

Introduction to Reinforcement Learning

Introduction to Artificial Intelligence - 2023 Summer

Aug 17, 2023
Thu 4 PM

Kwangwoon University MI:RU
Artificial Intelligence Study
Sungmin Yoon



Intro

In this course, you will learn

Part 1 – Quick Review of Probability Theorem

- Joint Probability
 - Joint probability distribution & Marginalization
- Expected Value as Integral
- Chain Rule

Part 2 – Reinforcement Learning

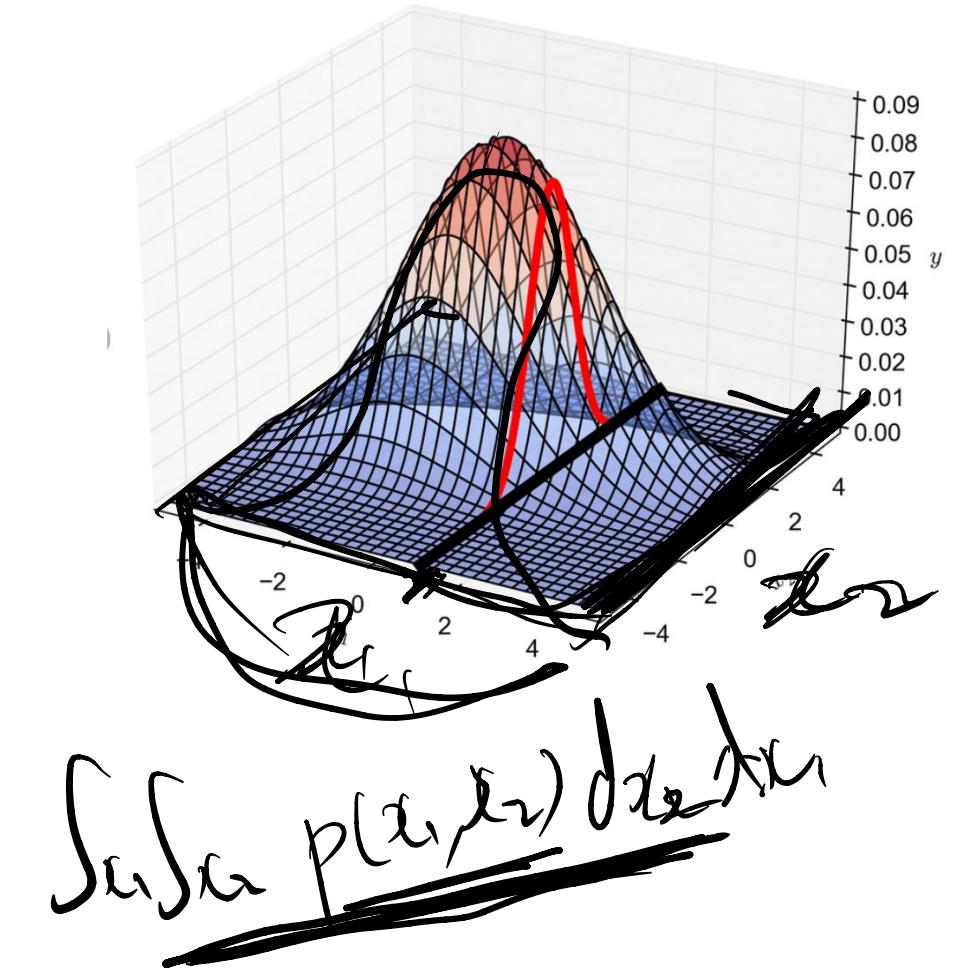
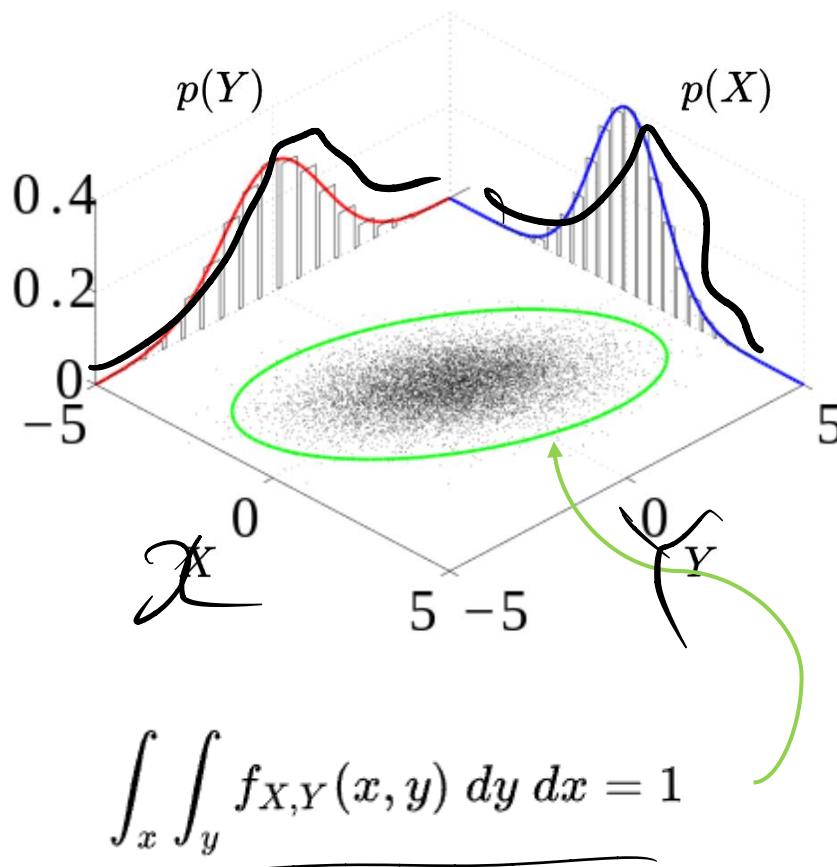
- Basic concepts & Terms
- Markov Property
- Bellman Equation

Part 3 – Policy based RL

- Calculation of Policy
 - Probability of Trajectory
 - Log-derivative trick
 - Expected gradient-log probability Lemma
- Advance of algorithms
 - (REINFORCE, A2C, A3C, PPO)

Probability Theory

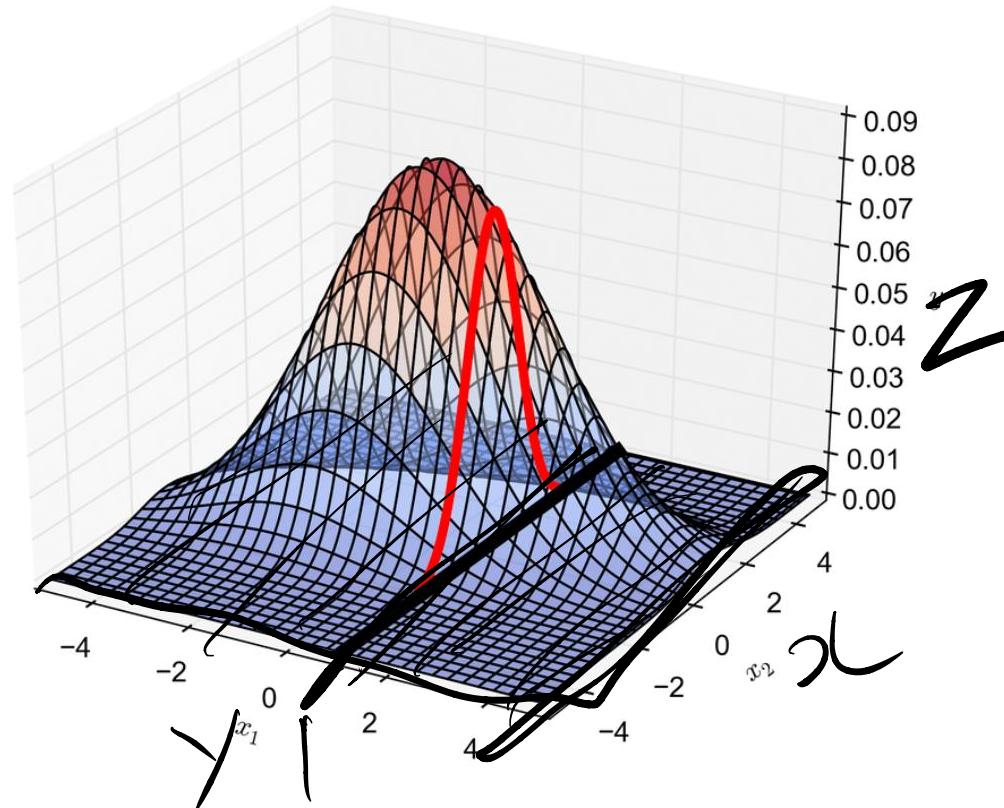
Joint Probability





Probability Theorem

Marginalization



$$\underline{P(X)} = \int_y P(X, Y = y) dy$$

Probability Theory

Expected Value as Integral



$$\int x p(x) dx = 1$$

$$X = [1, 2, 3, 4, 5, 6]$$

$$p(x) = \frac{1}{6} \text{ for each } X$$

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx.$$

$$\int x p(x) dx = E[X]$$

$$E[X] = \int x p(x) dx$$

$$\int p(x) dx = 3.5$$

$\frac{1}{6}(1+2+3+4+5+6)$

Probability Theorem

Chain Rule

Chain rule — Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space. Let $A_1, \dots, A_n \in \mathcal{A}$. Then

$$\mathbb{P}(A_1 \cap A_2 \cap \dots \cap A_n) = \mathbb{P}(A_1)\mathbb{P}(A_2 | A_1)\mathbb{P}(A_3 | A_1 \cap A_2) \cdot \dots \cdot \mathbb{P}(A_n | A_1 \cap \dots \cap A_{n-1})$$

