Previously: Q-Learning : Given $(s, a)$, predict $Q_{opt}(s, a)$
"regression"

Today: Policy Gradient : Given $s$, predict $a$
"classification"

Input:
Representation of states $\Big\}$ $\longrightarrow$ $\boxed{\begin{array}{c}\text{Model } \theta \\ \text{(e.g. neural} \\ \text{network)}\end{array}}$ $\longrightarrow$ probability distribution over actions

$$\pi_\theta (a \mid s)$$

probabilistic policy defined by model parameters $\theta$

How to train?
- Normal supervised learning requires knowing best action at given states as training data ✗
  Not known for any state
- Policy gradient : train $\pi_\theta (a \mid s)$ to achieve high total rewards

Want to <u>maximize</u> value of the policy $\pi_\theta$

$$V(\theta) = \sum_{\text{trajectories } z} P(z ; \theta) \cdot R(z)$$

Expected sum of rewards when using policy $\pi_\theta (a \mid s)$

Prob of $z$ happening

Trajectory = possible sequences of $[s_1, a_1, r_1, s_2, a_2, r_2, s_3, \ldots]$

Reward achieved
$$= \sum_{t=1}^{T} r_t$$

Plan: $V(\theta)$ is our training objective, maximize with gradient ascent

What is $\nabla_\theta V(\theta)$ ?

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

Sum over exponentially many trajectories — infeasible

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

$\rightleftharpoons \quad \nabla_\theta P(z;\theta) = P(z;\theta) \nabla_\theta \log P(z;\theta)$

Plug this in to $\nabla_\theta V(\theta)$:

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

Expected value of....    this quantity

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

Estimate this by sampling trajectories with $\pi_\theta(a|s)$ and taking average of $\nabla_\theta \log P(z;\theta) \cdot R(z)$

What is $\nabla_\theta \log P(z;\theta)$?

$$\log P(z;\theta) = \underbrace{\log P(s_1)}_{\text{start state}} + \underbrace{\log \pi_\theta(a_1|s_1)}_{\text{policy}} + \underbrace{\log T(s_1,a_1,s_2)}_{\text{transitions}}$$

$\{s_1, a_1, r_1, s_2, ...\}$    $+ \underbrace{\log \pi_\theta(a_2|s_2)}_{\text{policy}} + ...$

Don't depend on $\theta$
So $\nabla_\theta$ is $0$

So $\nabla_\theta \log P(z;\theta) = \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t)$

# Basic Policy Gradient Algorithm:

Initialize $\theta$ randomly

For each episode:

- Sample trajectory $z$ using $\pi_\theta(a|s)$
- Update:

$$\theta \leftarrow \theta + \eta \, R(z) \underbrace{\sum_{t=1}^{T} D_\theta \log \pi_\theta(a_t | s_t)}_{\approx \nabla_\theta V(\theta)}$$