

# CS-UY 6923: Lecture 14

## Reinforcement Learning

---

NYU Tandon School of Engineering, Prof. Christopher Musco

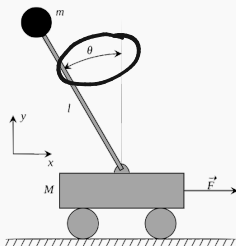
**Today:** Give flavor of the area and insight into one algorithm (Q-learning) which has been successful in recent years.

**Basic setup:**

- Agent interacts with environment over time  $1, \dots, t$ .
- Takes repeated sequence of **actions**,  $\underline{a}_1, \dots, \underline{a}_t$  which effect the environment.
- State of the environment over time denoted  $\underline{s}_1, \dots, \underline{s}_t$ .
- Earn **rewards**  $\underline{r}_1, \dots, \underline{r}_t$  depending on actions taken and states reached.
- Goal is to maximize reward over time.

## REINFORCEMENT LEARNING EXAMPLES

Classic inverted pendulum problem:

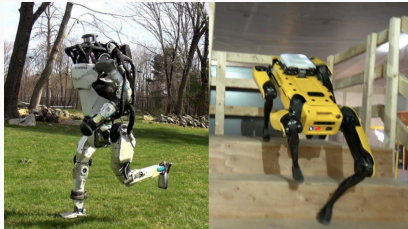


- **Agent**: Cart/software controlling cart.
- **Actions**: Move cart left or move right.
- **Reward**: 1 for every time step that  $|\theta| < 90^\circ$  (pendulum is upright)  $\textcircled{0}$  when  $|\theta| = 90^\circ$

# REINFORCEMENT LEARNING EXAMPLES

This problem has a long history in **Control Theory**. Other applications of classical control:

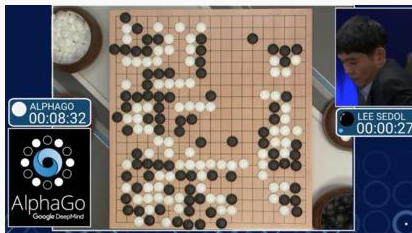
- Semi-autonomous vehicles (airplanes, helicopters, rockets, etc.)
- Industrial processes (e.g. controlling large chemical reactions)
- Robotics



control theory : reinforcement learning :: stats : machine learning

## REINFORCEMENT LEARNING EXAMPLES

Strategy games, like Go:



**State:** Position of all pieces on board.

**Actions:** Place new piece.

**Reward:** 1 if in winning position at time  $t$ . 0 otherwise.

This is a (sparse reward problem). Payoff only comes after many times steps, which makes the problem very challenging.

# REINFORCEMENT LEARNING EXAMPLES

Video games, like classic Atari games:



• **State:** Raw pixels on the screen (sometimes there is also hidden state which can't be observed by the player).

**Actions:** Actuate controller (up,down,left,right,click).

**Reward:** 1 if point scored at time  $t$ .

# MATHEMATICAL FRAMEWORK FOR RL

Model problem as a Markov Decision Process (MDP):  $3^{10}$

- $\mathcal{S}$ : Set of all possible states.  $|\mathcal{S}|$ .  $s \in \mathcal{S}$   $P_{\mathcal{S}}(s, a) = .5, .4, 1$
- $\mathcal{A}$ : Set of all possible actions.  $|\mathcal{A}|$ .  $-1, 0, 1$
- $\mathcal{R}$ : Set of possible rewards. Could have  $\mathcal{R} = \mathbb{R}$   $-1, -.9, \dots, .9, 1$
- Reward function  
 $R(s, a)$ :  $\mathcal{S} \times \mathcal{A} \rightarrow$  probability distribution over  $\mathcal{R}$ .  $r_t \sim R(s_t, a_t)$ .
- State transition function  
 $P(s, a)$ :  $\mathcal{S} \times \mathcal{A}$   $\rightarrow$  probability distribution over  $\mathcal{S}$ .  $s_{t+1} \sim P(s_t, a_t)$ .

Why is this called a Markov decision process? What does the term Markov refer to?

