

CSCE 636 Deep Learning

Lecture 21: Deep Reinforcement Learning (continued)

Anxiao (Andrew) Jiang

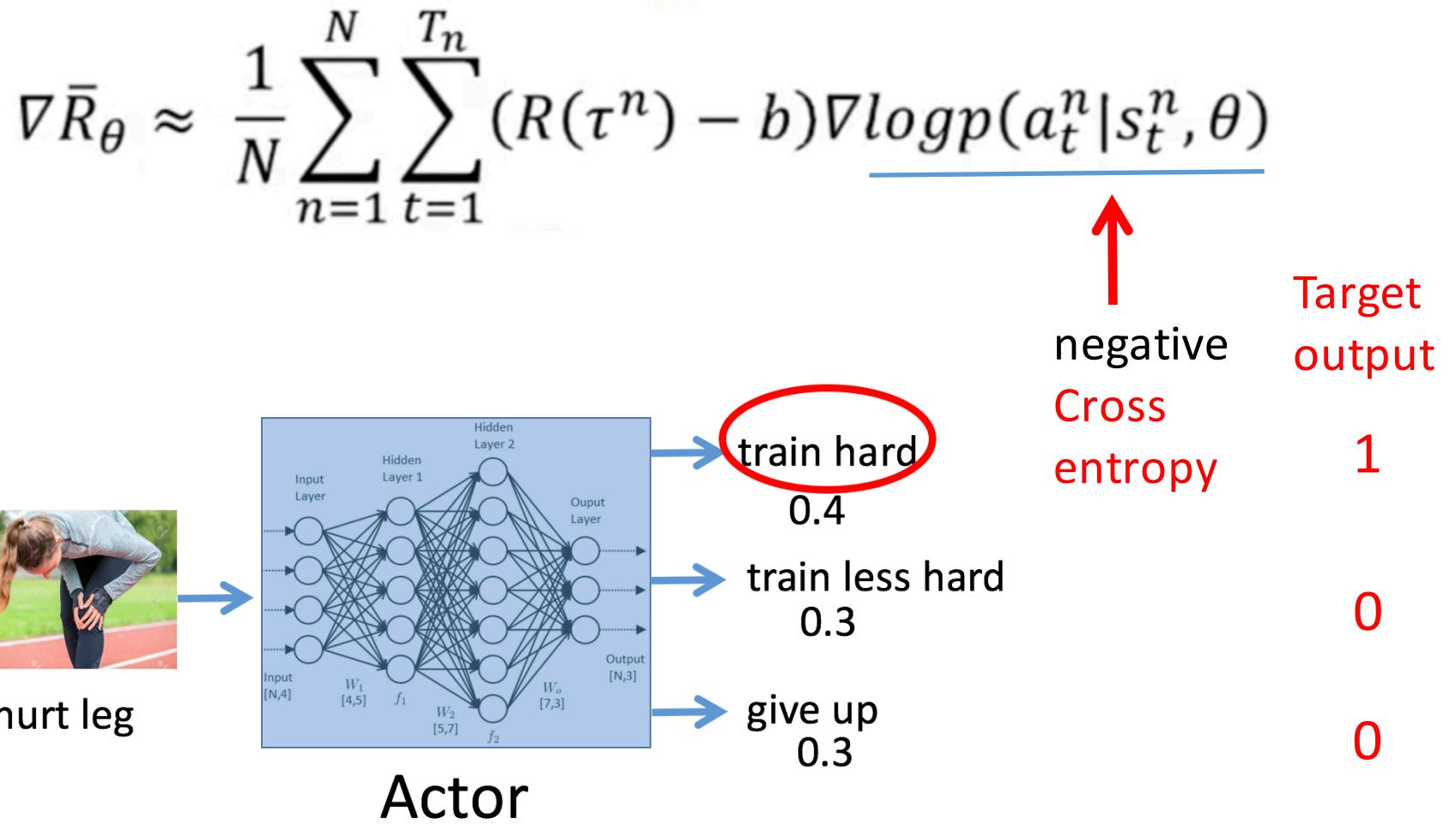
Based on the interesting lecture of Prof. Hung-yi Lee “Deep Reinforcement Learning”

<https://www.youtube.com/watch?v=z95ZYgPgXOY>

Policy-based Approach

(a.k.a. Policy Gradient Approach)

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p(a_t^n | s_t^n, \theta)$$

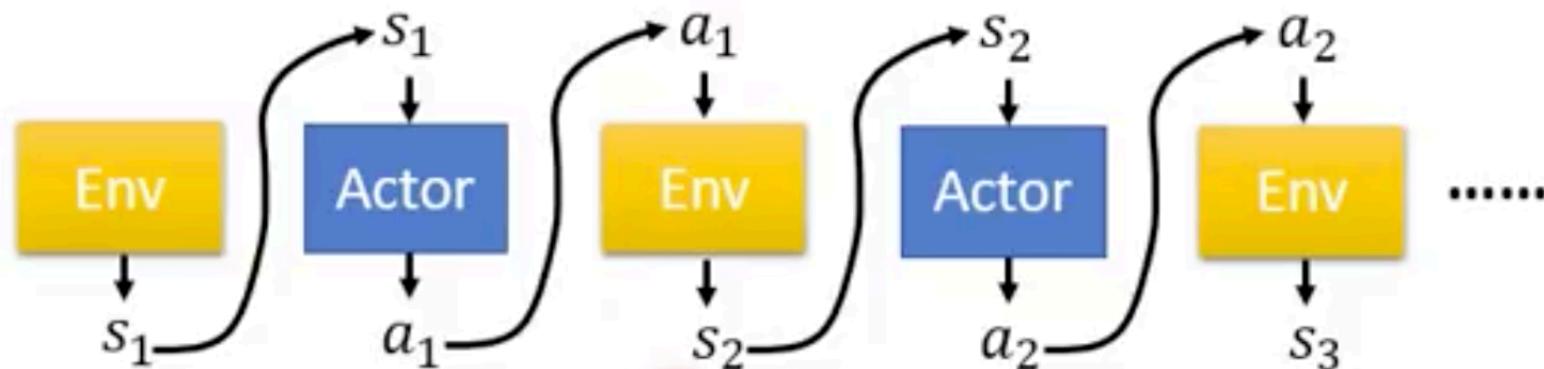


Proximal Policy Optimization(PPO)

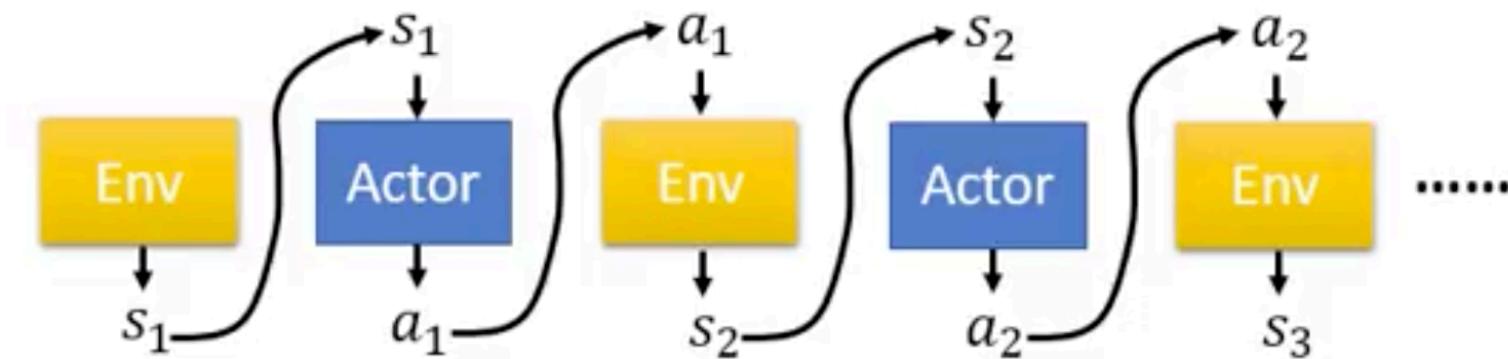
default reinforcement learning algorithm at OpenAI



Actor, Environment, Reward



Actor, Environment, Reward



Trajectory $\tau = \{s_1, a_1, s_2, a_2, \dots, s_T, a_T\}$

$p_\theta(\tau)$

$$= p(s_1)p_\theta(a_1|s_1)p(s_2|s_1, a_1)p_\theta(a_2|s_2)p(s_3|s_2, a_2) \dots$$

Policy Gradient

Given policy π_θ

$$\tau^1: \quad (s_1^1, a_1^1) \quad R(\tau^1)$$

$$(s_2^1, a_2^1) \quad R(\tau^1)$$

⋮

⋮

$$\tau^2: \quad (s_1^2, a_1^2) \quad R(\tau^2)$$

$$(s_2^2, a_2^2) \quad R(\tau^2)$$

⋮

⋮

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

Policy Gradient

Given policy π_θ

$$\tau^1: \quad (s_1^1, a_1^1) \quad R(\tau^1)$$

$$(s_2^1, a_2^1) \quad R(\tau^1)$$

⋮

⋮

$$\tau^2: \quad (s_1^2, a_1^2) \quad R(\tau^2)$$

$$(s_2^2, a_2^2) \quad R(\tau^2)$$

⋮

⋮

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta =$$

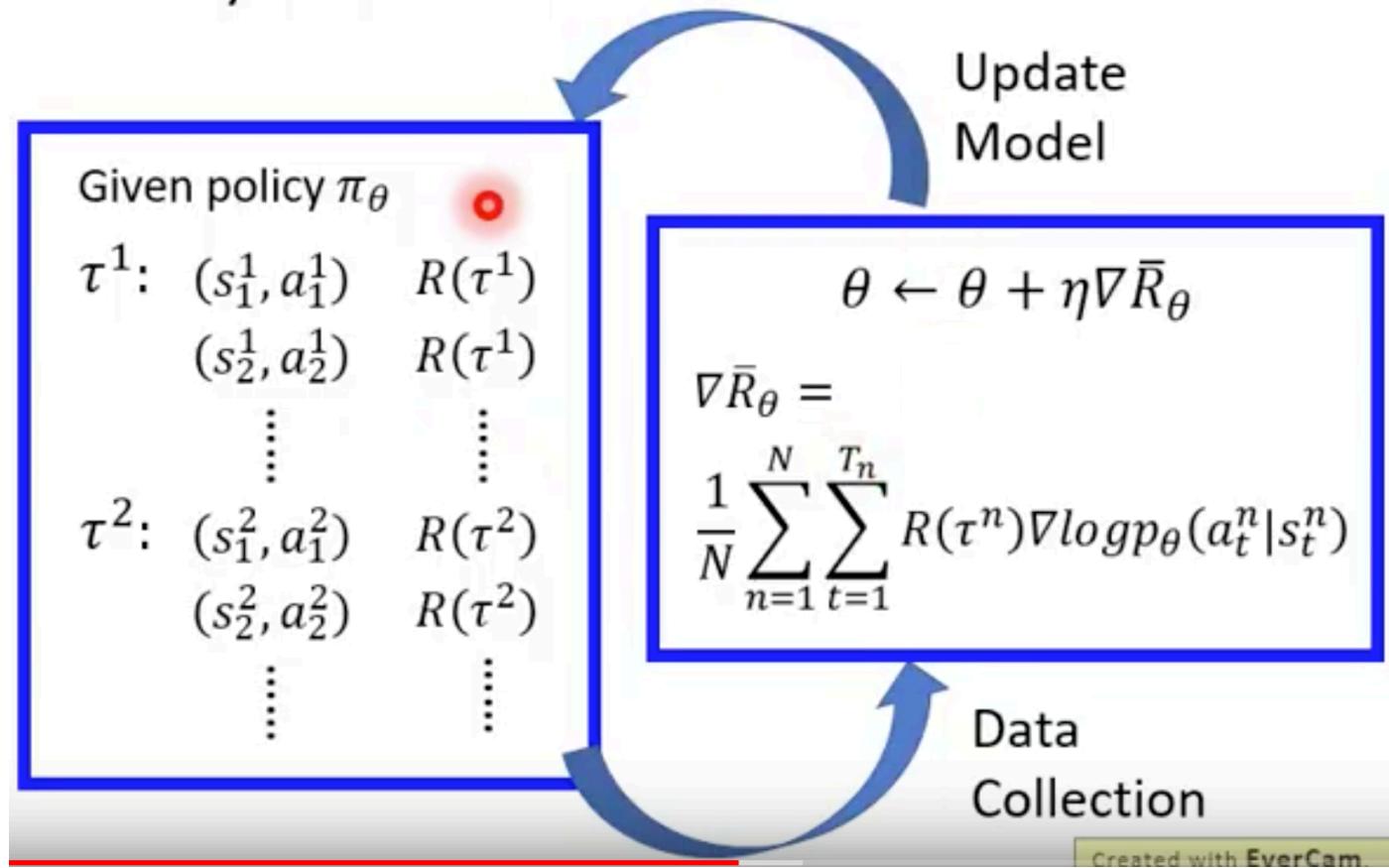
$$\frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$



