

Plan

Recap

Logistics

Q Learning

Deep Q Learning

Comparison to Policy
Gradients

Recap



State

Action

Reward

$$\min_{\theta} - \mathbb{E}[R(\tau)]$$

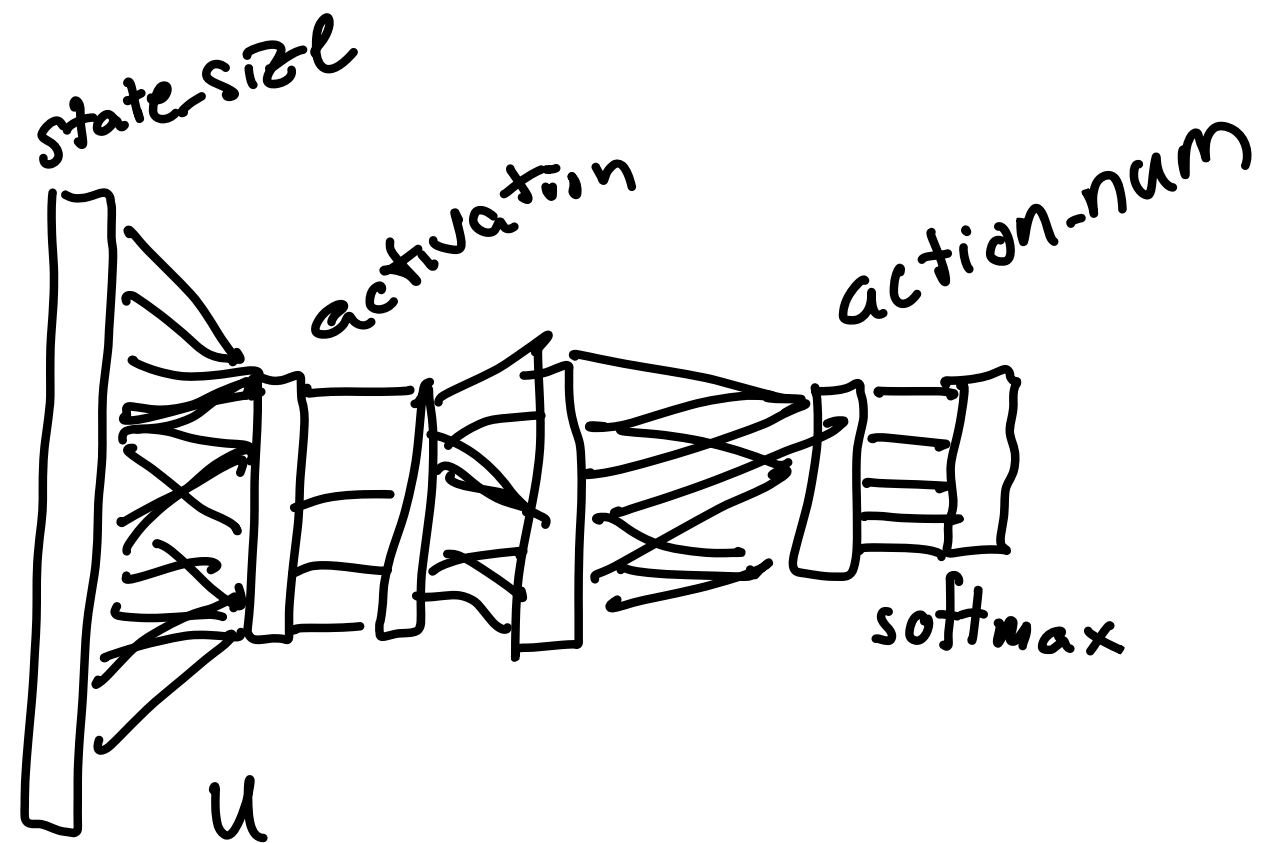
$\tau \sim \pi_{\theta}$

$\pi: \mathbb{R}^{\text{state-size}} \rightarrow \mathbb{R}^{\text{action-num}}$

REINFORCE

1. Sample trajectory τ according to π
2. Calculate $R(\tau)$
3. $\theta \leftarrow \theta + \eta R(\tau) \cdot \frac{\partial \log \pi(\tau)}{\partial \theta}$

$$\theta \leftarrow \theta + \eta \underbrace{\frac{\partial \mathbb{E}[R(\tau)]}{\partial \theta}}_{\mathbb{E}\left[R(\tau) \cdot \frac{\partial \log \pi(\tau)}{\partial \theta}\right]}$$



