# Reinforcement Learning

Designing, Visualizing and Understanding Deep Neural Networks

# CS W182/282A

Instructor: Sergey Levine UC Berkeley



### Recap: policy gradients



"reward to go"

#### Improving the policy gradient

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

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

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



#### 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?

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

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

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





fit  $V^{\pi}$ 

# 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: 
$$y_{i,t} = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[ r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t} \right] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \sum_{t' \in \mathbf{I}}^{T} (\underline{\mathbf{s}}_{i,t+1}) \sum_{\pi_{\theta}}^{T} [r((\mathbf{s}_{t',t}, \mathbf{a}_{t'})_{t'}] \approx r(\mathbf{s}_{i,t+1})$$
  
Monte Carlo target:  $y_{i,t} = \sum_{t'=t}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$  directly use previous fitted

directly use previous fitted 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:

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



### 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  $\infty$ ?  $\hat{V}^{\pi}_{\phi}$  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)

 $\gamma$  changes the MDP:



# 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  $\hat{V}^{\pi}_{\phi}$  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.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 

# 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  $\hat{V}^{\pi}_{\phi}$  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



# Can we use **just** a value function?

# Can we omit policy gradient completely?

 $A^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$ : how much better is  $\mathbf{a}_t$  than the average action according to  $\pi$ 

 $\arg \max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$ : best action from  $\mathbf{s}_t$ , if we then follow  $\pi$ 

at *least* as good as any  $\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)$ 

regardless of what  $\pi(\mathbf{a}_t|\mathbf{s}_t)$  is!







 $\pi \leftarrow \pi'$ 

# Policy iteration

High level idea:

policy iteration algorithm:

1. evaluate  $A^{\pi}(\mathbf{s}, \mathbf{a}) \longleftarrow$  how to do this? 2. set  $\pi \leftarrow \pi'$ 

$$\pi'(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg\max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

as before:  $A^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')] - V^{\pi}(\mathbf{s})$ let's evaluate  $V^{\pi}(\mathbf{s})!$ 



#### Dynamic programming

Let's assume we know  $p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ , and  $\mathbf{s}$  and  $\mathbf{a}$  are both discrete (and small)

0.2	<mark>0.3</mark>	0.4	0.3
0.3	0.3	<mark>0.</mark> 5	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5

16 states, 4 actions per state can store full  $V^{\pi}(\mathbf{s})$  in a table!  $\mathcal{T}$  is  $16 \times 16 \times 4$  tensor

bootstrapped update:  $V^{\pi}(\mathbf{s}) \leftarrow E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})}[V^{\pi}(\mathbf{s}')]]$ just use the current estimate here

$$\pi'(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg\max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases} \longrightarrow \text{ deterministic policy } \pi(\mathbf{s}) = \mathbf{a} \end{cases}$$

simplified:  $V^{\pi}(\mathbf{s}) \leftarrow r(\mathbf{s}, \pi(\mathbf{s})) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s}))}[V^{\pi}(\mathbf{s}')]$ 

# Policy iteration with dynamic programming



0.2	0.3	0.4	0.3
0.3	0.3	<mark>0.</mark> 5	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5

16 states, 4 actions per state can store full  $V^{\pi}(\mathbf{s})$  in a table!  $\mathcal{T}$  is  $16 \times 16 \times 4$  tensor

# Even simpler dynamic programming

$$\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

$$A^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')] - V^{\pi}(\mathbf{s})$$

 $\arg\max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \arg\max_{\mathbf{a}_t} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$ 

 $Q^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')]$  (a bit simpler)

skip the policy and compute values directly!

value iteration algorithm:

1. set 
$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$$
  
2. set  $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ 









# Fitted value iteration

how do we represent  $V(\mathbf{s})$ ?

big table, one entry for each discrete **s** neural net function  $V : S \to \mathbb{R}$ 



$$\mathcal{L}(\phi) = \frac{1}{2} \left\| V_{\phi}(\mathbf{s}) - \max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a}) \right\|^{2}$$

fitted value iteration algorithm:

1. set 
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$
  
2. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$ 

$$s = 0: V(s) = 0.2$$
  
 $s = 1: V(s) = 0.3$   
 $s = 2: V(s) = 0.5$ 



# curse of dimensionality

 $|\mathcal{S}| = (255^3)^{200 \times 200}$ (more than atoms in the universe)

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \mathbf{a})}[V^{\pi}(\mathbf{s}')]$$



# What if we don't know the transition dynamics?

fitted value iteration algorithm:

1. set 
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$
  
2. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$ 

need to know outcomes for different actions!

