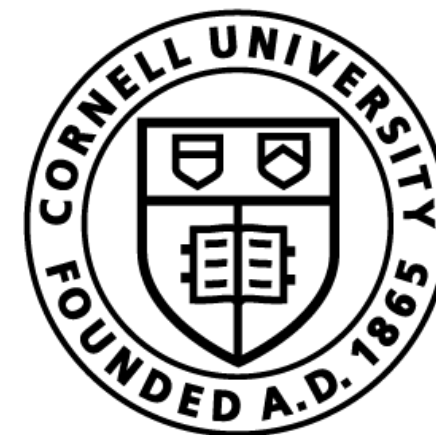


Nightmares of Policy Optimization

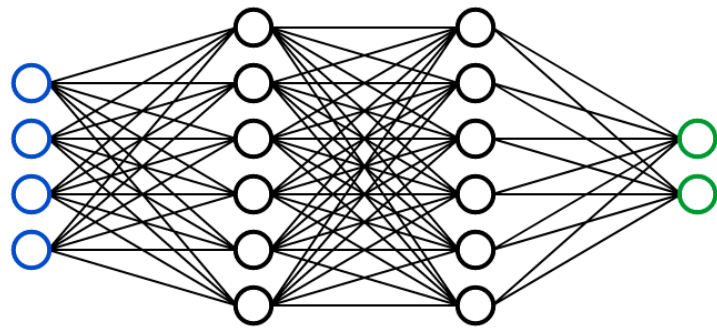
Sanjiban Choudhury



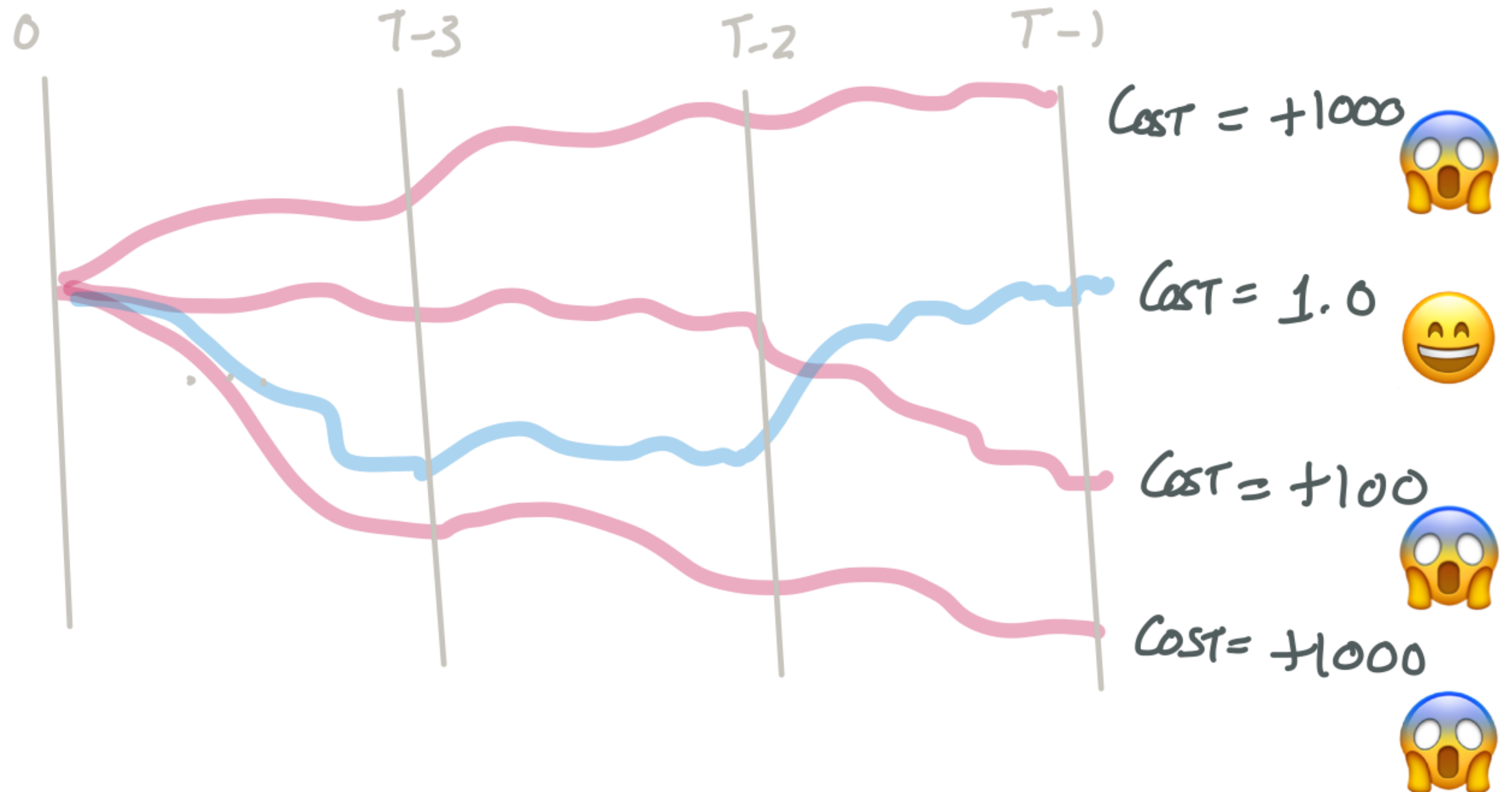
Cornell Bowers CIS
Computer Science

Can we just focus on finding a good policy?

$$\pi_{\theta} : S_t \rightarrow a_t$$



Learn a mapping from states to actions



Roll-out policies in the real-world to estimate value

We assumed black-box policies ...

BLACK BOX



DATE:
No
Pr
P

Form No. 1
THIS CASE ORIGINATED AT

REPORT MADE AT

TITLE

Mr. ARNOLD GIBBY

SYNOPSIS OF FACTS:

Subject employment as a worker chief for the

DETAILS:

BACKGROUND

Birth

Have we redacted too much?

SUBJECT:

I. [redacted] is being s[redacted] in the areas of [redacted] work of [redacted] and biological study at [redacted] The pr [redacted] until [redacted].

2. This [redacted] effects of [redacted] load into the [redacted] areas, [redacted] ors will continue to be D [redacted] and [redacted].

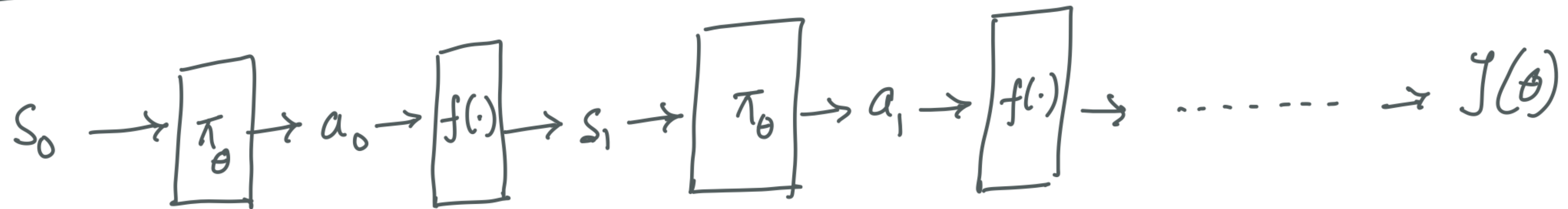
3. The [redacted] The [redacted] estimated to in the [redacted] required, under [redacted]

Black-box vs White-box vs Gray-box

BLACK BOX



WHITE BOX

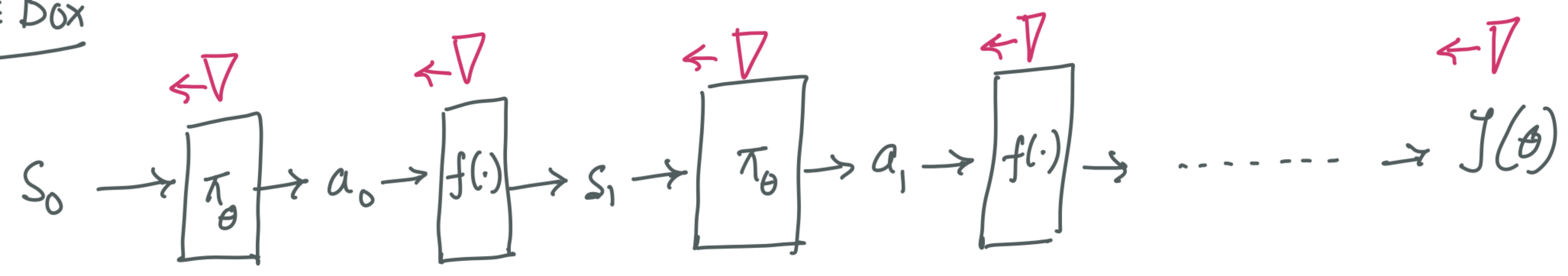


Black-box vs White-box vs Gray-box

Black Box



White Box



How can we take
gradients if we don't
know the dynamics?



