

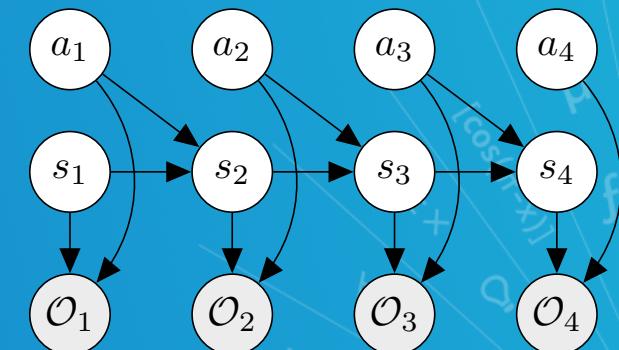
Probabilistic Graphical Models

Reinforcement Learning & Control
Through Inference in GM (part 2)

Maruan Al-Shedivat

Lecture 21, April 3, 2019

Reading: see class homepage



Logistics

Homework:

- HW3 is due tonight (April 3 @ 11:59 pm ET). HW4 will be out Monday, April 4.

Update on the project logistics:

- Poster session: NSH atrium, April 30 @ 2:30 – 5:30 pm.
- Presentation will count as 20% of the project grade.
- What is expected from your project posters:
 - Motivation of your project + clear description of your approach.
 - Preliminary results (good if you have most of the results ready, but totally okay to present only preliminary results for your method + baselines).
- Final report due Friday, May 10.



A note on materials used in this module

- Sutton & Barto. Reinforcement Learning: An Introduction. 2nd edition.
- David Silver's [UCL course](#) on reinforcement learning.
- Materials from UC Berkeley's [Deep RL course](#).
- Sergey Levine's [tutorial on RL and control as inference](#).
- Brian Ziebart's [PhD thesis](#) (maximum causal entropy models).

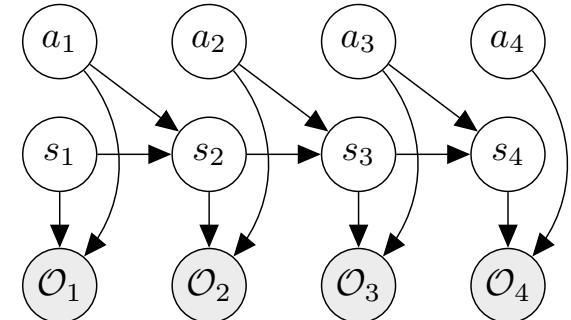




Plan

Part 1: Intro to RL and Control as Inference Framework

- ❑ Intro to Reinforcement Learning (RL)
- ❑ RL and Control as Inference: The GM framework
- ❑ Connections to variational inference



Part 2: Max-entropy RL Algorithms

- ❑ Recap and an inferential approach to inverse RL
- ❑ Classical Q-learning and policy gradient methods
- ❑ Soft Q-learning and soft policy gradients

Algorithm 1 Soft Actor-Critic

```
Initialize parameter vectors  $\psi, \bar{\psi}, \theta, \phi$ .  
for each iteration do  
    for each environment step do  
         $a_t \sim \pi_\phi(a_t | s_t)$   
         $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$   
         $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$   
    end for  
    for each gradient step do  
         $\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$   
         $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$  for  $i \in \{1, 2\}$   
         $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$   
         $\bar{\psi} \leftarrow \tau\psi + (1 - \tau)\bar{\psi}$   
    end for  
end for
```

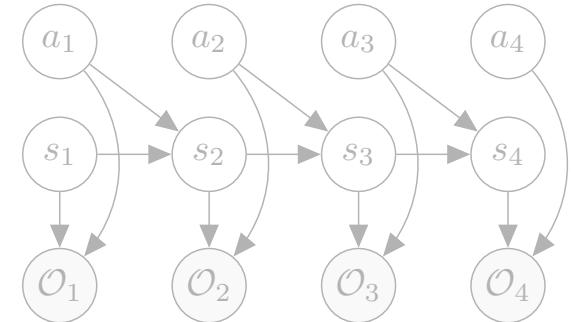




Plan

Part 1: Intro to RL and Control as Inference Framework

- ❑ Intro to Reinforcement Learning (RL)
- ❑ RL and Control as Inference: The GM framework
- ❑ Connections to variational inference



