
 - Schulman, Moritz, L., Jordan, Abbeel (2016). High-dimensional continuous control using generalized advantage estimation: batch-mode actor-critic with blended Monte Carlo and function approximator returns
 - Gu, Lillicrap, Ghahramani, Turner, L. (2017). Q-Prop: sample-efficient policygradient with an off-policy critic: policy gradient with Q-function control variate