

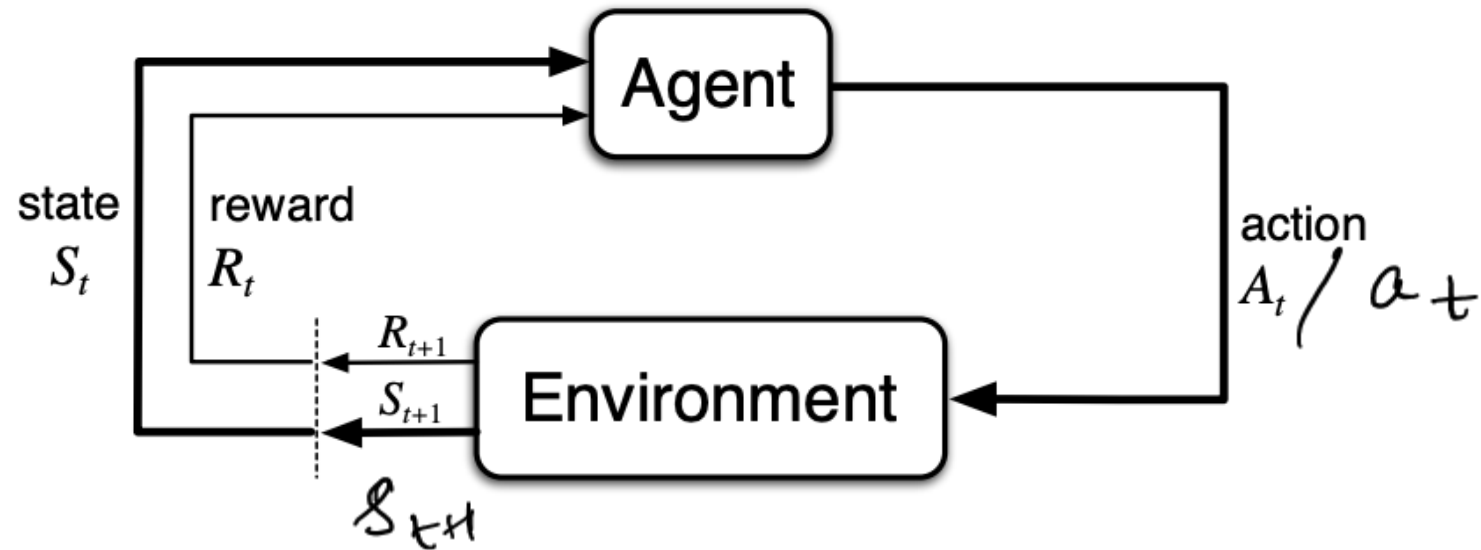
Policy optimization

$$Q_{\text{opt}}(s, a) \downarrow \arg \max_a Q_{\text{opt}}$$

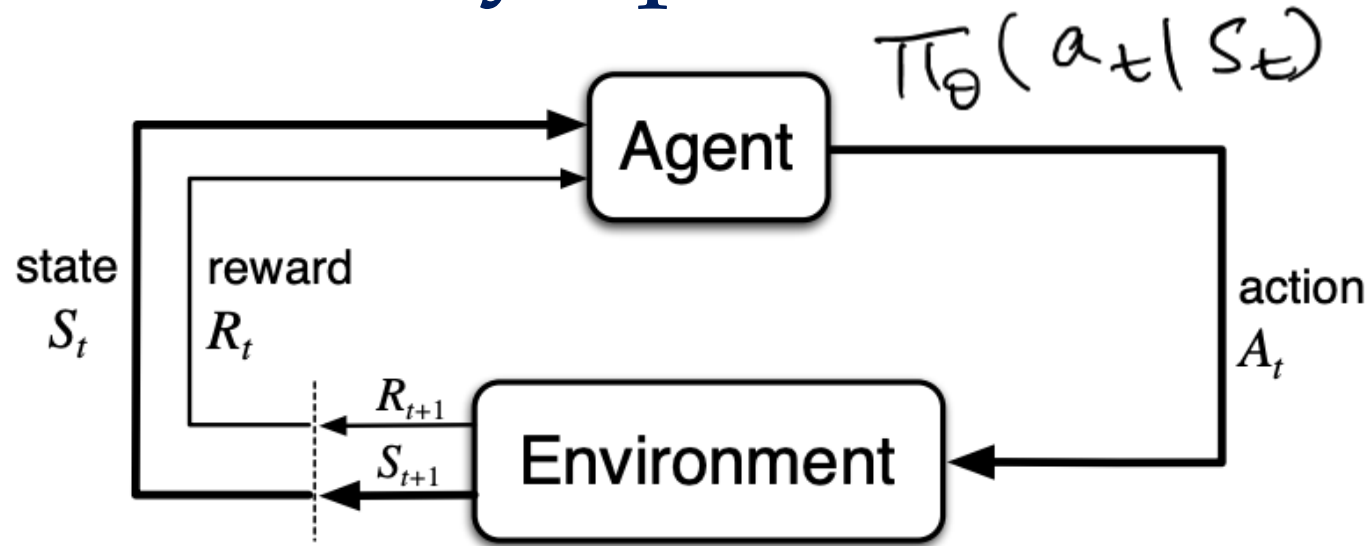
- Why even estimate Q-values? **Can we learn the policy directly?**
- Parameterize $\pi_{\theta}(a | s)$ as a neural network and directly train this network to maximize the rewards
- Should $\pi_{\theta}(a | s)$ be deterministic or stochastic?
 - The optimal policy is deterministic, but stochastic policies make the optimization process “smoother”
 - Also helps with exploration