Part 2: Max-entropy RL Algorithms

- ❑ Recap and an inferential approach to inverse RL
- ❑ Classical Q-learning and policy gradient methods
- ❑ Soft Q-learning and soft policy gradients

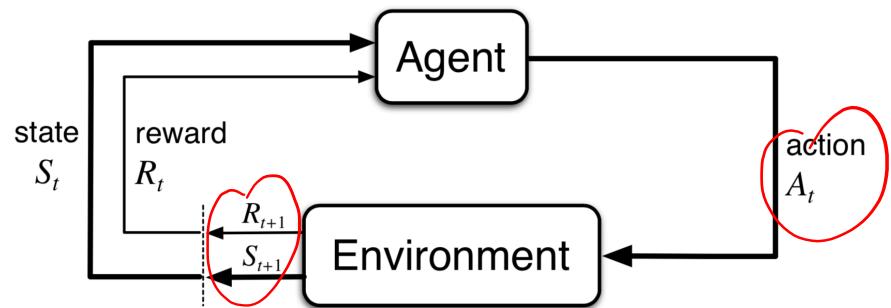
Algorithm 1 Soft Actor-Critic

```
Initialize parameter vectors  $\psi, \bar{\psi}, \theta, \phi$ .  
for each iteration do  
    for each environment step do  
         $a_t \sim \pi_\phi(a_t | s_t)$   
         $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$   
         $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$   
    end for  
    for each gradient step do  
         $\psi \leftarrow \psi - \lambda_V \hat{\nabla}_\psi J_V(\psi)$   
         $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$  for  $i \in \{1, 2\}$   
         $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$   
         $\bar{\psi} \leftarrow \tau\psi + (1 - \tau)\bar{\psi}$   
    end for  
end for
```





