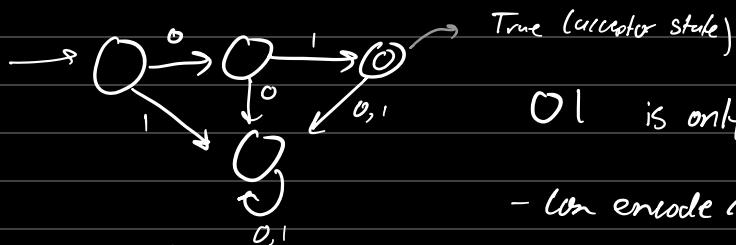


Reinforcement Learning

Finite State Machine / Deterministic Finite Automata

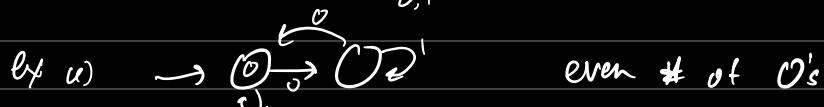
- Acceptors / recognizer \rightarrow return T/F based on final state in graph



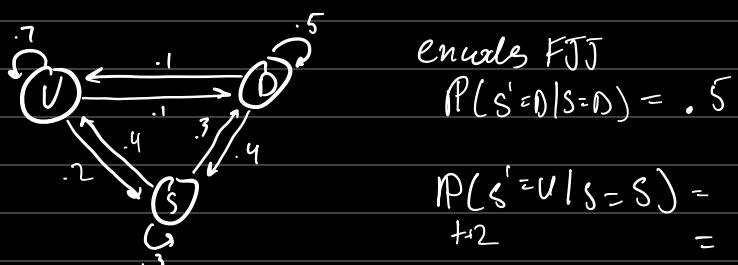
True (acceptor state)

01 is only accepted string

- can encode regex like $(01)^*$ or $(01)^*$



Markov Chain



$$P(S'=V|S=S) = \frac{P(VV) + P(SV) + P(DV)}{3} = .4(.7) + .3(.4) + .3(.1)$$

Adding inputs and Actions \rightarrow MDP need tuples S, a, S', R

$$T(s, a, s') = P(s'|s, a)$$

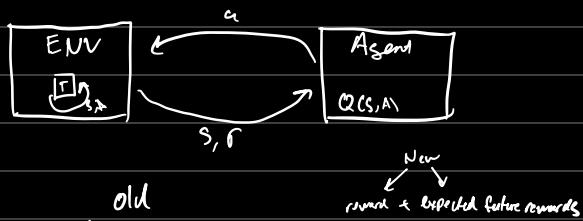
$$R(s), R(s, A), R(s, a, s')$$

$V(s)$ \approx goodness of state

use reinforcement learning to learn this

$Q(s, a)$ = total expected rewards for a discarde

$$\pi(s) = \arg \max_a Q(s, a)$$



Exploration exploitation \rightarrow epsilon greedy / decay rate

for project will train and test some except test does not learn and take no random actions

Alternative to epsilon greedy

- epsilon greedy ignores actions that are close to the best Q^* : $\{ \begin{array}{l} .49 \\ .48 \\ .01 \\ .02 \end{array} \}$
- ↳ Boltzmann (softmax): $P(a|s) = \frac{e^{Q_{s,a}}}{\sum_i e^{Q_{s,i}}} \text{ "normalized exponential"}$

Problems w/ QL

1. discrete state-action space

↳ never enough training data for large grid

→ can artificially generate more samples for training
make the table continuous: Deep Q
 $Q[s, a] \rightarrow Q(s, a)$ neural network

2. Why did I get this reward?

↳ sparse reward spread out over lots of actions

↳ chess ex: 100 moves → win or loss what caused it?

↳ credit assignment problem

3. Experimentation is SLOW!

↳ ex stock/cars only get data from trying → slow

4. Experimentation is Costly

↳ risk real \$ or damage

↳ train in simulation first

ϵ_θ

$$\alpha = .3$$

$$\gamma = .9$$

Sleep slope to cliff when
 $s_1, s_2, s_3, s_4 \rightarrow$ cliff

s_0	↑	s_1	s_2	s_3	s_4	↓
		.5	.3	.2	.4	.1

Go Right $\rightarrow \langle s_0, R, s_1 \rangle$

$$Q[s_0, R] = .7(.5) + .3(.09(.05)) = .431$$

Try Left $\rightarrow Q[s_0, L] = .318$

Try Left $\rightarrow Q[s_1, L] = .361$

Try Right $\rightarrow Q[s_1, R] = -29.79$

Going to take lots of iterations to figure out that all of actions beginning with right are bad
needs to hit cliff repeatedly and eventually error is propagated back to root cause

How to solve? \rightarrow hallucinate experiences \rightarrow "Dyna-Q"

Model + model free approach
calculate T & R

Slower than Q, but reduces real world interaction

Idea:

Loop forever:

1. Q-learning iteration as usual \rightarrow 1 real experience

2. estimate T & R from experience

$$T[S, A, S'] = \frac{T_c[S, a, s']}{\sum_i T_c[S, a, i]} \rightarrow \text{matrix of counts of transition has occurred}$$

\rightarrow normalize for all adjacent from current

Learning?

consciousness
representation
of state?

$$R[S, A] = (1-\alpha) R[S, A] + \alpha \overset{\text{this experience}}{r}$$

hyper parameter

3. Hallucinate n times: Sample from model

• \hat{s} is random state we've seen before \rightarrow weight by # times seen before?