## MATHEMATICAL FRAMEWORK FOR RL

**Goal:** Find a policy  $\Pi: \mathcal{S} \rightarrow \mathcal{A}$  from states to actions which maximize expected cumulative reward.

- Start is state  $s_0$ .

- For  $t = 0 \dots, T$

  - $r_t \sim R(s_t, \Pi(s_t))$ .

  - $s_{t+1} \sim P(s_t, \Pi(s_t))$ .

action choose to play at time  $t$

The **time horizon**  $T$  could be short (game with fixed number of steps), very long (stock investing), or infinite. Goal is to maximize:

$$\underline{\text{reward}(\Pi)} = \mathbb{E} \sum_{t=0}^T r_t$$

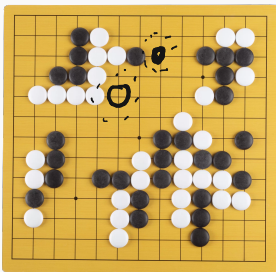
$[s_0, a_0, r_0], [s_1, a_1, r_1], \dots, [s_t, a_t, r_t]$  is called a trajectory of the MDP under policy  $\Pi$ .<sup>1</sup>

<sup>1</sup>It is not a priori clear that a fixed policy makes sense. Maybe we could get better reward by changing the policy over time. We will discuss this shortly.

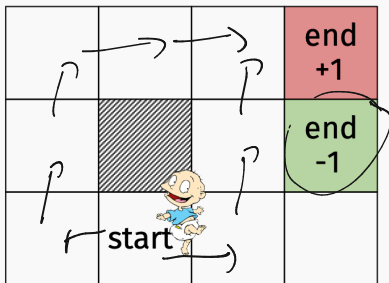


## FLEXIBILITY OF MDPs

- Can be used to model ~~time-varying environments~~. Just add time  $t$  to the state vector.
- Can be used to model games where actions have different effect if play in sequence (e.g. combo in a video game). Just add list of previous few actions to state.
- Can be used to model two-player games. Model adversary as part of the transition function.



## SIMPLE EXAMPLE: GRIDWORLD

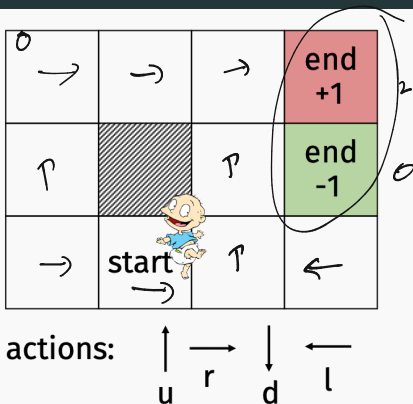


actions:  $\begin{matrix} \uparrow & \rightarrow \\ \text{u} & \text{r} \\ \downarrow & \leftarrow \\ \text{d} & \text{l} \end{matrix}$

- $r_t = -0.01$  if not at an end position.  $\pm 1$  if at end position.
- $P(s_t, a)$ : 70% of the time move in the direction indicated by  $a$ . 30% of the time move in a random direction.

**What is the optimal policy  $\Pi$ ?**

## SIMPLE EXAMPLE: GRIDWORLD



- $r_t = -0.5$  if not at an end position.  $\pm 1$  if at end position.
- $P(s_t, a)$ : 70% of the time move in the direction indicated by  $a$ . 30% of the time move in a random direction.

What is the optimal policy  $\Pi$ ?

## DISCOUNT FACTOR

$$r \geq 0$$

$$\mathbb{E} \sum_{t=0}^T r_t$$

For infinite or very long times horizon games (large  $T$ ), we often introduce a discount factor  $\gamma$  and seek instead to take actions which minimize:

$$\left( \mathbb{E} \sum_{t=0}^T \underline{\underline{\gamma^t r_t}} \right)$$

$$\gamma \leq 1$$

$$\gamma = .999$$

$$\gamma = .9$$

where  $r_t \sim R(s_t, \Pi(s_t))$  and  $s_{t+1} \sim P(s_t, \Pi(s_t))$  as before.

$\gamma \rightarrow 1$ : No discount. Standard MDP expected reward.

$\gamma \rightarrow 0$ : Care about short term reward more.