Recap: Control as Inference



Initial state

$$s_0 \sim p_0(s)$$

Transition

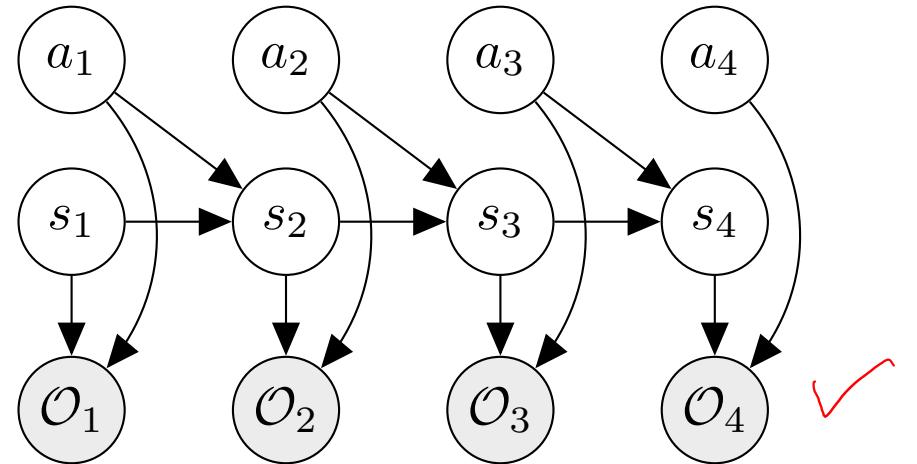
$$s_{t+1} \sim p(s_{t+1} | s_t, a_t)$$

Policy

$$a_t \sim \pi(a_t | s_t)$$

Reward

$$r_t = r(s_t, a_t)$$



Initial state

$$s_0 \sim p_0(s)$$

Transition

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

Policy

$$a_t \sim \pi(a_t | s_t)$$

Reward

$$r_t = r(s_t, a_t)$$

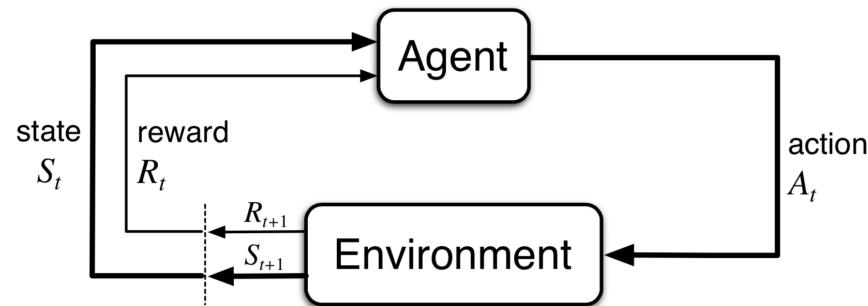
Optimality

$$p(\mathcal{O}_t = 1 | s_t, a_t) = \exp(r(s_t, a_t))$$





Recap: Control as Inference



Initial state

$$s_0 \sim p_0(s)$$

Transition

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

Policy

$$a_t \sim \pi(a_t | s_t)$$

Reward

$$r_t = r(s_t, a_t)$$

In the classical RL setup, we have:

$$V_\pi(s) := \mathbb{E}_\pi \left[\sum_{k=0}^T \gamma^k r_{t+k+1} \mid s_t = s \right]$$

$$Q_\pi(s, a) := \mathbb{E}_\pi \left[\sum_{k=0}^T \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right]$$

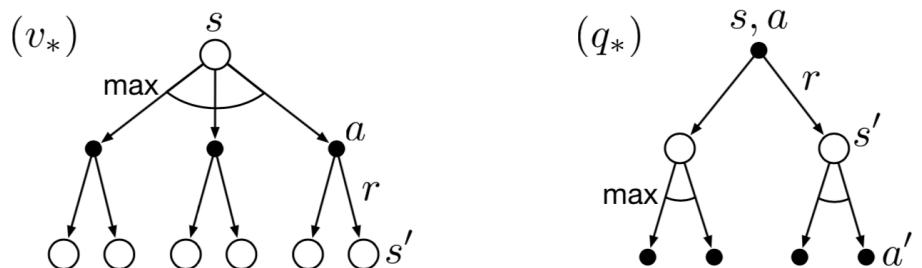


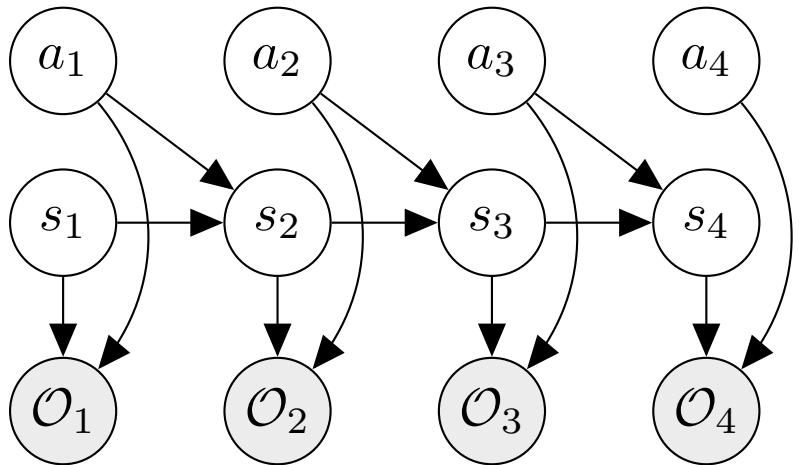
Figure 3.4: Backup diagrams for v_* and q_*

$$\pi_\star(a \mid s) = \delta \left(a = \arg \max_a Q_\star(s, a) \right)$$





Recap: Control as Inference



Initial state
Transition
Policy
Reward
Optimality

$$\begin{aligned}s_0 &\sim p_0(s) \\ s_{t+1} &\sim p(s_{t+1} | s_t, a_t) \\ a_t &\sim \pi(a_t | s_t) \\ r_t &= r(s_t, a_t) \\ p(O_t = 1 | s_t, a_t) &= \exp(r(s_t, a_t))\end{aligned}$$

Running inference in this GM allows us to compute:

$$p(\tau | \mathcal{O}_{1:T}) \propto \prod_{t=1}^T p(s_{t+1} | s_t, a_t)$$

$$\text{let } V_t(\mathbf{s}_t) = \log \beta_t(\mathbf{s}_t)$$

$$\text{let } Q_t(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t) \quad \underbrace{p(O_{t:T} | s_t, a_t)}$$

