

# CSC413 Tutorial11: Policy Gradient

by Sheng Jia

March 31st, 2020

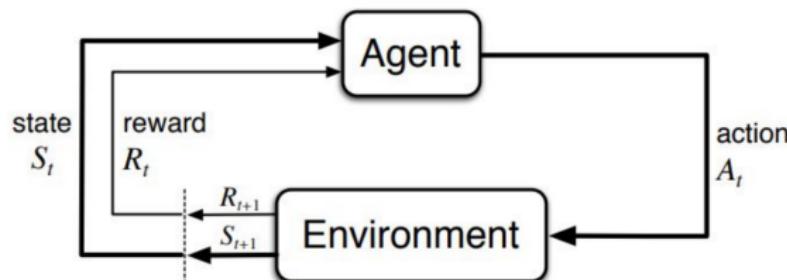
# Content

- State and Action
- Policy
- Trajectory and how to sample it
- Objective in Reinforcement Learning
- Policy optimization by policy gradient ascent
  - **Trajectory-based Policy Gradient Derivation**  
(Log-derivative trick. Exploit conditional independence)
  - **Break: Apply policy gradient for playing Dota2**
  - **Reward-to-go based Policy Gradient Derivation**  
(Exploit conditional independence. expected grad-log-prob equal 0)
  - **Reducing variance of policy gradient estimate by Baseline**  
( $\text{Var}(x - y)$  can be less than  $\text{Var}(x)$ . Expected grad-log-prob equal 0)
- Implementing Policy Gradient in Pytorch  
(Credit to the notebook from the last year's CSC421)

# Problem Setup

## State

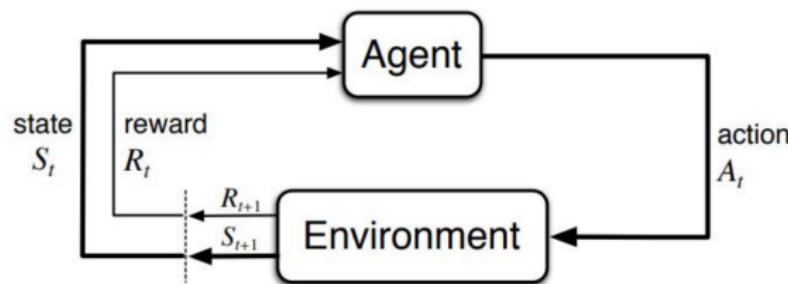
State  $s$  is the complete description of the task/environment from which the agent can make decisions for taking actions and receive rewards. Both state and action are indexed by the timestep as  $s_t, a_t$  during the agent-environment interaction.



# Problem Setup

## State

State  $s$  is the complete description of the task/environment from which the agent can make decisions for taking actions and receive rewards. Both state and action are indexed by the timestep as  $s_t, a_t$  during the agent-environment interaction.



State does not have to be the “physical location” of the agent.  
E.g.  $s_t$  : how many ppl infected with covid19 today.  
 $a_t$  : whether or not wash your hands now.

# Problem Setup

## Agent's Policy

- “Agent” is an abstract concept, but we can formulate how the agent behaves by, for example, a stochastic policy. This can be a conditional distribution that is parameterized by  $\theta$ .

$$p_{\theta}(a_t|s_t) = \pi_{\theta}(a_t|s_t) = \pi(a_t|s_t; \theta)$$

# Problem Setup

Different implementations of a stochastic policy

$$p_{\theta}(a_t|s_t) = \pi_{\theta}(a_t|s_t) = \pi(a_t|s_t; \theta)$$

Based on the problem, implement different types of stochastic policy

- If  $\mathcal{A}$  is discrete, but  $\mathcal{S}$  is continuous or too large (e.g. Atari), use a function approximator such as NN to map the state vector  $s$  to the distribution over actions using softmax for the output layer. i.e. The size of your network output will be  $\mathcal{A}$ , with each output denoting the probability of taking that action.

# Problem Setup

Different implementations of a stochastic policy

$$p_{\theta}(a_t|s_t) = \pi_{\theta}(a_t|s_t) = \pi(a_t|s_t; \theta)$$

Based on the problem, implement different types of stochastic policy

- If both  $\mathcal{S}$  and  $\mathcal{A}$  are continuous or too large (e.g. Robot control), map  $s$  to parameters associated with distributions such as  $\mu$  and  $\sigma^2$  for Gaussian distribution. Then sample the value, which we treat as the action, from this distribution under the mapped  $\mu$  and  $\sigma$ .  
(A simpler solution is to discretize continuous action space. e.g. OpenAI Dota2 bot [1])