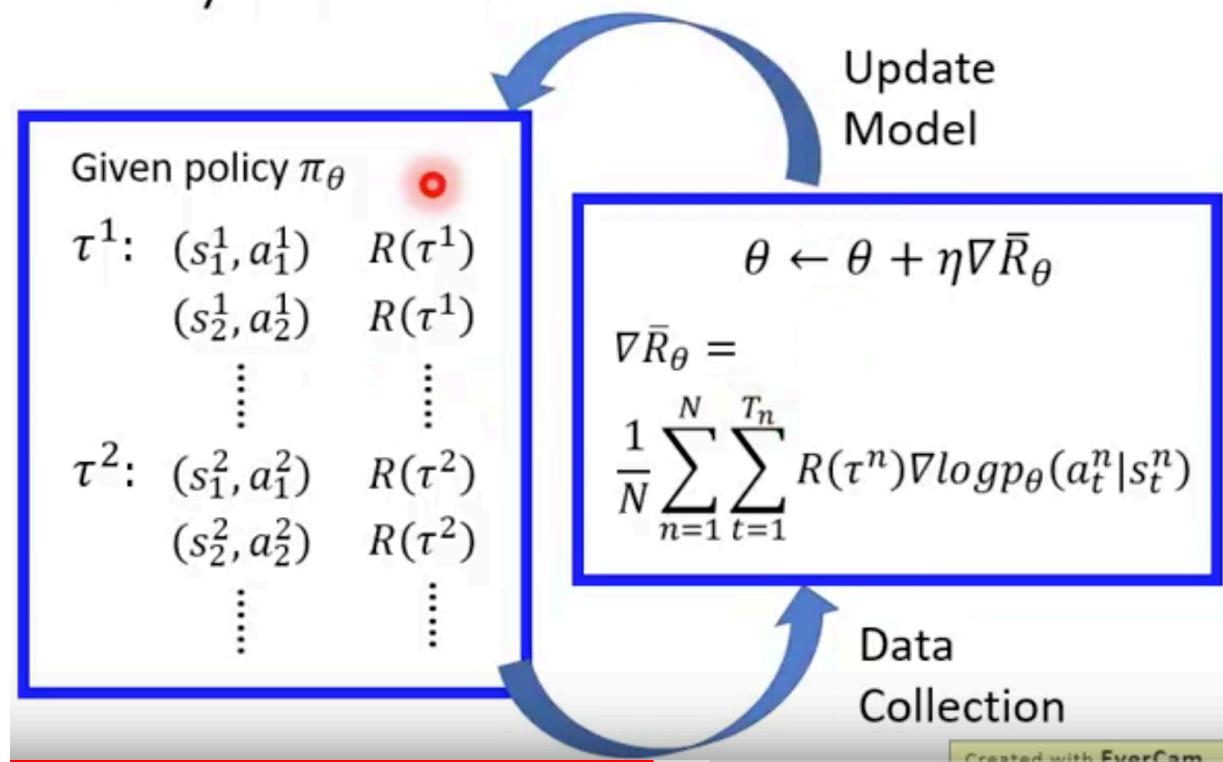
Data
Collection

Created with Easycam

Policy Gradient



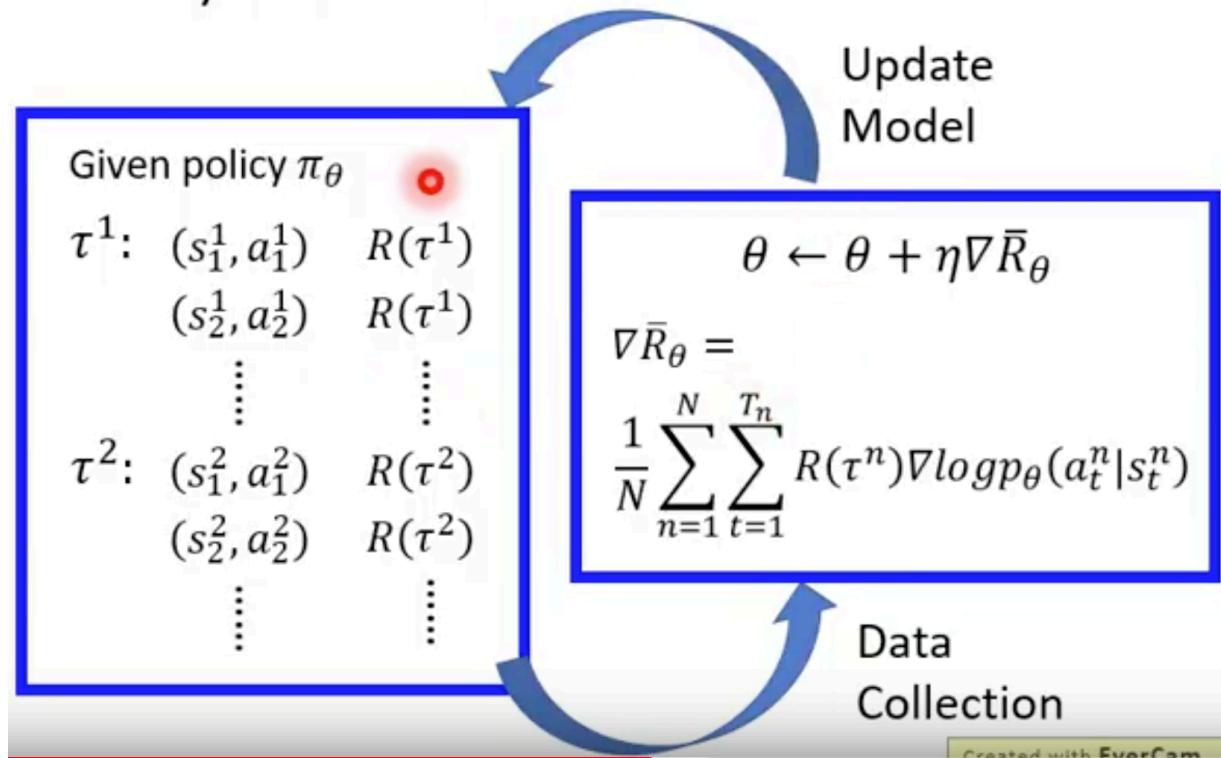
Policy Gradient



Rigorously speaking, we should use the above data only once (to update the NN's weights only once).

Then we need to collect data AGAIN!

Policy Gradient

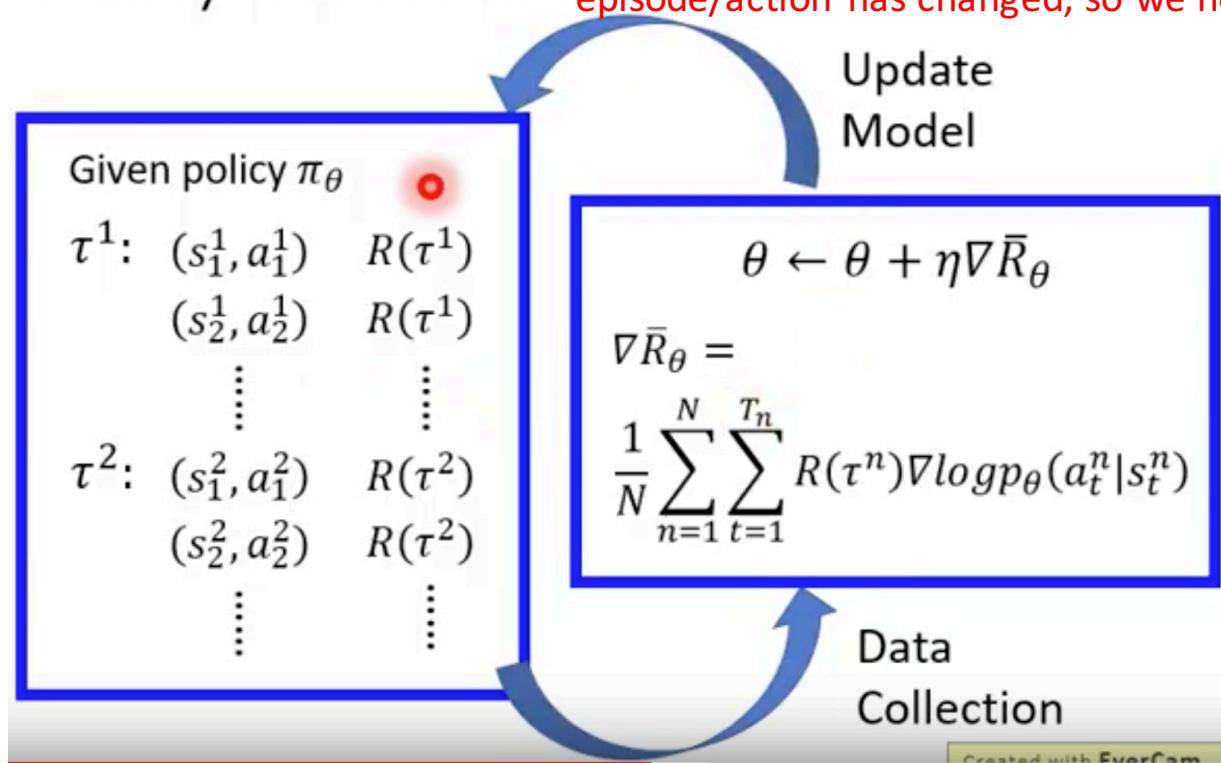


Rigorously speaking, we should use the above data only once (to update the NN's weights only once).

Policy Gradient

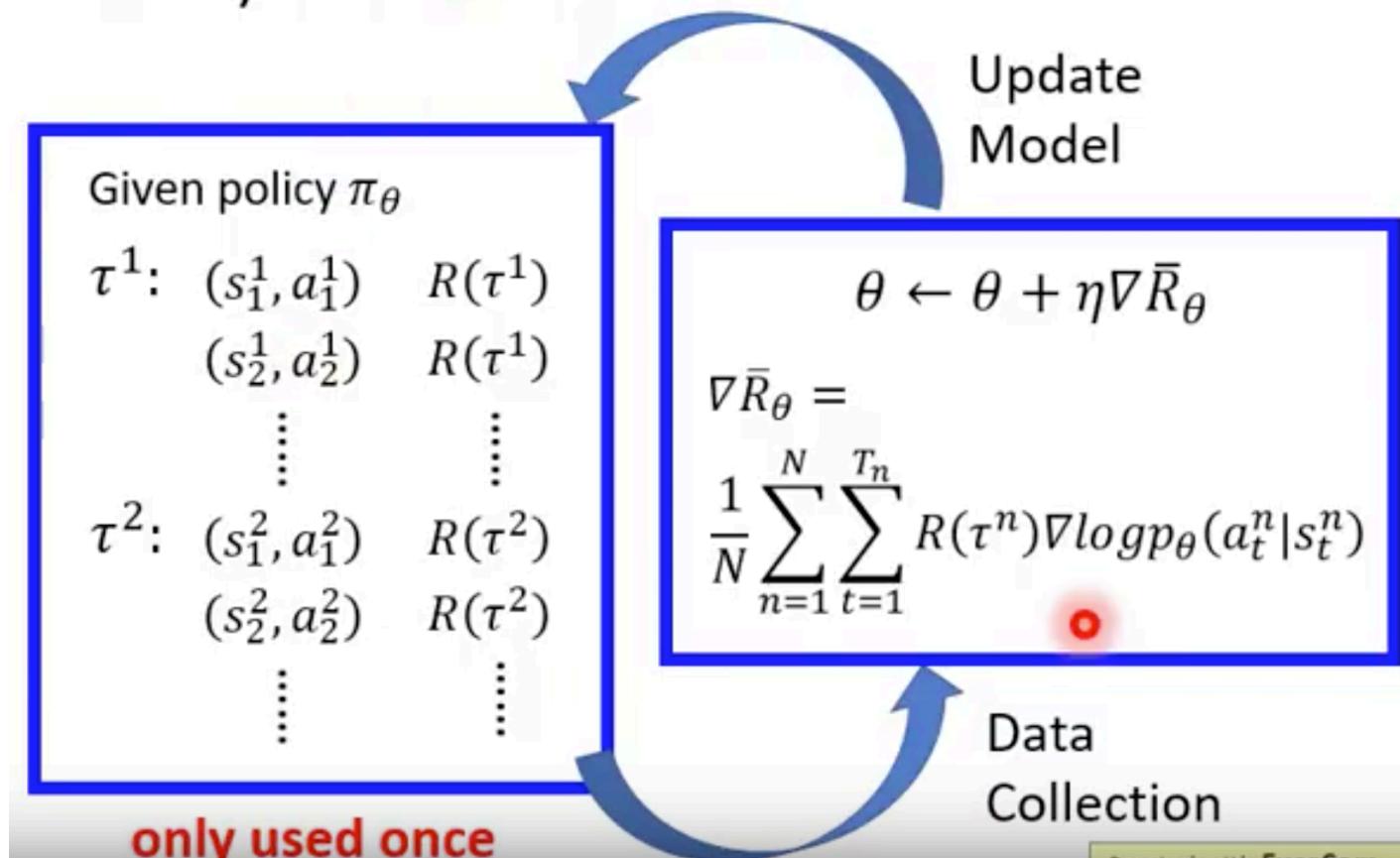
Then we need to collect data AGAIN!