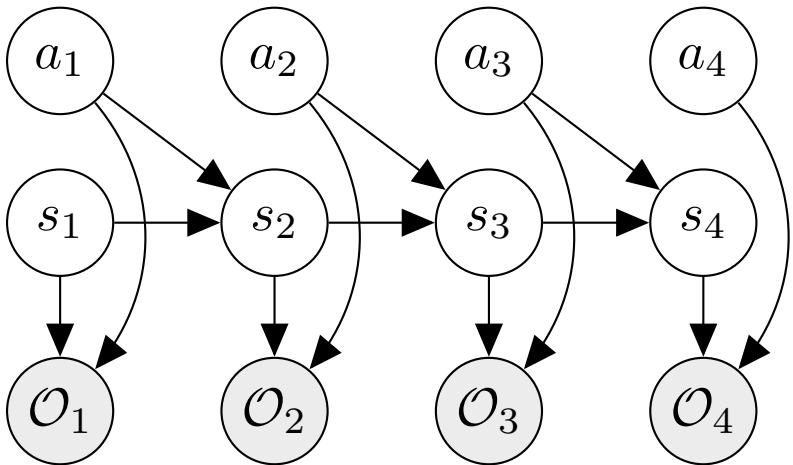
$$V(\mathbf{s}_t) = \log \int \exp(Q(\mathbf{s}_t, \mathbf{a}_t) + \log p(\mathbf{a}_t | \mathbf{s}_t)) \mathbf{a}_t$$

$$\underbrace{p(a_t | s_t, \mathcal{O}_{1:T})}_{\text{softmax}} = \exp(Q_t(s_t, a_t) - V_t(s_t))$$





Recap: Control as Inference



Initial state

$$s_0 \sim p_0(s)$$

Transition

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

Policy

$$a_t \sim \pi(a_t | s_t)$$

Reward

$$r_t = r(s_t, a_t)$$

Optimality

$$p(O_t = 1 | s_t, a_t) = \exp(r(s_t, a_t))$$

Which objective does inference optimize?

$$\begin{aligned}
 & - D_{\text{KL}} (\hat{p}(\tau) || p(\tau)) = \\
 & \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim \hat{p}(s_t, a_t)} [r(s_t, a_t)] + \\
 & \mathbb{E}_{(s_t) \sim \hat{p}(s_t)} [\mathcal{H}(\pi(a_t | s_t))]
 \end{aligned}$$

$\hat{p}(\tau | \mathcal{O}_1, \mathbf{f})$

- For deterministic dynamics, get it directly
- For stochastic dynamics, obtain it from the ELBO on the evidence



Policy Gradient Methods

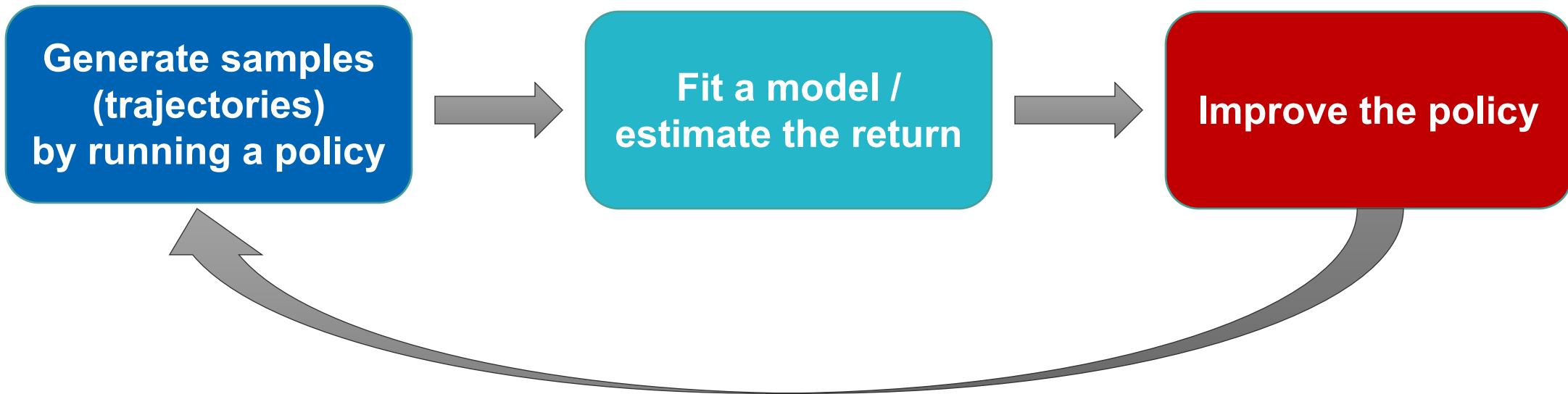
[Barto & Sutton's textbook; Sergey Levine's lectures]