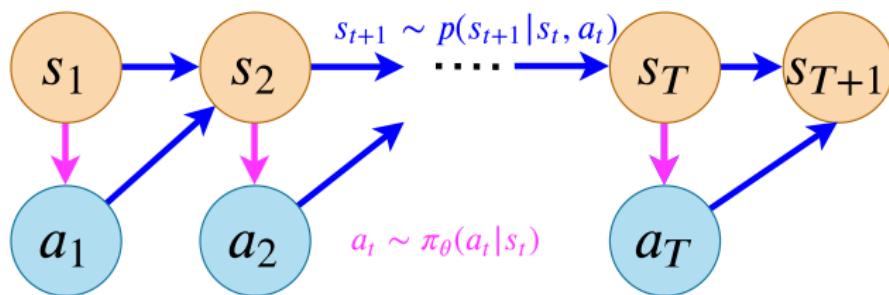
# Problem Setup

## Trajectory

- Trajectory is nothing but a set of random variables, and its distribution is a joint distribution over  $2T + 1$  r.v.:

$$\tau = (s_1, a_1, s_2, \dots, s_T, a_T, s_{T+1})$$

$$p(\tau; \theta) = p(s_1, a_1, s_2, \dots, s_T, a_T, s_{T+1}; \theta) = (\star)$$



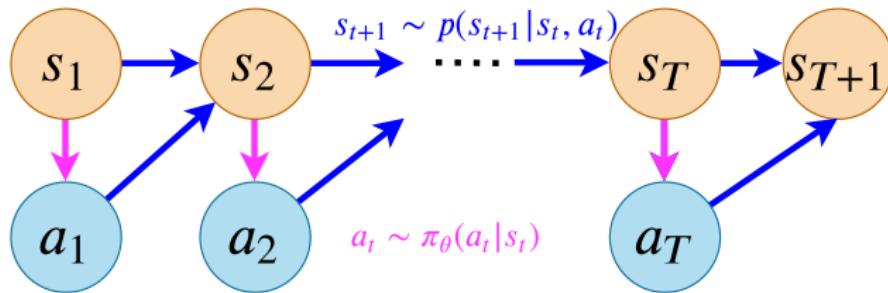
# Problem Setup

## Trajectory

- We can simplify using conditional independences from DAG:

$$(*) = \rho_0(s_1) \prod_{t=1}^T \pi_\theta(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

- Remark: we will use  $p(\tau; \theta)$  to denote that changing our policy parameters  $\theta$  induce a different trajectory distribution.

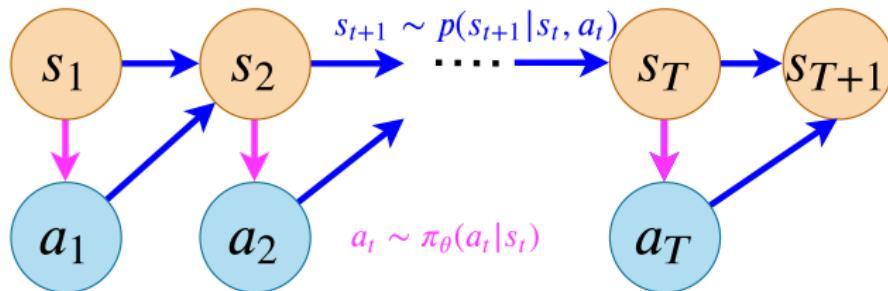


# Problem Setup

How to sample a trajectory (Run/Execute an agent)

- “Running/Executing the agent in a environment” means **ancestral sampling from this DAG**. (Sample the parent node and successively sample the child nodes.)

$$s_1 \sim \rho_0(s) \quad a_t \sim \pi_\theta(a_t|s_t) \quad s_{t+1} \sim p(s_{t+1}|s_t, a_t)$$



# Objective in Reinforcement Learning

## Reward, Return

- Consider reward  $r_t = R(s_t, a_t)$  as something that measures how well action  $a_t$  is in state  $s_t$ . This is computed by a blackbox function  $R(s_t, a_t)$  from the environment.

# Objective in Reinforcement Learning

## Reward, Return

- Consider reward  $r_t = R(s_t, a_t)$  as something that measures how well action  $a_t$  is in state  $s_t$ . This is computed by a blackbox function  $R(s_t, a_t)$  from the environment.
- Return is the cumulative reward for the trajectory  $\tau$ . (Consider finite-horizon undiscounted version in this tutorial)

$$R(\tau) = \sum_{t=1}^T R(s_t, a_t)$$

# Objective in Reinforcement Learning

## Reward, Return

- Consider reward  $r_t = R(s_t, a_t)$  as something that measures how well action  $a_t$  is in state  $s_t$ . This is computed by a blackbox function  $R(s_t, a_t)$  from the environment.
- Return is the cumulative reward for the trajectory  $\tau$ . (Consider finite-horizon undiscounted version in this tutorial)