~~$\mathbb{P}(A_1 \cap A_2 \cap \dots \cap A_n)$~~ = $\mathbb{P}(A_1) \prod_{j=2}^n \mathbb{P}(A_j | A_1 \cap \dots \cap A_{j-1}).$

Proof

The formula follows immediately by recursion

$$\begin{aligned} (1) \quad \mathbb{P}(A_1)\mathbb{P}(A_2 | A_1) &= \mathbb{P}(A_1 \cap A_2) \\ (2) \quad \mathbb{P}(A_1)\mathbb{P}(A_2 | A_1)\mathbb{P}(A_3 | A_1 \cap A_2) &= \mathbb{P}(A_1 \cap A_2) \mathbb{P}(A_3 | A_1 \cap A_2) \\ &\quad \text{---} \\ &= \mathbb{P}(A_1 \cap A_2 \cap A_3), \end{aligned}$$

where we used the definition of the conditional probability ~~$\mathbb{P}(A_3 | A_1 \cap A_2)$~~ .

Probability of trajectory

$$\mathbb{P}(s_0, a_0, s_1, \dots, a_T) = \frac{\mathbb{P}(s_0, a_0, s_1, \dots, a_T | s_0)}{\mathbb{P}(s_0) \mathbb{P}(a_0 | s_0) \mathbb{P}(s_1 - a_T | s_0, a_0)}$$
$$= \frac{\mathbb{P}(a_0 | s_0) \mathbb{P}(s_1 - a_T | s_0, a_0)}{\mathbb{P}(s_1 | a_0, s_0) \mathbb{P}(a_1 - a_T | s_1, a_0, s_0)}$$
$$= \frac{\mathbb{P}(a_1 | s_1, a_0, s_0) \mathbb{P}(s_2 - a_T | a_1, s_1, a_0, s_0)}{\mathbb{P}(s_2 | a_1, s_1, a_0, s_0)}$$

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$$

$$\mathbb{P}(A|B) = \mathbb{P}(B) \mathbb{P}(A|B)$$

$$\mathbb{P}(A|B,C) = \frac{\mathbb{P}(A|B,C)}{\mathbb{P}(B,C)}$$

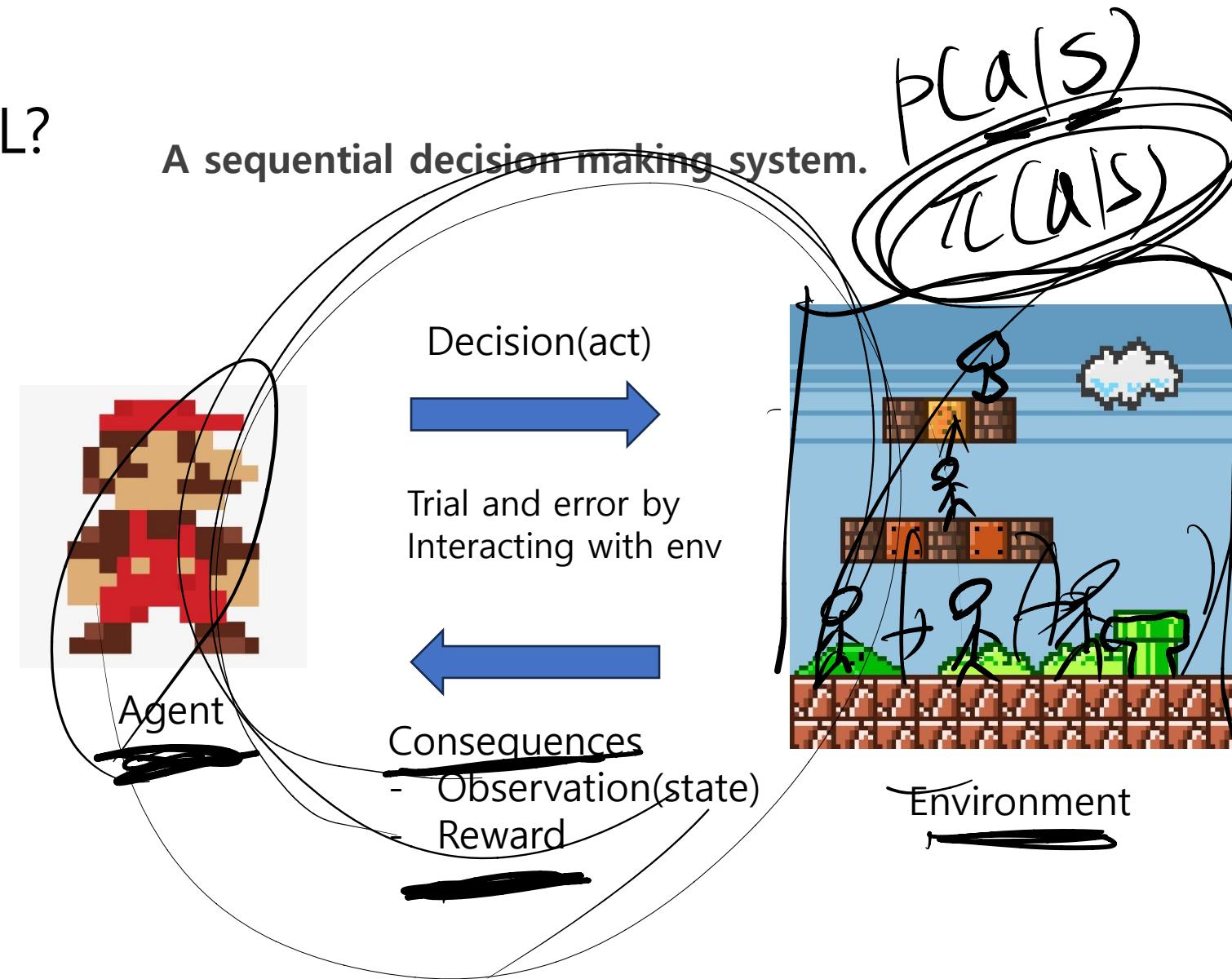
$$\mathbb{P}(A|B,C) = \mathbb{P}(A|B,C) \mathbb{P}(B,C)$$

$$\mathbb{P}(A|B,C,D) = \mathbb{P}(A|B,C,D) \mathbb{P}(B,C,D)$$

Reinforcement Learning

Key concept

What is RL?



Intuitive understanding

Emergence of Locomotion Behaviours in Rich Environments

Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne,
Yuval Tassa, Tom Erez, Ziyu Wang, S. M. Ali Eslami, Martin Riedmiller, David Silver
DeepMind

<https://arxiv.org/abs/1707.02286>
https://www.youtube.com/watch?v=hx_bgoTF7bs

- Example of misspecified rewards
 - [CoastRunner](#)
 - Reward : pass the checkpoint, get a power up item
 - [Escape Room](#)
 - Reward : -1 per second
 - Trained policy : suicide at start
 - [Half cheetah](#) (hard to escape local optima)
 - Reward : speed
 - Trained policy : Backflip and crawl