#### Back to policy iteration...

policy iteration: 1. evaluate  $Q^{\pi}(\mathbf{s}, \mathbf{a})$ 2. set  $\pi \leftarrow \pi'$   $\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$  policy evaluation:  $V^{\pi}(\mathbf{s}) \leftarrow r(\mathbf{s}, \pi(\mathbf{s})) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s}))}[V^{\pi}(\mathbf{s}')]$   $Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \mathbf{a})}[Q^{\pi}(\mathbf{s}', \pi(\mathbf{s}'))]$ can fit this using samples

# Can we do the "max" trick again?

policy iteration:

 $\begin{array}{c} \bullet \\ \bullet \\ 1. \text{ evaluate } V^{\pi}(\mathbf{s}) \\ 2. \text{ set } \pi \leftarrow \pi' \end{array}$ 

fitted value iteration algorithm:

1. set 
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$
  
2. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$ 

forget policy, compute value directly

can we do this with Q-values **also**, without knowing the transitions?

fitted Q iteration algorithm:

1. set  $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)] \leftarrow q = \text{approxiate } E[V(\mathbf{s}'_i)] \approx \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$ 2. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$  doesn't require simulation of actions!

+ works even for off-policy samples (unlike actor-critic)

+ only one network, no high-variance policy gradient

- no convergence guarantees for non-linear function approximation (more on this later)

### Fitted Q-iteration

full fitted Q-iteration algorithm:

1. collect dataset 
$$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$$
 using some policy  
2. set  $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   
3. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$ 

parameters

dataset size N, collection policy iterations Kgradient steps S



## Q-Learning

# Online Q-learning algorithms

full fitted Q-iteration algorithm:





off policy, so many choices here!

online Q iteration algorithm:

1. take some action 
$$\mathbf{a}_i$$
 and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ 

2. 
$$\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$$

3. 
$$\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$$

# Exploration with Q-learning

online Q iteration algorithm:

1. take some action 
$$\mathbf{a}_i$$
 and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$   
2.  $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   
3.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$ 

final policy:

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

why is this a bad idea for step 1?

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 - \epsilon \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ \epsilon/(|\mathcal{A}| - 1) \text{ otherwise} \end{cases}$$

"epsilon-greedy"

 $\pi(\mathbf{a}_t|\mathbf{s}_t) \propto \exp(Q_\phi(\mathbf{s}_t,\mathbf{a}_t))$ 

"Boltzmann exploration"

# What's wrong?

online Q iteration algorithm:

1. take some action 
$$\mathbf{a}_i$$
 and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$   
2.  $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   
3.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$   
isn't this just gradient descent? that converges, right?

Q-learning is *not* gradient descent!

$$\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{r}(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)])$$
  
no gradient through target value

# Correlated samples in online Q-learning

online Q iteration algorithm:

1. take some action  $\mathbf{a}_i$  and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$  - target value 2.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)])$ 

- sequential states are strongly correlated

- target value is always changing

### Replay buffers

online Q iteration algorithm:

1. take some action  $\mathbf{a}_i$  and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ 2.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)])$ 

full fitted Q-iteration algorithm:

1. collect dataset  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy  $K \times 2. \text{ set } \mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   $3. \text{ set } \phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$  special case with K = 1, and one gradient step

any policy will work! (with broad support) just load data from a buffer here still use one gradient step



### Replay buffers

Q-learning with a replay buffer:

▶ 1. sample a batch  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$  from  $\mathcal{B}$ 

+ samples are no longer correlated

• 2. 
$$\phi \leftarrow \phi - \alpha \sum_{i} \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_{i}, \mathbf{a}_{i})(Q_{\phi}(\mathbf{s}_{i}, \mathbf{a}_{i}) - [r(\mathbf{s}_{i}, \mathbf{a}_{i}) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_{i}, \mathbf{a}'_{i})])$$

+ multiple samples in the batch (low-variance gradient)

but where does the data come from?

need to periodically feed the replay buffer...



#### Putting it together

full Q-learning with replay buffer:

1. collect dataset  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy, add it to  $\mathcal{B}$ 

 $K \times \begin{cases} 2. \text{ sample a batch } (\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i) \text{ from } \mathcal{B} \\ 3. \phi \leftarrow \phi - \alpha \sum_i \frac{dQ_{\phi}}{d\phi} (\mathbf{s}_i, \mathbf{a}_i) (Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)]) \end{cases}$ 

K = 1 is common, though larger K more efficient



# What's wrong?

online Q iteration algorithm:

1. take some action 
$$\mathbf{a}_i$$
 and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$   
2.  $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   
3.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$   
use replay buffer

#### Q-learning is *not* gradient descent!

$$\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_{i}, \mathbf{a}_{i})(Q_{\phi}(\mathbf{s}_{i}, \mathbf{a}_{i}) - \mathbf{r}(\mathbf{s}_{i}, \mathbf{a}_{i}) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_{i}, \mathbf{a}'_{i})])$$
This is still a problem!