$$R(\tau) = \sum_{t=1}^T R(s_t, a_t)$$

**Return is also a random variable** because it is a function of  $2T$  random variables in the trajectory.

# Objective in Reinforcement Learning

## Expected Return

- As  $R(\tau)$  is random, the objective is to maximize the expected return  $\mathbb{E}[R(\tau)]$  w.r.t  $\theta$ . By the law of the unconscious statistician, we can write it as the expectation under  $\tau$  distribution  $p(\tau; \theta)$ :

$$\mathcal{J}(\theta) = \mathbb{E}[R(\tau)] = \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] = (\star)$$

# Objective in Reinforcement Learning

## Expected Return

- As  $R(\tau)$  is random, the objective is to maximize the expected return  $\mathbb{E}[R(\tau)]$  w.r.t  $\theta$ . By the law of the unconscious statistician, we can write it as the expectation under  $\tau$  distribution  $p(\tau; \theta)$ :

$$\mathcal{J}(\theta) = \mathbb{E}[R(\tau)] = \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] = (\star)$$

And by the ancestral sampling, we can further simplify:

$$(\star) = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_\theta(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \sum_{t=1}^T R(s_t, a_t) \right]$$

# Policy Optimization by Policy Gradient Ascent

A method to “skill up” the agent

## Policy Optimization by Policy Gradient Ascent

We can make a one-step optimization for the current policy  $\pi_{\theta_k}(a_t|s_t)$  to  $\pi_{\theta_{k+1}}(a_t|s_t)$  for maximizing  $\mathcal{J}(\theta)$  by gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} \mathcal{J}(\theta)|_{\theta_k}$$

# Policy Optimization by Policy Gradient Ascent

A method to “skill up” the agent

## Policy Optimization by Policy Gradient Ascent

We can make a one-step optimization for the current policy  $\pi_{\theta_k}(a_t|s_t)$  to  $\pi_{\theta_{k+1}}(a_t|s_t)$  for maximizing  $\mathcal{J}(\theta)$  by gradient ascent:

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} \mathcal{J}(\theta)|_{\theta_k}$$

## Gradient of the objective w.r.t policy

$$\nabla_{\theta} \mathcal{J}(\theta)|_{\theta_k} = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

# Policy Optimization by Policy Gradient Ascent

## Deriving policy gradient (Step1)

Step1 using log-derivative trick

$$\nabla_{\theta} \mathcal{J}(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)]$$

# Policy Optimization by Policy Gradient Ascent

## Deriving policy gradient (Step1)

Step1 using log-derivative trick

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] \\ &= \nabla_{\theta} \int p(\tau; \theta) R(\tau) d\tau\end{aligned}$$

# Policy Optimization by Policy Gradient Ascent

## Deriving policy gradient (Step1)

Step1 using log-derivative trick

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] \\ &= \nabla_{\theta} \int p(\tau; \theta) R(\tau) d\tau \\ &= \int \nabla_{\theta} p(\tau; \theta) R(\tau) d\tau\end{aligned}$$

# Policy Optimization by Policy Gradient Ascent

## Deriving policy gradient (Step1)

### Step1 using log-derivative trick

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] \\&= \nabla_{\theta} \int p(\tau; \theta) R(\tau) d\tau \\&= \int \nabla_{\theta} p(\tau; \theta) R(\tau) d\tau \\&= \int p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) R(\tau) d\tau \quad \because \nabla_{\theta} \log p(\tau; \theta) = \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)}\end{aligned}$$

# Policy Optimization by Policy Gradient Ascent

## Deriving policy gradient (Step1)

Step1 using log-derivative trick

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p(\tau; \theta)} [R(\tau)] \\&= \nabla_{\theta} \int p(\tau; \theta) R(\tau) d\tau \\&= \int \nabla_{\theta} p(\tau; \theta) R(\tau) d\tau \\&= \int p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta) R(\tau) d\tau \quad \because \nabla_{\theta} \log p(\tau; \theta) = \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} \\&= \mathbb{E}_{\tau \sim p(\tau; \theta)} [\nabla_{\theta} \log p(\tau; \theta) R(\tau)]\end{aligned}$$

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Step2)

Step2 using conditional independences

$$\mathbb{E}_{\tau \sim p(\tau; \theta)} [\nabla_{\theta} \log p(\tau; \theta) R(\tau)] \quad \text{Now use ancestral sampling}$$

$$= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta}(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \underbrace{\nabla_{\theta} \log (\rho_0(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t))}_{\textcircled{1}} \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Step2)