# VALUE FUNCTION

From now on assume  $T = \infty$ . We can do this without loss of generality by adding a time parameter to state and moving into an “end state” with no additional rewards once the time hits  $T$ .

**Value function:** Measures the expected return if we start in state  $s$  and follow policy  $\pi$ .

$$V^\pi(s) = \mathbb{E}_{\pi, s_0=s} \sum_{t \geq 0} \gamma^t r_t$$

Let  $\pi_s^*$  =  $\arg \max$   $V^\pi(s)$ . If we are in state  $s$ , at any point, we should always take action  $\pi_s^*(s)$ .



# VALUE FUNCTION

Value function:

$$V^\Pi(s) = \mathbb{E}_{\Pi, s_0=s} \sum_{t \geq 0} \gamma^t r_t$$

**Claim:** Let  $\Pi_s^* = \arg \max V^\Pi(s)$ . If we are in state  $s$ , at any point, we should always take action  $\Pi_s^*(s)$ .

**Proof:** Suppose we have already taken  $j - 1$  steps and seen trajectory  $[s_0, a_0, r_0], \dots, [s_{j-1}, a_{j-1}, r_{j-1}]$ . Then our expected reward is:

$$\begin{aligned} & \underline{r_0} + \underline{\gamma r_1} + \dots + \underline{\gamma^{j-1} r_{j-1}} + \mathbb{E}_\Pi \sum_{t \geq j} \gamma^t r_t \rightarrow \sum_{t \geq 0} \gamma^{t+j} r_{t+j} \\ & = \cancel{r_0} + \cancel{\gamma r_1} + \dots + \cancel{\gamma^{j-1} r_{j-1}} + \gamma^j \left( \mathbb{E}_\Pi \sum_{t \geq 0} \gamma^t r_{t+j} \right) \quad \begin{matrix} \downarrow \\ \gamma^t \gamma^j \end{matrix} \\ & = \underline{r_0 + \gamma r_1 + \dots + \gamma^{j-1} r_{j-1}} + \gamma^j \underline{V^\Pi(s_j)} \rightarrow \text{to maximize, choose } \Pi = \Pi_{s_j}^*(s_j) \end{aligned}$$

Value function:

 $\Pi$ 

$$V^\Pi(s) = \mathbb{E}_{\Pi, s_0=s} \sum_{t \geq 0} \gamma^t r_t$$

**Claim:** Let  $\Pi_s^* = \arg \max V^\Pi(s)$ . If we are in state  $s$ , at any point, we should always take action  $\Pi_s^*(s)$ .

So, there is a single optimal policy  $\Pi^*$  which simultaneously maximizes  $V^\Pi(s)$  for all  $s$ . I.e.  $\Pi_1^* = \Pi_2^* = \dots = \Pi_{|S|}^* = \Pi^*$ . We do not need to change the policy over time to maximize expected reward.

Goal in RL is to find this optimal policy  $\Pi^*$ .

## TWO SETTINGS

**Full information:** We know  $\mathcal{S}$ ,  $\mathcal{A}$ , the transition function  $P$  and reward function  $R$ . The optimal policy can  $\Pi^*$  can be found via dynamic programming. Sometimes called “planning” problem.

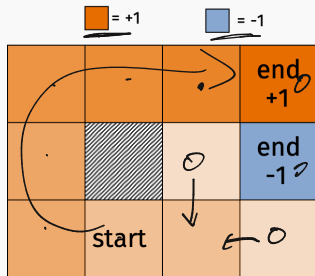
**Reinforcement Learning setting:** We do not know  $P$  or  $R$ , but we can repeatedly play the MDP, running whatever policy we like.



## VALUE ITERATION

$$V^{\pi}(s)$$

Let  $V^*(s) = V^{\pi^*}(s)$ . This function is equal to the expected future reward if we play optimally start in state  $s$ .



# VALUE ITERATION

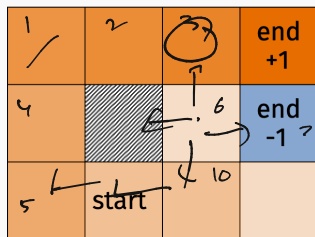
In the full information setting, if we knew  $V^*$  we can easily find the optimal policy  $\Pi^*$ .