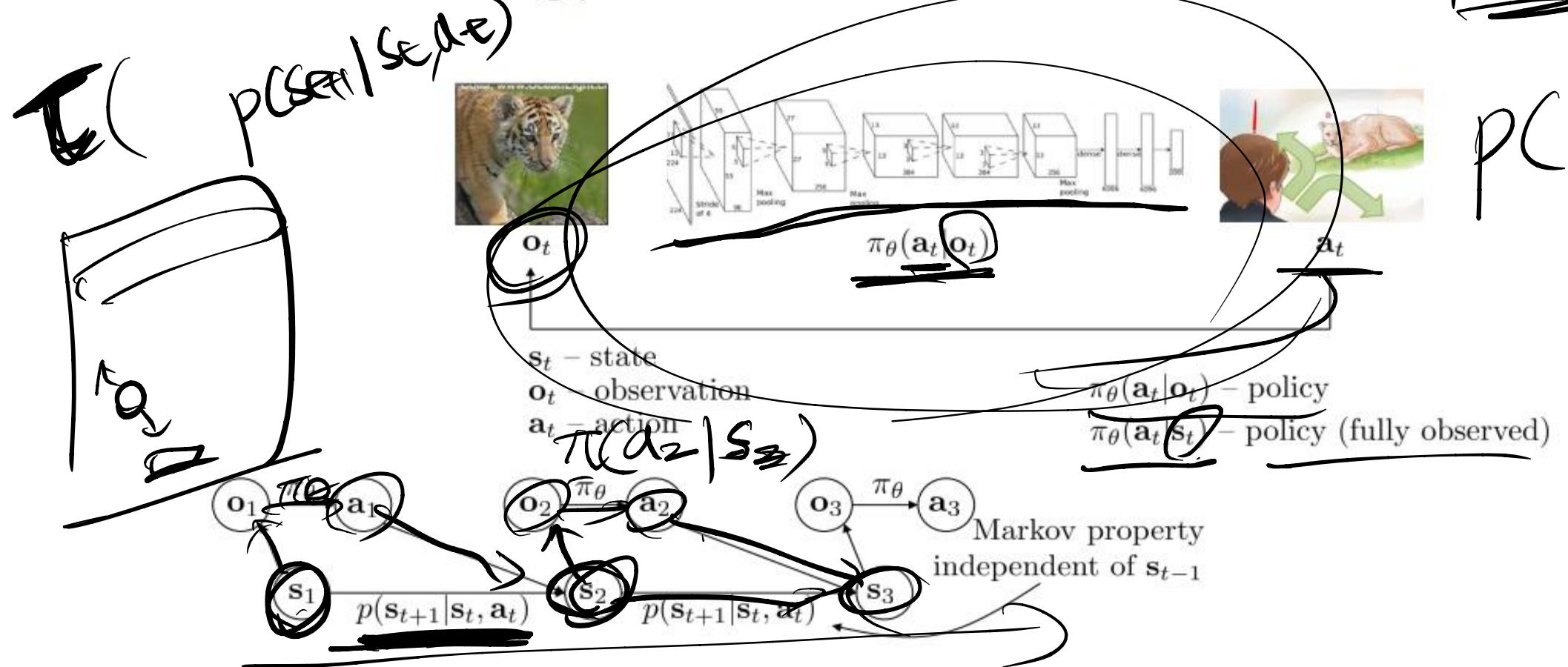


Reinforcement Learning

$p(a|s) \leftarrow p(s'|a)$
 (CalS) environment

$p(a_t|s_t)$

Terminology & notation



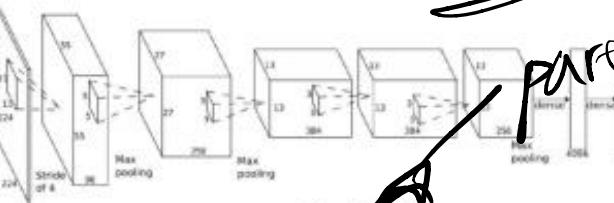
(Slide From CS285 (Sergey Levine) Lecture : <https://rail.eecs.berkeley.edu/deeprlcourse-fa21/>)

Reinforcement Learning

Reward functions



o_t



$\pi_\theta(a_t | s_t)$
 $\pi(a_t | s)$

which action is better or worse?

$r(s, a)$: reward function

tells us which states and actions are better



high reward



low reward

$$G_t = r^0 + \gamma V_{t+1} + \gamma^2 V_{t+2} + \dots + \gamma^\infty V_t$$

$\delta \in [0, 1]$

$R_t < \sum \gamma^t V_{\max}$

~~REINFORCEMENT LEARNING~~

(Slide From CS285 (Sergey Levine) Lecture : <https://rail.eecs.berkeley.edu/deeprlcourse-fa21/>)

Definitions

Markov chain

$$\mathcal{M} = \{\mathcal{S}, \mathcal{T}\}$$

\mathcal{S} – state space

states $s \in \mathcal{S}$ (discrete or continuous)

\mathcal{T} – transition operator

$$p(s_{t+1}|s_t)$$

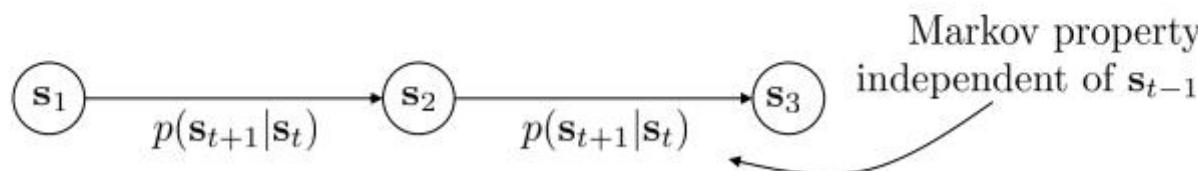
why “operator”?

$$\text{let } \mu_{t,i} = p(s_t = i)$$

$\vec{\mu}_t$ is a vector of probabilities

$$\text{let } \mathcal{T}_{i,j} = p(s_{t+1} = i | s_t = j)$$

$$\text{then } \vec{\mu}_{t+1} = \mathcal{T} \vec{\mu}_t$$



Andrey Markov

Definitions

Markov decision process

\mathcal{S} – state space

\mathcal{A} – action space

P^T

\mathcal{T} transition operator (now a tensor!)

let $\mu_{t,j} = p(s_t = j)$

let $\xi_{t,k} = p(a_t = k)$

let $\mathcal{T}_{i,j,k} = p(s_{t+1} = i | s_t = j, a_t = k)$

$$\mathcal{M} = \{S, A, T, r\}$$

states $s \in \mathcal{S}$ (discrete or continuous)

actions $a \in \mathcal{A}$ (discrete or continuous)

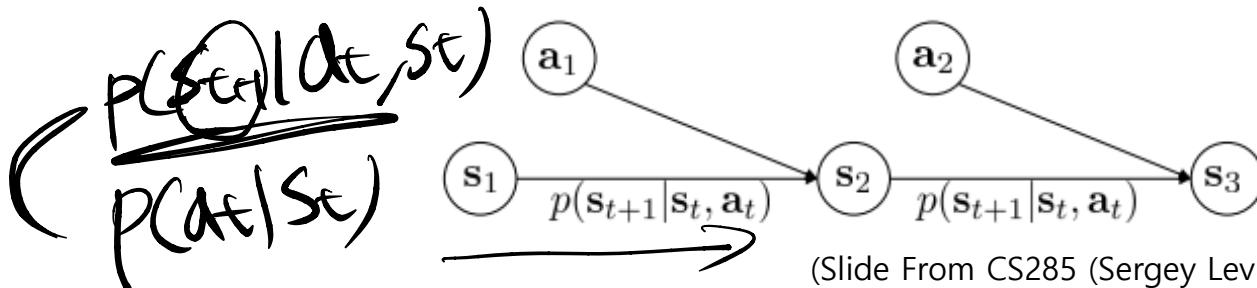
$p(s_{t+1}|a_t, s_t)$

$r_t(s_t, a_t)$



Richard Sutton

$$\mu_{t+1,i} = \sum_{j,k} \mathcal{T}_{i,j,k} \mu_{t,j} \xi_{t,k}$$



(Slide From CS285 (Sergey Levine) Lecture : <https://rail.eecs.berkeley.edu/deeprlcourse-fa21/>)

Definitions



Markov decision process

$$\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$$

\mathcal{S} – state space

states $s \in \mathcal{S}$ (discrete or continuous)

\mathcal{A}

\mathcal{A} – action space

actions $a \in \mathcal{A}$ (discrete or continuous)

\mathcal{T} – transition operator (now a tensor!)

r – reward function

$$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$$

$r(s_t, a_t)$ – reward



Richard Bellman

Bellman Equa.