Step2 using conditional independences

$$\mathbb{E}_{\tau \sim p(\tau; \theta)} [\nabla_{\theta} \log p(\tau; \theta) R(\tau)] \quad \text{Now use ancestral sampling}$$

$$= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta}(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \underbrace{\nabla_{\theta} \log (\rho_0(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t))}_{\textcircled{1}} \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

$$\text{where } \textcircled{1} = \nabla_{\theta} \left( \log \rho_0(s_1) + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \sum_{t=1}^T \log p(s_{t+1} | s_t, a_t) \right)$$

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Step2)

Step2 using conditional independences

$$\mathbb{E}_{\tau \sim p(\tau; \theta)} [\nabla_{\theta} \log p(\tau; \theta) R(\tau)] \quad \text{Now use ancestral sampling}$$

$$= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta}(a_t | s_t) \\ s_{t+1} \sim p(s_{t+1} | s_t, a_t)}} \left[ \underbrace{\nabla_{\theta} \log (\rho_0(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t))}_{\textcircled{1}} \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

$$\text{where } \textcircled{1} = \nabla_{\theta} \left( \log \rho_0(s_1) + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \sum_{t=1}^T \log p(s_{t+1} | s_t, a_t) \right)$$

$$= \cancel{\nabla_{\theta} \rho_0(s_1)} + \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) + \sum_{t=1}^T \nabla_{\theta} \log p(s_{t+1} | s_t, a_t) \xrightarrow{0}$$

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Final form)

Hence, the policy gradient w.r.t the current policy parameters is:

$$\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Final form)

Hence, the policy gradient w.r.t the current policy parameters is:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t^{(i)}|s_t^{(i)}) \left[ \sum_{t'=1}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right] \right]\end{aligned}$$

In practice, this gradient is estimated by executing the policy  $\pi_{\theta_k}$  in the environment  $N$  times ( $N$  times ancestral sampling).

# Policy Optimization by Policy Gradient Ascent

Deriving policy gradient (Final form)

Hence, the policy gradient w.r.t the current policy parameters is:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t^{(i)}|s_t^{(i)}) \left[ \sum_{t'=1}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right] \right]\end{aligned}$$

The log-derivative trick in step 1 allows for this type of gradient estimate of the expected value even though the thing inside expectation was a blackbox function using samples from the parameterized distribution.  
“score function estimator”

# Break: Apply policy gradient for playing Dota2

Successful application of policy optimization by policy gradient

- In Dota2, each team have five players controlling their unique agents. Players gather golds by killing monsters and enemies to buy items. The final objective is destroy an enemy structure called Ancient. OpenAI agents recently won against the best team in the world. [1]



# Break: Apply policy gradient for playing Dota2

## Observation (Input of the policy)

- State  $S$ : **16000-dimensional vector** with information such as the distances to the observed enemies. But **it is partially observable** because teams don't see the map far from the current locations even if they went there before. **LSTM is used to memorize previous states.**



# Break: Apply policy gradient for playing Dota2

## Action (Output of the policy)

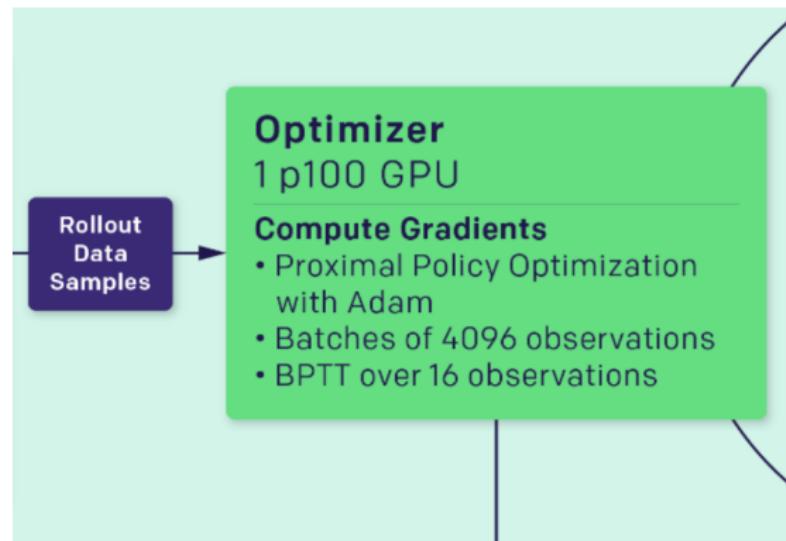
- Action  $\mathcal{A}$ : Continuous, but discretized into 8000-80000 actions.