The Likelihood Ratio Trick!



REINFORCE

Algorithm 20: The REINFORCE algorithm.

Start with an arbitrary initial policy π_θ

while *not converged* **do**

Run simulator with π_θ to collect $\{\zeta^{(i)}\}_{i=1}^N$

Compute estimated gradient

$$\tilde{\nabla}_\theta J = \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta \left(a_t^{(i)} | s_t^{(i)} \right) \right) R(\zeta^{(i)}) \right]$$

Update parameters $\theta \leftarrow \theta + \alpha \tilde{\nabla}_\theta J$

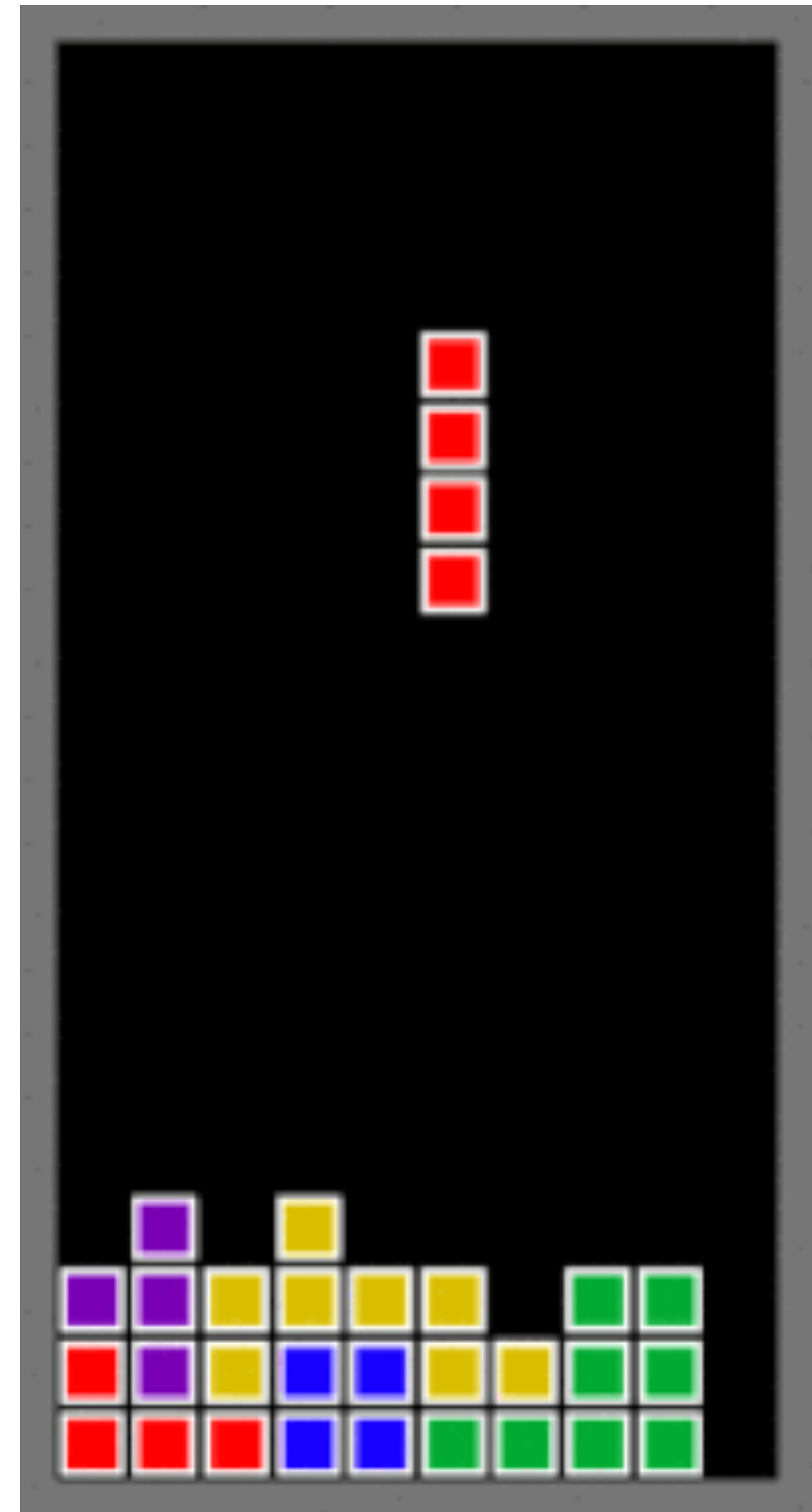
return π_θ

Tetris Policy

$$\pi_{\theta}(a|s) = \frac{\exp(\theta^{\top} f(s, a))}{\sum_{a'} \exp(\theta^{\top} f(s, a'))}$$

$f_1(s, a) = \#$ number of holes

$f_2(s, a) = \#$ max height

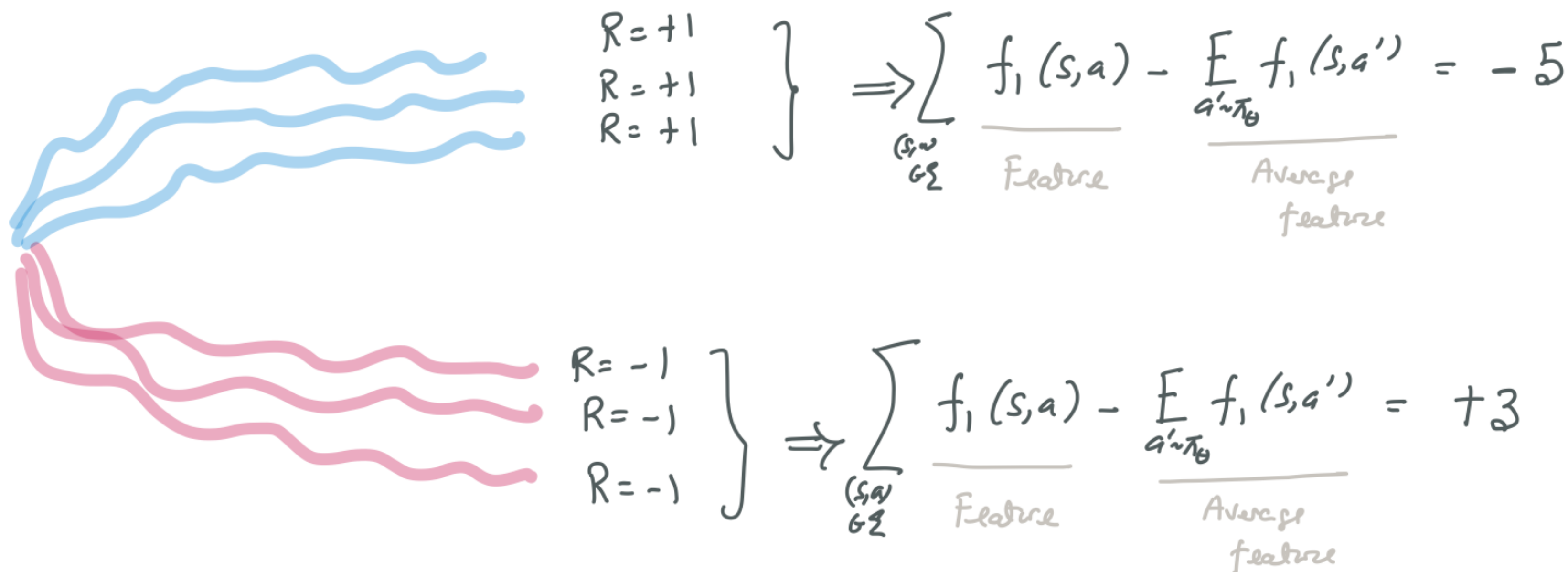


Chugging through the gradient ..

$$\begin{aligned}\nabla_{\theta} \log \pi_{\theta}(a|s) &= \nabla_{\theta} \left[\theta^{\top} f(s, a) - \log \sum_{a'} \exp \left(\theta^{\top} f(s, a') \right) \right] \\ &= f(s, a) - \frac{\sum_{a'} f(s, a') \exp \left(\theta^{\top} f(s, a') \right)}{\sum_{a'} \exp \left(\theta^{\top} f(s, a') \right)} \\ &= f(s, a) - \sum_{a'} f(s, a') \pi_{\theta} (a' | s) \\ &= f(s, a) - E_{\pi_{\theta}(a'|s)} [f(s, a')]\end{aligned}$$

Understanding the REINFORCE update

LET $f_1(s,a) = \# \text{ holes.}$



$$\begin{aligned} \Theta_1 &= \Theta_0 + \sum_{(s,a)} \left(\nabla_{\theta} \log \pi_{\theta}(a|s) \right) R(z) = \Theta_0 + \alpha \left(-5 \times (+1) + 3 \times (-1) \right) \\ &= \Theta_0 - \alpha 8 \quad (\text{Bump down this feature}) \end{aligned}$$

REINFORCE

Algorithm 20: The REINFORCE algorithm.