$V(s_\epsilon)$ $Q(s_\epsilon, a_\epsilon)$

Definitions

partially observed Markov decision process

\mathcal{S} – state space

\mathcal{A} – action space

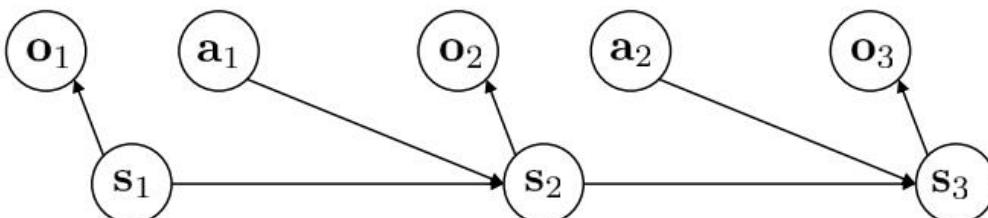
\mathcal{O} – observation space

\mathcal{T} – transition operator (like before)

\mathcal{E} – emission probability $p(o_t|s_t)$

r – reward function

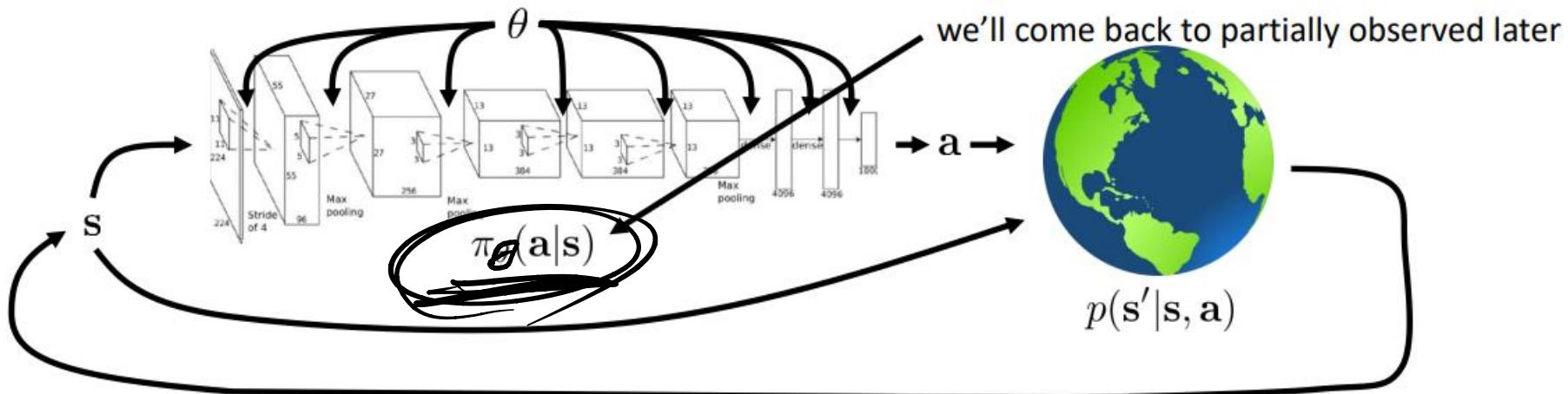
$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$



$$\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{E}, r\}$$



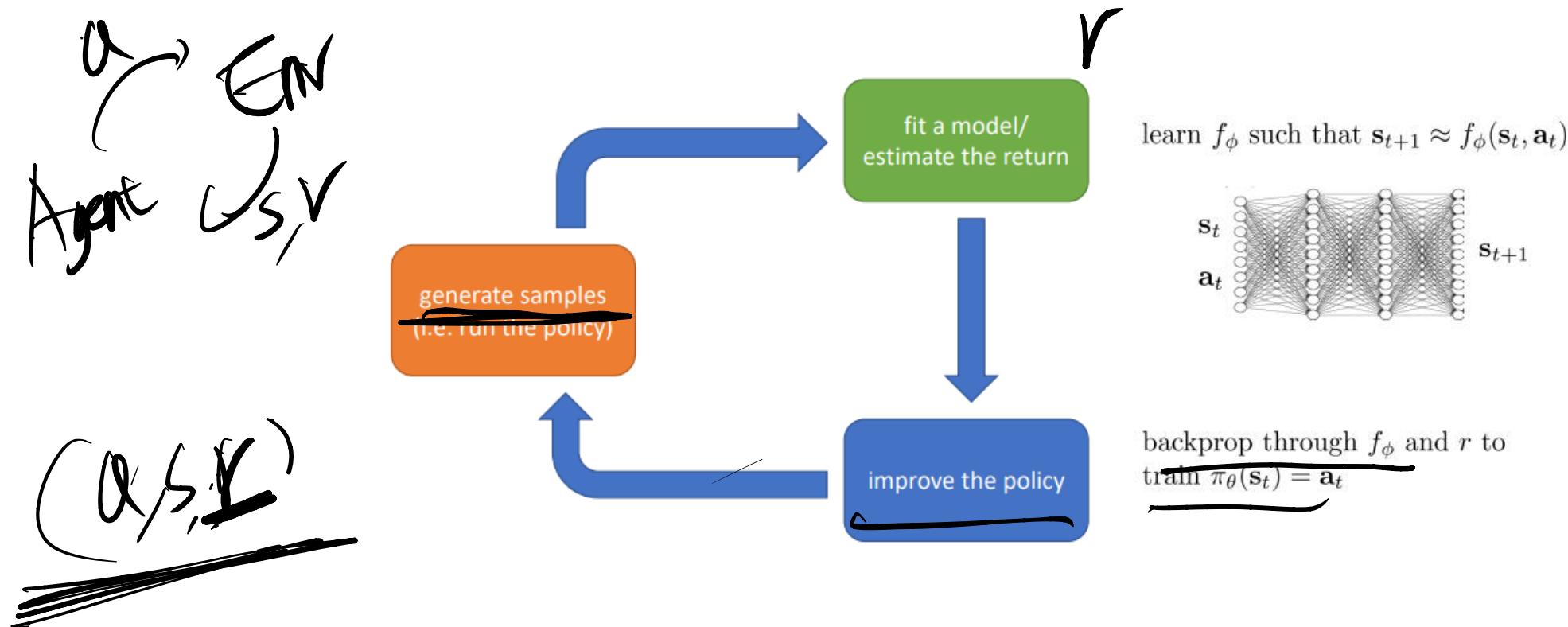
The goal of reinforcement learning



$$p_{\theta}(s_1, a_1, \dots, s_T, a_T) = \underbrace{p(s_1)}_{p_{\theta}(\tau)} \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(s_t, a_t) \right]$$

The anatomy of a reinforcement learning algorithm



~~value-based RL~~

Definition: Q-function = State-action function
Value

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_{\pi_\theta}[r(s_{t'}, a_{t'})|s_t, a_t]: \text{total reward from taking } a_t \text{ in } s_t$$

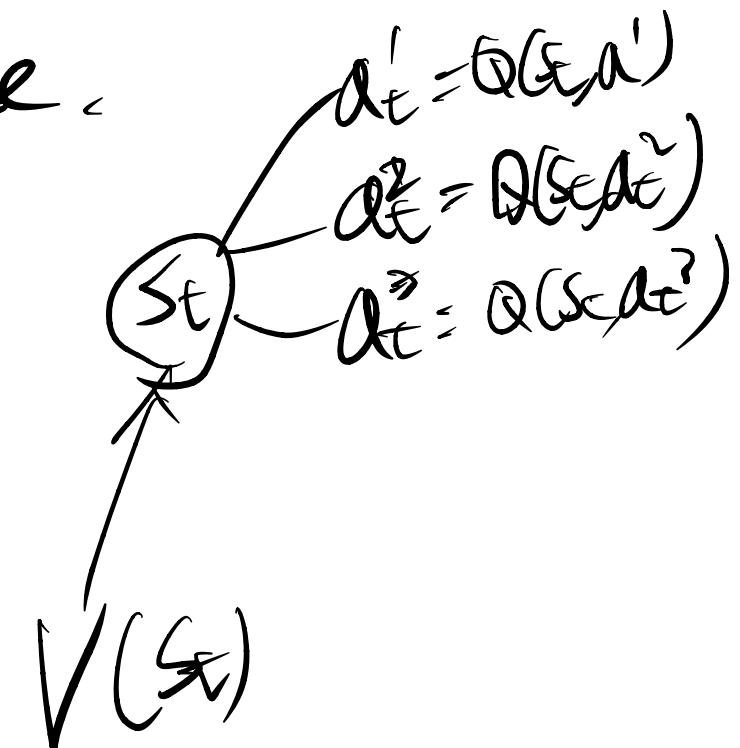
~~Definition: value function~~ State Value

$$V^\pi(s_t) = E_{a_t \sim \pi(a_t|s_t)}[Q^\pi(s_t, a_t)]$$

$E_{s_1 \sim p(s_1)}[V^\pi(s_1)]$ is the RL objective!

$$Q^\pi(s_t, a_t) = V_t + V_{t+1}^\pi(s_{t+1}, a_{t+1})$$

$$V^\pi_t = E_{a_t}[Q(a_t, s_t)]$$



Reinforcement Learning

Value Iteration : State-Value function V 를 반복적으론 계산 \Rightarrow 정책 도출

Bellman - Equation 을 활용해 V 업데이트 (DP)

