# RL

model-based

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

learn the actual model
of the world

model-free

don't model the world

## Q-learning

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

optimal reward at $s$ and
take action $a$

Q-learning with function
approximation

## policy gradient

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

one policy gradient
algorithm

## Q-learning



S1
S2  states
S3
S4

$a_1$ $a_2$ $a_3$
actions
$\hat{Q}(S_4, 2)$
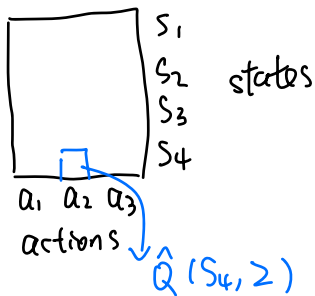
= estimate of $Q_{opt}(S_4, 2)$

we act with $\pi_{act}$ upon
observing $(s, a, r, s')$

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

$\approx 0.1$   where   $\hat{V}(s') = \max\limits_{a \in Action(s')} \hat{Q}(s', a)$

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

<span style="color:blue">↑</span>

<span style="color:blue">looks like gradient descent / ascent</span>

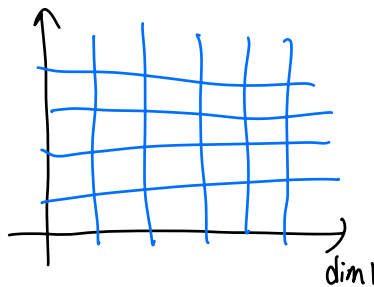<span style="color:blue">minimize $\text{Loss} = \frac{1}{2}(r + \gamma \hat{V}(s') - \hat{Q}(s,a))^2$</span>

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

$\Rightarrow$ Q-learning is doing gradient descent on squared loss

In real world, states are not discrete.

discretize by dividing each continuous dimension into buckets.

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



dim 1

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

first attempt linear model:
- constant feature function $\phi(s,a) \in \mathbb{R}^d$
- predict $\hat{Q}(s,a) = w^\top \phi(s,a)$ for $w \in \mathbb{R}^d$

To do Q-learning, minimize squared error

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

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

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

$$= w^T \phi(s',a)$$

second attempt: neural network

Deep Q Network (DQN)

$\hat{Q}(s,a)$ will be a neural network that maps $(s,a)$ to

estimate $Q_{opt}(s,a)$

$$\nabla_\theta \text{loss} = -(r + \gamma \hat{V}(s') - \hat{Q}_\theta(s,a)) \cdot \nabla_\theta \hat{Q}(s,a)$$

parameters
of NN

easy to compute

computable by
backpropagation

Example DQN architecture for video game
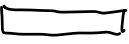represent state: last k frames
   − each frame is 84×84



feed to CNN  ↓

□  u(s) : vector encoding of state

3 actions: 1 vector per action

up ☐ $W_{up}$
down ☐ $W_{down}$
stay ☐ $W_{stay}$

predict

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

## Policy Gradient

$\Pi_\theta (a|s)$ probability distribution over
↑ actions given current state
params

$\begin{bmatrix} a_1 \\ a_2 \\ 0.7 \end{bmatrix}$ up
down
stay

How to train?

· normal classification: given correct y's for x's, maximize $p(y|x)$

· policy gradient: don't know best action in any state

Want to maximize value of policy $\Pi_\theta$:

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

expected total rewards when using $\Pi_\theta$

rewards for $z$: $\sum_{t=1}^T r_t$

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

$V(\theta)$ is for training,
    maximize with gradient descent

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 compute expected value over trajectories

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

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

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

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

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

$$+ \log \pi_\theta(a_2 | s_2) + \underbrace{\log T(s_2, a_2, s_3)}_{\text{independent of } \theta}$$

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

Final policy gradient algorithm

  initialize $\theta$ randomly
  for each episode:
     sample trajectory $z$ using $\pi_\theta(a|s)$
     $\theta \leftarrow \theta + \eta \, R(z) \cdot \sum_{t=1}^{T} \nabla \log \pi_\theta(a_t | s_t)$
              $\uparrow$
     gradient ascent on $V(\theta)$