Start with an arbitrary initial policy π_θ

while *not converged* **do**

Run simulator with π_θ to collect $\{\zeta^{(i)}\}_{i=1}^N$

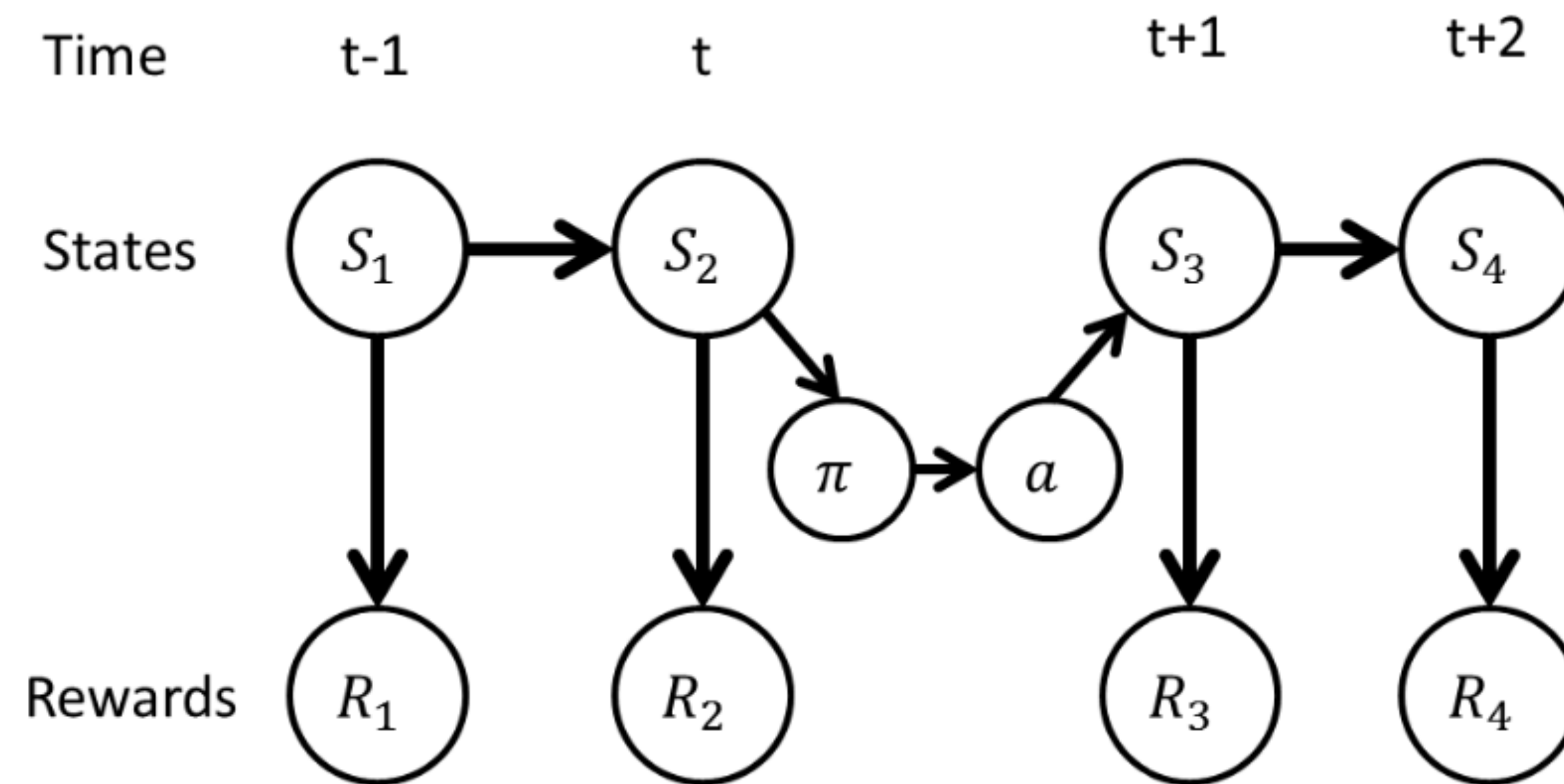
Compute estimated gradient

$$\tilde{\nabla}_\theta J = \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta \left(a_t^{(i)} | s_t^{(i)} \right) \right) R(\zeta^{(i)}) \right]$$

Update parameters $\theta \leftarrow \theta + \alpha \tilde{\nabla}_\theta J$

return π_θ

Causality: Can actions affect the past?



The Policy Gradient Theorem

$$\begin{aligned}\nabla_{\theta} J &= E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \left(\nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \left(\sum_{t'=0}^{t-1} r(s_{t'}, a_{t'}) + \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) \right) \right] \\ &= E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \left(\nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) \right],\end{aligned}$$

$$\nabla_{\theta} J = E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) Q^{\pi_{\theta}}(s_t, a_t) \right]$$

Life is good!

This solves
everything ...



The Three Nightmares of Policy Optimization



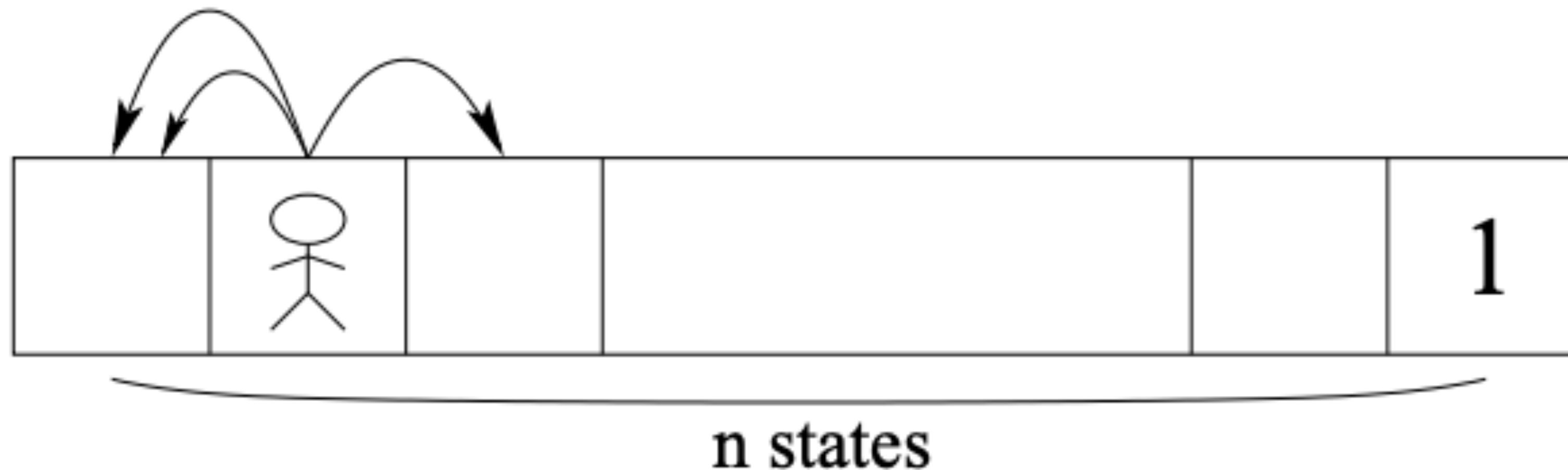
Nightmare 1:
Local Optima



Activity!