Policy Iteration : 정책 평가와 정책 개선

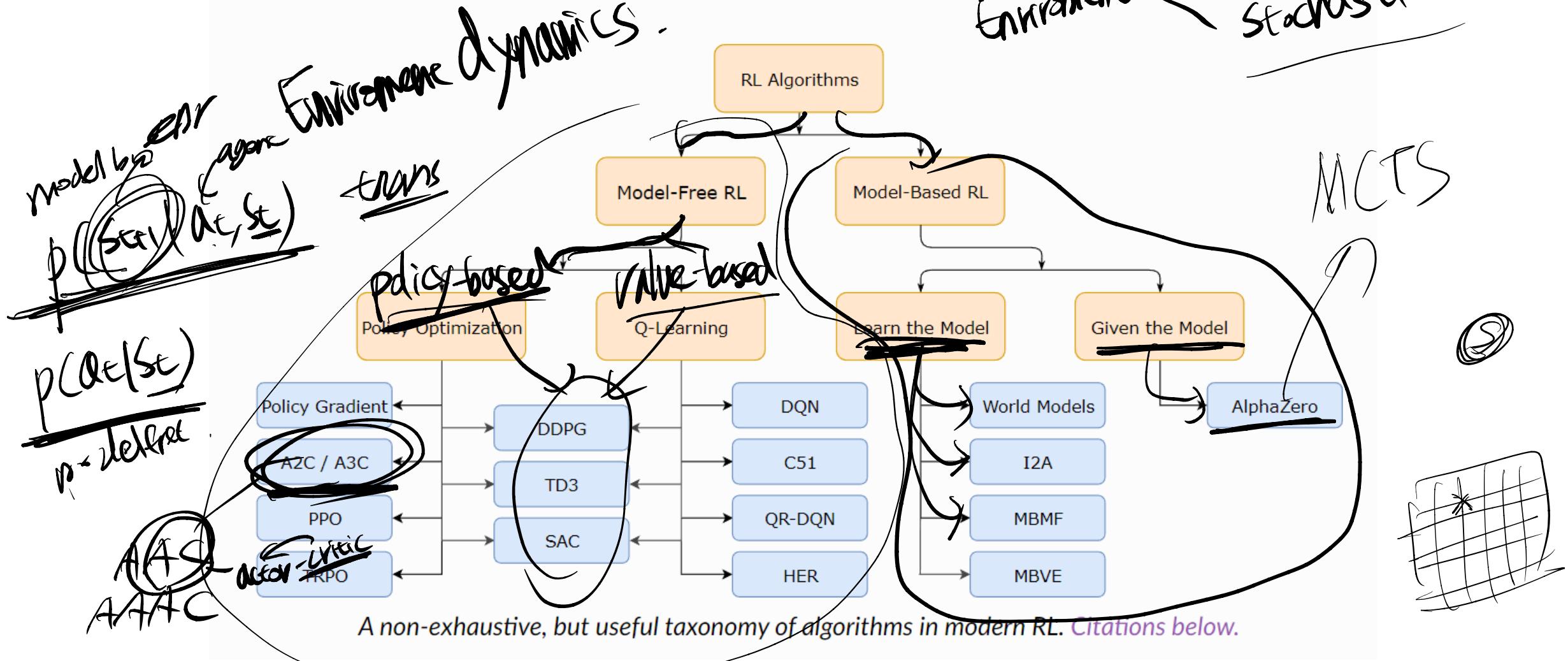
정책 평가 : 현재 정책에 따른 Value Function

정책 개선 : 기존 V 를 바탕으로 정책 update.