# Break: Apply policy gradient for playing Dota2

Policy optimization by policy gradient ascent

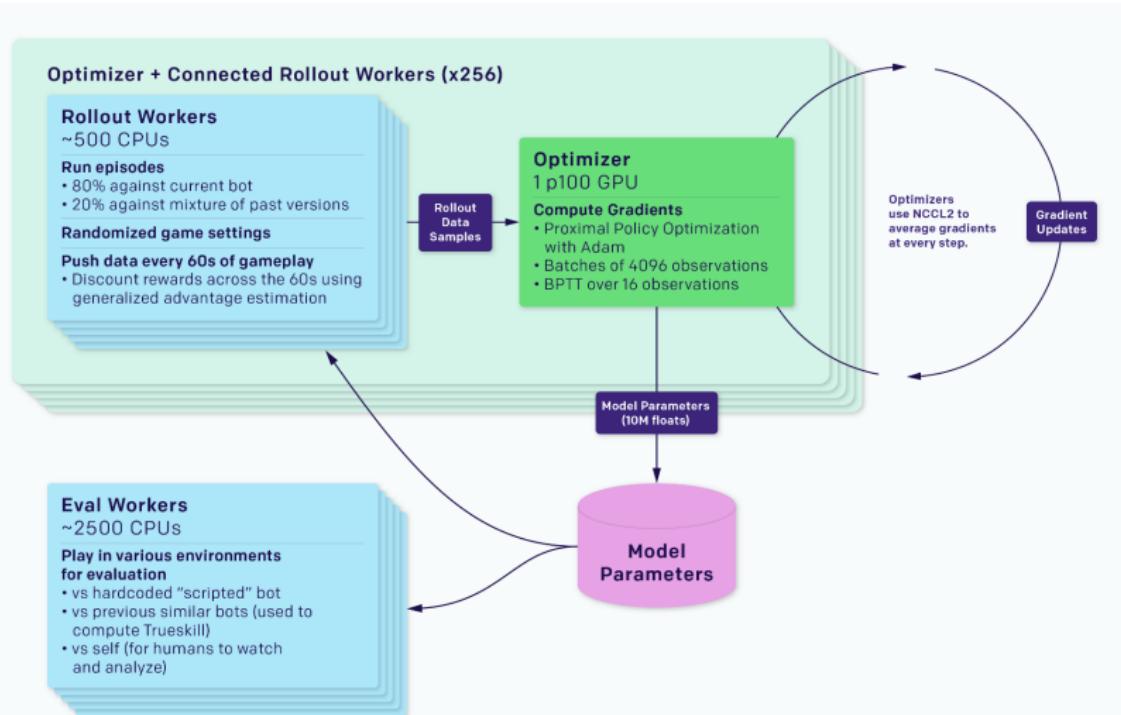
- Besides winning the game, intermediate rewards such as kill enemies are provided. PPO, an improved **policy gradient method**, is used to train the policy with **Adam** optimizer. [1]



# Break: Apply policy gradient for playing Dota2

## Large-scale engineering

- Rollouts against the current bot means self-play [1]



**Next: Deriving Reward-to-go Policy Gradient**

# Reward-to-go Policy Gradient

Intuitively, the rewards  $R(s_1, a_1), \dots R(s_{t-1}, a_{t-1})$  obtained before taking the action  $a_t$  should not tell how good action  $a_t$  is. This claim is saying:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right] \\ &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=t}^T R(s_{t'}, a_{t'}) \right] \right]\end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

Proved using the DAG structure and expected grad-log-prob equal 0:

$$\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

# Reward-to-go Policy Gradient (Proof)

Proved using the DAG structure and expected grad-log-prob equal 0:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right] \\ &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_t, s_{t'}, a_{t'}} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) R(s_{t'}, a_{t'})] \quad \text{by linearity}\end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

Proved using the DAG structure and expected grad-log-prob equal 0:

$$\begin{aligned}\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} &= \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right] \\ &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_t, s_{t'}, a_{t'}} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) R(s_{t'}, a_{t'})] \quad \text{by linearity} \\ &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_{t'}, a_{t'}} \left[ \mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) R(s_{t'}, a_{t'}) | s_{t'}, a_{t'}] \right] \quad \text{by iterated } \mathbb{E}\end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

Proved using the DAG structure and expected grad-log-prob equal 0:

$$\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=1}^T R(s_{t'}, a_{t'}) \right] \right]$$

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_t, s_{t'}, a_{t'}} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) R(s_{t'}, a_{t'})] \quad \text{by linearity}$$

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_{t'}, a_t} \left[ \mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) R(s_{t'}, a_{t'}) | s_{t'}, a_{t'}] \right] \quad \text{by iterated } \mathbb{E}$$