How do we learn in RL?

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

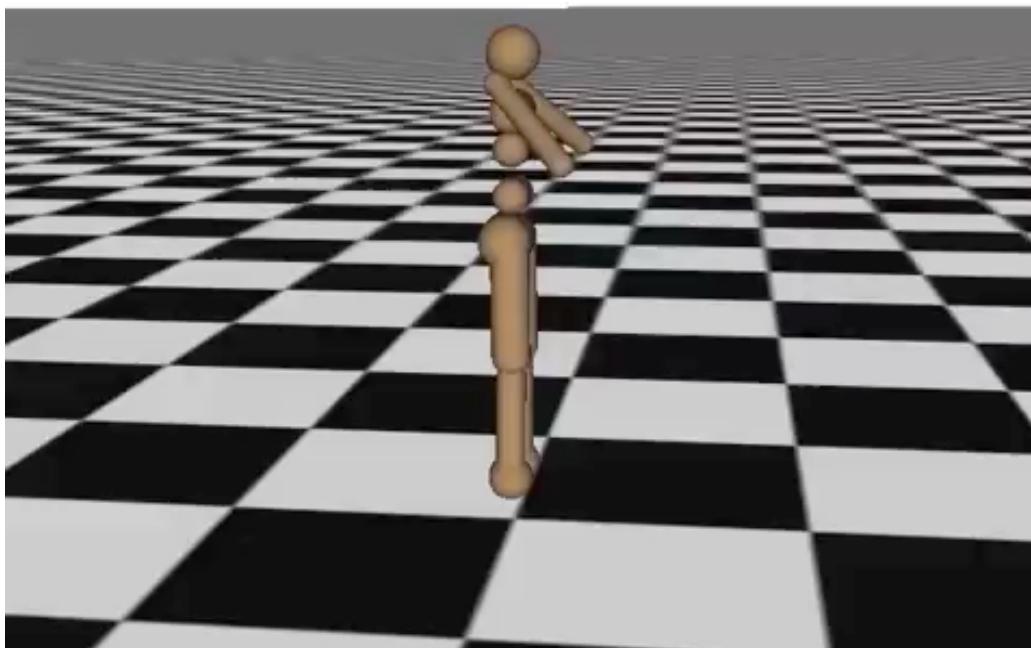
A red circle highlights the θ^* term, and a red bracket highlights the summation term $\sum_{t=1}^T r(s_t, a_t)$. A red arrow points from the highlighted bracket to the right side of the equation.





How do we learn in RL?

Iteration 0



Video from Schulman et al. (2016). High-dimensional continuous control using generalized advantage estimation.





Types of RL algorithms

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

- **Policy gradients:** directly optimize the above stochastic objective
- **Value-based:** estimate V-function or Q-function of the optimal policy (no explicit policy; the policy is derived from the value function)
- **Actor-critic:** estimate V-/Q-function under the current policy and use it to improve the policy
- **Model-based methods:** estimate the transition model $p(s_{t+1}|s_t, a_t)$ and...
 - Use it for planning (plug-in the model into the objective, optimize it w.r.t. a sequence of actions, pick the first action in the best sequence)
 - Use it to improve the policy (e.g., MCTS distillation in AlphaGo)