\Rightarrow Grid World example

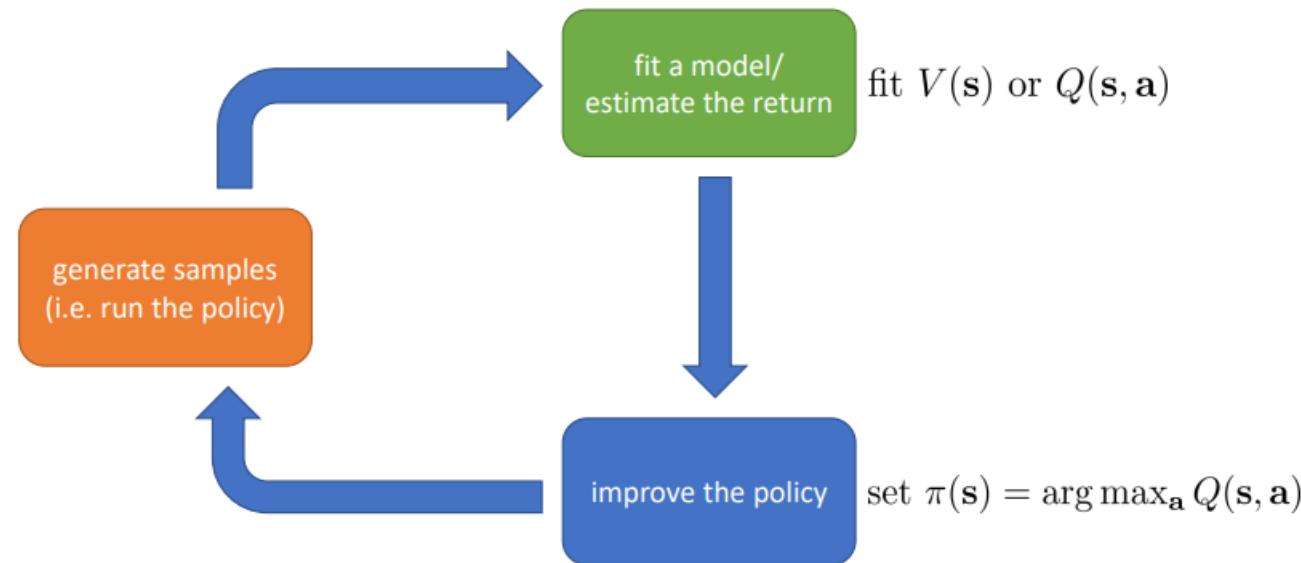
Reinforcement Learning

A Taxonomy of RL Algorithms



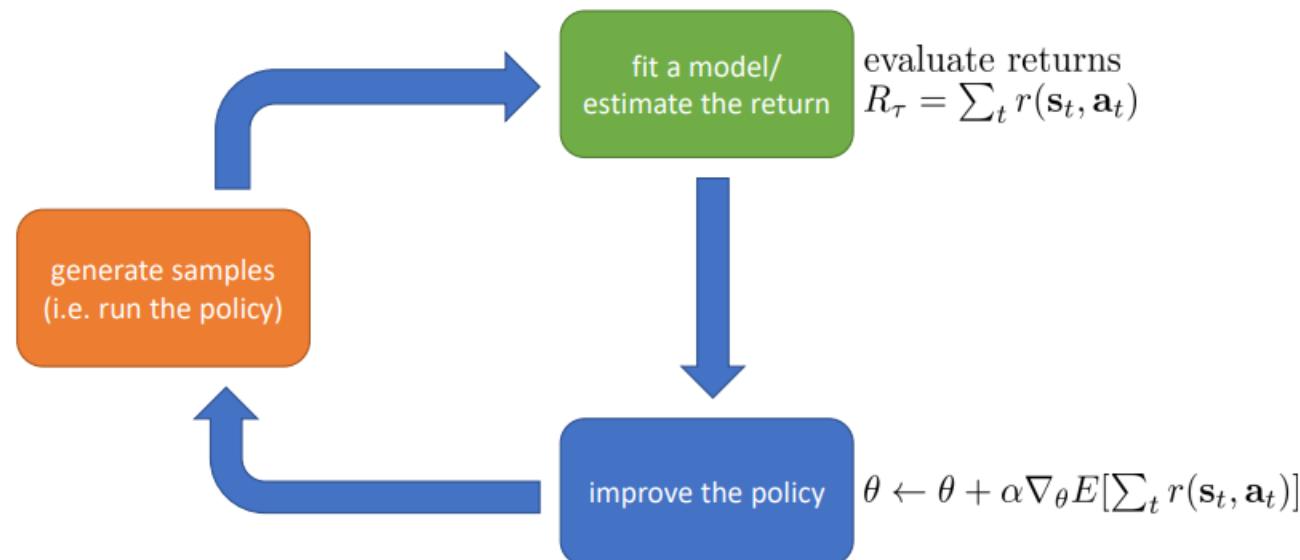


Value function based algorithms

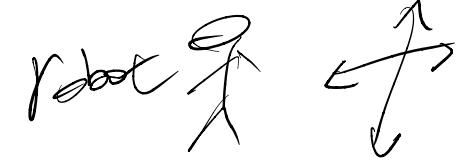




Direct policy gradients



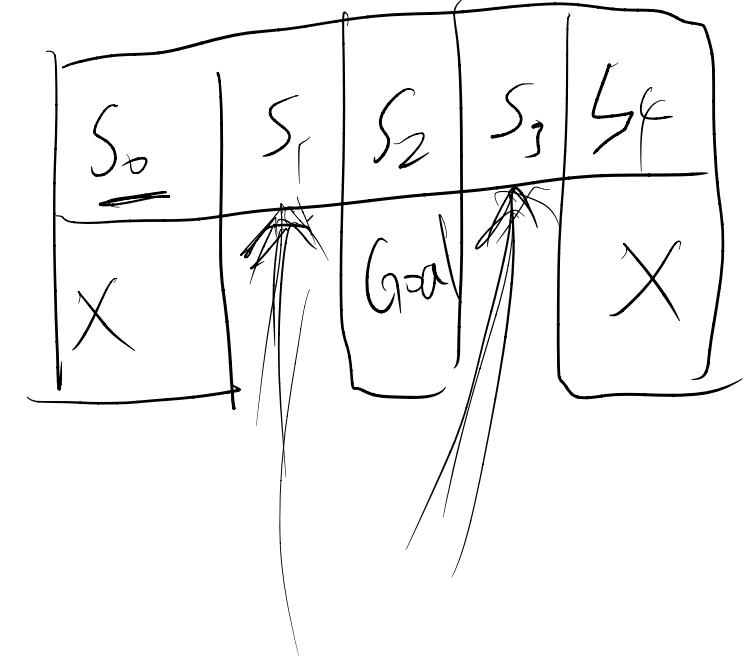
Reinforcement Learning



Value based VS Policy based

Continuous action

Stochastic policy
(especially in partially observable)



on Deterministic

S0 S1

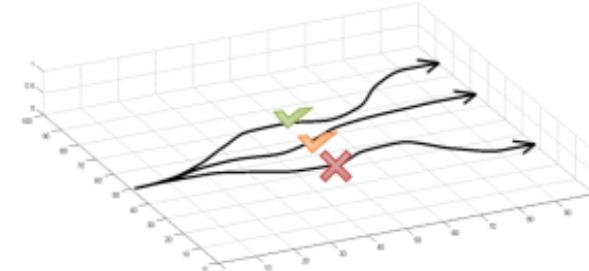
Goal state value



Evaluating the objective

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$J(\theta)$



$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

↑
sum over samples from π_{θ}

Direct policy differentiation

$$\theta^* = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[\underbrace{\sum_t r(\mathbf{s}_t, \mathbf{a}_t)}_{J(\theta)} \right]$$

a convenient identity

$$\underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} = p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} = \underline{\nabla_{\theta} p_{\theta}(\tau)}$$

$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau$$

$$\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)$$

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

$$\nabla_{\theta} J(\theta) = \int \underline{\nabla_{\theta} p_{\theta}(\tau)} r(\tau) d\tau = \int \underline{p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)} r(\tau) d\tau = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

$$\nabla_{\theta} p_{\theta}(c) = \frac{\nabla_{\theta} p_{\theta}(c)}{p_{\theta}(c)} p_{\theta}(c) = p_{\theta}(c) \log p_{\theta}(c)$$



Direct policy differentiation

$$\theta^* = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim p_\theta(\tau)}[r(\tau)]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_\theta(\tau)}[\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$$

$$\nabla_{\theta} \left[\cancel{\log p(s_1)} + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \cancel{\log p(s_{t+1} | s_t, a_t)} \right]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_\theta(\tau)} \left[\underbrace{\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right)}_{\text{Don't let the past distract you}} \right]$$

chain rule.

$$p_{\theta}(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod_{t=1}^T \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

log of both sides

$$\log p_{\theta}(\tau) = \log p(s_1) + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)$$

ev'g y'ur \rightarrow plz,



Evaluating the policy gradient

$$\text{recall: } J(\theta) = E_{\tau \sim p_\theta(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

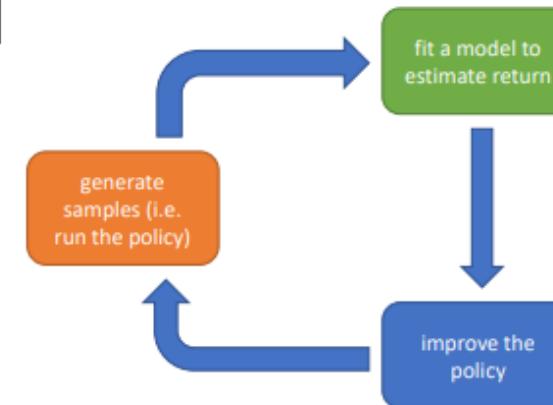
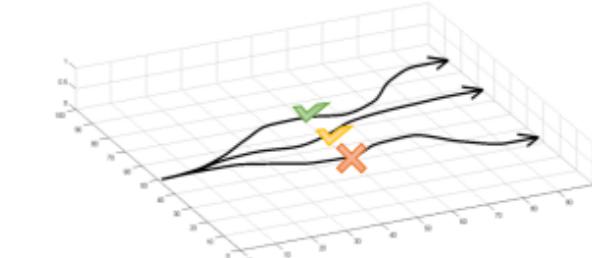
$$\nabla_\theta J(\theta) = E_{\tau \sim p_\theta(\tau)} \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

REINFORCE algorithm:

- 1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
- 2. $\nabla_\theta J(\theta) \approx \sum_i (\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
- 3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$





Reinforcement Learning

Example: Gaussian policies

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

example: $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = \mathcal{N}(f_{\text{neural network}}(\mathbf{s}_t); \Sigma)$

$$\log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = -\frac{1}{2} \|\mathbf{f}(\mathbf{s}_t) - \mathbf{a}_t\|_{\Sigma}^2 + \text{const}$$

$$\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) = -\frac{1}{2} \Sigma^{-1} (\mathbf{f}(\mathbf{s}_t) - \mathbf{a}_t) \frac{d\mathbf{f}}{d\theta}$$

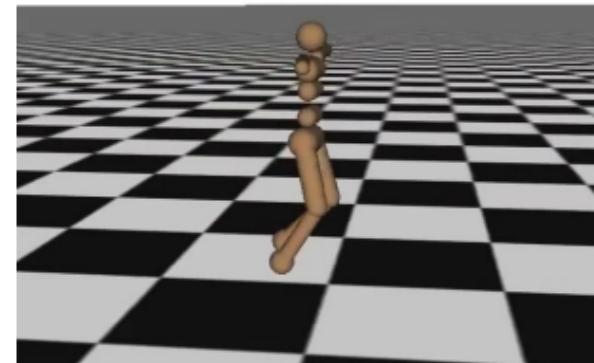
REINFORCE algorithm:

- 1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
- 2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
- 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

Gaussian distribution

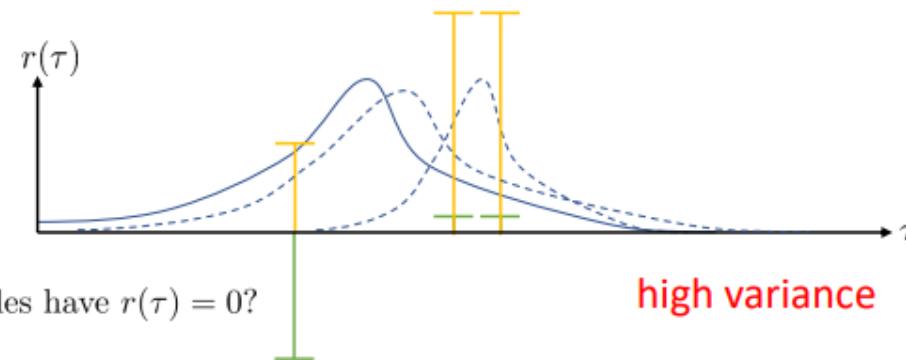
Iteration 2000





What is wrong with the policy gradient?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)$$



even worse: what if the two “good” samples have $r(\tau) = 0$?

Reducing variance

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Causality: policy at time t' cannot affect reward at time t when $t < t'$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \underbrace{\left(\sum_{t' \neq t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)}_{\text{return } \mathcal{R}}$$

“reward to go”

$$\hat{Q}_{i,t}$$

Baselines

a convenient identity

$$p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) = \nabla_\theta p_\theta(\tau)$$

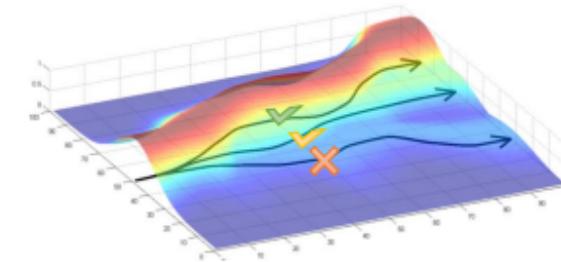
$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log p_\theta(\tau) [r(\tau) - b]$$

$$b = \underbrace{\frac{1}{N} \sum_{i=1}^N r(\tau)}$$

but... are we *allowed* to do that??

by wide *law*

★ $E[\nabla_\theta \log p_\theta(\tau)b] = \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) b d\tau = \int \nabla_\theta p_\theta(\tau) b d\tau = b \nabla_\theta \int p_\theta(\tau) d\tau = b \nabla_\theta 1 = 0$



subtracting a baseline is *unbiased* in expectation!

average reward is *not* the best baseline, but it's pretty good!