+1, -1, .01

■ = +1    ■ = -1

$$\sum_{s', r} \Pr(s', r | s, a) \cdot (r + \gamma V(s))$$

$$\begin{aligned} & P_r(3, +1) \cdot ( \quad ) \\ & + P_r(3, -1) \cdot ( \quad ) \\ & + \underline{P_r(3, .01)} \cdot (.01 + \underline{\quad}) \end{aligned}$$



$$\begin{aligned} & + P(3, +1) ( \quad ) \\ & + P(3, -1) (-1 + \underline{\quad}) \\ & \quad \quad \quad \underbrace{\quad}_{.15} \end{aligned}$$

$$\underline{\Pi^*(s)} = \arg \max_a \left( \sum_{s', r} \Pr(s', r | s, a) \underline{[r + \gamma V^*(s')] } \right)$$

# VALUE ITERATION

$$\hat{V}(s) \neq \max_a \sum \dots [r + \gamma \hat{V}(s')] ]$$

$V^*(s)$  satisfies what is called a Bellman equation:

$$\underline{V^*(s)} = \left( \max_a \sum_{s',r} \cdot \Pr(s', r | s, a) [r + \gamma \underline{V^*(s')}] \right)$$

$$r \sim B(s, a)$$

Run a fixed point iteration to find  $V^*$ :

- Start with initial guess  $V^0$

$$\underline{V^{i-1}} = \underline{V^*}$$

- For  $i = 1, \dots, z$ :

- For  $\underline{s} \in \underline{S}$ :

$$\underline{V^i(s)} = \max_a \sum_{s',r} \cdot \Pr(s', r | s, a) [r + \gamma \underline{V^{i-1}(s')}]$$

Can be shown to converge in roughly  $z = \frac{1}{1-\gamma}$  iterations. What is the computational cost of each iteration?

$$\gamma = .9 \quad \frac{1}{1-.9} = 10 \quad .99999 \quad \frac{1}{1-.99999} \rightarrow 10000$$

## TWO SETTINGS

**Full information:** We know  $\mathcal{S}$ ,  $\mathcal{A}$ , the transition function  $P$  and reward function  $R$ .

**Reinforcement Learning setting:** We do not know  $P$  or  $R$ , but we can repeatedly play the MDP, running whatever policy we like.

- **Model-based** RL methods essentially try to learn  $P$  and  $R$  very accurately and then find  $\Pi^*$  via dynamic programming. Require a lot of samples of the MDP.



How many parameters do we need to learn if we hope to learn

$P$  and  $R$ ?  $P(s, a) \rightarrow s' \quad (|\mathcal{S}| \cdot |\mathcal{S}| \cdot |\mathcal{A}| \cdot |\mathcal{R}|)$

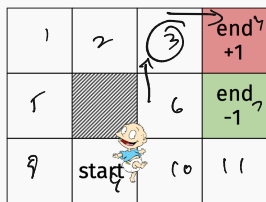
- **Model-free** RL methods try to learn  $\Pi^*$  without necessarily obtaining an accurate model of the world – i.e. without learning  $P$  and  $R$ .

## Q FUNCTION

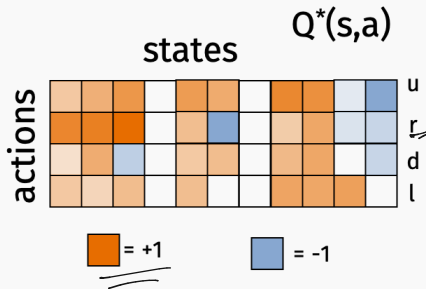
Another important function:  $V^\pi(s)$

- **Q-function:**  $Q^\pi(\underline{s}, \underline{a}) = \mathbb{E}_{\Pi, s_0=\underline{s}, a_0=\underline{a}} \sum_{t \geq 0} \gamma^t r_t$ . Measures the expected return if we start in state  $s$ , play action  $a$ , and then follow policy  $\Pi$ .

$$\left( \underline{Q^*(s, a)} \right) = \max_{\Pi} Q^\pi(\underline{s}, \underline{a}) = Q^{\Pi^*}(s, a).$$