• \hat{a} is random action from \hat{s} seen before

• $\hat{s}', r = \text{model}(\hat{s}, \hat{a})$

• Q update

Alternative: Store all experience in a dict and uniformly sample from it

Q-Learning (ipad was dead last class)

- Main take aways, can make epsilon depend on # of times the state has been seen

Assumptions

$$r_t \sim R : S \times A \times S' \rightarrow \mathbb{R} \quad (\text{fixed function})$$

$$\mathbb{E}[r^t | s, a, s'] = R_{sa}$$

$$S_{t+1} \sim p : S \times A \times S \rightarrow [0, 1] \quad (\text{fixed}), \quad \sum_s p_{sa} = 1$$

↳ Probabilities

discount learning

$\gamma \in [0, 1]$: domain dependent must be < 1 or else infinite loop to get reward

avg $(0, 1)$: average over $r_t \sim R$, need to have some sort of memory

Dependency problem

Core problem: none

↳ handles for defaults

Handle / ↳ handles for models

with Heedless in core.

model. client - self

$$Q = (1-\alpha) Q(s, a) + \alpha (r + \gamma \max Q(s', a'))$$

↳ overestimation bias "overly optimistic"

↳ looks for problems at beginning of training

- how to break this cycle? 2 learners

Double Q, 2 learners, Q^A, Q^B

· outlier experience only goes to one learner

↳ if Q^A overestimates, Q^B underestimates

- Algorithm: Store Q^A, Q^B

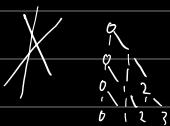
1. Select random baby let's say you set A

2. Use to get action / update $Q^A = (1-\alpha) Q^A(s, a) + \alpha (r + \gamma Q^B(s', \arg \max_a Q^A(s', a')))$

- need twice as much data → could use average of Q tables

- Roulette example, normal G learns \$22.63/226 after 100,000 trials
double Q learns 20 after 200 trials

dropping [ms] interactions RV



- problems w/ Q tables

- not enough data: non stationary } hard to fill discrete Q table
- need lots of steps

Vs?

Decision Tree

many
actions
more

↓
Learn type?

- Work arounds

A. limit actions

B. State neighborhoods

C. Dyna Q: data augmentation

↳ ex: for mouse, flip, tilt, etc. get more data

↳ could add noise but maintain trend → for legs?

D. ABIDES: Simulate w/ computers / agents to get more data

↳ Agent Based Interactive Discrete event simulation

E. train on multiple stockers to learn one

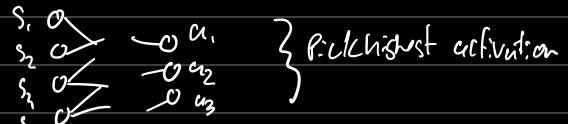
F. Make less states (bin)

- Neural Network → continuous state space $\mathbb{R} \rightarrow \mathbb{R}$

or

- Deep Q-Learning (DQ network)

- Q-table → Q-network



- Atari Example

- S: image of screen raw pixels

- partially observable

- perceptually Aliased - moving mouse - don't know direction different thing map to same mouse



- Actions: 8 directions, 1 button $\rightarrow 18$ (2 nothings)

- Reward: score don't know function

- Transition: don't know $\sum_{t=+}^T \gamma^{t-1} r_t$

- loss: error between estimated Q and true Q: GD

↳ $Q(s, a) \text{ vs. } r + \gamma \max_a (Q(s', a))$ part of update

↳ gradient descent has a built in "learning rate"

- Use convolutional neural network for image

→ is Block board on image that is selected?

- Problems

$$[r + \gamma \max_a (Q(s', a)) - Q(s, a)]^2$$

1. not sample efficient

2. consecutive samples is not very efficient - correlated

• not much time for action to have consequences, screen doesn't change much

3. Online learning drives NN to suboptimal results

• feedback loop from initial action - determines future states

Dyna

helps!

"experience
replay"

dream? mix

experiences - common and rare and simulate using mental model to get new experience. Model transition matrix as neural network?
reward/happiness as part of state?
instead of separate?

• Experience Replay

- store all exp in replay memory
- after each real exp.
 - Q-updates on mini-batches of experience randomly drawn from memory
 - ϵ -greedy for "u" as normal
 - ↳ run into memory issues with so many large states
 - can use circular buffer (overwrite oldest things)
 - ↳ Sampling uniformly might not be ideal
 - weight recent stuff more heavily or more impactful, or rare/common
 - ↳ downsample/reduce state data? \downarrow and out
 - ↳ Com stack multiple frames into one state
 - ↳ hard to train long term storage
 - ↳ can do Boltzmann exploration w/ softmax at final layer $\rightarrow P$ and sample

- Problem: Q learning can diverge \rightarrow extremely different values or train
- Gradient descent also can hit local optima
- ab state space invalidates the optimality proofs/guarantees
- also some problem w/ double Q \Rightarrow overestimation bias

↳ double deep Q learning

- training 2's really expensive

1. Q-network: selects action and gets update

eval only

2. target network: periodic slow copy of Q \rightarrow no training or action selection

↳ doesn't contain recent experiences - batch

$$[r + \max_a Q(s, a) - Q(s, a)]^2$$

↳ not fully decoupled action and value as from same network

$$[r + \gamma Q^*(s, \text{argmax}_a Q(s, a)) - Q(s, a)]^2$$

↳ also ask Q network what it would do next \Rightarrow makes all the actions

- Coach learner?

learns how to alter rewards from state to optimize learning

- Advantage learning

- actor-critic algorithms: 2 distinct learning entities
- Critic evaluation: advantage criterion - better or worse than expected
- A2C: Advantage Actor Critic
 - computes Q values (action probabilities) softmax
 - Critic network has no actions, one output neuron R - State Utility
 - ↳ does Value iteration / back prop.
 - $r + \gamma V(s') - V(s)$

- A3C: Asynchronous Advantage Actor Critic

