

WEEK 3: Reinforcement Learning

Example: Autonomous helicopter

Task: 10 times per second, decide how to move the control sticks.

How: Map state s \rightarrow action a

You could use supervised learning to learn the mapping $s \rightarrow a$. But it turns out the "correct" action is ambiguous.

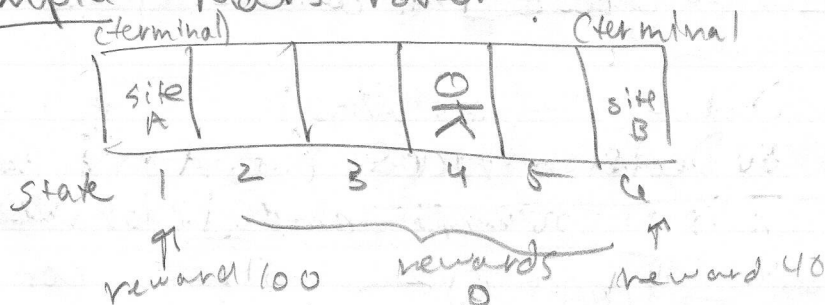
key input: reward function (like training a dog)

e.g. +1 for every second flying well,
-1000 for badly

Applications:

- controlling robots
- factory optimization
- playing games
- stock selling (e.g. to avoid washing the price)

Example: Mars rover!



at each step rover can go left or right
Go left: $(4, \leftarrow, 0, 3) = (s, a, R(s), s')$

How to know if a set of rewards is better than another? Return

Return for actions producing rewards R_1, R_2, \dots, R_n with discount factor γ :

$$\star \left[R = \sum_{i=1}^n \gamma^{i-1} R_i \right] \quad (\gamma < 1)$$

This prioritizes actions that give high rewards sooner (i.e. discount later rewards)

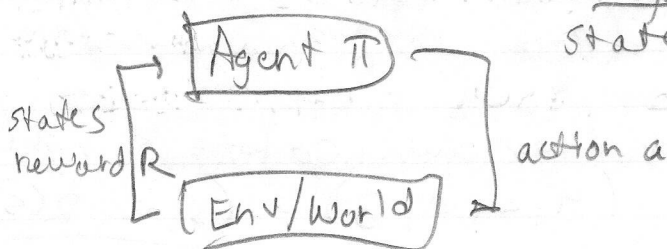
γ maps well to interest rate:
time value of money

Note that negative rewards would cause the algorithm to delay the rewards!

Formalization

Goal: define a policy $\pi: s \rightarrow a$ to maximize the return.

Markov Decision Process: future depends only on the current state, not history!



State-action value function

key quantity computed by algs. Learn

$Q: (s, a) \rightarrow$ return i.e. taking action a
 once w/ state s and the behavior
optimally after seems circular!

"Q function"
 (sometimes Q^*)

- Can compute this and then pick actions. (Recall there are multiple acceptable end states, so just pursuing the best outcome isn't necessarily the best return.)

- Let $\pi(s) = a$ where $Q(s, a) = \max_{a'} Q(s, a')$

Bellman Equation

let $S = \text{cur. state}$, $R(s) = \text{reward of cur. state}$

$a = \text{action taken in } S$

$S' = \text{state after taking } a$

$a' = \text{action taken in } S'$

$$\text{Then } \boxed{Q(s, a) = R(s) + \gamma \max_{a'} Q(s', a')}$$

(obvious)

Continuous State Spaces

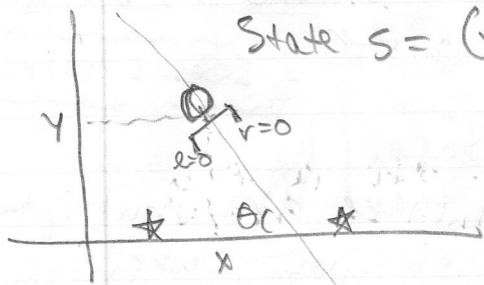
E.g. helicopter, state = $(x, y, z, \underbrace{\phi, \theta, \omega}_{\text{pitch}}, \underbrace{\dot{x}, \dot{y}, \dot{z}, \dot{\phi}, \dot{\theta}, \dot{\omega}}_{\text{velocity}})$

where $x, y, z, \phi, \theta, \omega: t \rightarrow \mathbb{R}$

Or, Lunar Lander:

Actions = { do nothing, fire left thruster, fire right thruster, fire main }

State $s = (x, y, \dot{x}, \dot{y}, \theta, \dot{\theta}, l, r)$



angle
whether left/right legs are on the ground, $\in \{0, 1\}$

Reward fn: at landing pad +100 -140
toward/away from addl.