Because when NN weights are updated, the probability of each episode/action has changed, so we need to sample again.



Rigorously speaking, we should use the above data only once (to update the NN's weights only once).

Policy Gradient

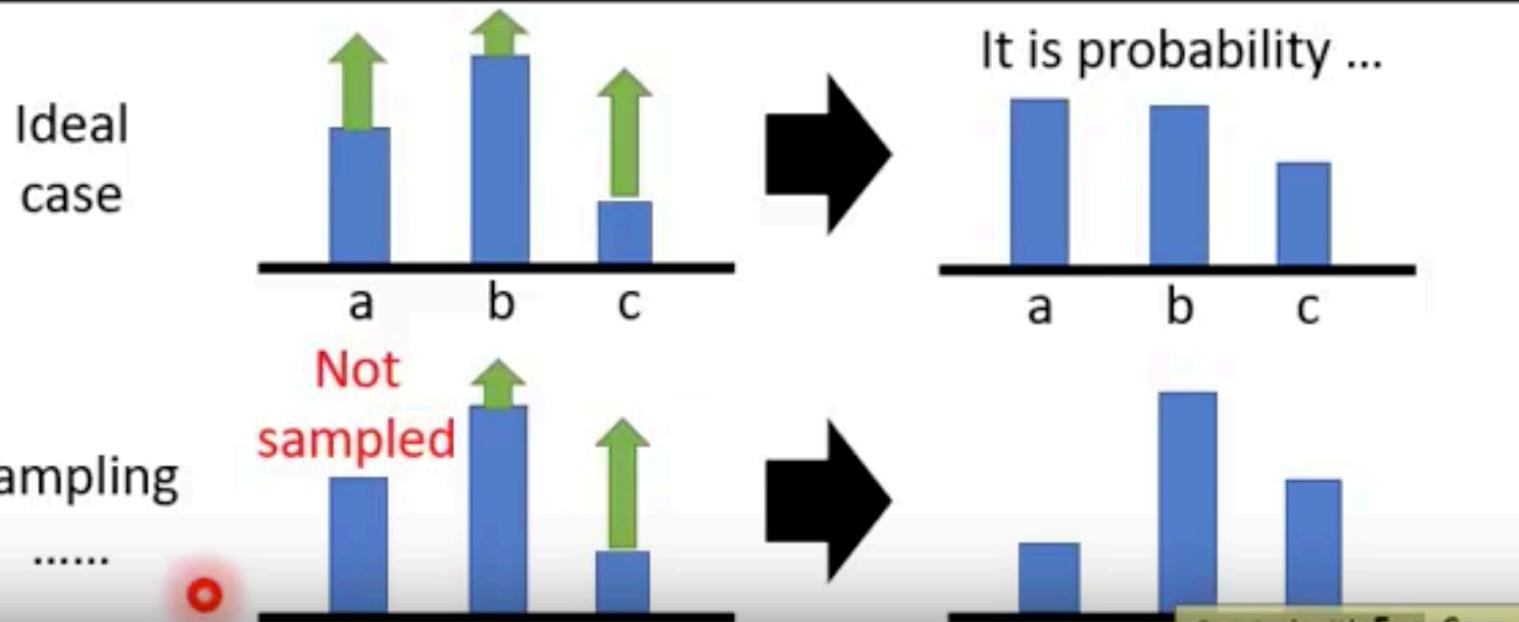


Very time consuming!

Tip 1: Add a Baseline

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta \quad \boxed{\text{It is possible that } R(\tau^n) \text{ is always positive.}}$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

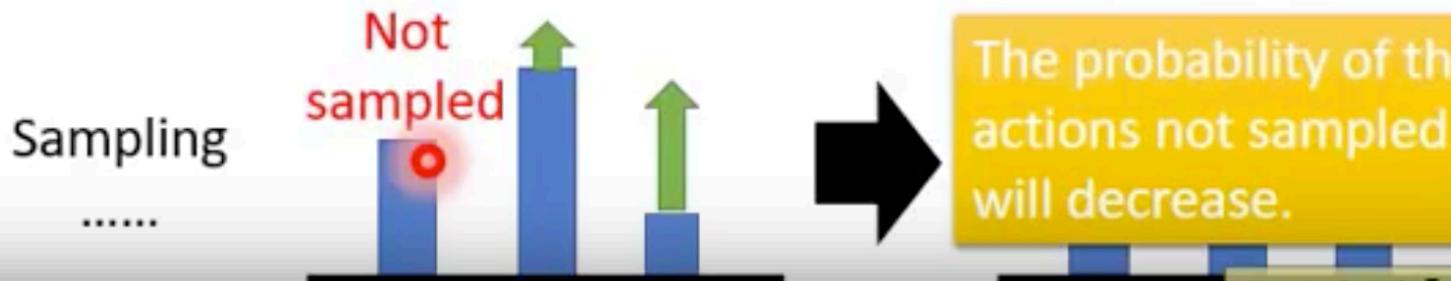
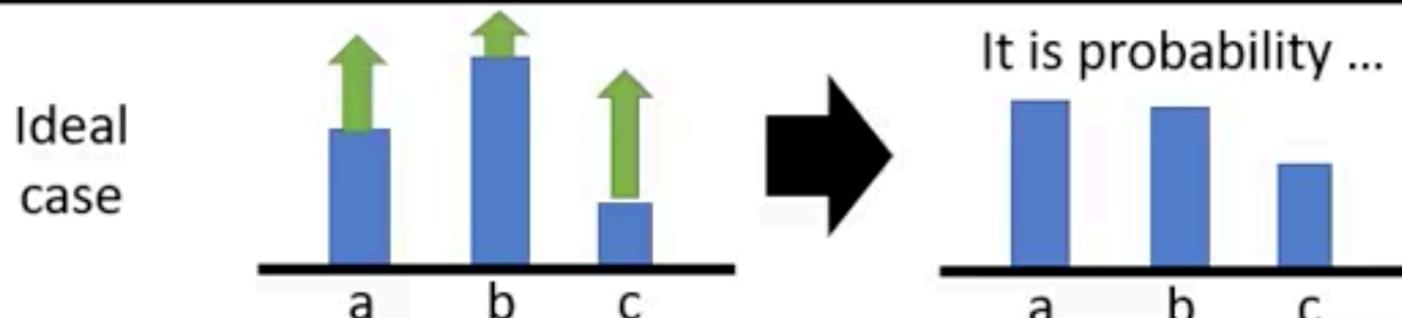


Tip 1: Add a Baseline

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta$$

It is possible that $R(\tau^n)$ is always positive.

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

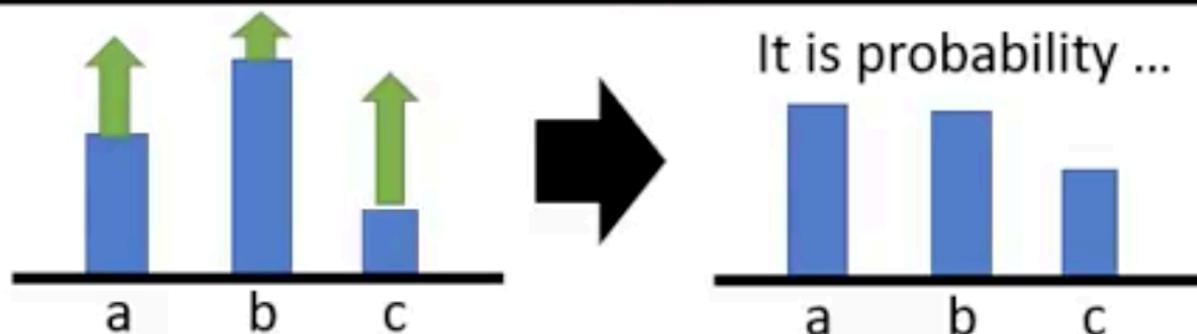


Tip 1: Add a Baseline

$$\theta \leftarrow \theta + \eta \nabla \bar{R}_\theta \quad \boxed{\text{It is possible that } R(\tau^n) \text{ is always positive.}}$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (R(\tau^n) - b) \nabla \log p_\theta(a_t^n | s_t^n) \quad b \approx E[R(\tau)]$$

Ideal
case



Sampling

.....



Tip 2: Assign Suitable Credit

$\times 3$	$\times 3$	$\times 3$	$\times -7$	$\times -7$	$\times -7$
(s_a, a_1)	(s_b, a_2)	(s_c, a_3)	(s_a, a_1)	(s_b, a_2)	(s_c, a_3)
+5	+0	-2	-5	+0	-2
$R = +3$			$R = -7$		

Tip 2: Assign Suitable Credit

$\times 3$	$\times -2$	$\times -2$	$\times -7$	$\times -2$	$\times -2$
(s_a, a_1)	(s_b, a_2)	(s_c, a_3)	(s_a, a_1)	(s_b, a_2)	(s_c, a_3)
+5	+0	-2	-5	+0	-2
$R = +3$			$R = -7$		

Let the weight for an action be the total reward since this action in the episode, instead of the total reward of this whole episode (before and after this action).

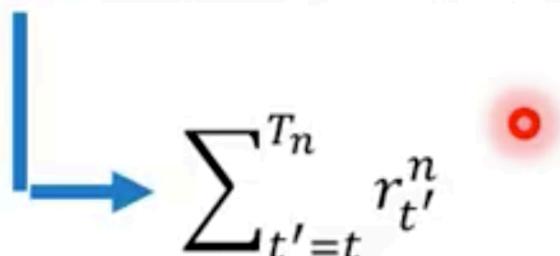
Tip 2: Assign Suitable Credit

$\times 3$	$\times -2$	$\times -2$	$\times -7$	$\times -2$	$\times -2$
(s_a, a_1)	(s_b, a_2)	(s_c, a_3)	(s_a, a_1)	(s_b, \bullet_2)	(s_c, a_3)
+5	+0	-2	-5	+0	-2
$R = +3$			$R = -7$		

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\cdot, \cdot)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$


 $\sum_{t'=t}^{T_n} r_{t'}^n$

Tip 2: Assign Suitable Credit

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\cdot)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$


Add discount factor

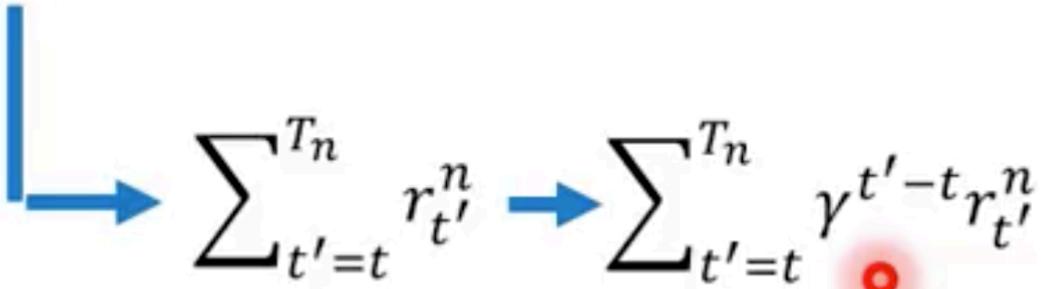
Tip 2: Assign Suitable Credit

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\tau^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$
$$\sum_{t'=t}^{T_n} r_{t'}^n \xrightarrow{\text{Add discount factor}} \sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n$$

Add discount factor

Tip 2: Assign Suitable Credit

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\epsilon^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$


Add discount factor $\gamma < 1$

Tip 2: Assign Suitable Credit

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\tau^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

Can be state-dependent



$$\sum_{t'=t}^{T_n} r_{t'}^n \rightarrow \sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n$$

Add discount factor $\gamma < 1$

(k)

Tip 2: Assign Suitable Credit

Advantage
Function

$$A^\theta(s_t, a_t)$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\tau^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

Can be state-dependent

$$\sum_{t'=t}^{T_n} r_{t'}^n \rightarrow \sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n$$

Add discount factor

$$\gamma < 1$$

Tip 2: Assign Suitable Credit

Advantage
Function

$$A^\theta(s_t, a_t)$$

How good it is if we take a_t other than other actions at s_t .
Estimated by “critic” (later)

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} (\cancel{R(\tau^n)} - b) \nabla \log p_\theta(a_t^n | s_t^n)$$

Can be state-dependent

$\sum_{t'=t}^{T_n} r_{t'}^n \rightarrow \sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n$

Add discount factor $\gamma < 1$

From on-policy to off-policy

Using the experience more than once

On-policy v.s. Off-policy

- On-policy: The agent learned and the agent interacting with the environment is the same.
- Off-policy: The agent learned and the agent interacting with the environment is different.