Types of RL algorithms

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

- **Policy gradients:** directly optimize the above stochastic objective
- **Value-based:** estimate V-function or Q-function of the optimal policy (no explicit policy; the policy is derived from the value function)
- **Actor-critic:** estimate V-/Q-function under the current policy and use it to improve the policy
- **Model-based methods:** estimate the transition model $p(s_{t+1}|s_t, a_t)$ and...
 - Use it for planning (plug-in the model into the objective, optimize it w.r.t. a sequence of actions, pick the first action in the best sequence)
 - Use it to improve the policy (e.g., MCTS distillation in AlphaGo)





Policy gradients

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$\underbrace{\qquad\qquad\qquad}_{J(\theta)}$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t})$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int r(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau$$

$s_1 = \text{env.init}()$

for $t = 1 \dots T$:

$a_t \sim \pi_{\theta}(a | s_t)$

$s_{t+1}, r_t = \text{env.step}(a_t)$





Policy gradients

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$\underbrace{\qquad\qquad\qquad}_{J(\theta)}$

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

" "

$$\nabla_{\theta} \log p_{\theta}(\tau)$$

↓

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t})$$

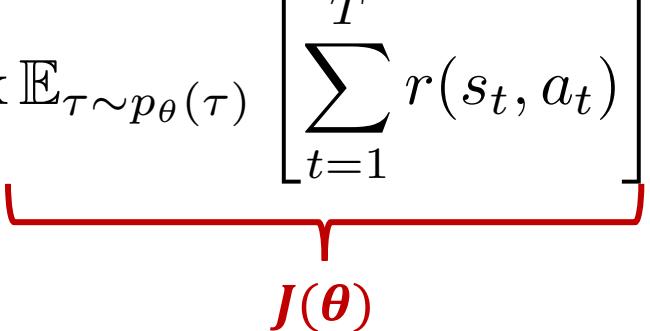
$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int r(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau = \int r(\tau) \underline{\underline{p_{\theta}(\tau)}} \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$





Policy gradients

$$\theta^* = \arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$


 $J(\theta)$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^T r(s_t, a_t) \right] \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t})$$

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau)] = \int r(\tau) \nabla_{\theta} p_{\theta}(\tau) d\tau = \int r(\tau) p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [r(\tau) \underline{\nabla_{\theta} \log p_{\theta}(\tau)}] \end{aligned}$$

For more on how to compute derivatives of stochastic objectives see:

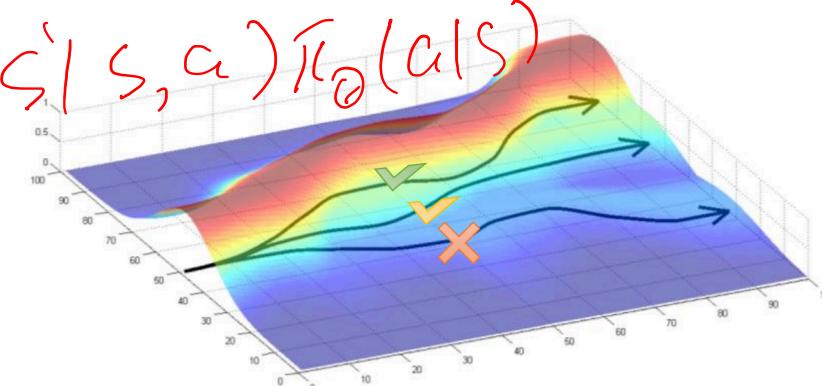
→ Schulman et al. (2015) Gradient estimation using stochastic computation graphs.

Foerster, Farquhar*, A.* et al. (2018) DiCE: The Infinitely Differentiable Monte-Carlo Estimator.





Policy gradients



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

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

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]$$

$$\nabla_\theta J_{ML}(\theta) = \mathbb{E}_{\tau \sim p_{\text{expert}}(\tau)} \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right] \left(\sum_{t=1}^T \gamma_e^{t-1} r_e(s_t, a_t) \right)$$





