

Lecture 24: RL

- HW #5 released
- Due 4/21

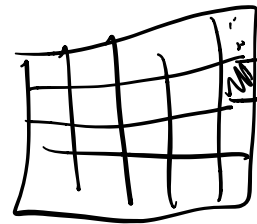
• Last time

- Dynamic programming
- Markov decision processes
- Value iteration

• Today

- Reinforced case {
- Value iteration, Q-iteration
 - Learn from data
 - Function approximation
- TD learning {
- Policy gradient

Jacob's made-up term



Value iteration

MDP w/ transition

probability
reward
discount

" $T(s', a, s)$ "

$P(s' | a, s)$

$R(s', a, s)$

γ

" $V^*(s)$ "

Value function $V(s)$ for optimal policy

$$V(s) = \max_a \sum_{s'} P(s'|a,s) (R(s',a,s) + \gamma V(s'))$$

problem: V defined in terms of itself

Solution: Add time component $V(s,t)$
 $V_t(s)$

go up instead of going down

$$\begin{cases} V_0(s) = 0 \\ V_{t+1}(s) = \max_a \sum_{s'} P(s'|a,s) (R(s',a,s) + \gamma V_t(s')) \end{cases}$$

Variant: Q-iteration

$V(s)$: optimal exp. reward from s

$Q(s,a)$: optimal exp. reward from s, a

$$V(s) = \max_a Q(s,a)$$

$$Q(s,a) = \sum_{s'} P(s'|s,a) (R(s',a,s) + \gamma V(s'))$$

$$= \sum_{s'} P(s'|s,a) (R(s',a,s) + \gamma \max_{a'} Q(s',a'))$$

Same algo as before ↓

- $Q_0(s, a) = 0$
- $Q_{t+1}(s, a) = \sum_{s'} \dots \left(\dots \max_{a'} Q_t(s', a') \right)$

What if $P(s'|a, s)$ unknown?

come from world

→ can still estimate $Q(s, a)$ from data

→ called Q-learning

$s_0, a_0, s_1, a_1, s_2, \dots$

trajectories

(come from policy π)

$\pi: S \rightarrow A$

Trajectory: $s_0, a_0, s_1, a_1, s_2, a_2, \dots$

After each state-action (s_t, a_t) , update:

$$Q(s_t, a_t) = (1-\alpha) Q_{\text{old}}(s_t, a_t) + \alpha \cdot (r_{t+1} + \gamma \cdot \max_{a'} Q_{\text{old}}(s_{t+1}, a'))$$

new value

old value

$R(s_{t+1}, a_t, s_t)$

expected value of Right-hand term

$$E_{s_{t+1}} \left[R(s_{t+1}, a_t, s_t) + \gamma \max_{a'} Q(s_{t+1}, a') \right]$$

$$= \sum_{s'} P(s' | a_t, s_t) \left(R(s', a_t, s_t) + \gamma \max_{a'} Q(s', a') \right)$$

exactly what we had
for Q-iteration

Convergence theorems

- require $\alpha \rightarrow 0$ | in practice, "step size" | fix some small step size
- other caveat: need to explore
 - analogy to multi-armed bandits: need to visit all states sufficiently often

- exploration policy
- induced policy from $Q(s, a)$ ←

↳ $Q_0(s, a) = 0$ rewards: all ≥ 0

each update: $Q(s_t, a_t) > 0$

or

Exploration

- some fraction of time, take random actions
- initialize $Q(s, a) = \infty$ some large value (optimistic value)
- ↳ similar idea to UCB

Large state spaces.

- Grid-world: ~ 20 states
- Starcraft, Dota
 - ↳ 200 units $\approx 10^5$ positions
 - $(10^5)^{200} = 10^{10000}$
- Each Q-learning update:
 - Only updates $Q(s, a)$ for single (s, a)
- Idea: Function approximation

$Q(s, a) \leftarrow$ parameterization of possible Q functions

QD update:

$$Q(s_t, a_t) = (1 - \alpha) Q_{old}(s_t, a_t) + \alpha (r_t + \max_{a'} Q_{old}(s_{t+1}, a'))$$

$$l = Q_{\text{old}}(s_t, a_t) + \alpha \left(r_t + \underbrace{\max_{a'} Q_{\text{old}}(s_{t+1}, a')}_{\text{try get new value}} - \underbrace{Q_{\text{old}}(s_t, a_t)}_{\text{old value}} \right)$$