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
-

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
 - Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.
-

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.

Old weights

New (updated)
weights

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)}[R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$, to train θ . θ' is fixed, so we can re-use the sample data.

Importance Sampling



$$E_{x \sim p}[f(x)]$$

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)}[R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.

Importance Sampling

$$E_{x \sim p}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N f(x^i)$$

x^i is sampled from $p(x)$

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)}[R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$, to train θ . θ' is fixed, so we can re-use the sample data.

Importance Sampling

$$E_{x \sim p}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N f(x^i)$$

x^i is sampled from $p(x)$
We only have x^i sampled
from $q(x)$

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.

Importance Sampling

$$E_{x \sim p}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N f(x^i)$$

x^i is sampled from $p(x)$

 We only have x^i sampled from $q(x)$

$$= \int f(x)p(x)dx = \int f(x) \frac{p(x)}{q(x)} q(x)dx = E_{x \sim q}[f(x) \frac{p(x)}{q(x)}]$$

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.

$$\nabla \bar{R}_\theta = E_{\tau \sim p_{\theta'}(\tau)} \left[\frac{p_\theta(\tau)}{p_{\theta'}(\tau)} R(\tau) \nabla \log p_\theta(\tau) \right]$$

- Sample the data from θ' .
- Use the data to train θ many times. until the two distributions are quite different.

On-policy → Off-policy

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

- Use π_θ to collect data. When θ is updated, we have to sample training data again.
- Goal: Using the sample from $\pi_{\theta'}$ to train θ . θ' is fixed, so we can re-use the sample data.

$$\nabla \bar{R}_\theta = E_{\tau \sim p_{\theta'}(\tau)} \left[\frac{p_\theta(\tau)}{p_{\theta'}(\tau)} R(\tau) \nabla \log p_\theta(\tau) \right]$$

- Sample the data from θ' .
- Use the data to train θ many times. until the two distributions are quite different.

When the two distributions are quite different, we need to use the NN with updated weights to collect data again.

On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x)\nabla \log f(x)$$

$$= E_{(s_t, a_t) \sim \pi_\theta} [A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n)]$$

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{P_\theta(s_t, a_t)}{P_{\theta'}(s_t, a_t)} A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$= E_{(s_t, a_t) \sim \pi_\theta} [A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n)]$$

$A^{\theta'}(s_t, a_t)$ This term is from
sampled data.

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{P_\theta(s_t, a_t)}{P_{\theta'}(s_t, a_t)} A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$= E_{(s_t, a_t) \sim \pi_\theta} [A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n)]$$

$$A^{\theta'}(s_t, a_t)$$

This term is from
sampled data.

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{P_\theta(s_t, a_t)}{P_{\theta'}(s_t, a_t)} A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_\theta(a_t | s_t)}{p_{\theta'}(a_t | s_t)} \frac{p_\theta(s_t)}{p_{\theta'}(s_t)} A^{\theta'}(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$= E_{(s_t, a_t) \sim \pi_\theta} [A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n)]$$

$$A^{\theta'}(s_t, a_t)$$

This term is from
sampled data.

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{P_\theta(s_t, a_t)}{P_{\theta'}(s_t, a_t)} A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_\theta(a_t | s_t)}{p_{\theta'}(a_t | s_t)} \frac{p_\theta(s_t)}{p_{\theta'}(s_t)} A^{\theta'}(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

On-policy → Off-policy

Gradient for update

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$= E_{(s_t, a_t) \sim \pi_\theta} [A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n)]$$

$$A^{\theta'}(s_t, a_t)$$

This term is from sampled data.

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{P_\theta(s_t, a_t)}{P_{\theta'}(s_t, a_t)} A^\theta(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

$$= E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_\theta(a_t | s_t)}{p_{\theta'}(a_t | s_t)} \frac{p_\theta(s_t)}{p_{\theta'}(s_t)} A^{\theta'}(s_t, a_t) \nabla \log p_\theta(a_t^n | s_t^n) \right]$$

Objective function:

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_\theta(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

Created with EverCam.

How to stop updating the NN's weights when the NN's updated weights are quite different from the old weights (with which the data were sampled)?

Add Constraint

PPO / TRPO

Proximal Policy Optimization (PPO)

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KL(\theta, \theta')$$

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

PPO / TRPO

Proximal Policy Optimization (PPO)

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta \underline{KL(\theta, \theta')}$$

KL divergence of the NN's two output probability vectors

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

KL divergence of two probability vectors

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

PPO / TRPO

Proximal Policy Optimization (PPO)

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KL(\theta, \theta')$$

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$J^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

TRPO (Trust Region Policy Optimization)

$$J_{TRPO}^{\theta'}(\theta) = E_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

$$KL(\theta, \theta') \leq \delta$$

Created with EverCam.
https://www.evercam.com

PPO / TRPO

θ cannot be very different from θ'
Constraint on behavior not parameters

Proximal Policy Optimization (PPO)

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KL(\theta, \theta')$$

$$\nabla f(x) = f(x) \nabla \log f(x)$$

$$J^{\theta'}(\theta) = \underset{\bullet}{E}_{(s_t, a_t) \sim \pi_{\theta'}} \left[\frac{p_{\theta}(a_t | s_t)}{p_{\theta'}(a_t | s_t)} A^{\theta'}(s_t, a_t) \right]$$

PPO algorithm

- Initial policy parameters θ^0
- In each iteration
 - Using θ^k to interact with the environment to collect $\{s_t, a_t\}$ and compute advantage $A^{\theta^k}(s_t, a_t)$
 - Find θ optimizing $J_{PPO}(\theta)$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t)$$

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

• Update parameters several times

PPO algorithm

- Initial policy parameters θ^0
- In each iteration
 - Using θ^k to interact with the environment to collect $\{s_t, a_t\}$ and compute advantage $A^{\theta^k}(s_t, a_t)$
 - Find θ optimizing $J_{PPO}(\theta)$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t)$$

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

Update parameters
several times

- If $KL(\theta, \theta^k) > KL_{max}$, increase β
- If $KL(\theta, \theta^k) < KL_{min}$, decrease β

Adaptive
KL Penalty

PPO algorithm

$$J_{PPO}^{\theta^k}(\theta) = J^{\theta^k}(\theta) - \beta KL(\theta, \theta^k)$$

$$J^{\theta^k}(\theta) \approx \sum_{(s_t, a_t)} \frac{p_\theta(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t)$$

PPO2 algorithm

$$\begin{aligned} J_{PPO2}^{\theta^k}(\theta) \approx & \sum_{(s_t, a_t)} \min \left(\frac{p_\theta(a_t | s_t)}{p_{\theta^k}(a_t | s_t)} A^{\theta^k}(s_t, a_t), \right. \\ & \left. \text{clip} \left(\frac{p_\theta(a_t | s_t)}{p_{\theta^k}(a_t | s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\theta^k}(s_t, a_t) \right) \end{aligned}$$

Q-Learning

Outline

Introduction of Q-Learning



Tips of Q-Learning

Q-Learning for Continuous Actions

Critic

- A critic does not directly determine the action.
- Given an actor π , it evaluates how good the actor is

Critic

- A critic does not directly determine the action.
- Given an actor π , it evaluates how good the actor is
- State value function $V^\pi(s)$
 - When using actor π , the *cumulated* reward expects to be obtained after visiting state s

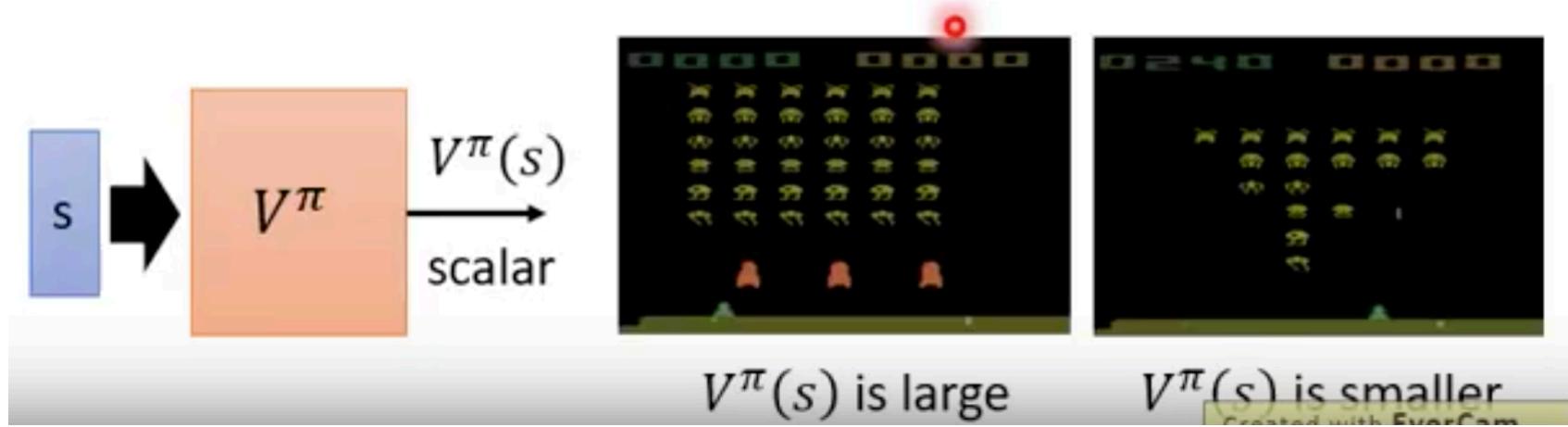
Critic

- A critic does not directly determine the action.
- Given an actor π , it evaluates how good the actor is
- State value function $V^\pi(s)$
 - When using actor π , the *cumulated reward expects to be obtained after visiting state s*



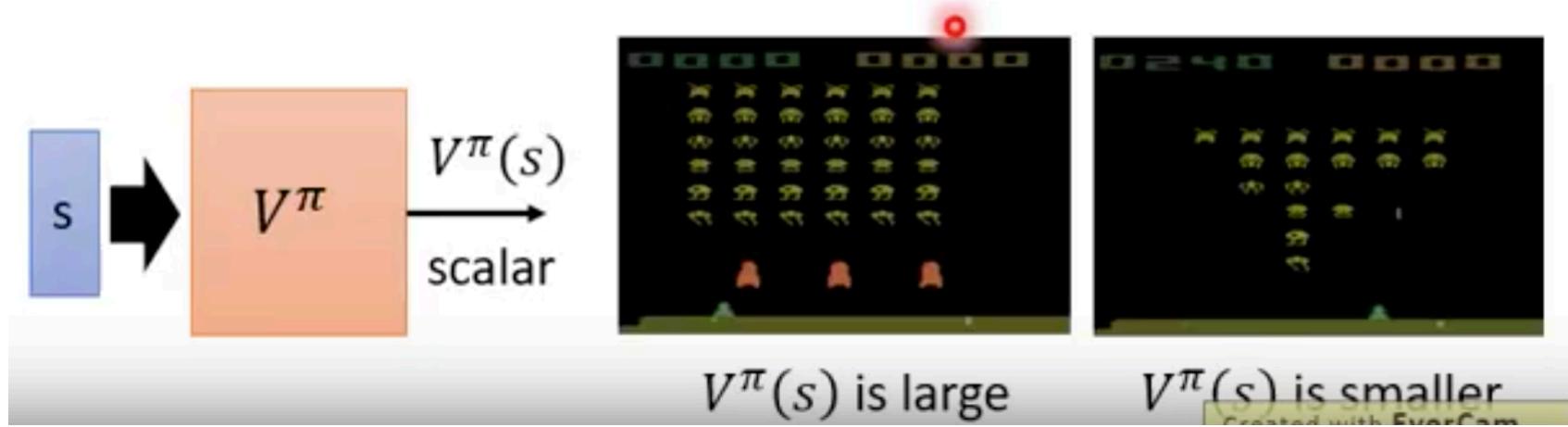
Critic

- A critic does not directly determine the action.
- Given an actor π , it evaluates how good the actor is
- State value function $V^\pi(s)$
 - When using actor π , the *cumulated* reward expects to be obtained after visiting state s



Critic

- A critic does not directly determine the action.
- Given an actor π , it evaluates how good the actor is
- State value function $V^\pi(s)$ **It depends on the actor**
 - When using actor π , the *cumulated* reward expects to be obtained after visiting state s