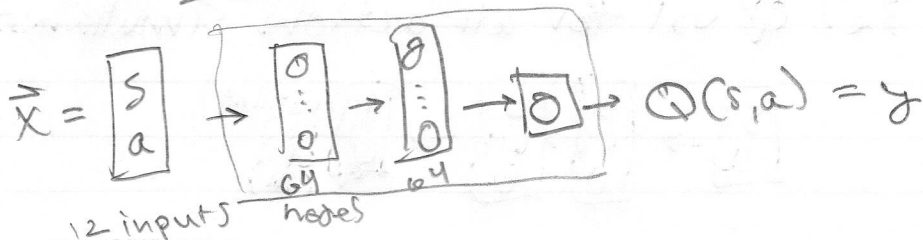
(but choice of reward fn determines what is learned incentives + costs)

crash	-140
soft land	+100
leg grounded	+10
fire main	-0.3
fire left/right	-0.03

$\gamma = 0.985$

learn policy π s.t. given s , $a = \pi(s)$ maximizes return.

How? Train NN to approximate Q . Deep RL



where s has 8 dims. and a is a one-hot encoding of $a \in \{\text{nothing}, L, R, \text{main}\}$.

→ Given state s , compute $Q(s, a) \forall a \in A$, and pick the a that maximizes.

Now use supervised learning to train the NN. But how to get training data?

→ Take a bunch of random actions in order to observe samples $(s, a, R(s), s')$. (Here, fire thrusters and observe outcomes.) And take Q random.

1. Init NN randomly as guess of Q .

2. Loop:

a. take action in lander, get $(s, a, R(s), s')$

b. store 10k most recent samples (replay buffer)

c. (occasionally) train NN for Q_{new}

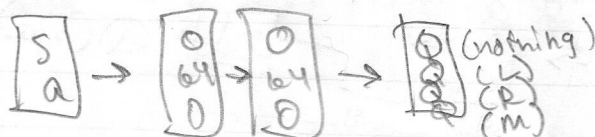
$$\vec{X} = (s, a) \quad y = R(s) + \gamma \max_a Q(s', a')$$

where Q is the previously trained Q

d. set $Q = Q_{\text{new}}$.

(Deep Q Network)

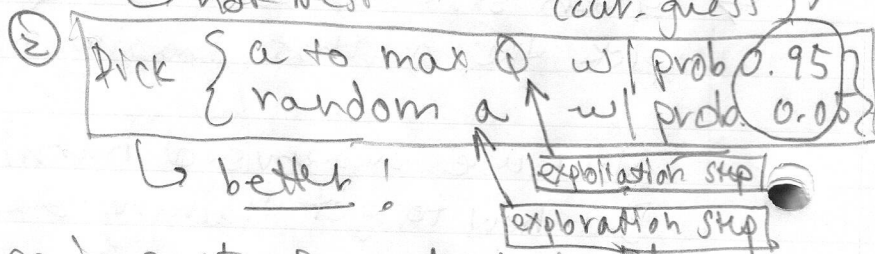
Improved architecture: output $Q(s, a) \forall a \in \mathcal{A}$,
 i.e. Q val for all actions simultaneously.



Further refinement: ϵ -greedy policy.

Recall step 2a: "take actions."

But what? ① Pick a to maximize $Q(s, a)$
 \hookrightarrow not best (cur. guess)



Why? Reduce impact of pathological random init. of Q !

Further: start with high ϵ and lower over time as Q improves.

Note: Deep RL can be more sensitive to hyperparameters.

Current State of RL

- Easier to get it working in sim/game than reality
- Far fewer applications than supervised/unsupervised
- Lots of exciting research directions.