input: s
output: $a \in A$

Reinforcement learning (recap)



Policy optimization



Goal

$$\max \sum_{t=0}^H R(s_t, a_t)$$

Handwritten notes:
→ MDP
could be discounted

Stochastic policy class

π_{θ} : distribution over actions given state

Why policy optimization?

- Can be simpler or faster to estimate optimal policy than Q or V
- If we compute V function, we still need to compute optimal policy by running Bellman update
- If we compute Q, we need to take argmax which can be challenging

$$\pi(s) = \arg\max_a Q_{\text{opt}}(s, a)$$

Likelihood ratio policy gradient

Notation: τ : state-action sequence, trajectory, roll-out

$s_0, a_0, s_1, a_1, s_2, \dots, s_H, a_H$

$R(\tau) : \sum_{t=0}^H R(s_t, a_t)$ $\pi_\theta(a|s)$ θ : determine a policy

$U(\theta)$: quality of policy corresponding to θ $\theta \rightarrow$ policy trajectory

$$= \mathbb{E} \left[\sum_{t=0}^H R(s_t, a_t) ; \theta \right] = \sum_{\tau} \underbrace{p(\tau; \theta)} R(\tau)$$

Goal: $\max_{\theta} V(\theta) = \max_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau)$ Expected reward

we need to compute $\nabla V(\theta)$ (gradient)

→ requires dynamics

$$\nabla_{\theta} V(\theta) = \nabla_{\theta} \left[\sum_{\tau} P(\tau; \theta) R(\tau) \right] = \sum_{\tau} \left[\nabla_{\theta} \underbrace{P(\tau; \theta) R(\tau)}_{\text{neural network}} \right]$$

For every gradient, we enumerate all trajectories

$$\sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau)$$

$$= \sum_{\tau} P(\tau; \theta) \left(\frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)} \cdot R(\tau) \right)$$

$$\nabla V(\theta) = \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)$$

unbiased estimate of gradient by sampling

→ $\mathbb{E} \left[\nabla_{\theta} \log P(\tau; \theta) R(\tau) ; \theta \right]$

↓

trajectories sampled from current policy
given by θ ($P(\tau; \theta)$)

Empirical estimate

Sample $\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(m)}$ according to π_θ

m paths under π_θ

$$\nabla V(\theta) \approx \frac{1}{m} \sum_{i=1}^m \nabla_\theta \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

R : need not be differentiable

$\log P(\tau^{(i)}; \theta)$: differentiable

Temporal decomposition

$$\begin{aligned}\nabla_{\theta} \log P(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[\prod_{t=0}^H P(S_{t+1}^{(i)} | S_t^{(i)}, a_t^{(i)}) \right] \\&= \nabla_{\theta} \left[\sum_{t=0}^H \log P(S_{t+1}^{(i)} | S_t^{(i)}, a_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(a_{t+1}^{(i)} | S_{t+1}^{(i)}) \right] \\&= \sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta}(a_{t+1}^{(i)} | S_{t+1}^{(i)})\end{aligned}$$