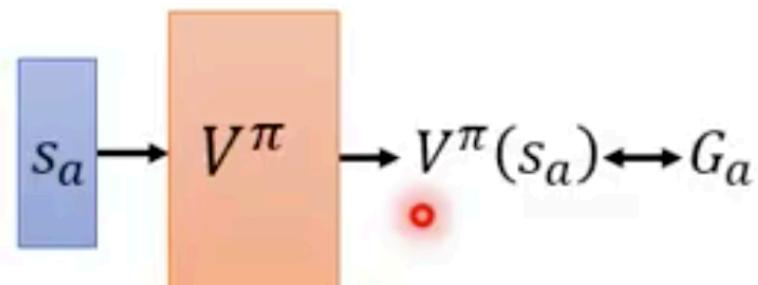


How to estimate $V^\pi(s)$

- **Monte-Carlo (MC) based approach**
 - The critic watches π playing the game

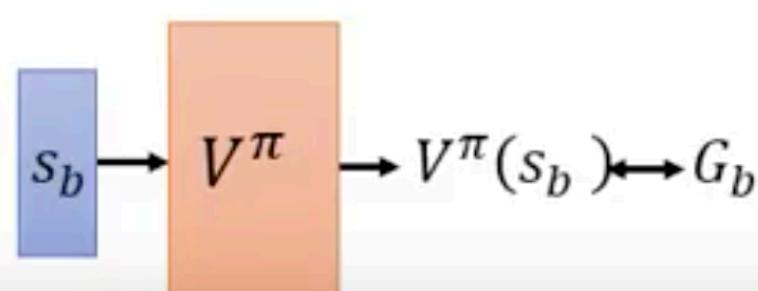
After seeing s_a ,

Until the end of the episode,
the cumulated reward is G_a



After seeing s_b ,

Until the end of the episode,
the cumulated reward is G_b

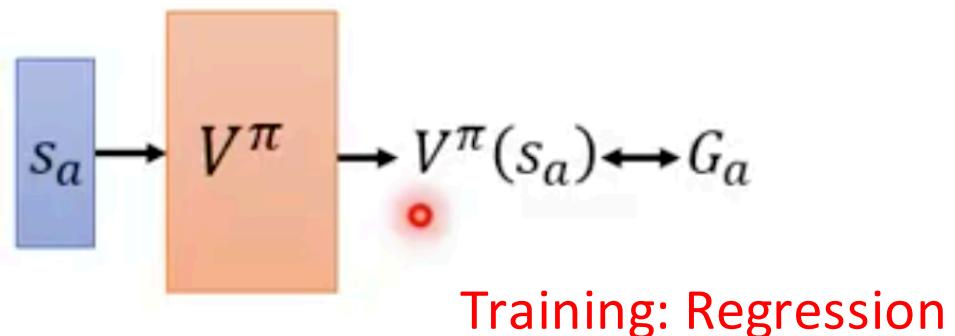


How to estimate $V^\pi(s)$

- **Monte-Carlo (MC) based approach**
 - The critic watches π playing the game

After seeing s_a ,

Until the end of the episode,
the cumulated reward is G_a



After seeing s_b ,

Until the end of the episode,
the cumulated reward is G_b

How to estimate $V^\pi(s)$

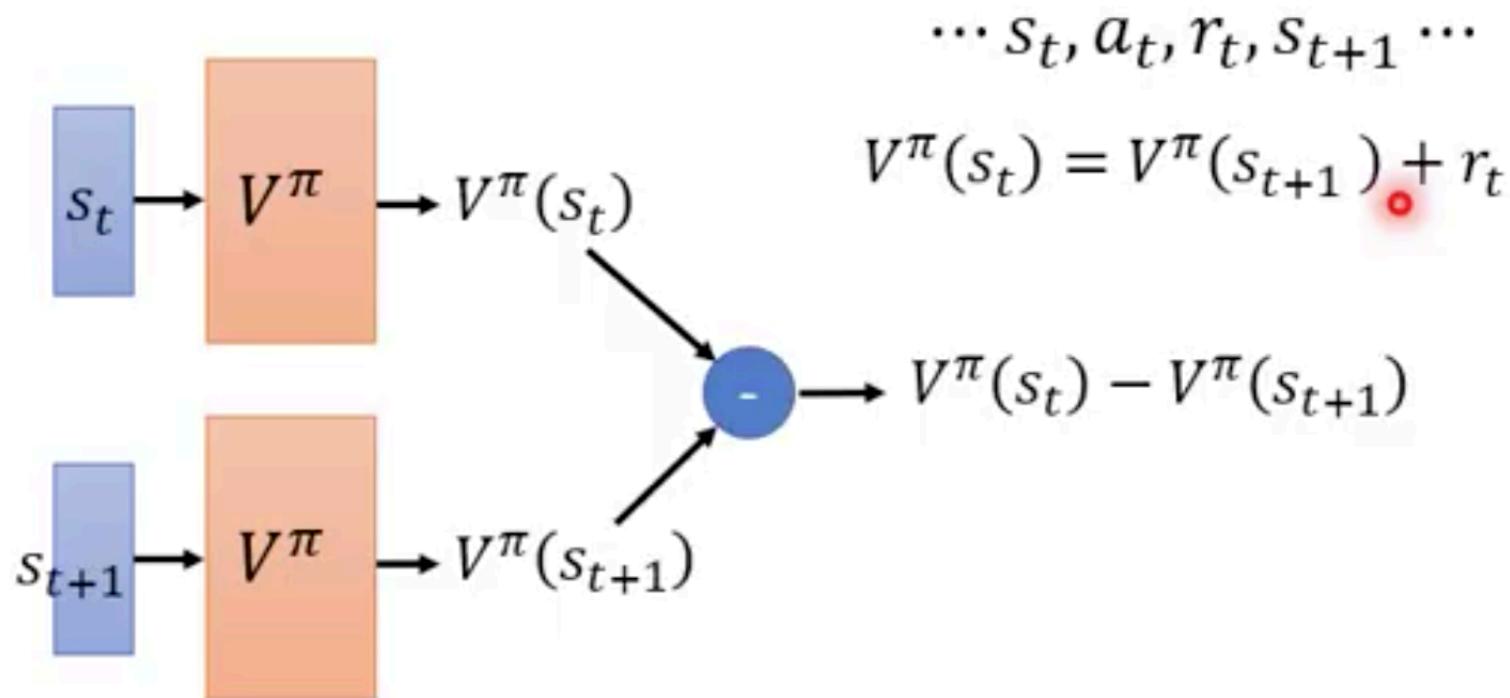
- **Temporal-difference (TD) approach**

$\cdots s_t, a_t, r_t, s_{t+1} \cdots$

$$V^\pi(s_t) = V^\pi(s_{t+1}) + r_t$$

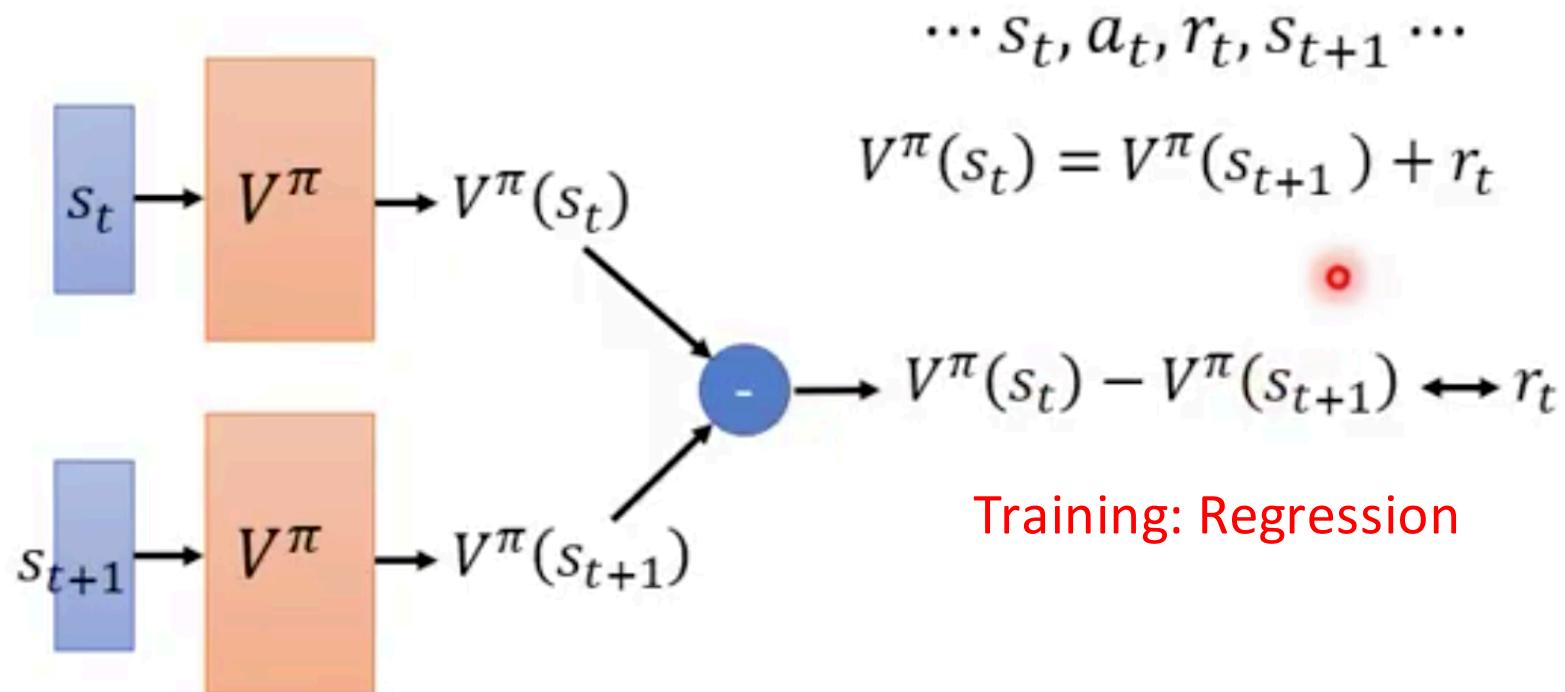
How to estimate $V^\pi(s)$

- **Temporal-difference (TD) approach**



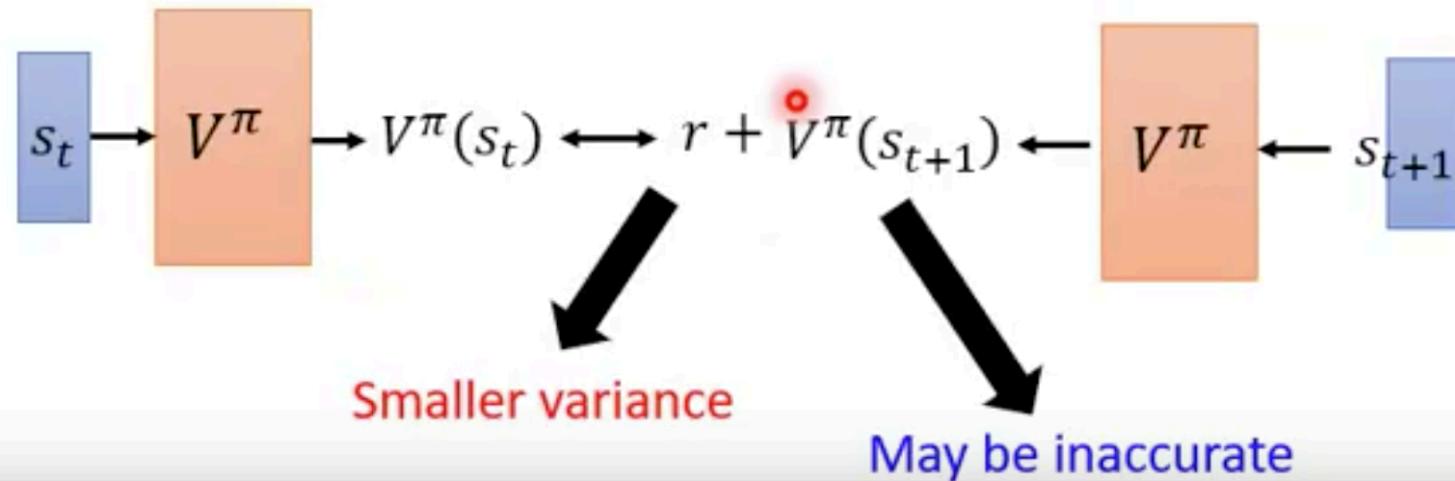
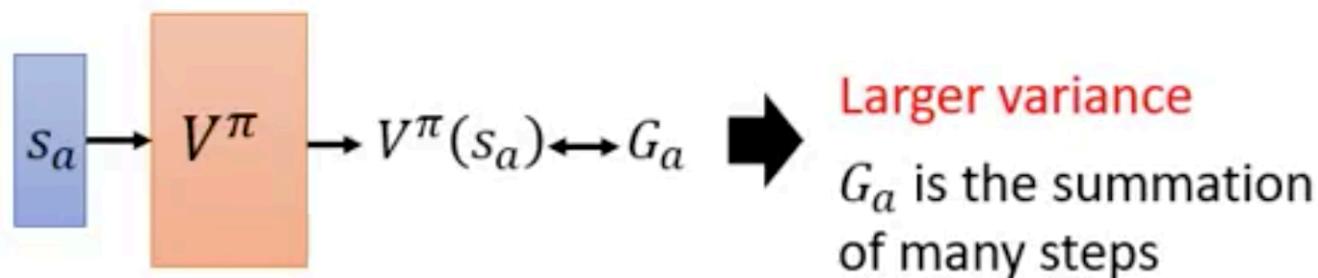
How to estimate $V^\pi(s)$