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_{t'}, a_{t'}} \left[ R(s_{t'}, a_{t'}) \underbrace{\mathbb{E}_{s_t, a_t} [\nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) | s_{t'}, a_{t'}]}_{(*) \text{ apply iterated expectation again}} \right]$$

# Reward-to-go Policy Gradient (Proof)

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_{t'}, a_{t'}} \left[ R(s_{t'}, a_{t'}) \mathbb{E}_{s_t} \left[ \mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_t] | s_{t'}, a_{t'} \right] \right]$$

# Reward-to-go Policy Gradient (Proof)

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_{t'}} \left[ R(s_{t'}, a_{t'}) \mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_{t'}, a_{t'}] \right]$$

The above step can be verified by:

$$(\star) = \int \int p(s_t, a_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t$$

# Reward-to-go Policy Gradient (Proof)

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_{t'}} \left[ R(s_{t'}, a_{t'}) \mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_{t'}, a_{t'}] \right]$$

The above step can be verified by:

$$\begin{aligned} (\star) &= \int \int p(s_t, a_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \\ &= \int \int p(a_t | s_t, s_{t'}, a_{t'}; \theta_k) p(s_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_t, a_{t'}} \left[ R(s_{t'}, a_{t'}) \mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_{t'}, a_{t'}] \right]$$

The above step can be verified by:

$$\begin{aligned} (*) &= \int \int p(s_t, a_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \\ &= \int \int p(a_t | s_t, s_{t'}, a_{t'}; \theta_k) p(s_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \\ &= \int p(s_t | s_{t'}, a_{t'}; \theta_k) \left( \int p(a_t | s_t, s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) da_t \right) ds_t \end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

$$= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}_{s_{t'}, a_{t'}} \left[ R(s_{t'}, a_{t'}) \mathbb{E}_{s_t} \left[ \underbrace{\mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_t]}_{(\diamond)} | s_{t'}, a_{t'} \right] \right]$$

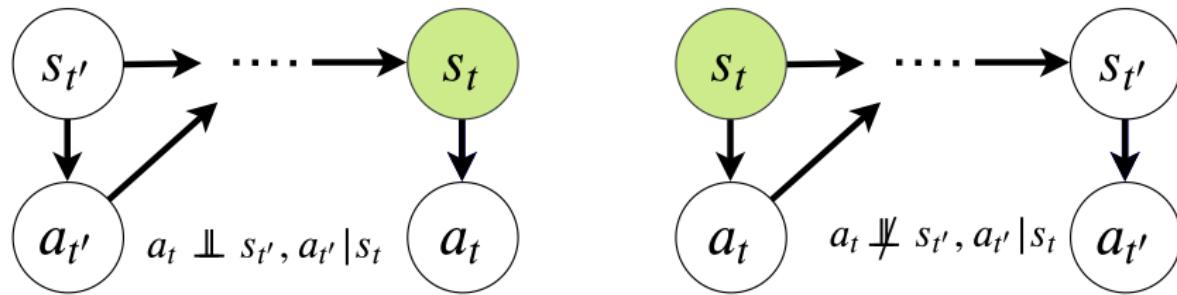
The above step can be verified by:

$$\begin{aligned} (\star) &= \int \int p(s_t, a_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \\ &= \int \int p(a_t | s_t, s_{t'}, a_{t'}; \theta_k) p(s_t | s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) ds_t da_t \\ &= \int p(s_t | s_{t'}, a_{t'}; \theta_k) \left( \int p(a_t | s_t, s_{t'}, a_{t'}; \theta_k) \nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) da_t \right) ds_t \\ &= \mathbb{E}_{s_t} \left[ \underbrace{\mathbb{E}_{a_t \sim p(a_t | s_t, s_{t'}, a_{t'}; \theta_k)} [\nabla_{\theta} \log \pi_{\theta_k}(a_t | s_t) | s_t]}_{(\diamond)} | s_{t'}, a_{t'} \right] \end{aligned}$$

# Reward-to-go Policy Gradient (Proof)

Final step using DAG structure and Expected grad-log-prob equal 0

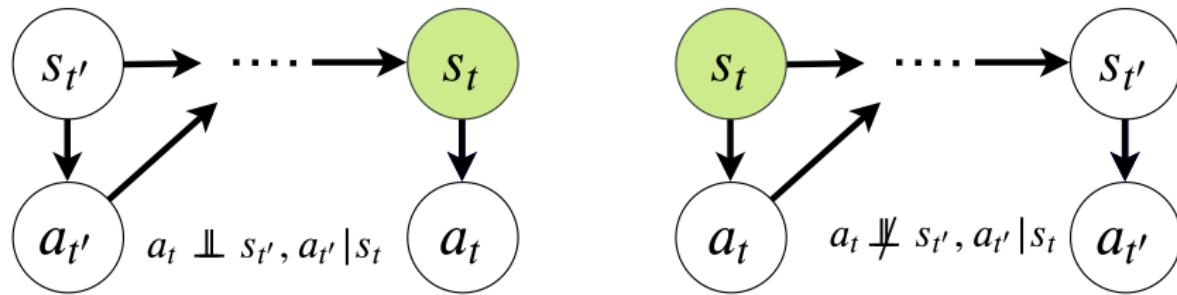
From the graphical model, we can observe the conditional independence when  $t' < t$ :



