

# 4/13/2023 Reinforcement Learning II

RL

Model-based

Learn  $T(s, a, s')$

old state  $\uparrow$  action  $\uparrow$  new state

and Reward  $(s, a, s')$

Learning the actual model of the world

Model-Free

Does not model the whole world

Q-Learning

Learn  $\hat{Q}(s, a)$   
to approximate  $Q_{opt}(s, a)$

optimal rewards  
starting at  $s$   
and taking action  $a$

Part 1 today Q-Learning  
with function approximation  
"Regression"

Policy Gradient

Directly predict best action  $a$  in a state  $s$

Part 2:  
one policy gradient algorithm  
"Classification"

## Q-Learning Recap



actions: 1 2 3

$\hat{Q}(s_4, 2)$

= our estimate of  $Q_{opt}(s_4, 2)$

$$\hat{Q}(s, a) \leftarrow (1-\eta) \hat{Q}(s, a) + \eta (r + \gamma \hat{V}(s'))$$

$\uparrow$   
 $\approx 0.1$

We act with policy  $\pi_{act}$  & upon observing  $(s, a, r, s')$ :

where  $\hat{V}(s') = \max_{a \in Actions(s')} \hat{Q}(s', a)$

$$\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \eta \left( \underbrace{r + \gamma \hat{V}(s')}_{\text{the target value}} - \underbrace{\hat{Q}(s,a)}_{\text{our prediction on input (s,a)}} \right)$$

Linear regression:

$$(y - w^T x)^2$$

$$-2 \cdot (y - w^T x) \cdot x$$

↑ This looks like gradient descent/ascent  
 It actually is where we minimize  

$$\text{Loss} = \frac{1}{2} (r + \gamma \hat{V}(s') - \hat{Q}(s,a))^2$$

Proof:  $\nabla_{\hat{Q}(s,a)} = \frac{1}{2} \cdot 2 (r + \gamma \hat{V}(s') - \hat{Q}(s,a)) \cdot -1$

So: Q-learning is just doing gradient descent on squared loss

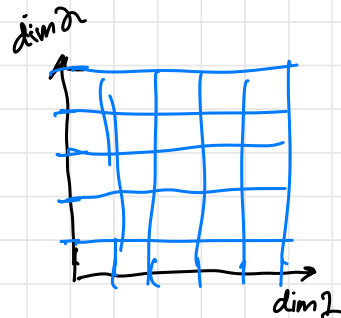
Going beyond tabular Q-learning

In real world, states are not discrete  
 - positions of joints in robot/objects

Can discretize by dividing each continuous dimension into buckets

$$\# \text{states} = (\# \text{buckets})^{\text{dimensions}}$$

Bad in high dimensions



Solution: Don't use a table, instead approximate  $Q_{opt}(s,a)$  with a learned model

First Attempt: Linear model

- Construct feature function  $\phi(s,a) \in \mathbb{R}^d$
- Predict  $\hat{Q}(s,a) = w^T \phi(s,a)$  for parameters  $w \in \mathbb{R}^d$

To do Q-Learning, minimize squared error

$$\text{Loss on } (s, a, r, s') = \frac{1}{2} (r + \gamma \hat{V}(s') - w^T \phi(s, a))^2$$

$$\nabla_w \text{loss} = \frac{1}{2} \cdot 2 \cdot (r + \gamma \hat{V}(s') - w^T \phi(s, a)) \cdot -\phi(s, a)$$

$$\hat{V}(s') \text{ is still } \max_{a \in \text{Actions}(s')} \underbrace{\hat{Q}(s', a)}_{= w^T \phi(s', a)}$$

Second Attempt: Neural network

Deep Q Network (DQN)

Idea:  $\hat{Q}_\theta(s, a)$  will be a neural network that maps  $(s, a)$  to estimate of  $Q_{opt}(s, a)$

Again, just from w/ gradient descent

$$\nabla_\theta \text{Loss} = - \underbrace{(r + \gamma \hat{V}(s') - \hat{Q}_\theta(s, a))}_{\text{Easy to compute}} \cdot \underbrace{\nabla_\theta \hat{Q}_\theta(s, a)}_{\text{Computable by backpropagation}}$$