- Temporal-difference (TD) approach



$$\text{Var}[kX] = k^2 \text{Var}[X]$$

MC v.s. TD

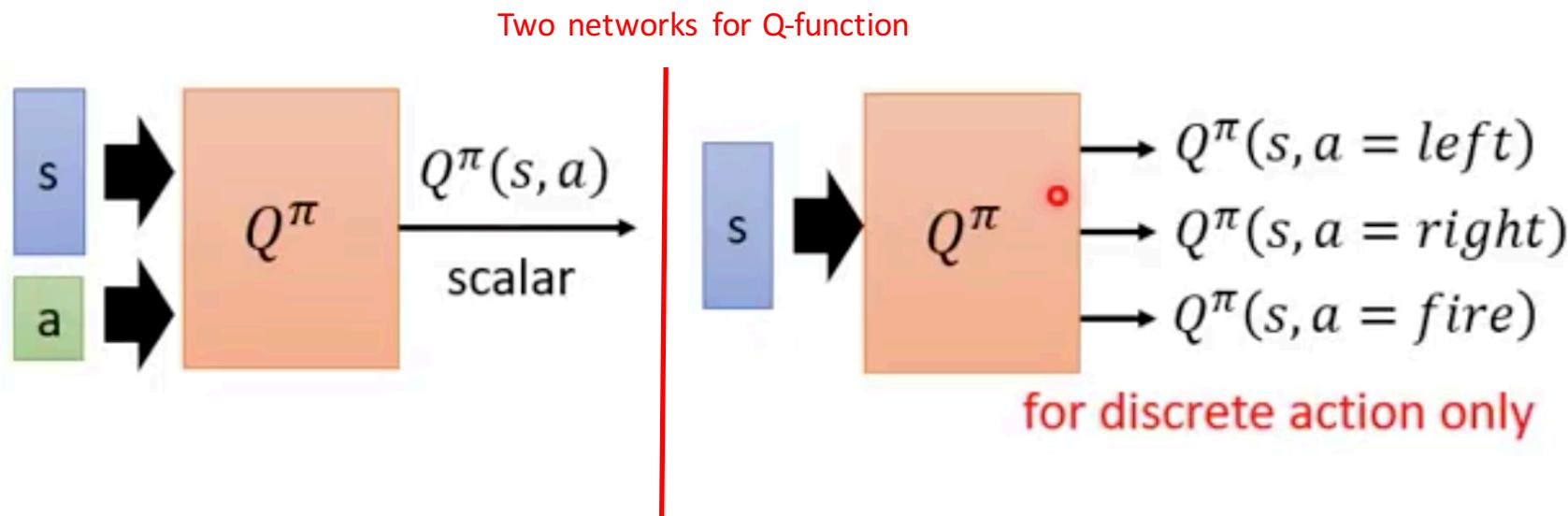


Another Critic

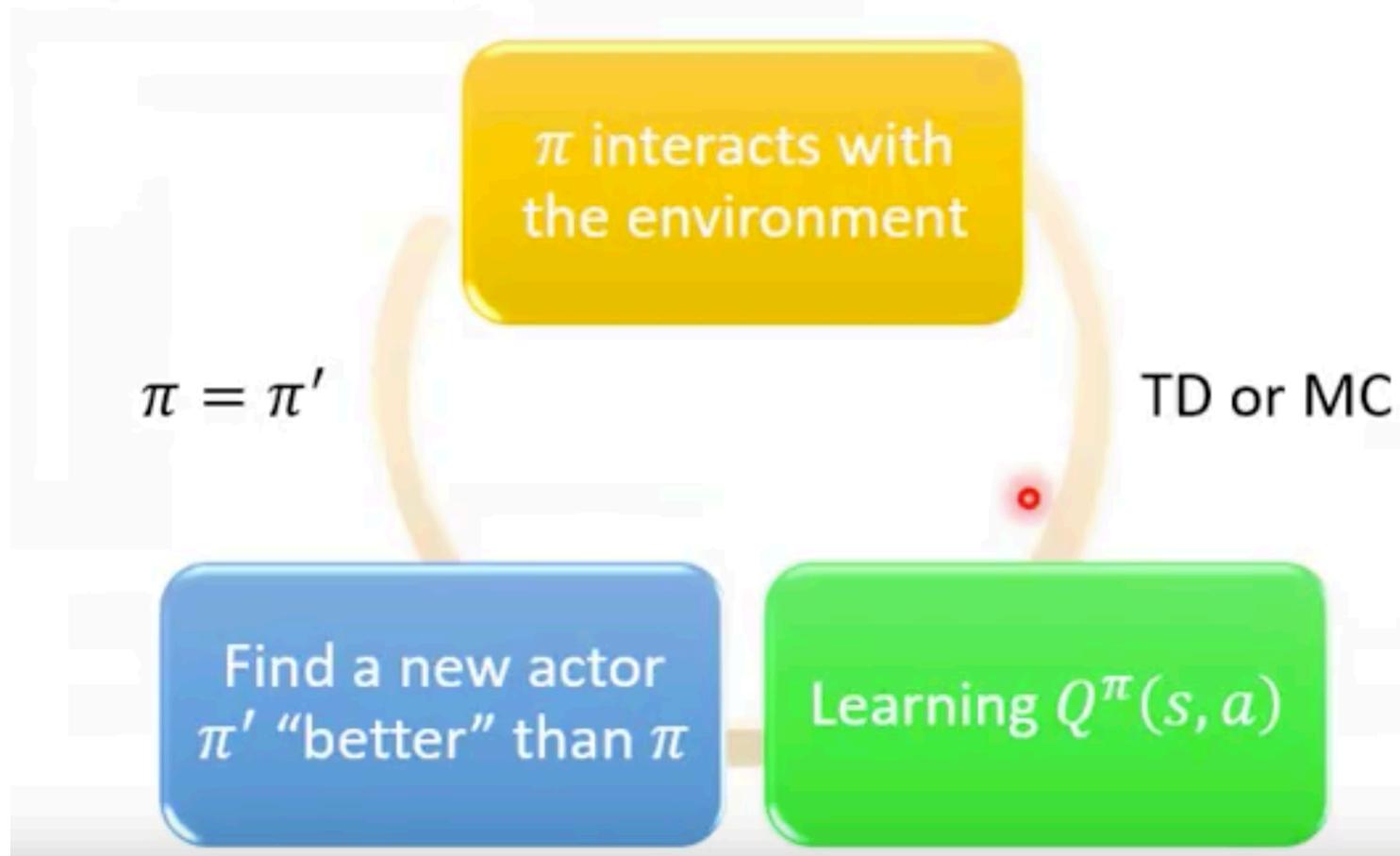
- State-action value function $Q^\pi(s, a)$
 - When using actor π , the *cumulated* reward expects to be obtained after taking **a** at state **s**

Another Critic

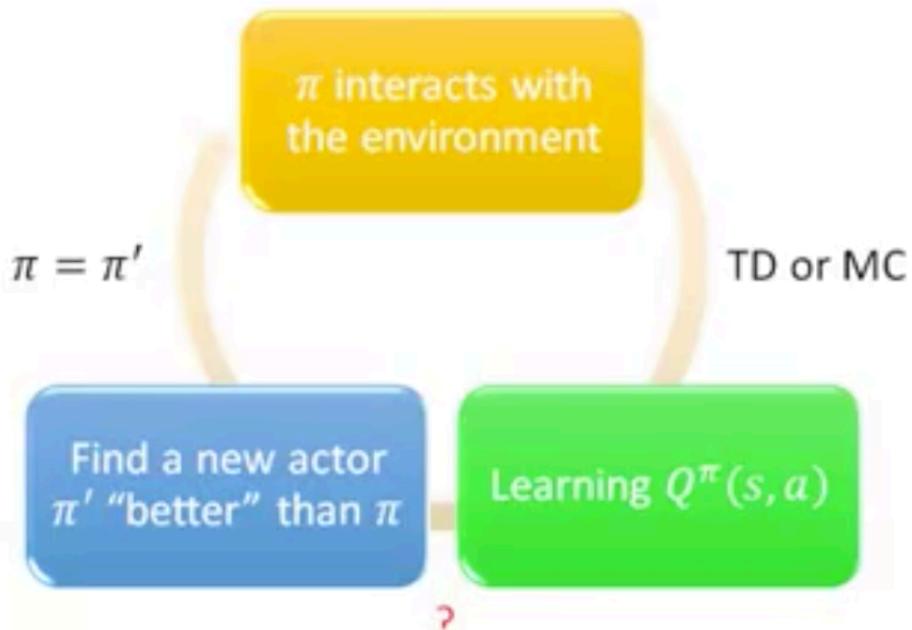
- State-action value function $Q^\pi(s, a)$
 - When using actor π , the *cumulated* reward expects to be obtained after taking a at state s



Another Way to use Critic: Q-Learning



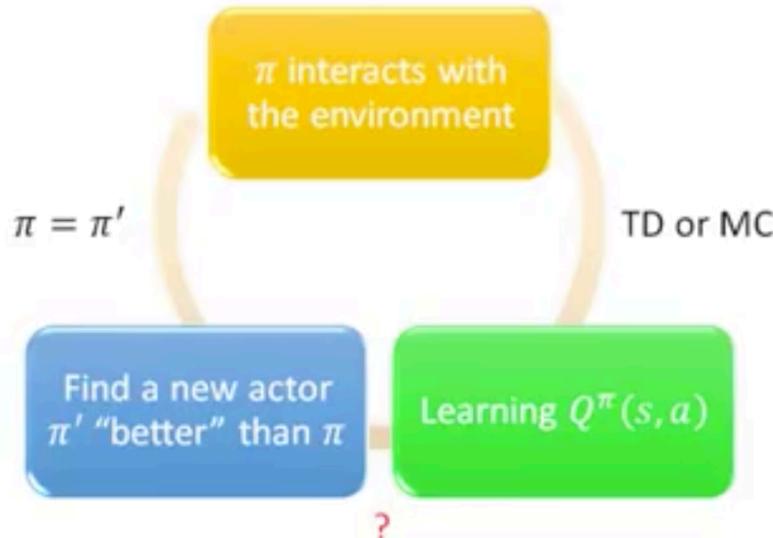
Q-Learning



- Given $Q^\pi(s, a)$, find a new actor π' “better” than π
 - “Better”: $V^{\pi'}(s) \geq V^\pi(s)$, for all state s

$$\pi'(s) = \arg \max_a Q^\pi(s, a)$$

Q-Learning

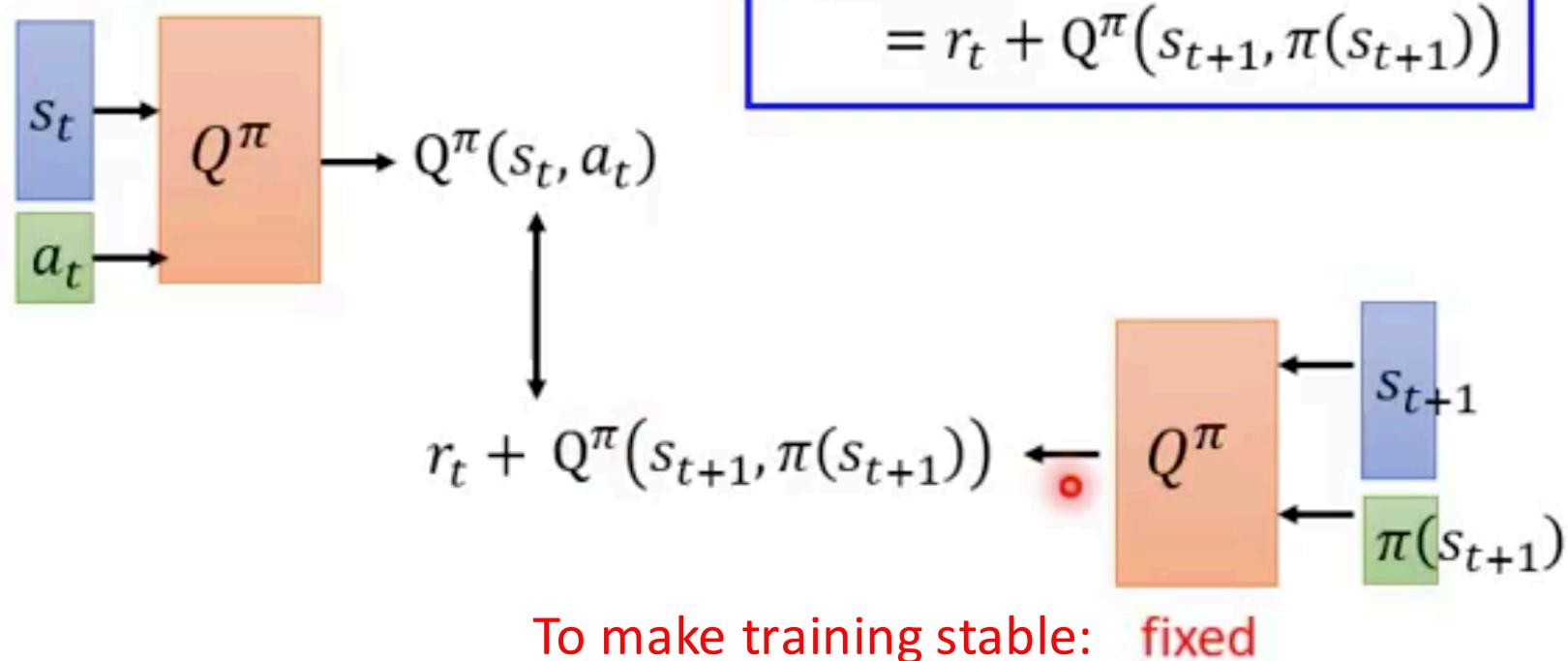


- Given $Q^\pi(s, a)$, find a new actor π' “better” than π
 - “Better”: $V^{\pi'}(s) \geq V^\pi(s)$, for all state s