Consider the following MDP



Let's say I picked actions uniformly.

How long would it take me to get to the state with reward=1?

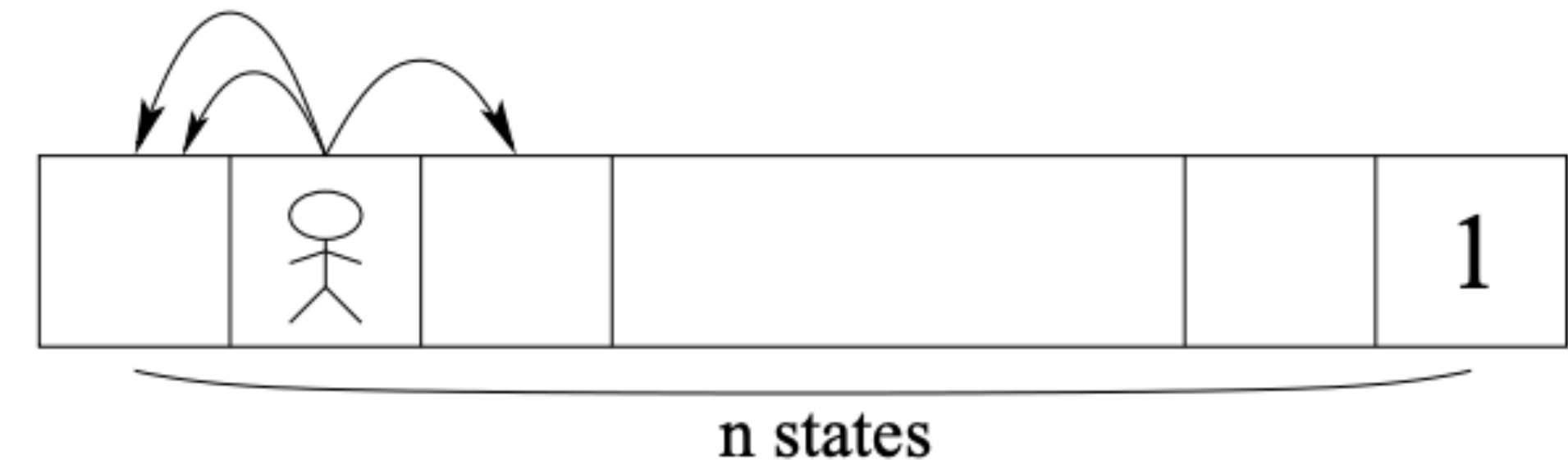
From Kakade and Langford

Think-Pair-Share

Think (30 sec): How long would it take me to get to the state with reward = 1? What does this imply if I run policy gradients?