11

### Q-Learning and Regression

full Q-learning with replay buffer:

1. collect dataset  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy, add it to  $\mathcal{B}$ 

 $K \times \begin{cases} 2. \text{ sample a batch } (\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i) \text{ from } \mathcal{B} \\ 3. \phi \leftarrow \phi - \alpha \sum_i \frac{dQ_{\phi}}{d\phi} (\mathbf{s}_i, \mathbf{a}_i) (Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)]) \end{cases}$ 

#### one gradient step, moving target

full fitted Q-iteration algorithm:

1. collect dataset 
$$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$$
 using some policy  
2. set  $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$   
3. set  $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i ||Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i||^2$ 

#### perfectly well-defined, stable regression

#### Q-Learning with target networks

Q-learning with replay buffer and target network:

→ 1. save target network parameters: 
$$\phi' \leftarrow \phi$$

▶ 2. collect dataset  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy, add it to  $\mathcal{B}$ 

$$\times$$
  $\rightarrow$  3. sample a batch  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$  from  $\mathcal{B}$ 

4. 
$$\phi \leftarrow \phi - \alpha \sum_{i} \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_{i}, \mathbf{a}_{i})(Q_{\phi}(\mathbf{s}_{i}, \mathbf{a}_{i}) - [r(\mathbf{s}_{i}, \mathbf{a}_{i}) + \gamma \max_{\mathbf{a}'} Q_{\phi'}(\mathbf{s}'_{i}, \mathbf{a}'_{i})]$$

targets don't change in inner loop!

### "Classic" deep Q-learning algorithm (DQN)

Q-learning with replay buffer and target network:

1. save target network parameters: 
$$\phi' \leftarrow \phi$$
  
2. collect dataset  $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$  using some policy, add it to  $\mathcal{B}$   
 $N \times \mathbf{k} \times$ 

"classic" deep Q-learning algorithm:

1. take some action 
$$\mathbf{a}_i$$
 and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ , add it to  $\mathcal{B}$   
2. sample mini-batch  $\{\mathbf{s}_j, \mathbf{a}_j, \mathbf{s}'_j, r_j\}$  from  $\mathcal{B}$  uniformly  
3. compute  $y_j = r_j + \gamma \max_{\mathbf{a}'_j} Q_{\phi'}(\mathbf{s}'_j, \mathbf{a}'_j)$  using *target* network  $Q_{\phi'}$   
4.  $\phi \leftarrow \phi - \alpha \sum_j \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_j, \mathbf{a}_j)(Q_{\phi}(\mathbf{s}_j, \mathbf{a}_j) - y_j)$   
5. update  $\phi'$ : copy  $\phi$  every  $N$  steps

### Representing the Q-function



#### more common with continuous actions



more common with discrete actions

#### Back to actor-critic

online Q iteration algorithm: 1. take some action  $\mathbf{a}_i$  and observe  $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ 2.  $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$ 3.  $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$ 

with **continuous actions**, this is very inconvenient (but not impossible)

**Idea:** use actor-critic, but with Q-functions (to train off-policy)

# Simple practical tips for Q-learning

- Q-learning takes some care to stabilize
  - Test on easy, reliable tasks first, make sure your implementation is correct



Figure: From T. Schaul, J. Quan, I. Antonoglou, and D. Silver. "Prioritized experience replay". *arXiv preprint arXiv:1511.05952* (2015), Figure 7

- Large replay buffers help improve stability
  - Looks more like fitted Q-iteration
- It takes time, be patient might be no better than random for a while
- Start with high exploration (epsilon) and gradually reduce

# Q-learning with convolutional networks

- "Human-level control through deep reinforcement learning," Mnih et al. '13
- Q-learning with convolutional networks
- Uses replay buffer and target network
- One-step backup
- One gradient step
- Can be improved a lot with double Q-learning (and other tricks)





Kalashnikov, Irpan, Pastor, Ibarz, Herzong, Jang, Quillen, Holly, Kalakrishnan, Vanhoucke, Levine. QT-Opt: Scalable Deep Reinforcement Learning of Vision-Based Robotic Manipulation Skills

minimize 
$$\sum_{i} (Q(\mathbf{s}_i, \mathbf{a}_i) - [r(\mathbf{s}_i, \mathbf{a}_i) + \max_{\mathbf{a}'_i} Q(\mathbf{s}'_i, \mathbf{a}'_i)])^2$$

# Q-learning suggested readings

- Classic papers
  - Watkins. (1989). Learning from delayed rewards: introduces Q-learning
  - Riedmiller. (2005). Neural fitted Q-iteration: batch-mode Q-learning with neural networks
- Deep reinforcement learning Q-learning papers
  - Lange, Riedmiller. (2010). Deep auto-encoder neural networks in reinforcement learning: early image-based Q-learning method using autoencoders to construct embeddings
  - Mnih et al. (2013). Human-level control through deep reinforcement learning: Qlearning with convolutional networks for playing Atari.
  - Van Hasselt, Guez, Silver. (2015). Deep reinforcement learning with double Q-learning: a very effective trick to improve performance of deep Q-learning.
  - Lillicrap et al. (2016). Continuous control with deep reinforcement learning: continuous Q-learning with actor network for approximate maximization.
  - Gu, Lillicrap, Stuskever, L. (2016). Continuous deep Q-learning with model-based acceleration: continuous Q-learning with action-quadratic value functions.
  - Wang, Schaul, Hessel, van Hasselt, Lanctot, de Freitas (2016). Dueling network architectures for deep reinforcement learning: separates value and advantage estimation in Q-function.