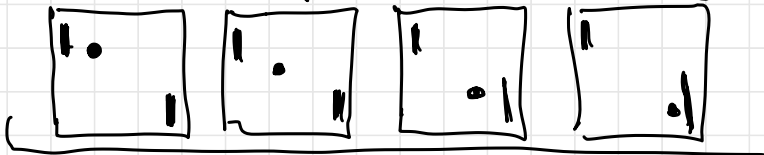
$\uparrow$   
All parameters of network

Example DQN architecture for video games

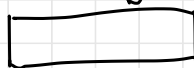
Represent state: Last  $k$  frames

- Each frame is  $84 \times 84$  (for Atari)

- Use last 4 frames  $\rightarrow 84 \times 84 \times 4$  state representation

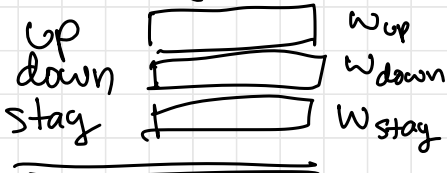


↓ Feed to CNN



$u(s)$ : vector encoding of the state

3 actions: 1 vector per action



Product:

$$Q(s, up) = w_{up}^T u(s)$$

$$Q(s, down) = w_{down}^T u(s)$$

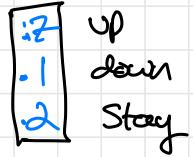
$$Q(s, stay) = w_{stay}^T u(s)$$

### Policy Gradient

$$\pi_{\theta}(a|s)$$

↑  
parameters

probability distribution over actions given current state



### How do we train?

- Normal classification: Given correct  $y$ 's for a bunch of  $x$ 's maximize  $P(y|x)$
- Policy gradient: Nobody tells you best action in any state what training algorithm to use?

Want to maximize value of policy  $\pi_{\theta}$ :

$$V(\theta) = \sum_{\text{trajectories } z} P(z; \theta) \cdot R(z)$$

Expected total rewards when using policy  $\pi_{\theta}$

Prob. of  $z$  happening when using  $\pi_{\theta}$

Reward for  $z = \sum_{t=1} r_t$

$$z = [s_1, a_1, \underline{r_1}, s_2, a_2, \underline{r_2}, s_3, \dots]$$

$V(\theta)$  is our training objective maximize with gradient ascent

What is  $\nabla_{\theta} V(\theta)$ ?

$$\nabla_{\theta} V(\theta) = \sum_{z} \nabla P(z; \theta) \cdot R(z)$$

Sum over exponentially many  $z$ 's — infeasible  
 Hope to instead compute an expected value over trajectories

Key trick:  $\nabla_{\theta} \log P(z; \theta) = \frac{1}{P(z; \theta)} \nabla P(z; \theta)$

$$\nabla P(z; \theta) = P(z; \theta) \cdot \nabla \log P(z; \theta)$$

$$\nabla_{\theta} V(\theta) = \sum_z P(z; \theta) \underbrace{\nabla \log P(z; \theta)}_{\text{this quantity}} \cdot R(z)$$

Expected Value of...

$$z = [s_1, a_1, r_1, s_2, a_2, r_2, s_3, \dots]$$

$$\log P(z; \theta) = \underbrace{\log P(s_1)}_{\substack{\text{Start state} \\ \text{prob.} \\ \text{in MDP}}} + \underbrace{\log \pi_{\theta}(a_1 | s_1)}_{\substack{\text{Prob of} \\ \text{taking } a_1 \\ \text{in} \\ \text{state } s_1}} + \underbrace{\log T(s_1, a_1, s_2)}_{\substack{\text{transition prob} \\ \text{of MDP}}}$$

$$+ \log \pi_{\theta}(a_2 | s_2) + \log T(s_2, a_2, s_3) + \dots$$

Do not depend on  $\theta$

$$\nabla \log P(z; \theta) = \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Final policy gradient algorithm:

Initialize  $\theta$  randomly

For each episode:

Sample trajectory  $z$ , using  $\pi_{\theta}(a | s)$

$$\theta \leftarrow \theta + \eta R(z) \cdot \sum_{t=1}^T \nabla \log \pi_{\theta}(a_t | s_t)$$

↑  
Gradient ascent  
on  $V(\theta)$