Pair: Find a partner

Share (45 sec): Partners exchange ideas



Problem: Lack of exploration

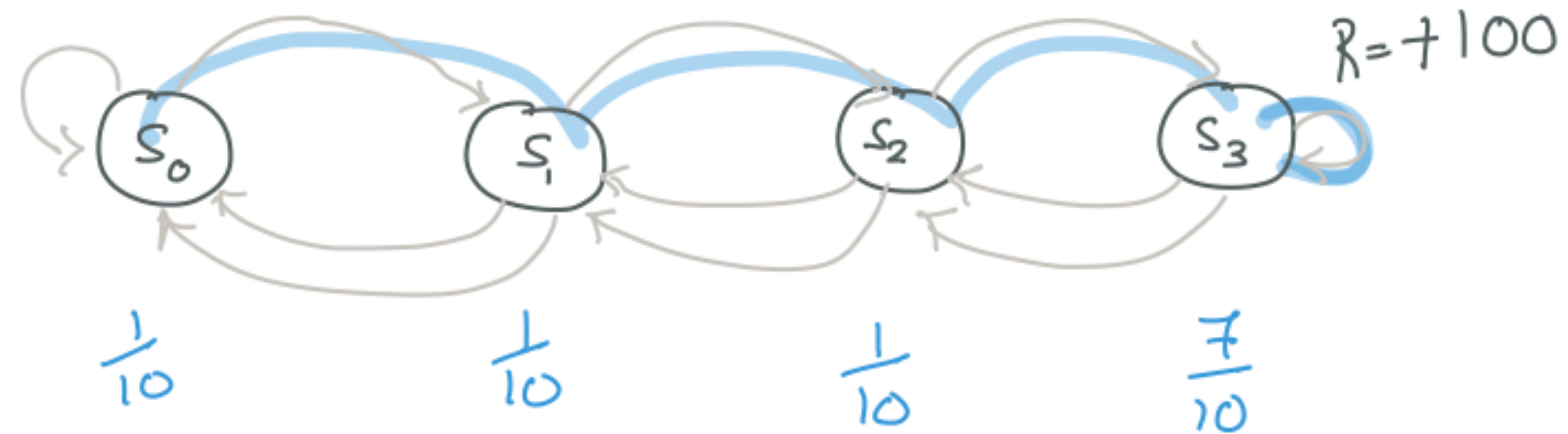


Optimal policy

Problem: Lack of exploration

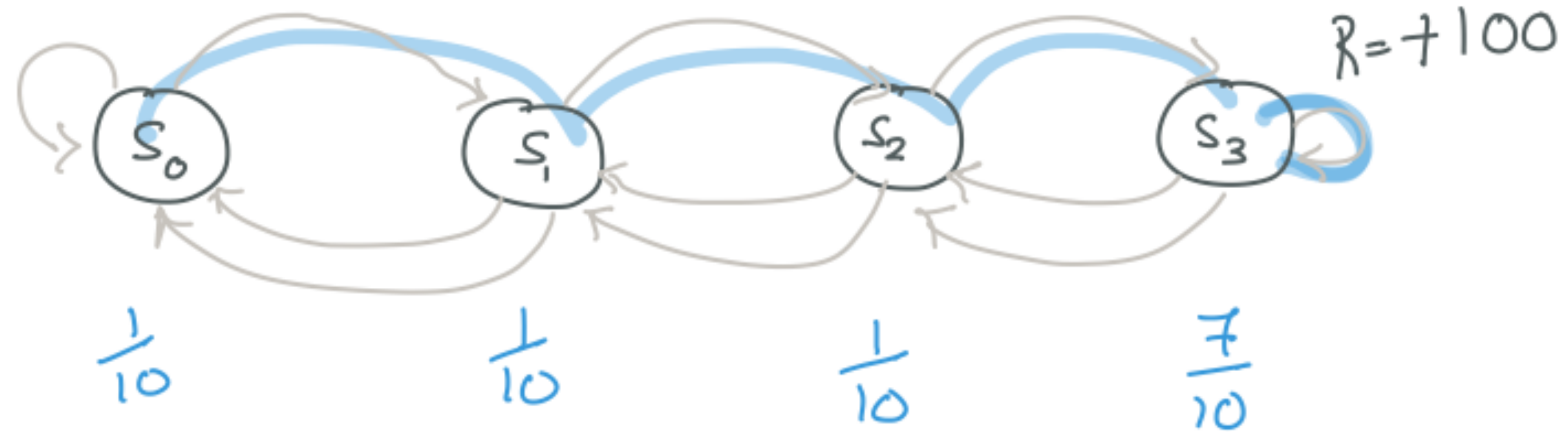


Optimal policy

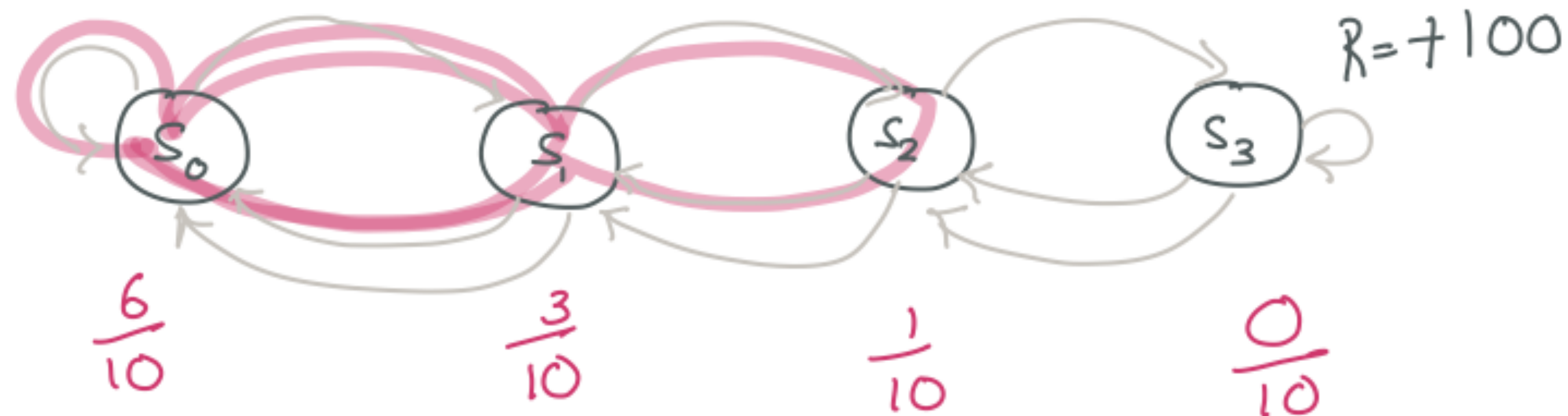


Problem: Lack of exploration

Optimal policy

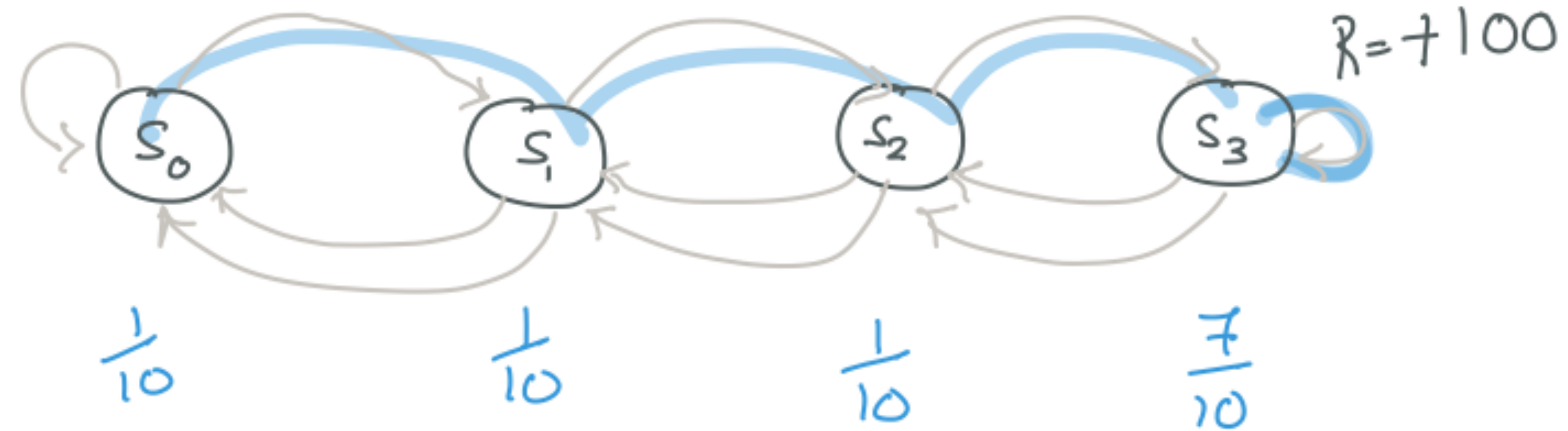


Random Policy starting from s_0

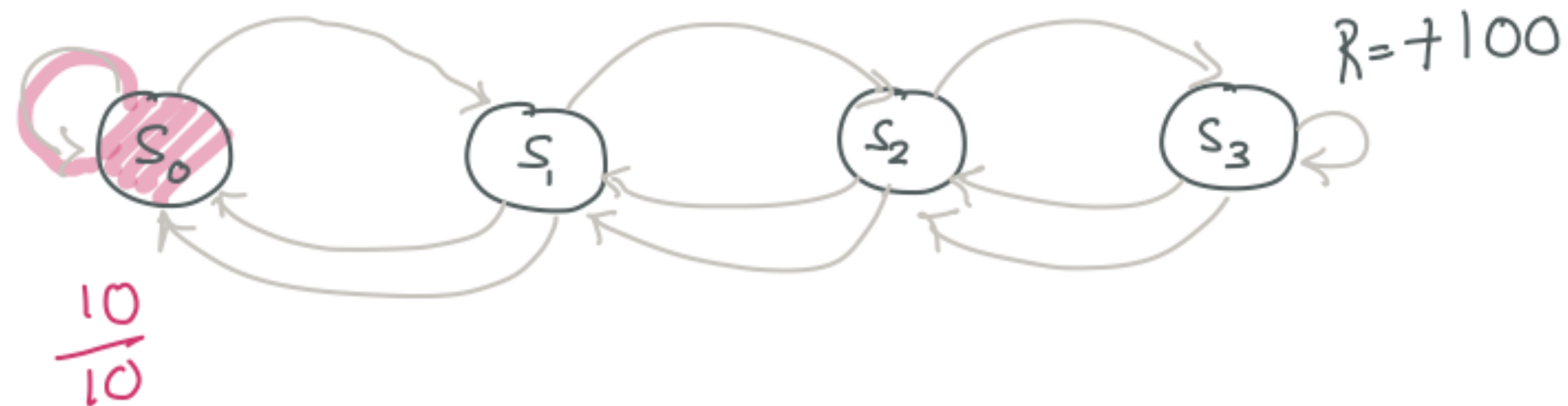


Problem: Lack of exploration

Optimal policy

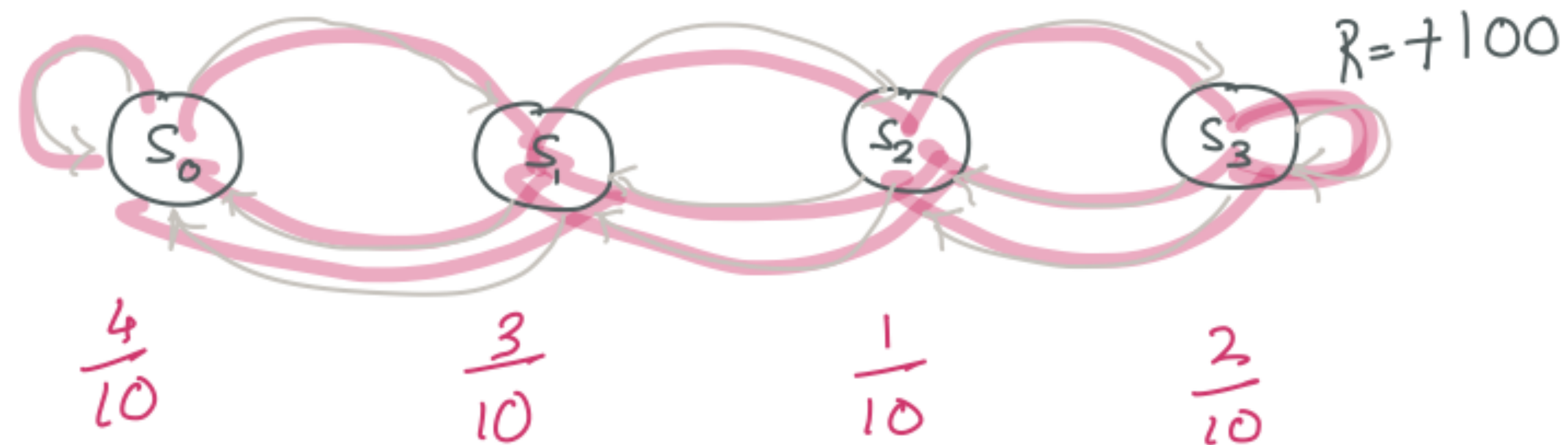


AFTER MANY ROUNDS OF POLICY ITERATION

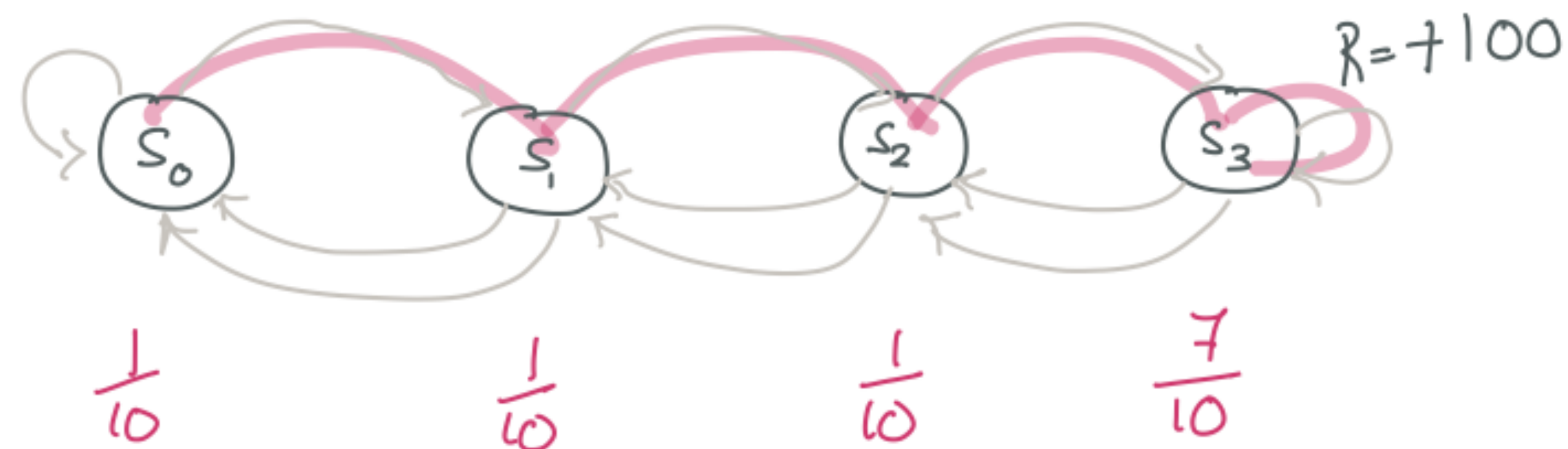


Solution: Demand improvement from *all states*

CHOOSING AN UNIFORM START STATE DISTRIBUTION



AFTER POLICY ITERATION



Key Idea: Use a good “restart” distribution

Choose a restart distribution $\mu(s)$ instead of start state distribution

Try your best to “cover” states the expert will visit

Suffer at most a penalty of $\left\| \frac{d_{\pi^*}}{\mu} \right\|_{\infty}$

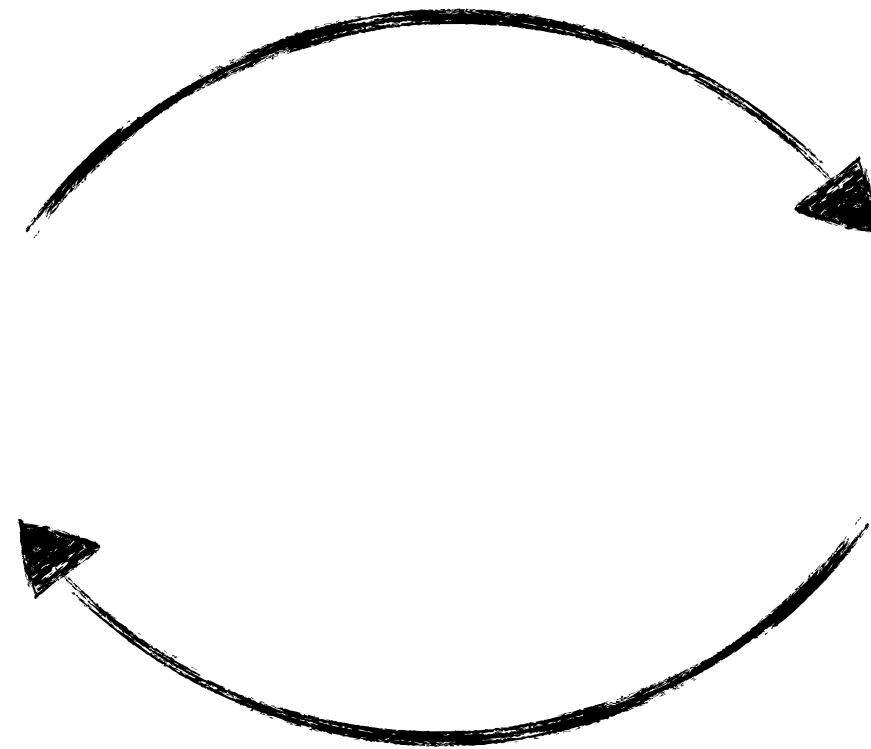
Nightmare 2:
Distribution Shift



Approximate Policy Iteration

Estimate advantage

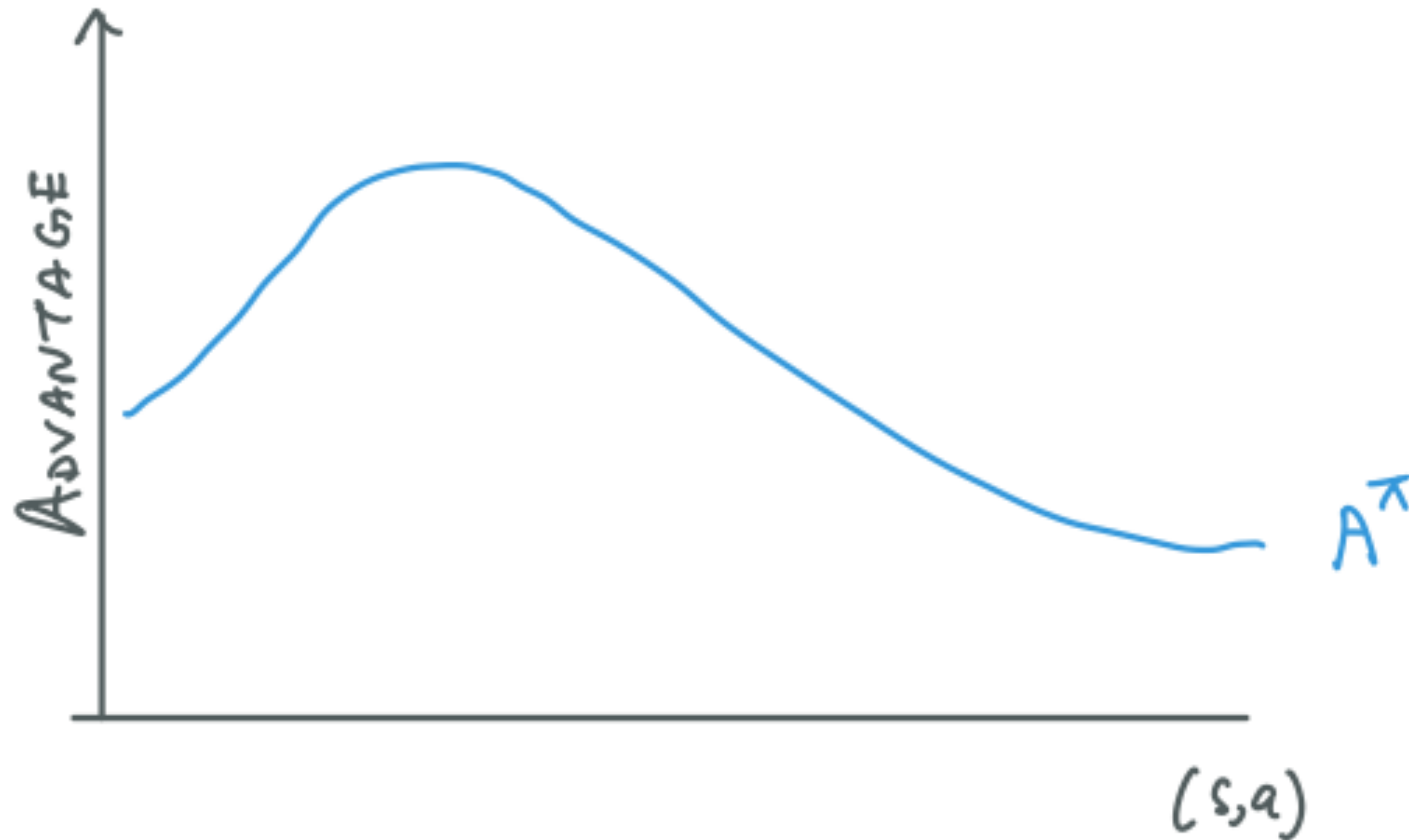
$$A^{\pi}(s, a)$$



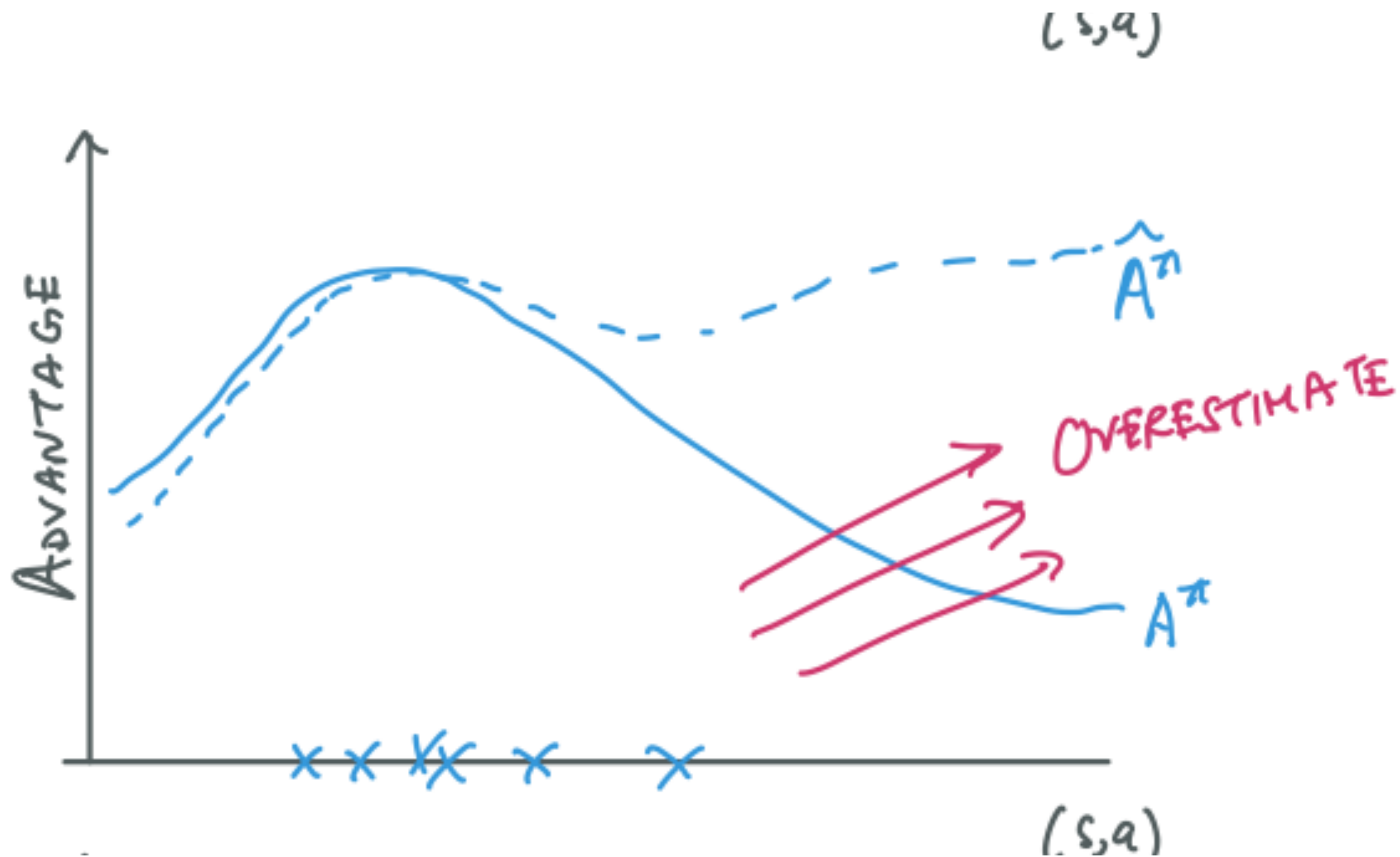
Greedily improve policy

$$\pi' = \arg \min_{\pi'} A^{\pi}(s, \pi'(s))$$

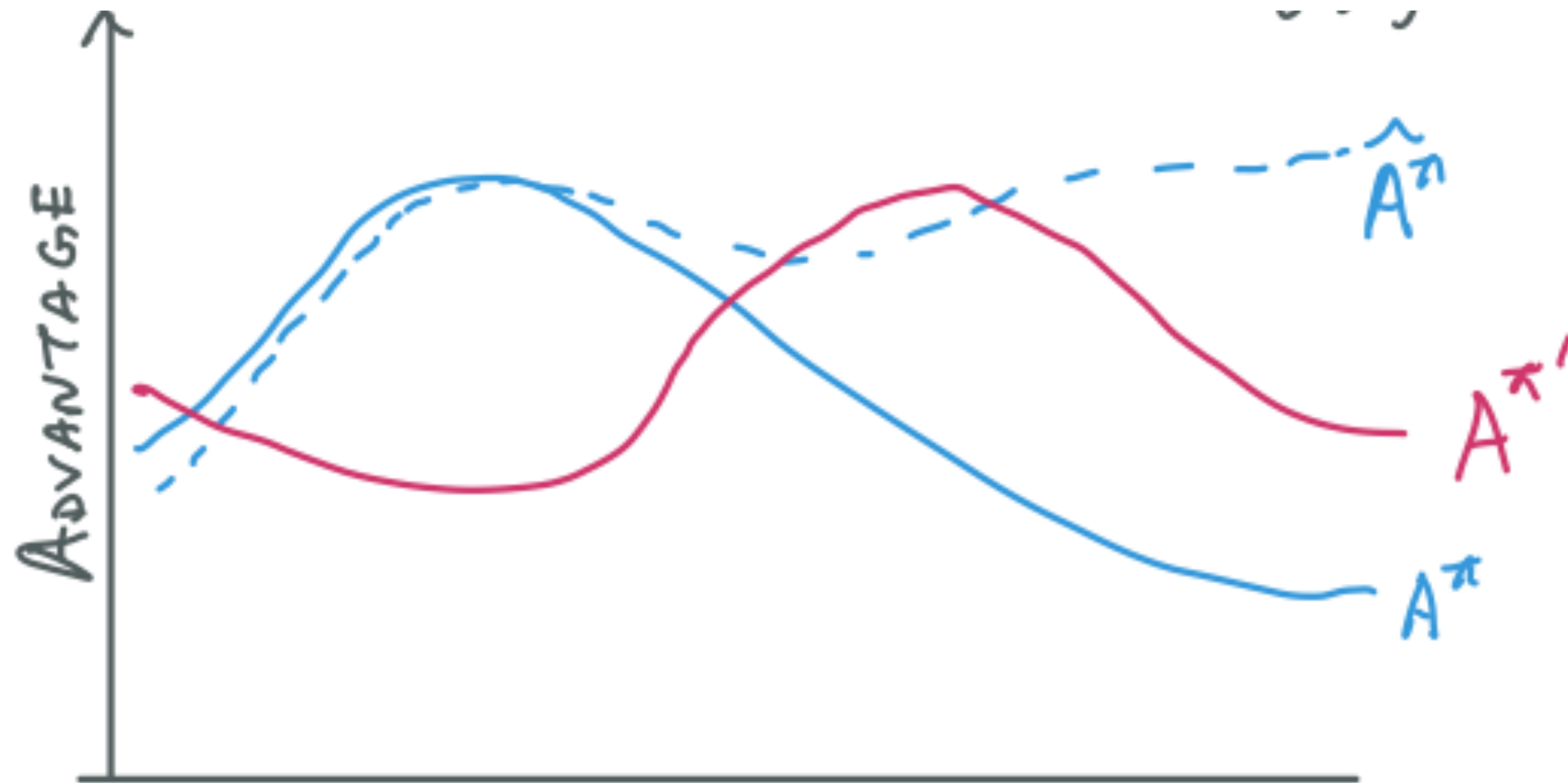
The problem of distribution shift



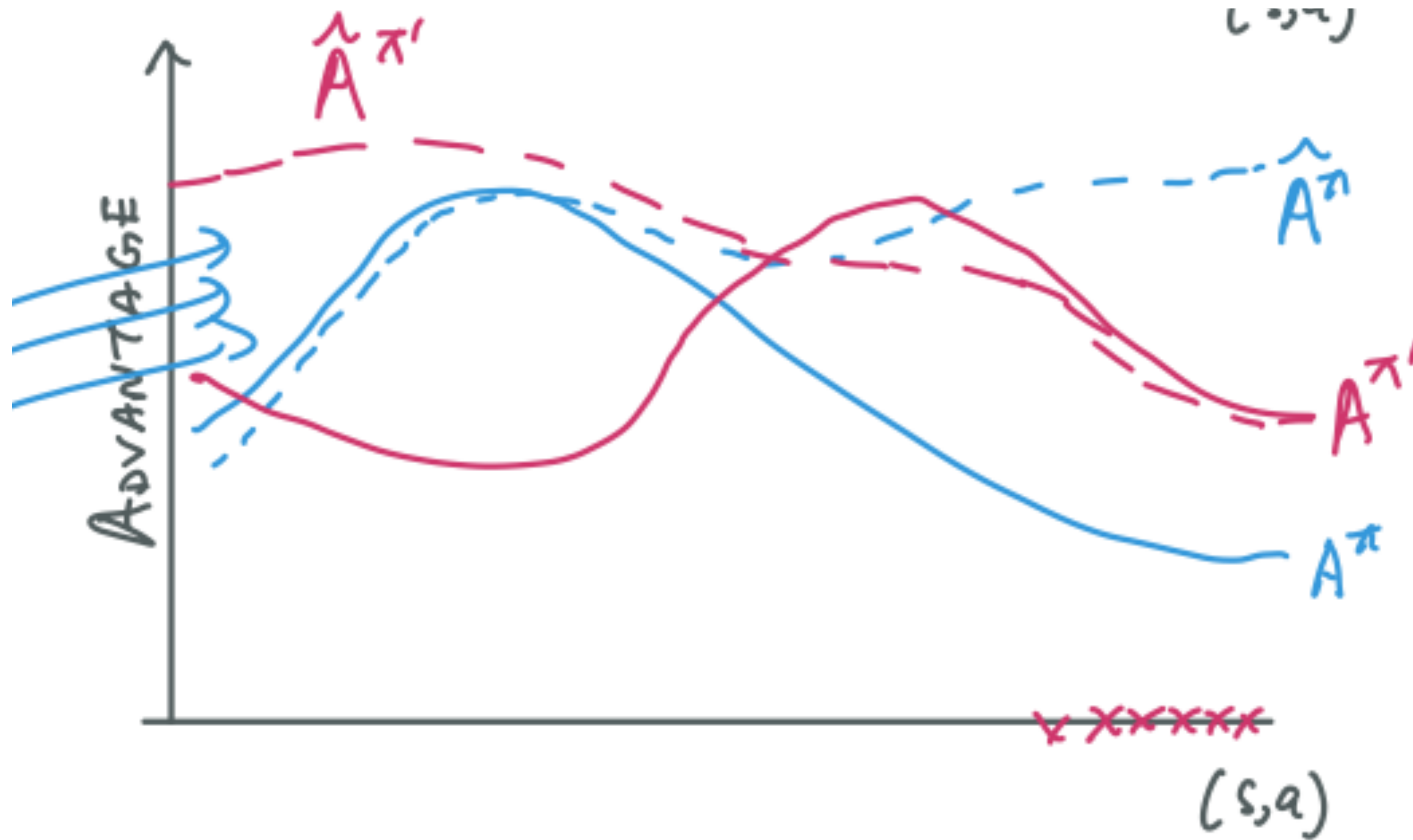
The problem of distribution shift



The problem of distribution shift



The problem of distribution shift



How does distribution shift manifest?

The true performance difference