Logistics

Grades

↳ By point

↳ Participate and turn
everything → B or higher

↳ Improvement 😊

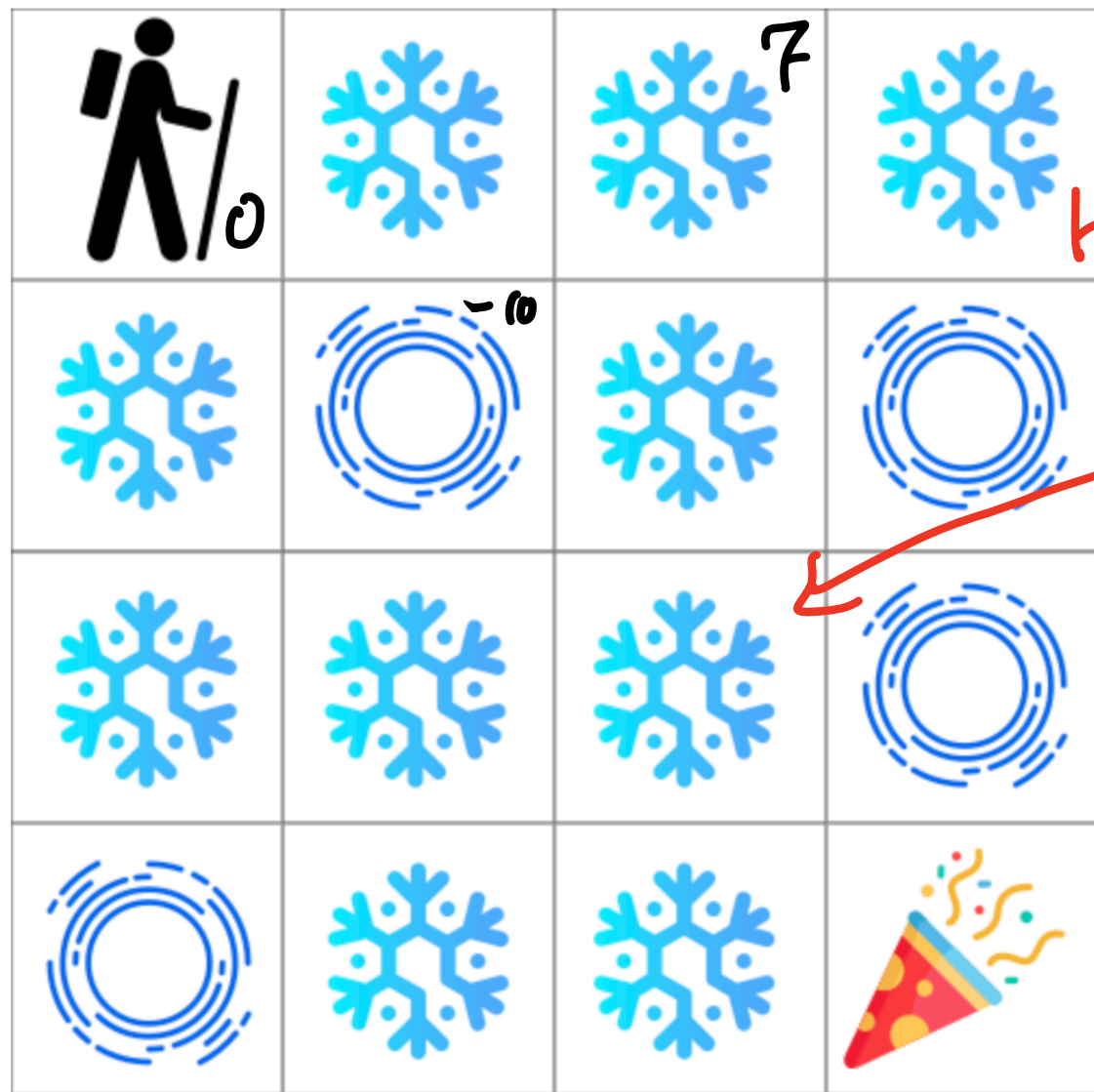
Form 23 / 25

Proposal

Project

Q Learning

Goal: Go to good states



How good is this state?