Handwritten annotations:

- dynamics model (env)* (under the first product term)
- policy* (under the second sum term)
- policy* (under the final sum term)

$\nabla U(\theta)$: unbiased estimate $= \hat{g}$

$\tau^{(1)}, \tau^{(2)} \dots \tau^{(m)}$

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta} (a_{t+1}^{(i)} | s_{t+1}^{(i)}) R(\pi^{(i)})$$

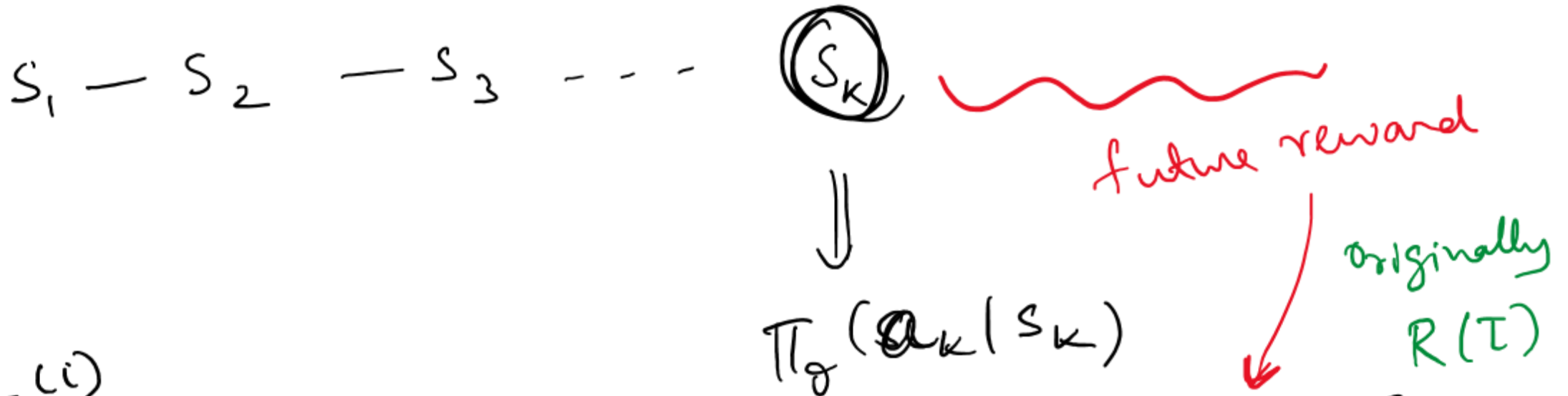
$$E[\hat{g}] = \nabla U(\theta)$$

Intuition

→ if $R(\tau)$ is high, ~~we~~ want τ to be more likely under θ

→ if $R(\tau)$ is low, τ should be less likely

Decompose to states and actions



For $\tau^{(i)}$:

$$\sum_{t=0}^H \pi_\theta \log \pi_\theta(a_t^{(i)} | s_t^{(i)}) \left[\sum_{k=t}^H R(s_k^{(i)}, a_k^{(i)}) \right]$$

Likelihood ratio gradient estimate

Variance reduction: baseline

→ if $R(\tau)$ is high, make τ more likely

high on average

For $\tau^{(i)}$

$$\sum_{t=0}^H \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \left[\sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)}) - b(s_k^{(i)}) \right]$$

can be anything!

Baselines

→ adding baseline lowers variance, keeps bias unchanged
baseline doesn't depend on probs of action

Large family of methods to come up with baselines
⇒ "actor critic methods"
→ baseline

baseline: $V^{\pi}(S_k^{(i)})$: value function from before

Practical update

How to estimate $V_{\phi}^{\pi}(s)$? \rightarrow a neural network

Collect $\tau^{(1)}, \tau^{(2)} \dots \tau^{(m)}$

Regress againsts empirical estimate

$\phi_{i:H}$

$$\phi \leftarrow \arg \min_{\phi} \left(\frac{1}{m} \sum_{i=1}^m \sum_{t=0}^{H-1} \left(V_{\phi}^{\pi}(s_t^{(i)}) - \underbrace{\left(\sum_{k=t}^{H-1} R(s_k^{(i)}, a_k^{(i)}) \right)}_{\text{target}} \right)^2 \right)$$