$$\begin{array}{cc} V^{\pi'}(s) & - & V^{\pi}(s) & = & \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi}^t} A^{\pi}(s, \pi'(s)) \\ \text{(New)} & & \text{(Old)} & & \end{array}$$

What our estimator currently approximates

$$\frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi}^t} A^{\pi}(s, \pi'(s))$$

Be stable

Slowly change
policies

Keep d_{π}^t close to $d_{\pi'}^t$



Idea 1: Conservative Policy Iteration (CPI)

$$\pi' = (1 - \alpha)\pi + \alpha\pi_{greedy}$$

Mix in old policy and greedy policy

Can prove that performance difference is bounded by

$$V^{\pi'}(s) - V^{\pi}(s) \geq \alpha A_{greedy} - 2\alpha^2 \frac{\gamma}{1 - \gamma}$$

How much greedy policy
improves based on
estimate

How much distribution
shift hurts!

Approximately Optimal Approximate Reinforcement Learning

Sham Kakade

Gatsby Computational Neuroscience Unit, UCL, London WC1N 3AR, UK

SHAM@GATSBY.UCL.AC.UK

John Langford

Computer Science Department, Carnegie-Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15217

JCL@CS.CMU.EDU

Idea 2: Update distributions slowly

CPI requires keeping around *all* the policies you have seen thus far,
which is not scalable ...

Instead can we change policies slowly?

Does this simply mean do gradient descent with a small step size?

Nightmare 2:

~~Distribution Shift~~
Correlated Features



Activity!



What happens if we have correlated features?

Parameterization 1: $f_1 = \#$ of Holes after the placement, $f_2 = \text{Height}$ after the placement. We use θ to denote the parameter for this parameterization.

Parameterization 2: $g_1 = \dots = g_{100} = \#$ of Holes after the placement, $g_{101} = \text{Height}$ after the placement. We use ϕ to denote the parameter for this parameterization