Analyzing variance

can we write down the variance?

$$\text{Var}[x] = E[x^2] - E[x]^2$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]$$

$$\text{Var} = E_{\tau \sim p_{\theta}(\tau)} [(\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b))^2] - \underbrace{E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) (r(\tau) - b)]^2}$$

this bit is just $E_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]$
(baselines are unbiased in expectation)

$$\begin{aligned} \frac{d\text{Var}}{db} &= \frac{d}{db} E[g(\tau)^2 (r(\tau) - b)^2] = \frac{d}{db} (E[\cancel{g(\tau)^2 r(\tau)^2}] - 2E[g(\tau)^2 r(\tau)b] + b^2 E[g(\tau)^2]) \\ &= -2E[g(\tau)^2 r(\tau)] + 2bE[g(\tau)^2] = 0 \end{aligned}$$

$$b = \frac{E[g(\tau)^2 r(\tau)]}{E[g(\tau)^2]} \quad \leftarrow \quad \text{This is just expected reward, but weighted by gradient magnitudes!}$$

Policy gradient is on-policy

$$\theta^* = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim p_\theta(\tau)}[r(\tau)]$$

$$\nabla_{\theta} J(\theta) = \underline{E_{\tau \sim p_\theta(\tau)}[\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)]}$$



this is trouble...

- Neural networks change only a little bit with each gradient step
- On-policy learning can be extremely inefficient!

REINFORCE algorithm:

- 
1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
 2. $\nabla_{\theta} J(\theta) \approx \sum_i (\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i)) (\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i))$
 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

can't just skip this!



Recap: policy gradients

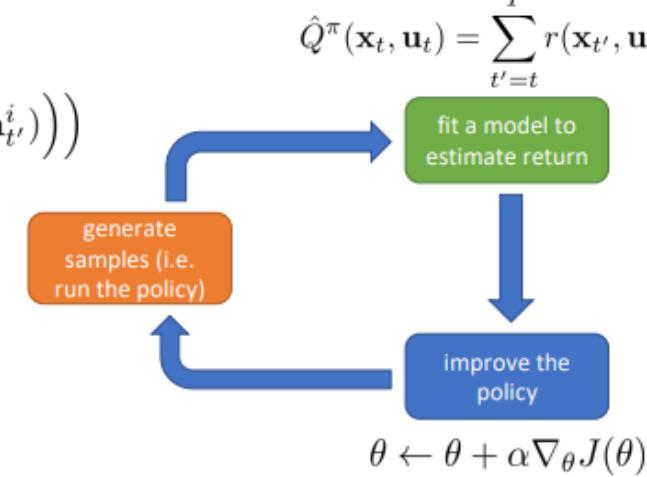
REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy)
2. $\nabla_\theta J(\theta) \approx \sum_i \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i) \left(\sum_{t'=t}^T r(\mathbf{s}_{t'}^i, \mathbf{a}_{t'}^i) \right) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}^\pi$$

“reward to go”

$$\hat{Q}^\pi(\mathbf{x}_t, \mathbf{u}_t) = \sum_{t'=t}^T r(\mathbf{x}_{t'}, \mathbf{u}_{t'})$$





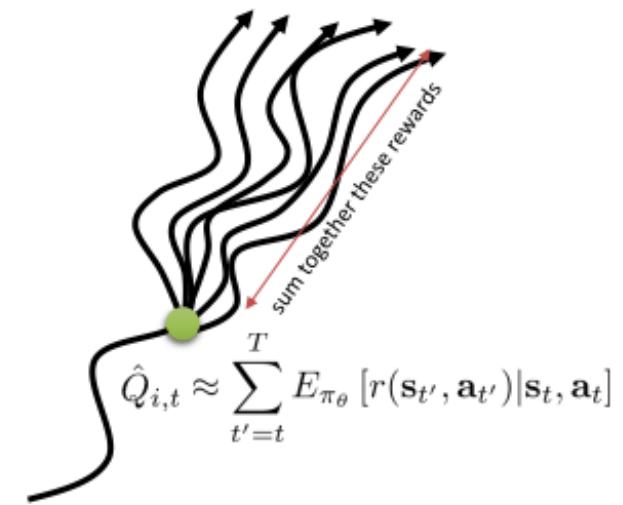
What about the baseline?

$Q(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: true expected reward-to-go

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) (Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - V(\mathbf{s}_{i,t}))$$

$$b_t = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad \text{average what?}$$

$$V(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)]$$



State & state-action value functions

$Q^\pi(s_t, a_t) = \sum_{t'=t}^T E_{\pi_\theta}[r(s_{t'}, a_{t'})|s_t, a_t]$: total reward from taking a_t in s_t

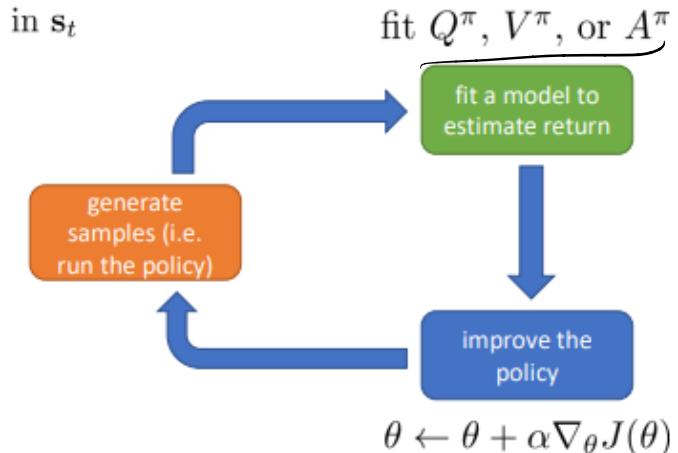
$V^\pi(s_t) = E_{a_t \sim \pi_\theta(a_t|s_t)}[Q^\pi(s_t, a_t)]$: total reward from s_t

$A^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t)$: how much better a_t is

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

the better this estimate, the lower the variance

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\underbrace{\sum_{t'=1}^T r(s_{i,t'}, a_{i,t'}) - b}_{\text{unbiased, but high variance single-sample estimate}} \right)$$



Value function fitting

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta}[r(\mathbf{s}_{t'}, \mathbf{a}_{t'})|\mathbf{s}_t, \mathbf{a}_t]$$

$$V^\pi(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)}[Q^\pi(\mathbf{s}_t, \mathbf{a}_t)]$$

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) - V^\pi(\mathbf{s}_t)$$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t}) A^\pi(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

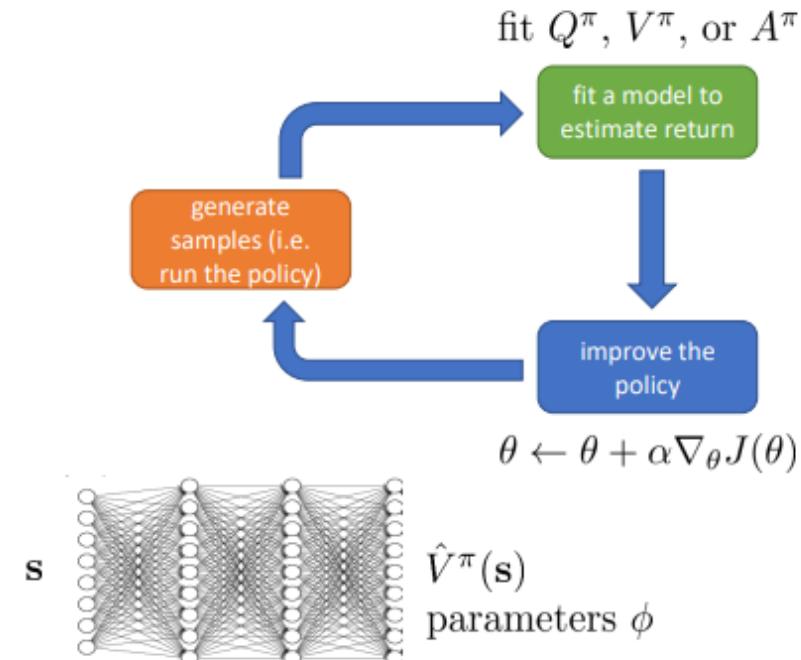
fit what to what?

