# 11/14/2023: Reinforcement Learning Contd.

$s^1$
$s^2$
:
states :
$s^{|s|}$

$a^1 \ a^2 \ a^3 \ a^4 \ a^5$

At every timestep:
   Start at some state $s$
   Take an action $a$
   Get a reward $r$
   Transition to new state $s'$

This is "1 training example" for
Q learning

## Tabular Q-Learning:
   each $\hat{Q}(S,a)$ is 1 parameter to learn
Total # of params is

$$\boxed{\# \text{States} \times \# \text{actions}}$$ ← Can be very very large

## Q-Learning Update Rule:
Input: $(s, a, r, s')$
Update

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

Our guess of
Q-value for action
we just took

learning
rate
(e.g. 0.1)

old
guess

immediate
reward

anticipated
future
reward

Estimate of total
future reward

Where $\hat{V}(s) = \begin{cases} \max\limits_{a \in \text{Actions}(s)} \hat{Q}(S,a) & \text{if NOT IsEnd}(s) \\ 0 & \text{else} \end{cases}$

Our guess of how good a state is
based on $\hat{Q}$

One more consideration: How to choose actions
during training?

Obvious answer (wrong):
At state s, choose $a = \text{argmax}_{a'} \hat{Q}(s, a')$

- Optimal If $\hat{Q}$ values are accurate
- Bad idea early in training!

Suppose we do action $a$ in state s once,
receive large reward.

$\Rightarrow \hat{Q}(s,a)$ will update to be large (larger than other actions at same state)

$\Rightarrow$ Naive policy always chooses $a$ in state s
forever

All exploitation,    no exploration

use knowledge     try different things to
we've learned     see what's best

Simple Solution: $\varepsilon$ - greedy
At each timestep: At state s
  - with probability $1-\varepsilon$, choose $\text{argmax}_a \hat{Q}(s,a)$
                                               (Exploitation)

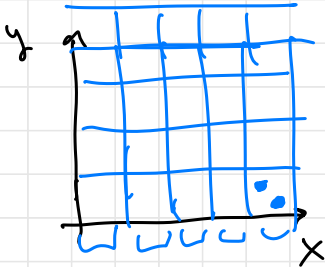  - with probability $\varepsilon$, choose random action
       (exploration)

Usually choose small but nonzero $\varepsilon$ during training
                          (eg $\varepsilon = 0.1$)

At test time, use $\varepsilon = 0$

# How to deal with very large state spaces?

Option #1: Discretize the state
- States might be continuous (eg. location)
- Divide a continuous dimension into buckets



Discretize $xy$ plane to $5 \times 5$ grid = 25 states

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

Bad in high dimensions

Option #2: Q-learning with Linear function approximation

Idea: Q-learning is kind of like regression:
Input: $(s,a)$, output = Q-value

So let's learn a linear model:
1. Need feature function $\phi(s,a) \in \mathbb{R}^d$
2. Learn parameter vector $w \in \mathbb{R}^d$
   to predict $\hat{Q}(s,a) = w^T \phi(s,a)$

How to learn $w$? Revisit Q-Learning update rule

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

$$= \hat{Q}(s,a) + \eta \left( \underbrace{r + \gamma \hat{V}(s')}_{\text{"target"}} - \underbrace{\hat{Q}(s,a)}_{\text{current prediction}} \right)$$

Review: linear regression:
$$\nabla_w (w^T x - y)^2$$
$$= 2(w^T x - y) \cdot x$$

$$= \hat{Q}(s,a) - \eta \underbrace{\left( \hat{Q}(s,a) - r - \gamma \hat{V}(s') \right)}_{\text{Basically the linear regression gradient}}$$

So, for Q-learning w/ function approx:

minimize squared error between

$$\underbrace{\hat{Q}(s,a)}_{\text{Prediction}} \quad \text{and} \quad \underbrace{r + \gamma \hat{V}(s')}_{\text{"target"}}$$

$$\text{Loss}(w) = \frac{1}{2}\left(r + \gamma \hat{V}(s') - \underbrace{w^T \phi(sa)}_{= \hat{Q}(s,a)}\right)^2$$

Gradient:

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

Update Rule:

$$w \leftarrow w - \eta \nabla_w \text{loss}(w)$$
$$= w + \eta\left(r + \gamma \hat{V}(s') - w^T \phi(sa)\right)\phi(s,a)$$

## Option 3: Deep Q Network (DQN)

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

Let $\Theta$ be parameters of network:

$$\text{Loss}(\Theta) = \frac{1}{2}\left(\underbrace{r + \gamma \hat{V}(s')}_{\text{"target"}} - \underbrace{\hat{Q}_\Theta(s,a)}_{\text{prediction}}\right)^2$$

$$\nabla_\Theta \text{loss}(\Theta) = \frac{1}{2} \cdot 2 \cdot \underbrace{\left(r + \gamma \hat{V}(s') - \hat{Q}_\Theta(s,a)\right)}_{\text{Compute directly}} \cdot \underbrace{- \nabla_\Theta \hat{Q}_\Theta(s,a)}_{\substack{\text{Compute by} \\ \text{backpropagation}}}$$
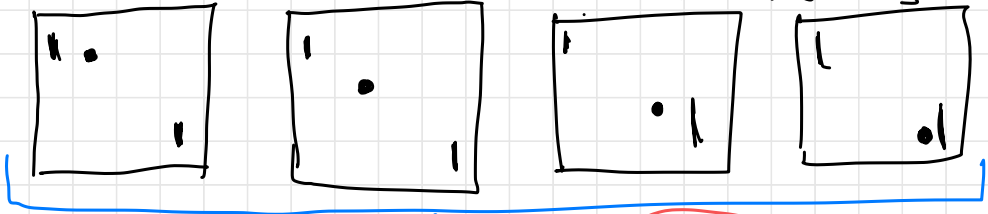
Run gradient descent to update $\Theta$

# Example DQN to play Pong

Represent state of game with last k frames
- Each frame is $84 \times 84$ image
- Set $k = 4 \rightarrow$ Input to DQN is $84 \times 84 \times 4$
  block of numbers



√ Feed to (CNN), generates a
vector $U(s)$

□ $U(s)$

Has lots of parameters

3 actions: Learn 1 vector 1 per action

up:     □  $W_{up}$
down:   □  $W_{down}$
stay:   □  $W_{stay}$

parameters

Predictions:
$$\hat{Q}(s, up) = W_{up}^T \, U(s)$$
$$\hat{Q}(s, down) = W_{down}^T \, U(s)$$
$$\hat{Q}(s, stay) = W_{stay}^T \, U(s)$$

## Taxonomy of RL methods

**Model-Based RL**
Learns transitions + Rewards

**Model-Free RL**
Dont try to directly learn transition & reward probabilities

**Q-Learning**
Learn $\hat{Q}(s,a)$
Infer good policy by maximizing $\hat{Q}(s,a)$

**Policy Gradient**
Directly learn a policy
≈ classifier to predict action given state