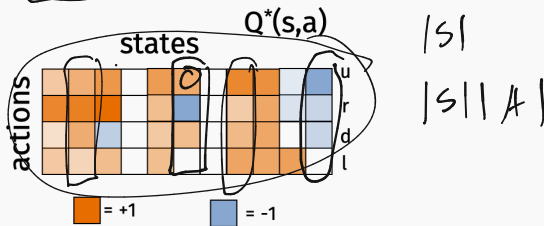


actions:  $\begin{matrix} \uparrow & \rightarrow & \downarrow & \leftarrow \\ u & r & d & l \end{matrix}$



$$\underline{Q^*(s, a)} = \max_{\pi} \mathbb{E}_{\pi, s_0=s, a_0=a} \sum_{t \geq 0} \gamma^t r_t.$$

If we knew the function  $Q^*$ , we would immediately know an optimal policy. Whenever we're in state  $s$ , we should always play action  $a^* = \arg \max_a Q^*(s, a)$ .



$Q$  has more parameters than  $V$ , but you can use it to determine an optimal policy without knowing transition probabilities.

## BELLMAN EQUATION

$Q^*$  also satisfies a Bellman equation:

$$\underbrace{Q^*(s, a)} = \underbrace{\mathbb{E}[R(s, a)]} + \underbrace{\gamma \mathbb{E}_{s' \sim P(s, a)}} \max_{a'} \underbrace{Q^*(s', a')}.$$

Bellman equation:

$$r \sim R(s, A)$$

$$\mathbb{E}[r] = \mathbb{E}[R(s, A)]$$

$$Q^*(s, a) = \mathbb{E}[R(s, a)] + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q^*(s', a')$$

Again use fixed point iteration to find  $Q^*$ . Let  $Q^{i-1}$  be our current guess for  $Q^*$  and suppose we are at some state  $s, a$ .

For  $i = 1, \dots, 2$   $r \sim R(s, A)$

$s' \sim P(s, a)$

$$Q^i(s, a) = \mathbb{E}[R(s, a)] + \gamma \mathbb{E}_{s' \sim P(s, a)} \max_{a'} Q^{i-1}(s', a')$$

In reality, drop expectations and use a learning rate  $\alpha$

$\alpha = .01$

$$Q^i(s, a) = (1 - \alpha) Q^i(s, a) + \alpha \left( \overset{r}{\cancel{R(s, a)}} + \gamma \max_{a'} Q^{i-1}(s', a') \right)$$

$$Q^i(s, a) = r + \gamma \max_{a'} Q^{i-1}(s', a')$$



How do we choose states  $s$  and  $a$  to make the update for? In principal you can do anything you want! E.g. choose some policy  $\Pi$  and run:

- Initialize  $Q^0$  (e.g. all zeros)
  - Start at  $s$ , play action  $a = \Pi(s)$ , observe reward  $R(s, a)$ .
  - For  $i = 1, \dots, Z$ 
    - $Q^i(s, a) = (1 - \alpha)Q^i(s, a) + \alpha (r + \gamma \max_{a'} Q^{i-1}(s', a'))$
    - $s \leftarrow P(s, a)$
    - $a \leftarrow \Pi(s)$
- (restart if we reach a terminating state)

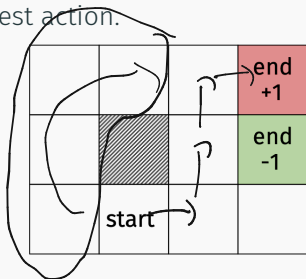
Q-learning is considered an off-policy RL method because it runs a policy  $\Pi$  that is not necessarily related to its current guess for an optimal policy, which in this case would be  $\Pi(s) = \max_a Q^i(s, a)$  at time  $i$ .

## EXPLORATION VS. EXPLOITATION

For small enough  $\alpha$ , Q-learning converges to  $Q^*$  as long as we follow a policy  $\Pi$  that visits every start  $(s, a)$  with non-zero probability.

Mild condition, but exact choice of  $\Pi$  matters for convergence rate.