New update:

$$\Theta = \Theta_{\text{old}} + \alpha \left(r_t + \max_{a'} Q_{\Theta_{\text{old}}}(s_{t+1}, a') - Q_{\Theta_{\text{old}}}(s_t, a_t) \right) \cdot \nabla_{\Theta} Q_{\Theta_{\text{old}}}(s_t, a_t)$$

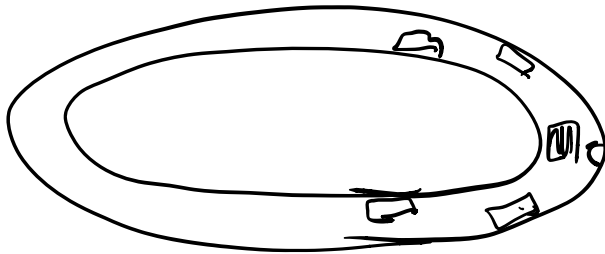
Special case: one-hot encoding

$\Theta_{s,a}$ for each state-action
 $Q(s,a) = \Theta_{s,a}$

1 in s_t, a_t
 0 everywhere else

Environment

↳ racecar



- our car $\rightarrow x_0$
- other cars $\rightarrow x_1, \dots, x_K$
- gas remaining $\rightarrow g$
- time left $\rightarrow t$

$$s = (x_0, x_1, \dots, x_K, g, t)$$

features we should care about:

- velocity, heading
 - distance to closest car
 - angle closest car
 - gas
 - brake
 - ...
- features

$$Q_{\theta}(s, a) = \theta_1 (\text{velocity}) + \theta_2 (\text{distance}) + \theta_3 (\text{angle}) + \dots$$

New update:

$$\theta = \theta_{\text{old}} + \alpha \left(r_t + \max_{a'} Q_{\theta_{\text{old}}}(s_{t+1}, a') - Q_{\theta_{\text{old}}}(s_t, a_t) \right) \cdot \nabla_{\theta} Q_{\theta_{\text{old}}}(s_t, a_t)$$

"Intuitive derivation"

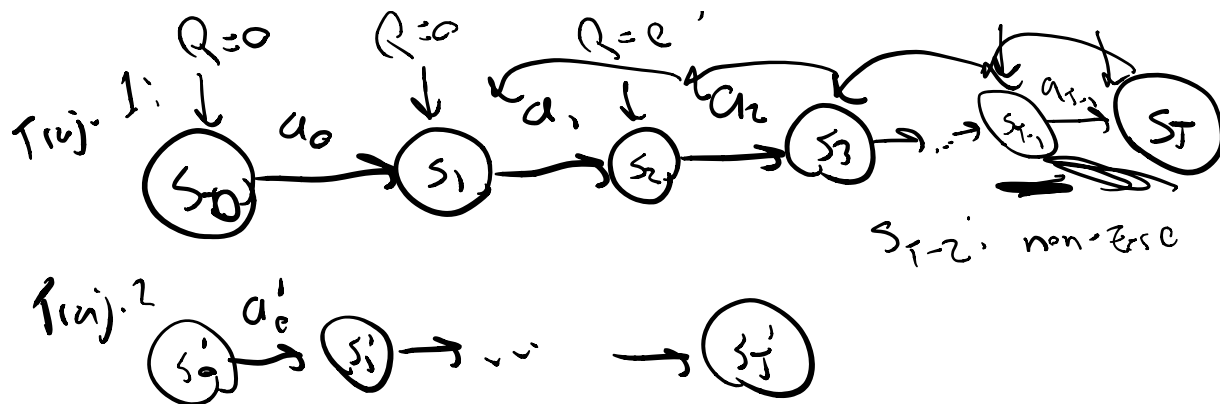
Prediction task: ~~max~~ (target - value)²

$$\nabla_{\theta} \left(\left(r_t + \max_{a'} Q_{\theta_{\text{old}}}(s_{t+1}, a') - Q_{\theta_{\text{old}}}(s_t, a_t) \right)^2 \right)$$

