# Policy Gradient [2]

## DOROZHKO Anton

Novosibirsk State University

June 6, 2020

[2]SpinningUP RL research : Part 3

# Outline

1. Recap Q-learning

2. Policy Gradient

# Q-learning

$$MDP \ (S, A, \cancel{P}, \cancel{R}, \gamma)$$

$$goal \quad \max_{\pi_\theta} \mathbb{E}[G_t]$$

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

$$\downarrow$$

$$\pi(a|s) = \underset{a}{argmax} \ \hat{Q}(s, a) \quad | \quad \varepsilon\text{-greedy}$$

# Action-Value approximation

target          NN

- Find $w$ that minimizes loss between $Q^\pi(s,a)$ and $\hat{Q}(s,a;w)$
- MSE

$$J(w) = \mathbb{E}_\pi[(Q^\pi(s,a) - \hat{Q}(s,a;w))^2]$$

- SGD samples the gradient

$$-\frac{1}{2}\nabla_w J(w) = \mathbb{E}_\pi[(Q^\pi(s,a) - \hat{Q}(s,a;w))\nabla_w \hat{Q}^\pi(s,a;w)]$$

$$\Delta w = -\frac{1}{2}\alpha \nabla_w J(w)$$

Novosibirsk
State
University
*THE REAL SCIENCE

# Incremental Model-Free Control

$r_0 + \gamma r_1 + \ldots + \gamma^i r_i$

- MC

200 apm
5 minutes
1000 steps

50s kps

$0.9^{50} \sim 5 \cdot 10^{-3}$

$\gamma = 0.9$

$\gamma^i \approx 0$

$$\Delta w = \alpha(G_t - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

- SARSA

$$\Delta w = \alpha(r + \gamma \hat{Q}(s', a'; w) - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

- Q-learning

$G_t = r_t + \gamma R_{t+1} + \ldots$

$$\Delta w = \alpha(r + \gamma \max_a \hat{Q}(s', a'; w) - \hat{Q}^\pi(s_t, a_t; w))\nabla_w \hat{Q}(s_t, a_t; w)$$

$\gamma = 0$

$\gamma \to 1$

$Q \sim r_t$

$Q \approx r_t + r_{t+1} + \ldots r$

N* Novosibirsk
State
University
*THE REAL SCIENCE

# What if ?

$$\hat{Q}(s, a; w) \longrightarrow \pi$$

What if we can directly optimize our expected reward w.r.t. the parameters of our policy ?

$$\pi_\theta(a|s) \qquad J_\theta(\pi_\theta) = \mathbb{E}[G_t]$$

$$\nabla_\theta J_\theta(\pi)$$

# Policy Optimization

- Stochastic, parametrized policy $\pi_\theta \rightsquigarrow f(\theta) : S \rightarrow \Delta A$
- maximize expected return

  $NN : x(S) \rightarrow$ logits over actions

$$J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[G_t(\tau)]$$

$\tau$ trajectory

$(S_0, a_0, s_1 a_1)$

- SGD with **policy gradient**

  $r_0, r_1 \ldots$

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_\theta)|_{\theta_k}$$
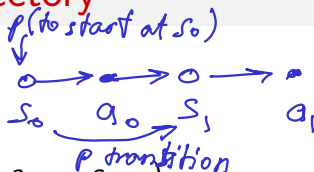
# How to compute policy gradient

- Derive the analitycal gradient
- Compute sample estimate

# 1) Probability of a Trajectory

$p($ to start at $s_0)$

$s_0 \quad a_0 \quad s_1 \quad a_1$

$p$ transition

Trajectory

$$\tau = (s_0, a_0, ..., s_{T+1})$$

Probability of a Trajectory

Markov property

$$P(\tau|\theta) = \rho_0(s_0) \prod_{t=0}^{T} P(s_{t+1}|s_t, a_t) \pi_\theta(a_t, s_t)$$

Take log

$$log P(\tau|\theta) =$$

$\log \rho_0(s_0) + \sum_t \log P(s_{t+1}|s_t, a_t) + \sum_t \log \pi_\theta$

# 2) The Log-Derivative Trick

$$(\log x)' = \frac{1}{x}$$

$$\nabla_x \log f(x) = \frac{1}{f(x)} \nabla_x f(x)$$

easy → (underbraced term) ?hard

$$\nabla_\theta log P(\tau | theta) = \frac{\nabla_\theta P(\tau|\theta)}{P(\tau|\theta)}$$

Rearrange

$$f(x) \nabla_x \log f(x)$$

$$\nabla_\theta P(\tau|\theta) = P(\tau|\theta) \nabla_\theta log P(\tau|\theta)$$

