4/16/2024: Deep Q-Learning, Policy gradient

## 3rd Q-Learning variant = Deep Q-Learning

Idea: $\hat{Q}(s,a)$ = output of a neural network whose input is $s$ & $a$

Let $\Theta$ be parameters of model.

$$\text{Loss}(\Theta) = \frac{1}{2}\left( \underbrace{\hat{Q}_\Theta(s,a)}_{\substack{\text{prediction} \\ \text{from} \\ \text{NN with} \\ \text{parameters } \Theta}} - \underbrace{\left[ r + \gamma \hat{V}(s') \right]}_{\text{"target"}} \right)^{\boxed{2}}$$

↑ minimize squared difference

$$\nabla_\Theta \text{loss}(\Theta) = \frac{1}{2} \cdot 2 \cdot \underbrace{\left( \hat{Q}_\Theta(s,a) - [r + \gamma \hat{V}(s')] \right)}_{\text{compute directly}} \cdot \underbrace{\nabla_\Theta \hat{Q}_\Theta(s,a)}_{\substack{\text{computed by} \\ \text{backprop}}}$$
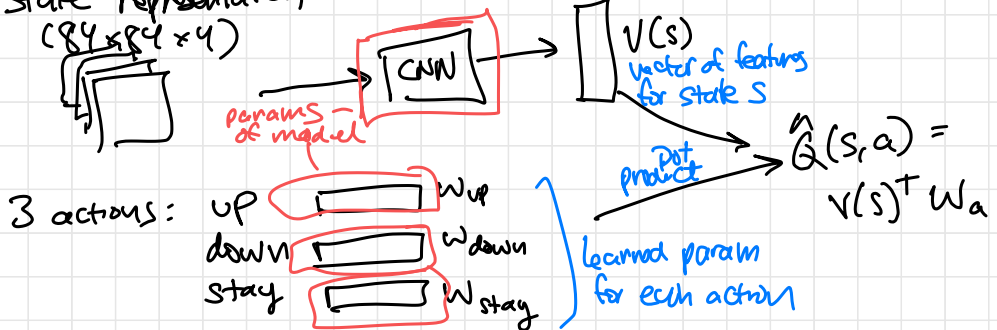
Update $\Theta$ with gradient descent when we receive new rewards $r$ & new state $s'$

## Example DQN for Pong
Represent state as last $k$ frames
- Each frame is $84 \times 84$
- Set $k = 4$
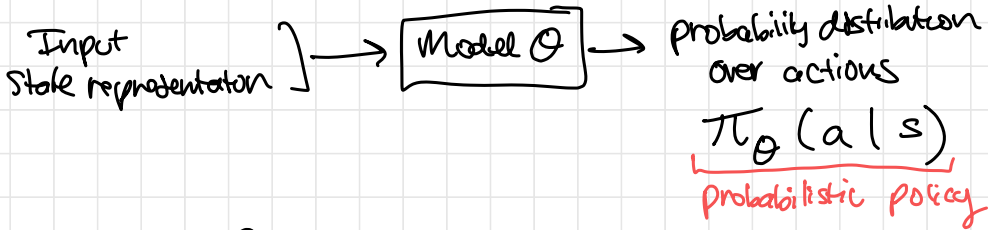- Input is $84 \times 84 \times 4$ block of numbers

state representation for $s$
$(84 \times 84 \times 4)$



params of model

CNN → $V(s)$ vector of features for state $s$

3 actions: up $\quad$ $W_{up}$
down $\quad$ $W_{down}$
stay $\quad$ $W_{stay}$

learned param for each action

dot product

$\hat{Q}(s,a) = V(s)^T W_a$

# Policy Gradient methods

- Q-Learning: Given $(s, a)$, predict $Q(s,a) \in \mathbb{R}$
  "Regression"

- Policy Gradients: Given $s$, predict optimal $a$
  "classification"

Input
State representation $\Big] \longrightarrow$ Model $\theta$ $\longrightarrow$ probability distribution over actions

$$\underbrace{\pi_\theta (a \mid s)}_{\text{Probabilistic policy}}$$

## How to train?

- Cant do supervised learning b/c no supervised training data
- Solution: Train $\pi_\theta (a\mid s)$ to achieve large expected total reward

Want to maximize:

$$\underbrace{V(\theta)}_{\substack{\text{Expected sum of rewards} \\ \text{when using} \\ \text{policy } \pi_\theta(a\mid s)}} = \sum_{\substack{\text{trajectories} \\ Z \\ \uparrow \\ \text{all possible sequences} \\ [s_1, a_1, r_1, s_2, a_2, r_2, s_3 \cdots]}} \underbrace{P(z; \theta)}_{\substack{\text{prob. that} \\ z \text{ happens}}} \cdot \underbrace{R(z)}_{\substack{\text{total reward} \\ \text{of } z \\ = \sum_{t=1}^{T} r_t}}$$

Plan: maximize $V(\theta)$, use gradient ascent.
Need to compute $\nabla_\theta V(\theta)$

Useful trick:

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

$$\iff \nabla_\theta P(z;\theta) = P(z;\theta) \cdot \nabla_\theta \log P(z;\theta)$$

$$\nabla_\Theta V(\Theta) = \sum_{\substack{\text{traj's} \\ z}} P(z;\theta) \cdot \nabla \log P(z;\theta) \cdot R(z)$$

this

expected value ... of

$$= \mathbb{E}_{z \sim P_\theta} \left[ \nabla \log P(z;\theta) \cdot R(z) \right]$$

Approximate this with sampled trajectories & computing the mean

computing this

$$\log P(z;\Theta) = \boxed{\log P(s_1)} + \boxed{\log \pi(a_1|s_1)} + \boxed{\log T(s_1, a_1, s_2)}$$
$$+ \boxed{\log \pi(a_2|s_2)} + \boxed{\log T(s_2, a_2, s_3)}$$
$$+ \dots$$

✗ unknown to us
∴ does not depend on $\Theta$
so $\nabla_\Theta$ is $\bigcirc$

compute $\nabla_\Theta$ with backprop