Monte-carlo estimation of V^π

→ from prev slide

Bootstrapped estimate

$$\phi_{i+1} \leftarrow \min \sum_{s, a, s', r}$$

$$\| \underbrace{r + V_{\phi_i}^\pi(s') - V_{\phi_i}^\pi(s)}_{\text{target}} \|_2^2 + \lambda \| \phi - \phi_i \|_2^2$$

monte-carlo
↔ SARSA

Algorithm 2 “Vanilla” policy gradient algorithm

Initialize policy parameter θ , baseline b

for iteration=1,2,... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute

the return $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$, and
the advantage estimate $\hat{A}_t = R_t - b(s_t)$.

Re-fit the baseline, by minimizing $\|b(s_t) - R_t\|^2$,
summed over all trajectories and timesteps.

Update the policy, using a policy gradient estimate \hat{g} ,

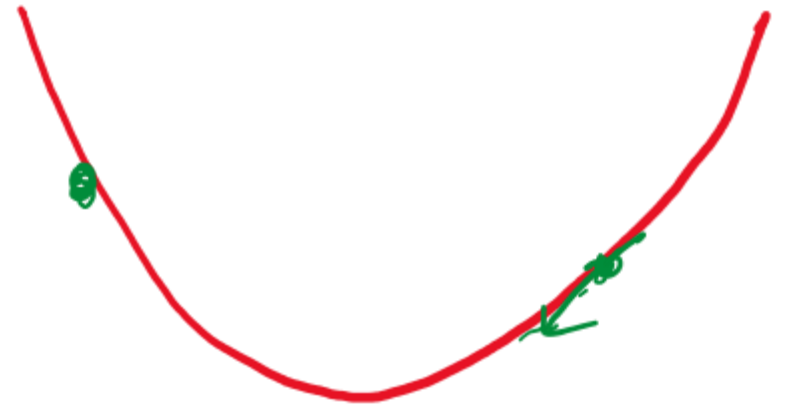
which is a sum of terms $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$

end for

lower
noise

ϕ is updated
value fn
network
policy is updated

Step sizing



- Cannot move too far or too less

- How to re-formulate policy gradient to allow a natural notion of “how far to move”?

rewards ^{gradients} come from current policy

$J(\theta)$: $\nabla J(\theta)$ from policy grad equation

Surrogate loss interpretation

$$V(\theta) = \sum_{\tau} P(\tau|\theta) R(\tau)$$

θ_{old} : fixed policy (current)

collect trajectories using θ_{old} , we estimating
expected reward of a different θ
"importance sampling"

$$V(\theta) = \mathbb{E}_{\tau \sim \theta_{old}} \left[\underbrace{\frac{P(\tau|\theta)}{P(\tau|\theta_{old})}}_{\gamma} R(\tau) \right]$$

γ

= importance weights

π_1	0.2	0.8
	A	B
	5	4

π_2	0.8	0.2
	A	B
	5	4

Trust-region Policy Optimization

TRPO:

$$\max_{\pi} L(\pi) = \mathbb{E}_{\pi_{\text{old}}} \left[\frac{\pi(a|s)}{\pi_{\text{old}}(a|s)} \underbrace{A^{\pi_{\text{old}}}(s,a)}_{\text{advantage rather than reward}} \right]$$

st. constraint that

$$\underbrace{K L(\pi \parallel \pi_{\text{old}})}_{\text{distance b/w prob distributions}} \leq \epsilon$$

original TRPO: ϵ for optimization

KL divergence : independent of dynamics

Constrained optimization is hard

TRPO

$$\max_{\theta} \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} \hat{A}_t \right]$$

st. $KL(\pi_{\theta_{old}} \parallel \pi_{\theta}) < \epsilon$

PPO: v1

$$\max_{\theta} \left[\hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} \hat{A}_t \right] - \beta KL(\pi_{old} \parallel \pi_{\theta}) \right]$$

PPO v2

$$\eta_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$$

$$\eta_t(\theta_{old}) = 1$$

"Don't deviate too much from 1"

$$L(\theta) = \mathbb{E}_t^{\text{CLIP}} \left[\min \left(\eta_t(\theta) \hat{A}_t, \text{clip} \left(\eta_t(\theta), 1-\epsilon, 1+\epsilon \right) \hat{A}_t \right) \right]$$