fixed
change-able

Temporal difference learning

TD(0) $Q(s, a) \Rightarrow \theta$ Q non-zero



Often only see reward at end
 "credit assignment" problem

regular Q-learning:
 need T updates for trajectory of
 length T

$TD(\lambda)$: propagate rewards backward
 \hookrightarrow TD-Gammon

Policy gradient

- Q-learning $Q_\theta(s, a)$
- directly learn policy $\pi_\theta(a|s)$

probability distⁿ
 of action | state

$$\max_{\theta} \mathbb{E}_{\pi_\theta} \left[R(s_0, a_0, s_1) + \gamma R(s_1, a_1, s_2) + \dots \right]$$

what gradient

log-derivative trick

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}} [R]$$

$$= \mathbb{E}_{\pi_{\theta}} \left[R \cdot \frac{\nabla_{\theta} \log \pi_{\theta}}{\pi_{\theta}} \right]$$

cancel $\frac{1}{\pi_{\theta}}$ $\nabla \pi_{\theta}$

Logistic regression

$(x_1, y_1), \dots, (x_n, y_n)$

Stochastic GD:

(x_i, y_i)

$$\Theta = \Theta_{old} + \alpha \nabla_{\Theta}$$

↑
model parameters

label ↓ features ↓

$$\frac{\log(1 + \exp(-y_i \Theta^T \phi(x_i)))}{\text{logistic loss } \Theta_{old}}$$

$Q(s, a)$: expected reward for taking action a in state s , then following optimal policy

approximate this by current guess based on $\max_{a'} Q(s', a')$ in new state

Small State space case EN or S'

$$Q(s, a) \approx \sum_{s'} P(s'|s, a) (R(s', a, s) + \gamma \max_{a'} Q(s', a'))$$

samples
immediate reward
looking ahead

100 times where

$$S_t = s, a_t = a$$

$$s_{t+1}^{(1)}, \dots, s_{t+1}^{(100)}$$

$$\frac{1}{100} \sum_{i=1}^{100} (R(s_{t+1}^{(i)}, a_t, s_t) + \gamma \dots)$$

hoping that $Q(s', a')$

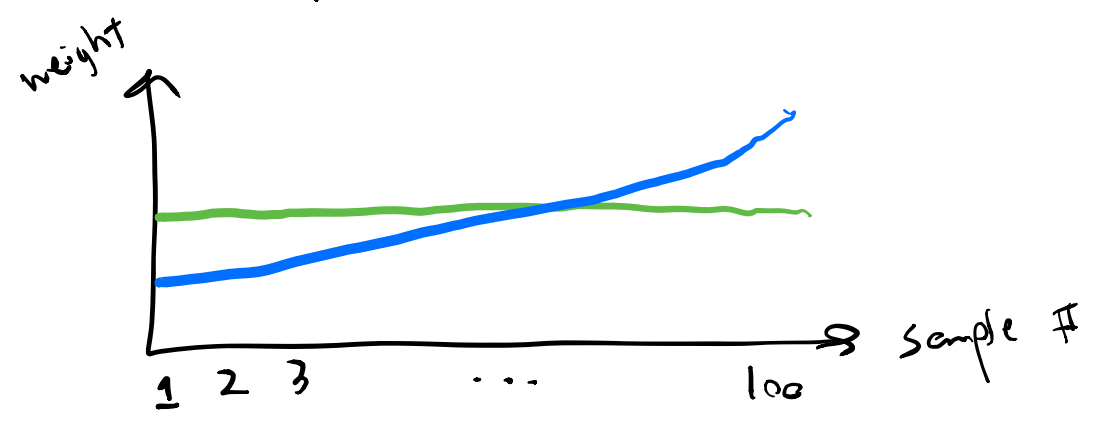
$$\sum_{s'} P(s'|s, a) (R(s', a, s) + \gamma \max_{a'} \sum_{s''} P(s''|s', a') \cdot (R(s'', a', s') + \dots))$$