$$\pi'(s) = \arg \max_a Q^\pi(s, a)$$

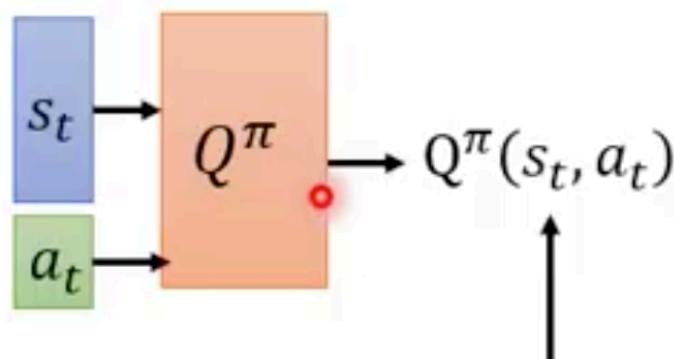
- π' does not have extra parameters. It depends on Q
- Not suitable for continuous action a (solve it later)

Target Network



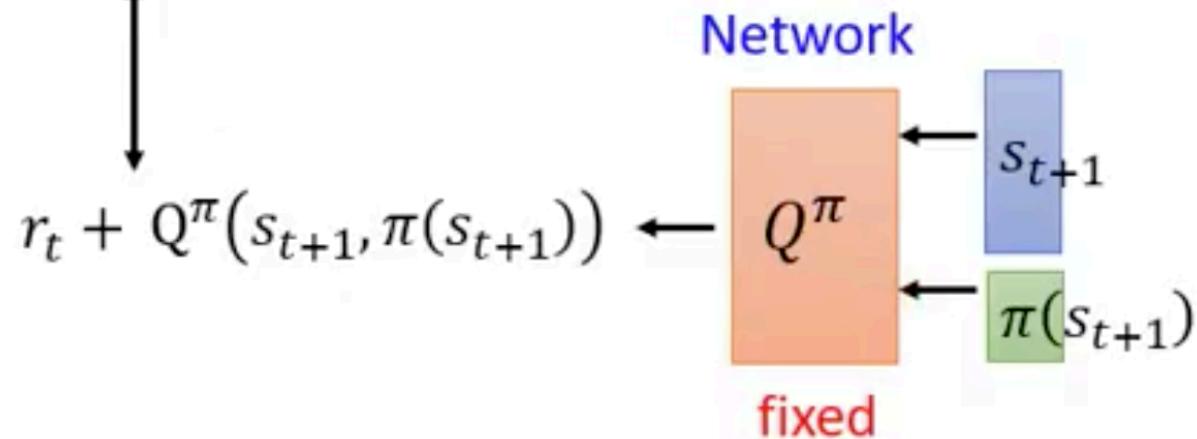
To make training stable: fixed

Target Network



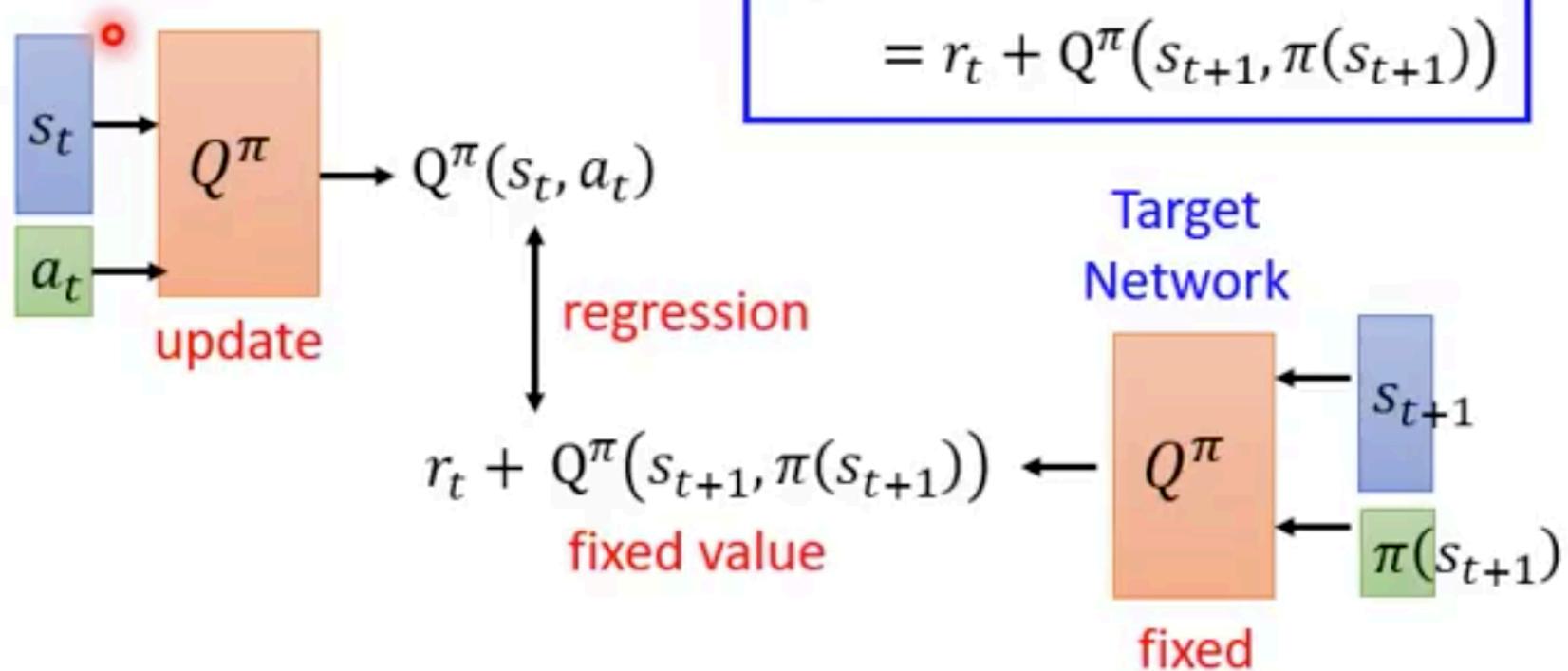
$$\begin{aligned} & \cdots s_t, a_t, r_t, s_{t+1} \cdots \\ & Q^\pi(s_t, a_t) \\ & = r_t + Q^\pi(s_{t+1}, \pi(s_{t+1})) \end{aligned}$$

Target
Network

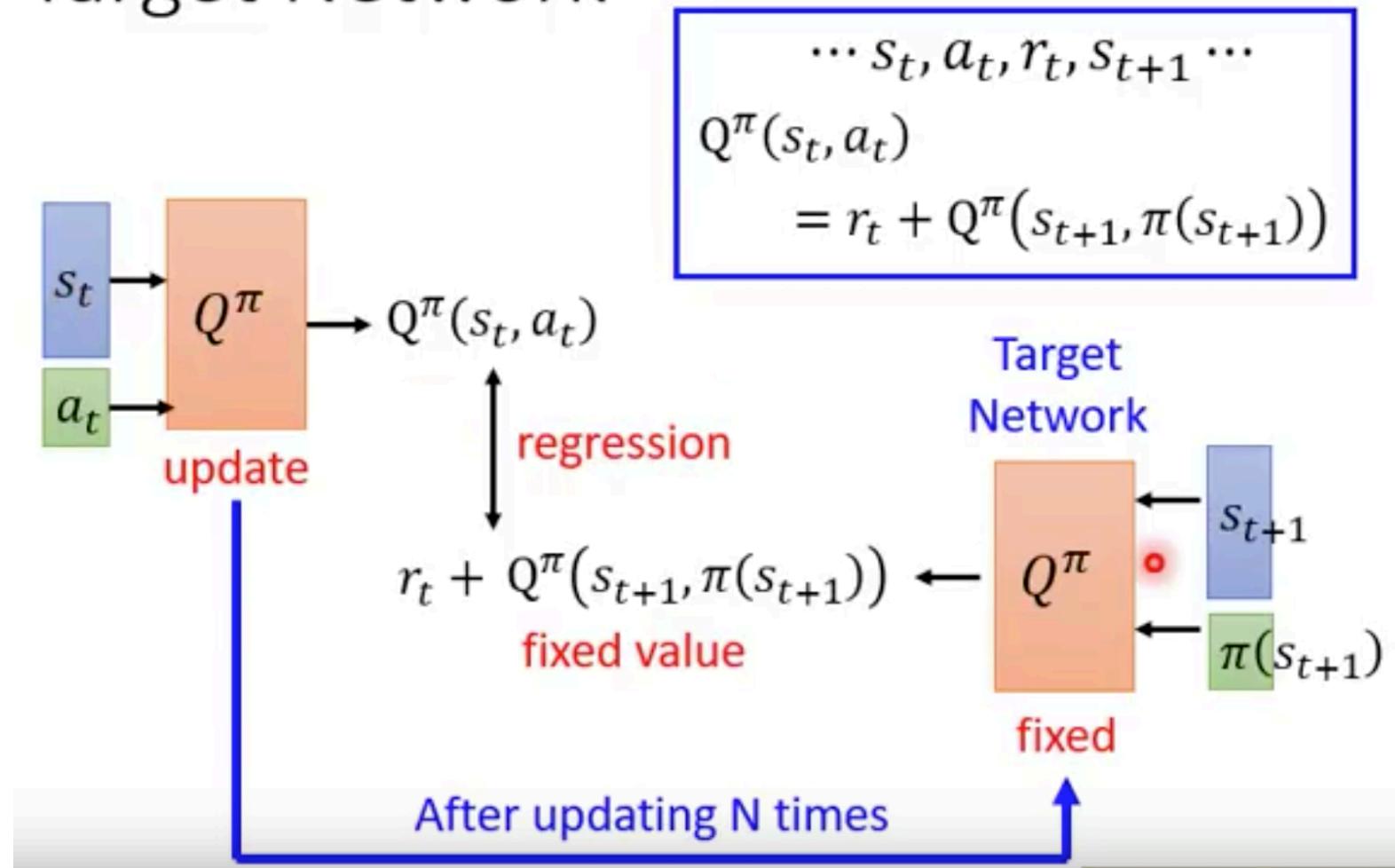


fixed

Target Network



Target Network



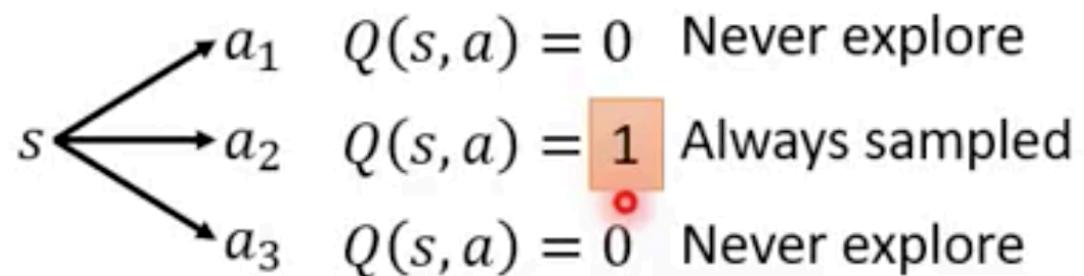
Exploration

- The policy is based on Q-function

$$\pi'(s) = \arg \max_a Q^\pi(s, a)$$

This is not a good way to collect data,
because if an action a is not sampled,
it will not be chosen by the neural network.