- **Random:** At state  $s$ , choose a random action  $a$ .
- **Greedy:** At state  $s$ , choose  $\arg \max_a Q^i(s, a)$ . I.e. the current guess for the best action.

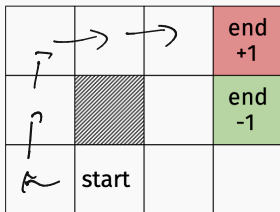
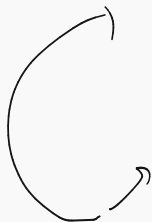


**Random** can be wasteful. Spend time improving parts of  $Q$  that aren't relevant to optimal play. **Greedy** can cause you to zero in on a locally optimal policy without learning new strategies.

## EXPLORATION VS. EXPLOITATION

Possible choices for  $\Pi$ :

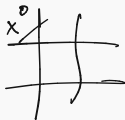
- **Random:** At state  $s$ , choose a random action  $a$ .
- **Greedy:** At state  $s$ , choose  $\arg \max_a Q^i(s, a)$ . I.e. the current guess for the best action.
- **$\epsilon$ -Greedy:** At state  $s$ , choose  $\arg \max_a Q^i(s, a)$  with probability  $1 - \epsilon$  and a random action with probability  $\epsilon$ .



Exploration-exploitation tradeoff. Increasing  $\epsilon$  = more exploration.

## CENTRAL ISSUE IN MODERN REINFORCEMENT LEARNING

**Another issue:** Even writing down  $Q^*$  is intractable... This is a function over  $|S||A|$  possible inputs. Even for relatively simple games,  $|S|$  is gigantic...



Back of the envelope calculations:

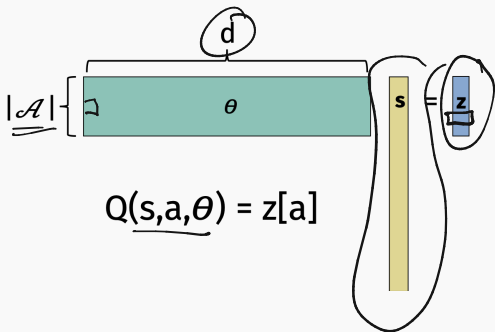
- Tic-tac-toe:  $3^{(3 \times 3)} \approx (20,000)$
- Chess:  $\approx 10^{43} < 28^{64}$  (due to Claude Shannon).
- Go:  $3^{(19 \times 19)} \approx 10^{171}$ .
- Atari:  $128^{(210 \times 160)} \approx 10^{71,000}$ .

Number of atoms in the universe:  $\approx 10^{82}$

## MACHINE LEARNING APPROACH

Learn a **simpler** function  $Q(s, a, \theta) \approx Q^*(s, a)$  parameterized by a small number of parameters  $\theta$ .

**Example:** Suppose our state can be represented by a vector in  $\mathbb{R}^d$  and our action  $a$  by an integer in  $1, \dots, |\mathcal{A}|$ . We could use a linear function where  $\theta$  is a small matrix:

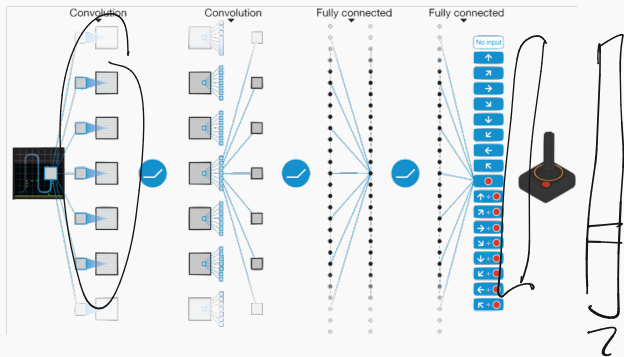


$|S|$

## MACHINE LEARNING APPROACH

Learn a **simpler** function  $Q(s, a, \theta) \approx Q^*(s, a)$  parameterized by a small number of parameters  $\theta$

**Example:** Could also use a (deep) neural network.



DeepMind: “Human-level control through deep reinforcement learning”, Nature 2015.