Q^π, V^π, A^π ?

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \underbrace{\sum_{t'=t+1}^T E_{\pi_\theta}[r(\mathbf{s}_{t'}, \mathbf{a}_{t'})|\mathbf{s}_t, \mathbf{a}_t]}_{\hat{V}^\pi(\mathbf{s}_t)}$$

$$A^\pi(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^\pi(\mathbf{s}_{t+1}) - \hat{V}^\pi(\mathbf{s}_t)$$

let's just fit $V^\pi(\mathbf{s})$!





Reinforcement Learning

Policy evaluation

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$

$$J(\theta) = E_{\mathbf{s}_1 \sim p(\mathbf{s}_1)} [V^\pi(\mathbf{s}_1)]$$

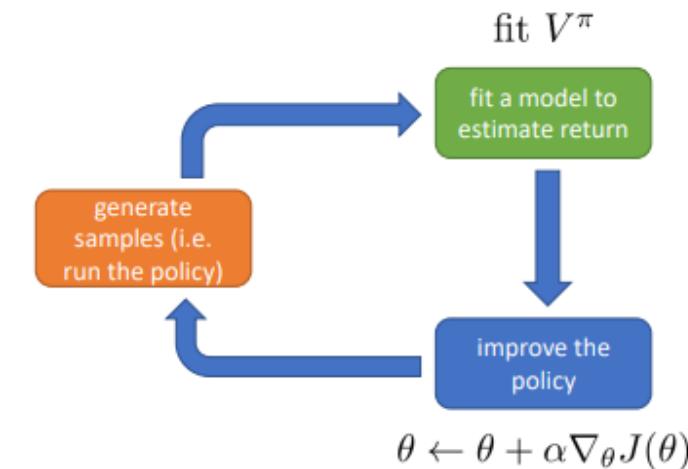
how can we perform policy evaluation?

Monte Carlo policy evaluation (this is what policy gradient does)

$$V^\pi(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

$$V^\pi(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

(requires us to reset the simulator)





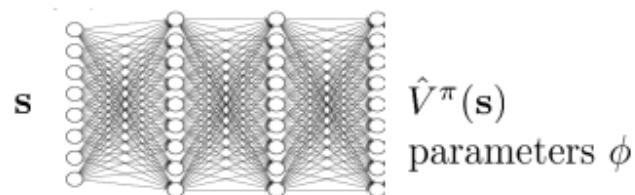
An actor-critic algorithm

batch actor-critic algorithm:

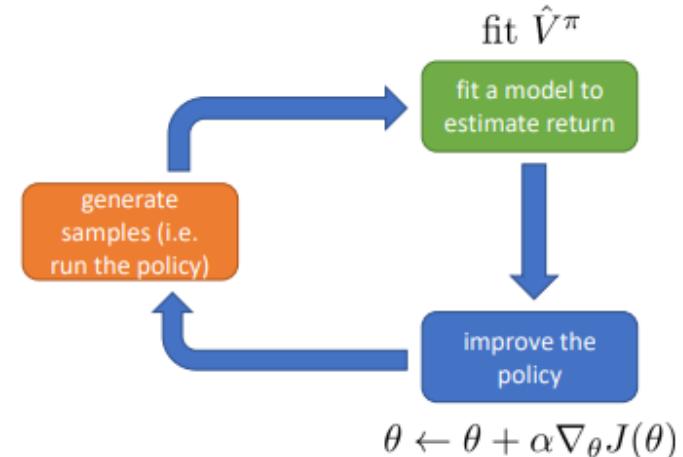
1. sample $\{\mathbf{s}_i, \mathbf{a}_i\}$ from $\pi_\theta(\mathbf{a}|\mathbf{s})$ (run it on the robot)
2. fit $\hat{V}_\phi^\pi(\mathbf{s})$ to sampled reward sums
3. evaluate $\hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \hat{V}_\phi^\pi(\mathbf{s}'_i) - \hat{V}_\phi^\pi(\mathbf{s}_i)$
4. $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^\pi(\mathbf{s}_i, \mathbf{a}_i)$
5. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$$y_{i,t} \approx \sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$$

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$



$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t]$$



Aside: discount factors

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$$

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \left\| \hat{V}_\phi^\pi(\mathbf{s}_i) - y_i \right\|^2$$

what if T (episode length) is ∞ ?

\hat{V}_ϕ^π can get infinitely large in many cases



episodic tasks



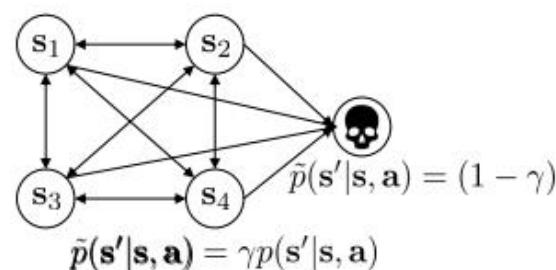
continuous/cyclical tasks

simple trick: better to get rewards sooner than later

$$y_{i,t} \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1})$$

↑
discount factor $\gamma \in [0, 1]$ (0.99 works well)

γ changes the MDP:

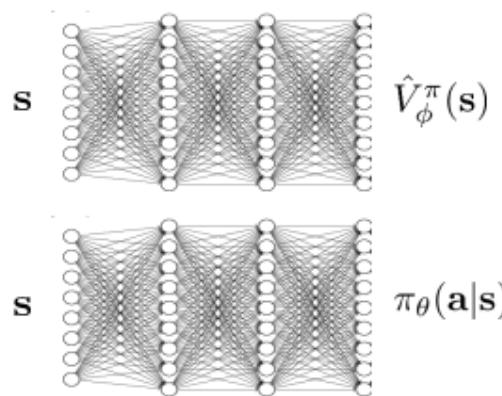


Architecture design

online actor-critic algorithm:

- 1. take action $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
- 2. update \hat{V}_ϕ^π using target $r + \gamma \hat{V}_\phi^\pi(\mathbf{s}')$
- 3. evaluate $\hat{A}^\pi(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}') - \hat{V}_\phi^\pi(\mathbf{s})$
- 4. $\nabla_\theta J(\theta) \approx \nabla_\theta \log \pi_\theta(\mathbf{a}|\mathbf{s}) \hat{A}^\pi(\mathbf{s}, \mathbf{a})$
- 5. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

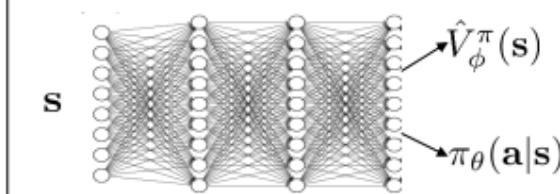
two network design



+ simple & stable

- no shared features between actor & critic

shared network design



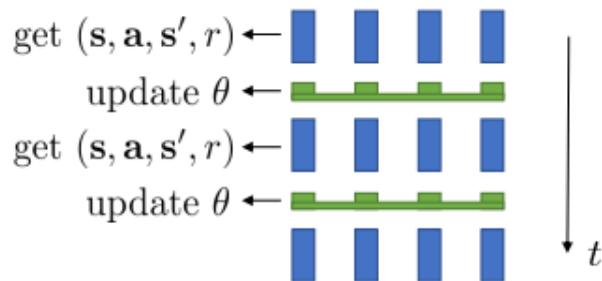


Online actor-critic in practice

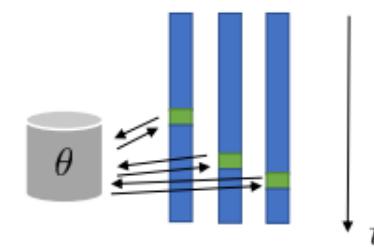
online actor-critic algorithm:

1. take action $\mathbf{a} \sim \pi_\theta(\mathbf{a}|\mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$
2. update \hat{V}_ϕ^π using target $r + \gamma \hat{V}_\phi^\pi(\mathbf{s}') \leftarrow$ works best with a batch (e.g., parallel workers)
3. evaluate $\hat{A}^\pi(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \hat{V}_\phi^\pi(\mathbf{s}') - \hat{V}_\phi^\pi(\mathbf{s})$
4. $\nabla_\theta J(\theta) \approx \nabla_\theta \log \pi_\theta(\mathbf{a}|\mathbf{s}) \hat{A}^\pi(\mathbf{s}, \mathbf{a}) \leftarrow$
5. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

synchronized parallel actor-critic



asynchronous parallel actor-critic





Reinforcement Learning



Reinforcement Learning



Reinforcement Learning