$$\alpha = 0.01$$

100 samples

... ale: weight $\alpha = 0.01$

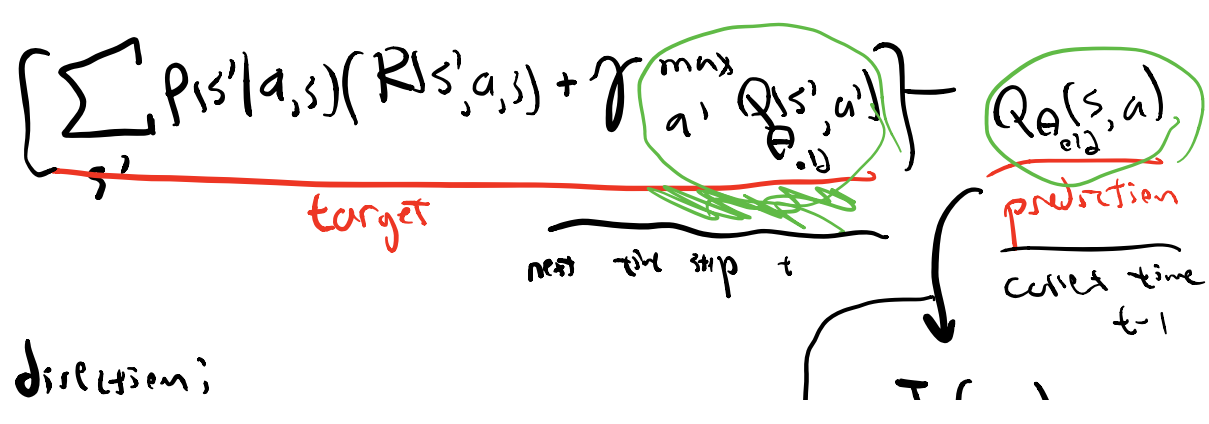
Most recent sample: α
 2nd most recent sample: weight $\alpha \cdot (1-\alpha) = 0.01 \cdot 0.99$
 3rd most recent: $\alpha \cdot (1-\alpha)^2$
 1st sample: $\alpha \cdot (1-\alpha)^{99}$

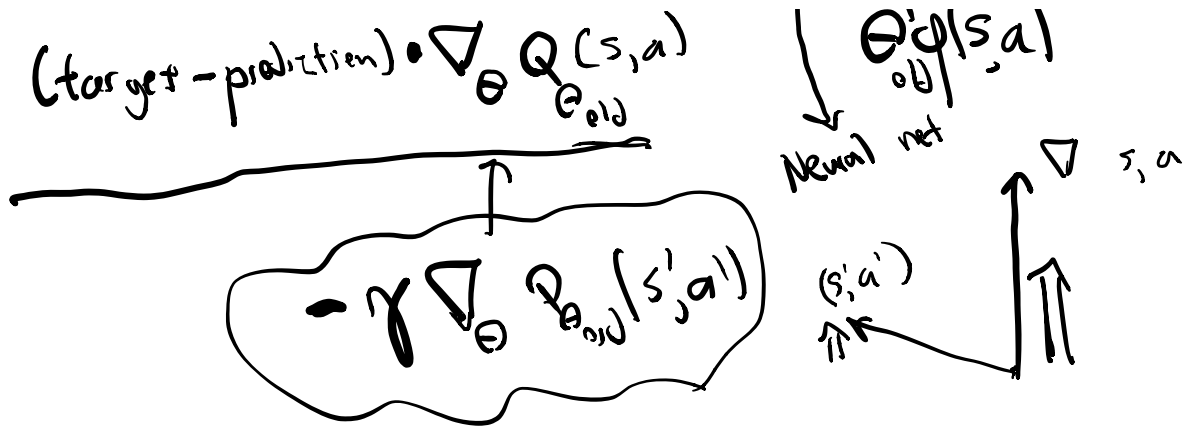


$$Q(s, a) \approx \sum_{s'} P(s' | a, s) \cdot (R(s', a, s) + \gamma \max_{a'} Q(s', a'))$$

$$Q_{\theta}(s, a) \approx \sum_{s'} P(s' | a, s) \cdot (R(s', a, s) + \gamma \max_{a'} Q_{\theta}(s', a'))$$

Find θ such that \approx true





$$\theta = \theta_{old} + \alpha \cdot \frac{(\text{target} - \text{prediction}) \cdot \nabla_{\theta} Q_{\theta_{old}}(s, a)}{Q_{\theta_{old}}(s, a)}$$

$$\frac{\sum_{s'} P(s'|a, s) \cdot (\dots)}{\text{Hpare w/ sample}}$$

$$\underline{R(s_{t+1}, a_t, s_t)} \quad \max_{a'} Q_{\theta_{old}}(s_{t+1}, a')$$

s_{t+1} is a sample from $P(s'|a_t, s_t)$