# Reward-to-go Policy Gradient (Proof)

Final step using DAG structure and Expected grad-log-prob equal 0

From the graphical model, we can observe the conditional independence when  $t' < t$ :

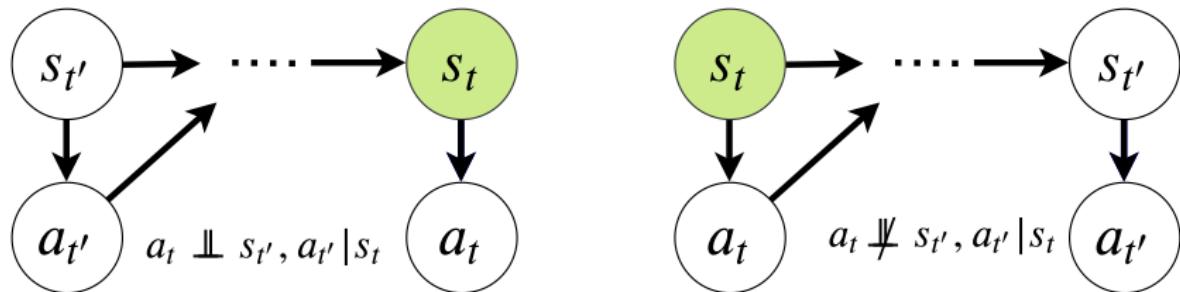


Hence, if  $t' < t$ ,  $p(a_t | s_t, s_{t'}, a_{t'}; \theta) = p(a_t | s_t; \theta) = \pi_\theta(a_t | s_t)$

# Reward-to-go Policy Gradient (Proof)

Final step using DAG structure and Expected grad-log-prob equal 0

From the graphical model, we can observe the conditional independence when  $t' < t$ :



Hence, if  $t' < t$ ,  $p(a_t | s_t, s_{t'}, a_{t'}; \theta) = p(a_t | s_t; \theta) = \pi_\theta(a_t | s_t)$

$$\begin{aligned} (\diamond) &= \mathbb{E}_{a_t \sim \pi_\theta(a_t | s_t)} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_t] = \int \pi_\theta(a_t | s_t) \nabla_\theta \log \pi_\theta(a_t | s_t) da_t \\ &= \int \nabla_\theta \pi_\theta(a_t | s_t) da_t = \nabla_\theta \int \pi_\theta(a_t | s_t) da_t = \nabla_\theta 1 = \mathbf{0} \end{aligned}$$

# Reward-to-go Policy Gradient

Hence, all the reward terms for  $t' < t$  will naturally disappear when taking the expectation over  $\tau = (s_1, a_1, \dots, s_T, a_T, s_{T+1})$

## Reward-to-go Policy Gradient

$$\nabla_{\theta} \mathcal{J}(\theta) |_{\theta_k} = \mathbb{E}_{\substack{s_1 \sim \rho_0(s) \\ a_t \sim \pi_{\theta_k}(a_t|s_t) \\ s_{t+1} \sim p(s_{t+1}|s_t, a_t)}} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t) \left[ \sum_{t'=t}^T R(s_{t'}, a_{t'}) \right] \right]$$

## Reducing variance of policy gradient estimate by Baseline

Gradient of the objective was an expectation, so we can only compute the gradient estimate (which is a random variable) from sampled trajectories:

$$\hat{\mathbf{g}} = \hat{\nabla}_{\theta} \mathcal{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \sum_{t'=t}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right] \right]$$

As  $\hat{\mathbf{g}}$  is random, we can talk about **bias** and **variance**. It is easy to see that this estimator is unbiased,  $\mathbb{E}[\hat{\mathbf{g}}] = \mathbf{g} = \nabla_{\theta} \mathcal{J}(\theta)$ .