Then, for Parameterization 1, we have,

$$\theta^\top f(x, a) = \theta_1 \times \# \text{ of Holes}(x, a) + \theta_2 \times \text{Height}(x, a).$$

While for Parameterization 2, we have,

$$\phi^\top g = \left(\sum_{i=1}^{100} \phi_i \right) \times \# \text{ of Holes}(x, a) + \phi_{101} \times \text{Height}(x, a).$$

Think-Pair-Share

Think (30 sec): What would happen if we ran policy gradient with Feature Set 1 vs Feature Set 2? How can we fix it?

Pair: Find a partner

Then, for Parameterization 1, we have,

$$\theta^\top f(x, a) = \theta_1 \times \# \text{ of Holes}(x, a) + \theta_2 \times \text{Height}(x, a).$$

While for Parameterization 2, we have,

$$\phi^\top g = \left(\sum_{i=1}^{100} \phi_i \right) \times \# \text{ of Holes}(x, a) + \phi_{101} \times \text{Height}(x, a).$$

Share (45 sec): Partners exchange ideas

Gradient Descent as Steepest Descent

Gradient Descent is simply Steepest Descent with L2 norm

$$\max_{\Delta\theta} J(\theta + \Delta\theta) \quad s.t. \quad \|\Delta\theta\| \leq \epsilon$$

An alternative norm: KL Divergence! Gives rise to Fisher Information Matrix

$$G(\theta) = E_{p_\theta} \left[\nabla_\theta \log(p_\theta) \nabla_\theta \log(p_\theta)^\top \right] \quad \Delta\theta = \frac{1}{2\lambda} \tilde{G}^{-1}(\theta) \tilde{\nabla}_\theta J.$$

Natural Gradient Descent

Estimate Fisher Information Matrix

$$\tilde{G}(\theta) = \frac{1}{N} \sum_{i=1}^N \left[\nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \nabla_{\theta} \log \pi_{\theta}(a_i | s_i)^{\top} \right]$$

Parameter update:

$$\Delta\theta = \frac{1}{2\lambda} \tilde{G}^{-1}(\theta) \tilde{\nabla}_{\theta} J.$$