If  $Q(s, a, \theta)$  is a good approximation to  $Q^*(s, a)$  then we have an approximately optimal policy:  $\tilde{\pi}^*(s) = \arg \max_a Q(s, a, \theta)$ .

- Start in state  $s_0$ .
- For  $t = 1, 2, \dots$ 
  - $a^* = \arg \max_a Q(s, a, \theta)$
  - $s_t \sim P(s_{t-1}, a^*)$



How do we find an optimal  $\theta$ ? If we knew  $Q^*(s, a)$  could use supervised learning, but the true  $Q$  function is infeasible to compute.

Find  $\theta$  which satisfies the Bellman equation:

$$\begin{aligned}
 \underline{Q^*(s, a)} &= \underline{\mathbb{E}_{s' \sim P(s, a)}} \left[ \underline{R(s, a)} + \underline{\gamma \max_{a'} Q^*(s', a')} \right] \\
 \left\{ \underline{Q(s, a, \theta)} \approx \underline{\mathbb{E}_{s' \sim P(s, a)}} \left[ \underline{R(s, a)} + \underline{\gamma \max_{a'} Q(s, a, \theta)} \right] \right\}
 \end{aligned}$$

Should be true for all  $a, s$ . Should also be true for  $\underline{a, s \sim \mathcal{D}}$  for any distribution  $\mathcal{D}$ :

$$\underline{\mathbb{E}_{s, a \sim \mathcal{D}} Q(s, a, \theta)} \approx \underline{\mathbb{E}_{s, a \sim \mathcal{D}} \mathbb{E}_{s' \sim P(s, a)} \left[ \underline{R(s, a)} + \underline{\gamma \max_{a'} Q(s, a, \theta)} \right]}.$$

Loss function:

$$\underline{L(\theta)} = \underline{\mathbb{E}_{s, a \sim \mathcal{D}}} (y - Q(s, a, \theta))^2$$

where  $y = \mathbb{E}_{s' \sim P(s, a)} [R(s, a) + \gamma \max_{a'} Q(s', a', \theta)]$ .



## Q-LEARNING W/ FUNCTION APPROXIMATION

Minimize loss with **gradient descent**:

$$\underline{\nabla L(\theta)} = \underline{\mathbb{E}_{s,a \sim \mathcal{D}}} [\underline{-2 \nabla Q(s, a, \theta)} \cdot \underline{[y - Q(s, a, \theta)]]}$$

In practice use stochastic gradient:

$$\underline{\nabla L(\theta, s, a)} = -2 \cdot \nabla Q(s, a, \theta) \cdot \left[ R(s, a) + \gamma \underline{\max_{a'} Q(s', a', \theta)} - \underline{Q(s, a, \theta)} \right]$$

- Initialize  $\theta_0$
- For  $i = 0, 1, 2, \dots$ 
  - Run policy  $\Pi$  to obtain  $s, a$  and  $s' \sim P(s, a)$
  - Set  $\theta_{i+1} = \theta_i - \eta \cdot \nabla L(\theta_i, s, a)$

$\eta$  is a learning rate parameter.

Again, the choice of  $\Pi$  matters a lot. **Random play** can be wastefully, putting effort into approximating  $Q^*$  well in parts of the state-action space that don't actually matter for optimal play.  $\epsilon$ -greedy approach is much more common:

- Initialize  $s_0$ .
- For  $t = 0, 1, 2, \dots$ ,
  - $a_i = \begin{cases} \arg \max_a Q(s_t, a, \theta_{curr}) & \text{with probability } (1 - \epsilon) \\ \text{random action} & \text{with probability } \epsilon \end{cases}$

Lots of other details we don't have time for! References:

- Original DeepMind Atari paper: 2015  
<https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>,  
which is very readable.
- Stanford lecture video:  
<https://www.youtube.com/watch?v=lvoHnicueoE> and  
slides: [http://cs231n.stanford.edu/slides/2017/  
cs231n\\_2017\\_lecture14.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture14.pdf)

Important concept we did not cover: **experience replay**.



<https://www.youtube.com/watch?v=V1eYniJ0Rnk>

- Don't forget about the last problem set!
- I will release a study document for the exam and also schedule and extra office hours for next week.