Exploration

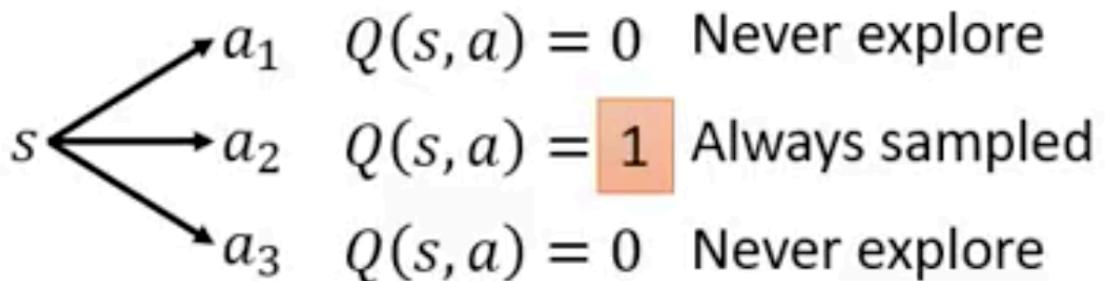


- The policy is based on Q-function

$$\pi'(s) = \arg \max_a Q(s, a)$$

This is not a good way
for data collection.

Exploration



- The policy is based on Q-function

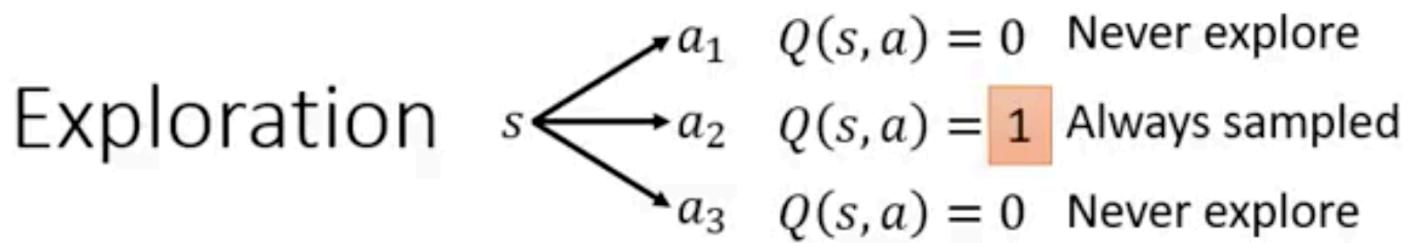
$$a = \arg \max_a Q(s, a)$$

This is not a good way
for data collection.

Epsilon Greedy

ε would decay during learning

$$\pi'(s) = \begin{cases} \arg \max_a Q(s, a), & \text{with probability } 1 - \varepsilon \\ \text{random}, & \text{otherwise} \end{cases}$$



- The policy is based on Q-function

$$a = \arg \max_a Q(s, a)$$

This is not a good way
for data collection.

Epsilon Greedy ε would decay during learning

$$a = \begin{cases} \arg \max_a Q(s, a), & \text{with probability } 1 - \varepsilon \\ \text{random,} & \text{otherwise} \end{cases}$$

Boltzmann Exploration

$$P(a|s) = \frac{\exp(Q(s, a))}{\sum_a \exp(Q(s, a))}$$

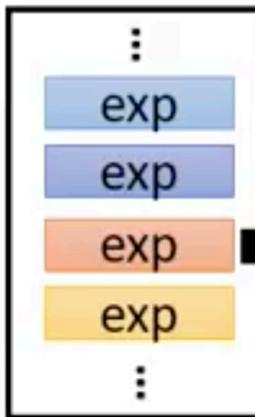
Replay Buffer

Put the experience into buffer.

π interacts with
the environment

$\pi = \pi'$

Buffer



s_t, a_t, r_t, s_{t+1}

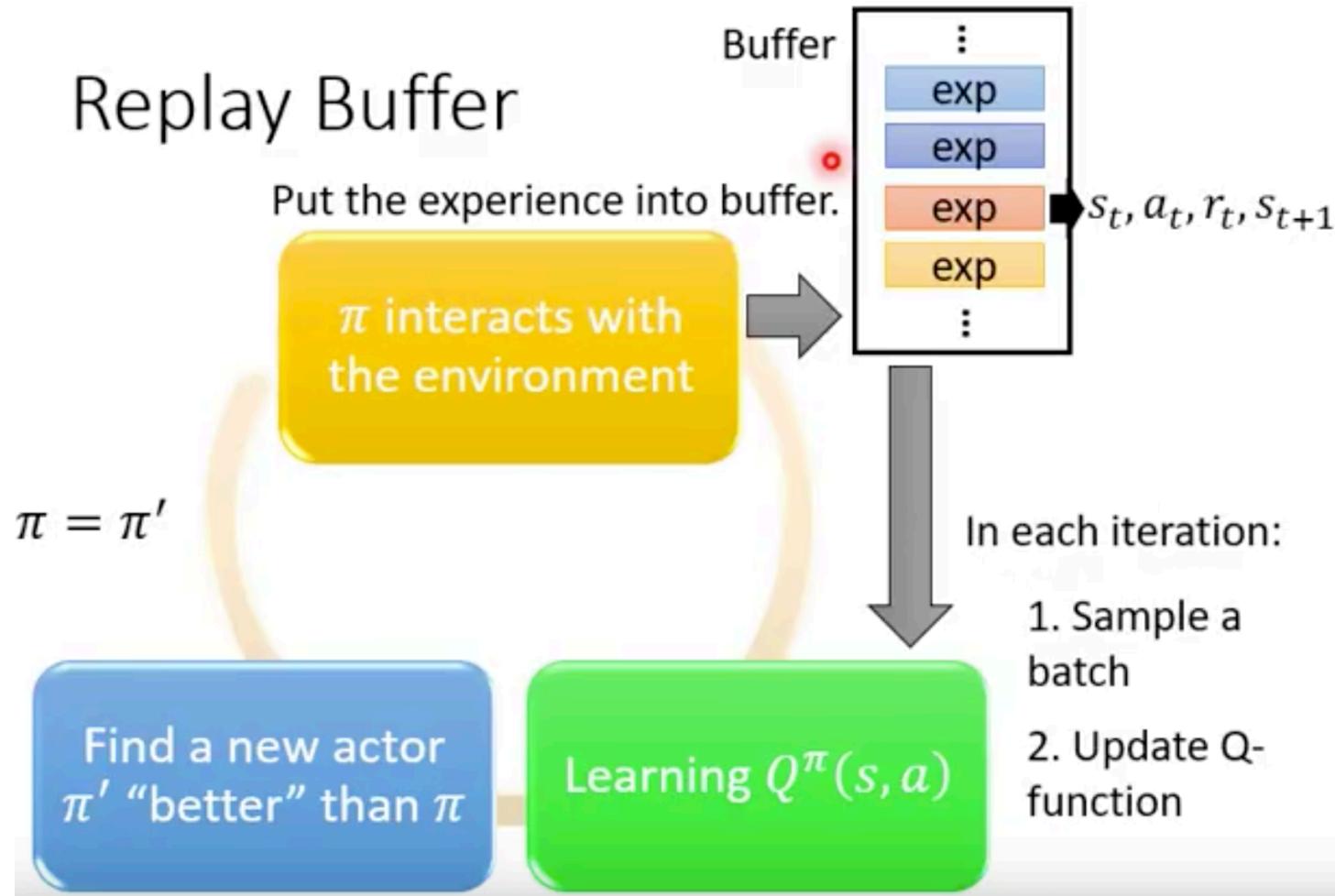
The experience in the
buffer comes from
different policies.

Drop the old experience
if the buffer is full.

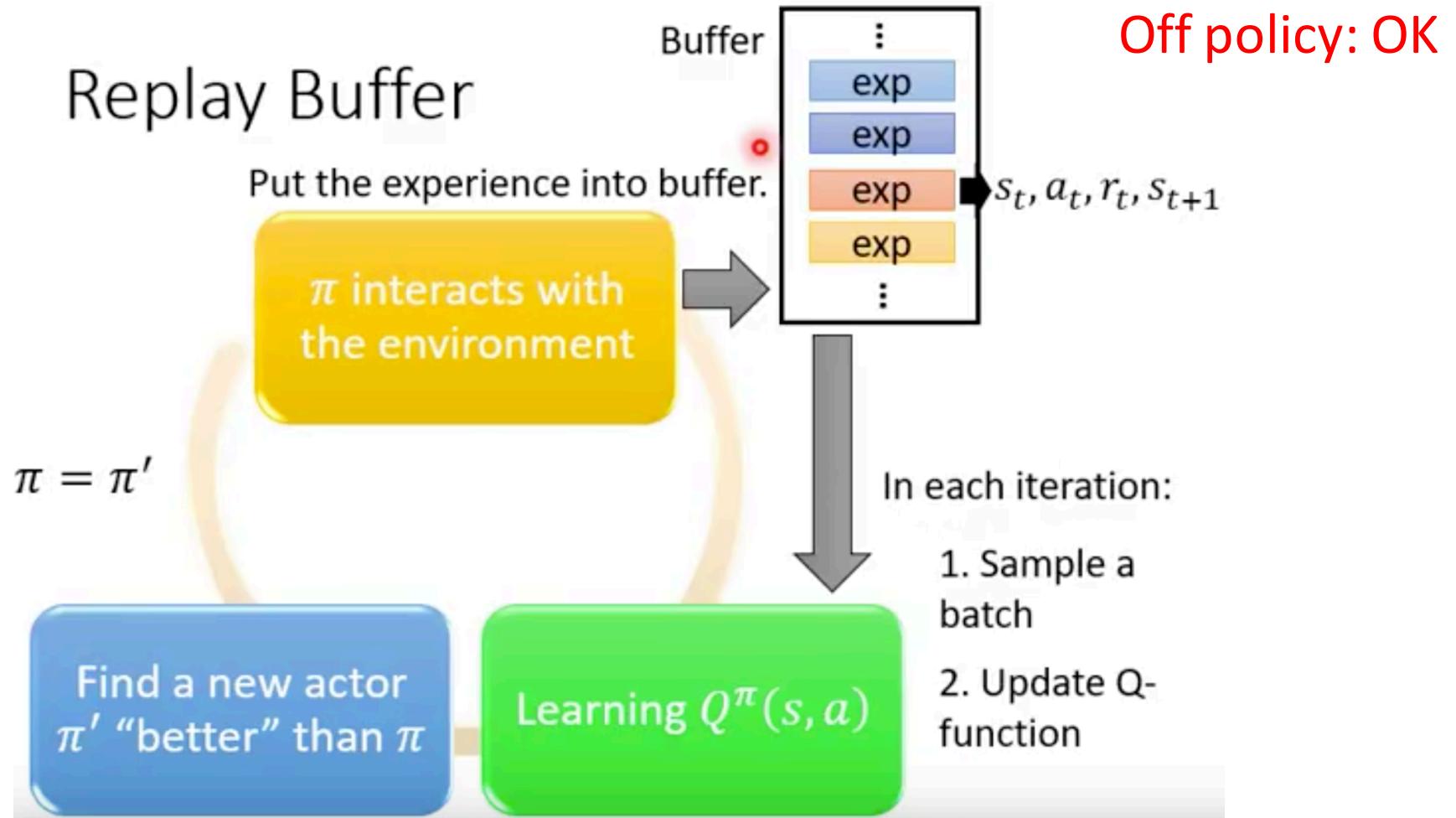
Find a new actor
 π' "better" than π

Learning $Q^\pi(s, a)$

Replay Buffer



Replay Buffer



Typical Q-Learning Algorithm

- Initialize Q-function Q , target Q-function $\hat{Q} = Q$
- In each episode
 - For each time step t
 - Given state s_t , take action a_t based on Q (epsilon greedy)
 - Obtain reward r_t , and reach new state s_{t+1}
 - Store (s_t, a_t, r_t, s_{t+1}) into butter
 - Sample (s_i, a_i, r_i, s_{i+1}) from butter (usually a batch)
 - Target $y = r_i + \max_a \hat{Q}(s_{i+1}, a)$
 - Update the parameters of Q to make $Q(s_i, a_i)$ close to y (regression)
 - Every C steps reset $\hat{Q} = Q$

Additional slide: some early demos of projects

- Music Genre Classification (Joshua Crockett),
<https://www.youtube.com/watch?v=OvO67VXRK7s&feature=youtu.be>
- Stock Price Prediction (Tsao Yuan Chang),
https://www.youtube.com/watch?v=T_9-PG96P7k&feature=youtu.be
- Traffic detection (Jackson Delametter),
https://www.youtube.com/watch?v=bU_lE6iKNzQ&feature=youtu.be
- Hand Gesture Recognition (Syed Ali Hasnain),
<https://www.youtube.com/watch?v=KCxmaGc2U8Q&feature=youtu.be>