What makes a state good?

$$\tau = (\Delta_0, a_0, \Delta_1, a_1, \Delta_2, a_2, \dots)$$

"gain" $r_t = r(\Delta_t, a_t)$

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

$$(\Delta_t, a_t, \Delta_{t+1}, a_{t+1}, \dots)$$

$$0 < \gamma < 1 \quad \gamma \approx .99$$

value of state Δ is

$$E[G_t \text{ when } \Delta_t = \Delta]$$

$$\gamma \sim \pi$$

↑ expected when discounted reward we start in Δ and follow policy π

We need to know how
to get to good states.

Need action-value function

$Q^\pi(s, a)$ = the expected
long term reward
of taking action
 a in state s
then following
policy π

$$Q^\pi(s, a) = \mathbb{E}_\pi [G_t | \Delta_t = s, a_t = a]$$

$$\text{best action} = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$$

Bellman Equation

$$Q^\pi(s, a) = \mathbb{E}[G_t | \Delta_t = s, a_t = a]$$

$$Q(\Delta_t, a_t) = \mathbb{E}\left[\sum_{i=0}^{\infty} \gamma^i r_{t+i} \mid \Delta_t = s, a_t = a\right]$$

$$= \mathbb{E}\left[r_t + \sum_{i=1}^{\infty} \gamma^i r_{t+i} \mid \Delta_t = s, a_t = a\right]$$

$$= r_t + \gamma \mathbb{E}\left[\sum_{i=1}^{\infty} r_{t+i} \gamma^{i-1} \mid \Delta_{t+1} = s', a_{t+1} \text{ is best action in state } s'\right]$$

$\underbrace{\sum_{i=0}^{\infty} r_{t+i} \gamma^i}_{\leftarrow} \leftarrow r_{t+1} \cdot 1 + r_{t+2} \cdot \gamma + r_{t+3} \gamma^2 + \dots$

deterministic
 $r_t = r(\Delta_t, a_t)$

s' is state
if we take a
in state s

$$= r_t + \gamma \mathbb{E}[G_{t+1} | \Delta_{t+1} = s', a_{t+1} \text{ is best}]$$

$$= r_t + \gamma \max_{a'} \mathbb{E}[G_{t+1} | \Delta_{t+1} = s', a_{t+1} = a'] = r_t + \gamma \max_{a'} Q(s', a')$$

Bellman Equation (s, a, r, s')

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) \approx \overset{\text{want}}{r} + \gamma \max_{a'} Q(s', a')$$

$$\mathcal{L} = \frac{1}{2} \left(Q(s, a) - \underbrace{\left(r + \gamma \max_{a'} Q(s', a') \right)}_{\text{fixed } Q} \right)^2$$

$$\frac{\partial \mathcal{L}}{\partial Q(s, a)} = \left(Q(s, a) - \overset{\text{target}}{r + \gamma \max_{a'} Q(s', a')} \right) \frac{\partial Q(s, a)}{\partial Q(s, a)}$$

Exploration vs Exploitation

Initialize Q (s, a, s', r)

Repeat until done:

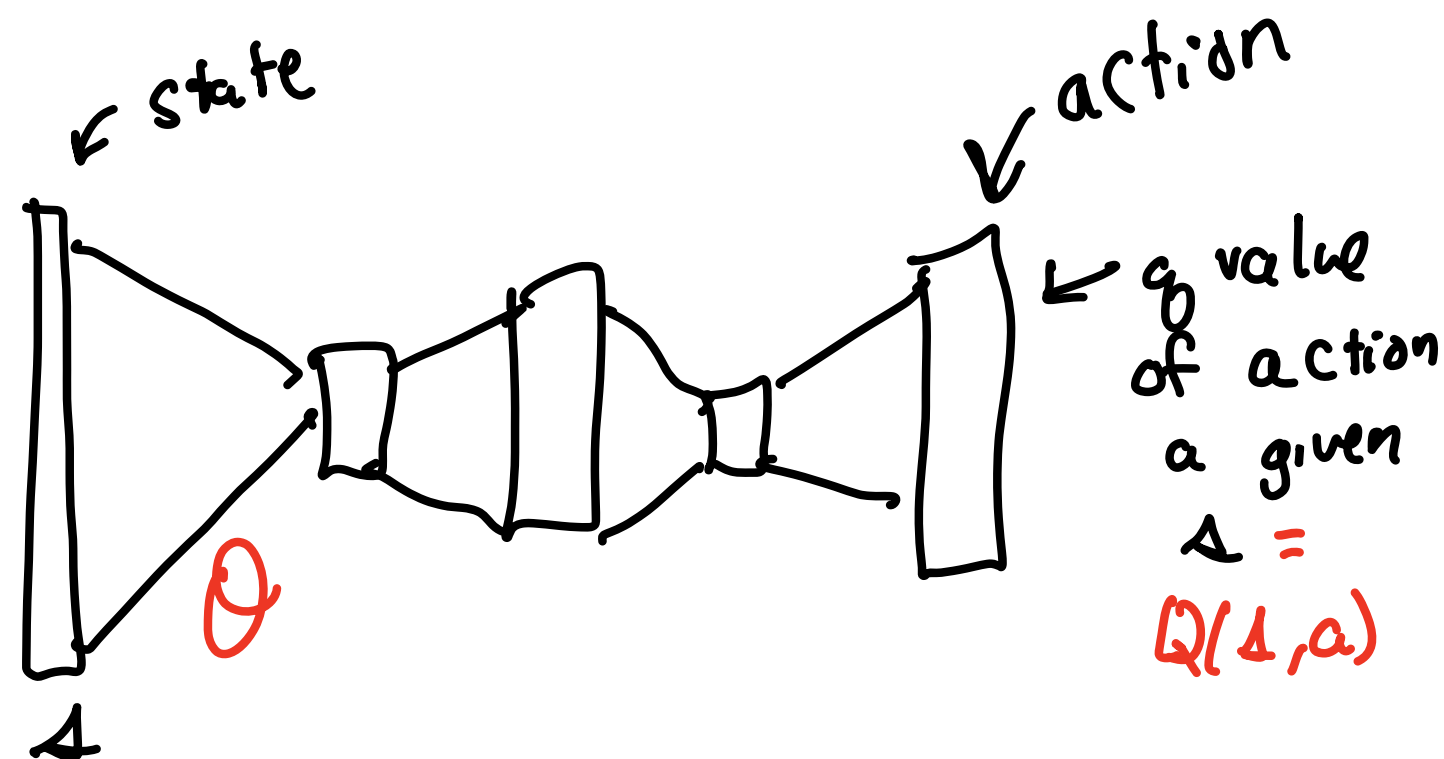
1. Choose $\underset{a}{\operatorname{argmax}} Q(s, a)$ with prob. $1-\epsilon$, random else
2. Take a and observe s', r
3. $Q(s, a) \leftarrow \eta (Q(s, a) - \text{target})$

$$\theta \leftarrow \eta \frac{\partial \mathcal{L}}{\partial \theta} = \eta (Q(s, a) - \text{target}) \frac{\partial Q(s, a)}{\partial \theta}$$

$$\text{target} = r + \gamma \max_{a'} Q(s', a')$$

$$Q = \begin{matrix} & \text{\# actions} \\ \begin{matrix} \text{\# states} \\ 3^{361} \end{matrix} & \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \\ & 361 \end{matrix}$$

$Q =$ neural net $\ddot{\smile}$



Policy

- no intermediate value
- Rock paper scissors need distribution

Deep Q Learning

- Estimate value (helpful in general)
- More sample efficient