## Reducing variance of policy gradient estimate by Baseline

Gradient of the objective was an expectation, so we can only compute the gradient estimate (which is a random variable) from sampled trajectories:

$$\hat{\mathbf{g}} = \hat{\nabla}_{\theta} \mathcal{J}(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \sum_{t'=t}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right] \right]$$

As  $\hat{\mathbf{g}}$  is random, we can talk about **bias** and **variance**. It is easy to see that this estimator is unbiased,  $\mathbb{E}[\hat{\mathbf{g}}] = \mathbf{g} = \nabla_{\theta} \mathcal{J}(\theta)$ . Consider

$$\hat{\mathbf{g}}' = \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \sum_{t'=t}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) - V_{\pi_{\theta}}(s_t^{(i)}) \right] \right]$$

where  $V_{\pi_{\theta}}(s_t)$  is random since  $s_t$  is random in this context.

# Reducing variance of policy gradient estimate by Baseline

$$\begin{aligned}\hat{\mathbf{g}}' &= \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \sum_{t'=t}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) - V_{\pi_{\theta}}(s_t^{(i)}) \right] \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ \sum_{t'=t}^T R(s_{t'}^{(i)}, a_{t'}^{(i)}) \right] \right] \\ &\quad - \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ V_{\pi_{\theta}}(s_t^{(i)}) \right] \right] \\ &= \hat{\mathbf{g}} - \mathbf{f}\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathbf{f}] &= \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[ V_{\pi_{\theta}}(s_t^{(i)}) \right] \right] \right] \\ &= \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) [V_{\pi_{\theta}}(s_t)] \right] \right] \quad \tau_i \text{ i.i.d}\end{aligned}$$

## Reducing variance of policy gradient estimate by Baseline

Similar to the derivation in Reward-to-go PG, the expected grad-log-prob equal 0 is also useful here.

$$= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_t] \right] \quad \text{out from inner } \mathbb{E}$$

## Reducing variance of policy gradient estimate by Baseline

Similar to the derivation in Reward-to-go PG, the expected grad-log-prob equal 0 is also useful here.

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_t] \right] \quad \text{out from inner } \mathbb{E} \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \int \pi_\theta(a_t | s_t) \nabla_\theta \log \pi_\theta(a_t | s_t) da_t \right] \end{aligned}$$

# Reducing variance of policy gradient estimate by Baseline

Similar to the derivation in Reward-to-go PG, the expected grad-log-prob equal 0 is also useful here.

$$\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \mathbb{E}_{a_t \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a_t | s_t) | s_t] \right] \quad \text{out from inner } \mathbb{E} \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \int \pi_\theta(a_t | s_t) \nabla_\theta \log \pi_\theta(a_t | s_t) da_t \right] \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} \left[ V_{\pi_\theta}(s_t) \int \nabla \pi_\theta(a_t | s_t) da_t \right] = \dots \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \mathbb{E}_{s_t} [V_{\pi_\theta}(s_t) \mathbf{0}] \\ &= 0 \end{aligned}$$

## Reducing variance of policy gradient estimate by Baseline

$$\mathbb{E} [\hat{\mathbf{g}}'] = \mathbb{E} [\hat{\mathbf{g}} - \mathbf{f}] = \mathbb{E} [\hat{\mathbf{g}}] - \mathbb{E} [\mathbf{f}] = \mathbf{g} + 0 = \nabla_{\theta} \mathcal{J}(\theta)$$

Hence, this new gradient estimate  $\hat{\mathbf{g}'}$  is unbiased so we can use it for policy gradient ascent.

## Reducing variance of policy gradient estimate by Baseline

$$\mathbb{E} [\hat{\mathbf{g}}'] = \mathbb{E} [\hat{\mathbf{g}} - \mathbf{f}] = \mathbb{E} [\hat{\mathbf{g}}] - \mathbb{E} [\mathbf{f}] = \mathbf{g} + 0 = \nabla_{\theta} \mathcal{J}(\theta)$$

Hence, this new gradient estimate  $\hat{\mathbf{g}'}$  is unbiased so we can use it for policy gradient ascent. **But the point is that we want to decrease the variance by:**

$$\text{Var}(\hat{\mathbf{g}'}) = \text{Var}(\hat{\mathbf{g}}) + \text{Var}(\mathbf{f}) - 2\text{Cov}(\hat{\mathbf{g}}, \mathbf{f}) \leq \text{Var}(\hat{\mathbf{g}})$$

$$\text{if } \text{Cov}(\hat{\mathbf{g}}, \mathbf{f}) \geq \frac{1}{2}\text{Var}(\mathbf{f})$$

In practice, we do see strong positive correlations between  $\hat{\mathbf{g}}$  and  $\mathbf{f}$  because the empirical rewards for  $(s_t, a_t, \dots)$  and the value function evaluation for the sampled state  $s_t$  do positively correlate.

# Demo in PyTorch

## (Credit to last year CSC421 RL tutorial)

# Reference



Christopher Berner et al. “Dota 2 with Large Scale Deep Reinforcement Learning”. In: *arXiv preprint arXiv:1912.06680* (2019).