Novosibirsk State University
THE REAL SCIENCE

# 3) Gradients of Environment Functions

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \mathop{\mathbb{E}}_{\tau \sim \pi_\theta}\left[R(\tau)\right]$$

$$= \nabla_\theta \int R(\tau) \cdot P(\tau | \theta)$$

- $\rho_0(s_0)$, $P(s_{t+1}|s_t, a_t)$ and $G_t(\tau)$
- no dependence on $\theta$
- Gradients w.r.t $\theta$ are **zero**

$$\mathop{\mathbb{E}}_{x \sim p} f(x) = \int f(x) \cdot p(x)\, dx$$

N* Novosibirsk
State
University
*THE REAL SCIENCE

# Grad-Log-Prob of Trajectory

$$\nabla_\theta \log P(\tau|\theta) = \nabla_\theta \log \rho_0(s_0) + \sum_{t=0}^{T} \left( \nabla_\theta \log P(s_{t+1}|s_t, a_t) \right.$$

$$\left. + \nabla_\theta \log \pi_\theta(a_t|s_t) \right)$$

$$= \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t).$$

$$\nabla_\theta J(\pi_\theta) = \ldots$$

$$\nabla_\theta \, \mathbb{E}_{\tau \sim \pi_\theta} \left[ R(\tau) \right] = \nabla_\theta \int R(\tau) \cdot P(\tau | \theta)$$

$$= \int R(\tau) \, \nabla_\theta P(\tau | \theta)$$

$$= \int R(\tau) \, P(\tau | \theta) \, \nabla_\theta \log P(\tau | \theta)$$

$$= \mathbb{E}_{\tau \sim \pi_\theta} \left[ R(\tau) \, \nabla_\theta \log P(\tau | \theta) \right]$$

$$= \mathbb{E}_{\tau \sim \pi} \left[ R(\tau) \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta (a_t | s_t) \right]$$

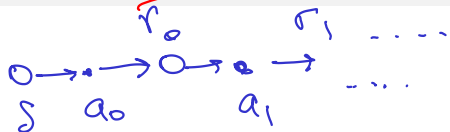# Basic Policy Gradient

**❗ Derivation for Basic Policy Gradient**

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \operatorname*{E}_{\tau \sim \pi_\theta}[R(\tau)] \qquad \text{def}$$

$$= \nabla_\theta \int_\tau P(\tau|\theta)R(\tau) \qquad \text{Expand expectation}$$

$$= \int_\tau \nabla_\theta P(\tau|\theta)R(\tau) \qquad \text{Bring gradient under integral}$$

$$= \int_\tau P(\tau|\theta)\nabla_\theta \log P(\tau|\theta)R(\tau) \qquad \text{Log-derivative trick}$$

$$= \operatorname*{E}_{\tau \sim \pi_\theta}[\nabla_\theta \log P(\tau|\theta)R(\tau)] \qquad \text{Return to expectation form}$$

$$\therefore \nabla_\theta J(\pi_\theta) = \operatorname*{E}_{\tau \sim \pi_\theta}\left[\sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t)R(\tau)\right] \qquad \text{Expression for grad-log-prob}$$

policy gradient          REINFORCE

Novosibirsk State University

*THE REAL SCIENCE

# How to compute expectation ?

$$R(\tau) = r_0 + \delta r_1$$

$$\underbrace{r_0 \qquad r_1 \quad \cdots}$$

$$O \rightarrow \cdots \rightarrow O \rightarrow \bullet \rightarrow$$
$$s \quad a_0 \qquad a_1$$

$$\mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t)R(\tau)]$$

$$\hat{g} = \frac{1}{N}\sum_{i=1}^{N}\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t)R(\tau) \qquad \nabla_\theta J$$

$$N \text{ games } \text{ with } \pi_\theta \qquad \theta_{k+1} = \theta_k - \alpha \nabla_\theta J$$

# Basic Policy Gradient in words

$$Q_w(s,a) : S \times A \to \mathbb{R}$$

$$\pi_\theta(a|s) : S \to \Delta A$$

- Define policy as parametrized function
- Write the gradient of the loss function
- Play games collect trajectories (1 or more)
- Compute cumulative discounted rewards for each $t$ $\quad R(\tau)$
- Compute gradients $\quad \nabla_\theta J(\pi_\theta)$

after each step of SGD
- Update parameters of your policy $\quad \theta_{k+1} = \theta_k - \alpha \nabla_\theta J$
- **Comment: throw out trajectories**

THE REAL SCIENCE

Novosibirsk State University

# Loss is not a normal loss

$$MSE \; loss$$
$$\frac{1}{N} \Sigma (f(x) - y)^2$$

1. Data is dependent on the parameters $\pi_\theta$
2. Is not connected to performance $J_{(\pi_\theta)}$ after one step
3. So, minimization of this "loss" itself for several steps on the same batch will not give better rewards.