Modern variants known as TRPO, PPO

Nightmare 3:

Variance



What happens when Q values for all rollouts are similar?

$$\nabla_{\theta} J = E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) Q^{\pi_{\theta}}(s_t, a_t) \right]$$

Recall that one of the reasons for the high variance is that the algorithm does not know how well the trajectories perform compared to other trajectories. Therefore, by introducing a baseline for the total reward (or reward to go), we can update the policy based on how well the policy performs compared to a baseline

Solution: Subtract a baseline!

$$\nabla_{\theta} J = E_{d^{\pi_{\theta}}(s)} E_{\pi_{\theta}(a|s)} [\nabla_{\theta} \log(\pi_{\theta}(a|s)) (Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s))].$$

We can prove that this does not change the gradient

But turns Q values into advantage (which is lower magnitude)

tl;dr

The Policy Gradient Theorem

$$\begin{aligned}\nabla_{\theta} J &= E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \left(\nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \left(\sum_{t'=0}^{t-1} r(s_{t'}, a_{t'}) + \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) \right) \right] \\ &= E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \left(\nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right) \right],\end{aligned}$$

$$\nabla_{\theta} J = E_{p(\xi|\theta)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) Q^{\pi_{\theta}}(s_t, a_t) \right]$$

15

The Three Nightmares of Policy Optimization



1. Local Optima: Use Exploration Distribution
2. Distribution Shift: *Natural* Gradient Descent
3. High Variance: Subtract baseline