N* Novosibirsk
State
University
*THE REAL SCIENCE

# Consider only future rewards



$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t)R(\tau)]$$

# Consider only future rewards

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t) R(\tau)]$$

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t) \sum_{t'=t}^{T} R(s_{t'}, a_{t'}, s_{t'+1})]$$

reward-to-go

---

[20]SpinningUP RL research : Part 3

### Lemma

$$\mathbb{E}_{x \sim P_\theta}[\nabla_\theta log P_\theta(x)] = 0$$

Proof: distributions are normalized $\rightarrow$ take gradient $\rightarrow$ log derivative trick

$$\int_x P_\theta(x) = 1$$

$$\nabla_\theta \int_x P_\theta(x) = \nabla_\theta 1 = 0$$

$$\nabla_\theta \int_x P_\theta(x) = \int_x P_\theta(x) \nabla_\theta log P_\theta(x) = \underbrace{\mathbb{E}_{x \sim P_\theta}[\nabla_\theta log P_\theta(x)]}$$
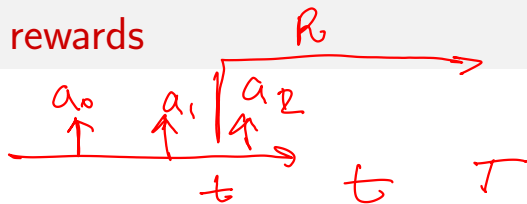
log deriv
trick

# Proof of go-to-reward usage

$$\mathbb{E}_\tau \left[ \sum \nabla_\theta \log \pi_\theta \sum_{t'=1}^{T} R \right] =$$

$$\mathbb{E}_\tau \left[ \sum_t \nabla_\theta \log \pi_\theta \left[ \sum_{t'=1}^{t-1} R + \sum_{t'=t}^{T} R \right] \right]$$

$$\nabla_\theta J_\theta(\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t) \sum_{t'=t}^{T} R(s_{t'}, a_{t'}, s_{t'+1})]$$

$$= \sum_t \mathbb{E}_\tau \nabla_\theta \log \pi_\theta \left[ \sum_{t'=1}^{t-1} R + \right.$$

$$0$$

doesn't depend on
$\tau$ after $t$

# Consider only future rewards



$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t) \sum_{t'=t}^{T} R(s_{t'}, a_{t'}, s_{t'+1})]$$

Why it's better ?

Recall the sum of random variables with 0 mean and some variance.

$$\sum_{t=0}^{T} R \sim r.v \quad Var(R)$$

# Improvement: Baseline

Consequence of lemma

For $\forall b(s)$:

const w.r.t. $a_t$

$$\mathbb{E}_{a_t \sim \pi_\theta}[\nabla_\theta log \pi_\theta(a_t|s_t)b(s_t)] = 0$$

This allows:

Var

$\frac{\partial Var}{\partial b} = 0$

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[\sum_{t=0}^{T} \nabla_\theta log \pi_\theta(a_t|s_t) \left( \sum_{t'=t}^{T} R(s_{t'}, a_{t'}, s_{t'+1}) - b(s_t) \right)]$$

$\hat{Q}$

We can choose $b$ to reduce variance. In practice $b(s_t) = V^\pi(s_t)$ value function works well.

# Basic Policy Gradient

**Algorithm 1** Vanilla Policy Gradient Algorithm

1: Input: initial policy parameters $\theta_0$, initial value function parameters $\phi_0$
2: **for** $k = 0, 1, 2, \ldots$ **do**
3:   Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
4:   Compute rewards-to-go $\hat{R}_t$.
5:   Compute advantage estimates, $\hat{A}_t$ (using any method of advantage estimation) based on the current value function $V_{\phi_k}$.
6:   Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t)|_{\theta_k} \hat{A}_t.$$

7:   Compute policy update, either using standard gradient ascent,

$$\theta_{k+1} = \theta_k + \alpha_k \hat{g}_k,$$

  or via another gradient ascent algorithm like Adam.
8:   Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg\min_\phi \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \left( V_\phi(s_t) - \hat{R}_t \right)^2,$$

  typically via some gradient descent algorithm.
9: **end for**

*(handwritten annotations)* play; $t \to T$; $Q(s,a) = V(s) + A(a)$; const $\forall a$; 101, 102, 99; +100, +1, +2, -1

# Policy Gradient with automatic differentiation

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_\theta log \pi_\theta(a_{i,t}|s_{i,t}) \hat{G}_{i,t}$$

How can we compute policy gradient with autodiff ?
Create a function s.t. it's gradient is policy gradient.
You will implement "pseudo-loss"

$$J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} log \pi_\theta(a_{i,t}|s_{i,t}) \hat{G}_{i,t}$$

cross entropy (discrete) or squared error (gaussian)

# It is good if you can ...

1. Define Markov Decision Process
2. Describe some task in terms of MDP
3. Understand value and action-value functions
4. Be able to apply Policy Iteration and Value Iteration
5. Understand model-based vs model-free
6. Be able to apply MC, TD(0) algorithms to Q-function
7. Be able to apply Q-learning, SARSA, Expected value SARSA
8. Describe the "Exploitation / Exploration" dilemma
9. Understand $\epsilon$-Greedy and "Optimism in the face of uncertainty" (UCB) ideas
10. Understand the idea of Policy Optimization (Policy Gradient)