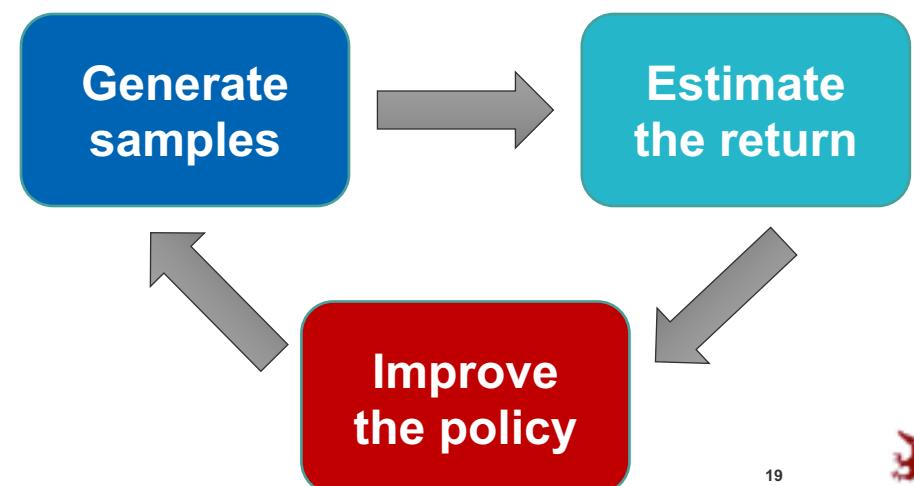
Policy gradients

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\underbrace{\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right)}_{\text{Generate samples}} \underbrace{\left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)}_{\text{Estimate the return}} \right] \quad \checkmark$$

REINFORCE algorithm:

1. sample $\{\tau_i\}_{i=1}^N$ under $\pi_{\theta}(a_t | s_t)$

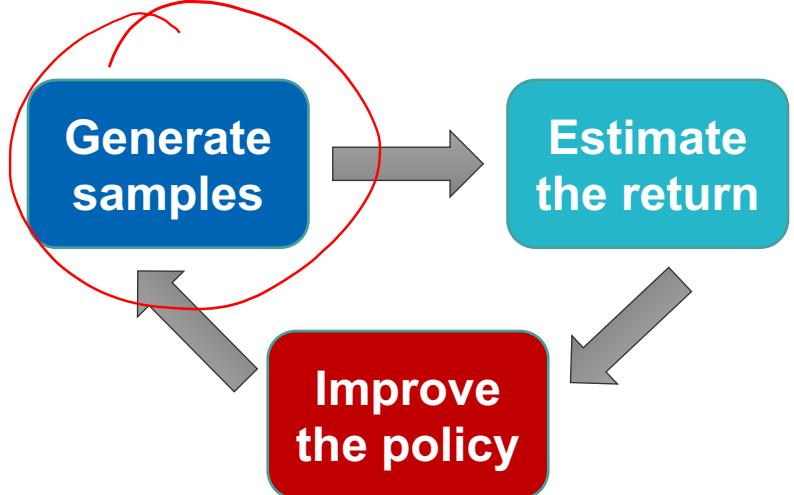




Policy gradients: Summary

REINFORCE algorithm:

1. sample $\{\tau_i\}_{i=1}^N$ under $\pi_\theta(a_t | s_t)$
2. $\hat{J}(\theta) = \sum_i \left(\sum_t \underbrace{\log \pi_\theta(a_{i,t} | s_{i,t})}_{\text{action causality}} \right) \left(\sum_t \underbrace{r(s_{i,t}, a_{i,t})}_{\text{rewards}} \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta \hat{J}(\theta)$



- Represent a policy with a parametric function (e.g., neural net) and learn it by optimizing the REINFORCE objective
- Relationship between PG objective and MLE objective: rewards reweight the samples
- REINFORCE gradients are often extremely high-variance → make use of *action causality* + value estimators to reduce the variance
(look up Sutton & Barto's textbook, Ch. 13 or check out a deep RL course)



Q-learning

[Barto & Sutton's textbook; Sergey Levine's lectures]



Can we get rid of the dependence on the policy?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right]$$

- Recall that if we have access to an optimal $Q_{\star}(s, a)$, it gives us right away a corresponding optimal greedy (deterministic) policy:

$$\pi_{\star}(a | s) = \delta \left(a = \arg \max_a Q_{\star}(s, a) \right)$$

- If we don't have access to an optimal $\underline{Q_{\star}(s, a)}$, we can still try:

$$\underline{a_t} = \arg \max_a [Q_{\pi}(a, s_t)]$$





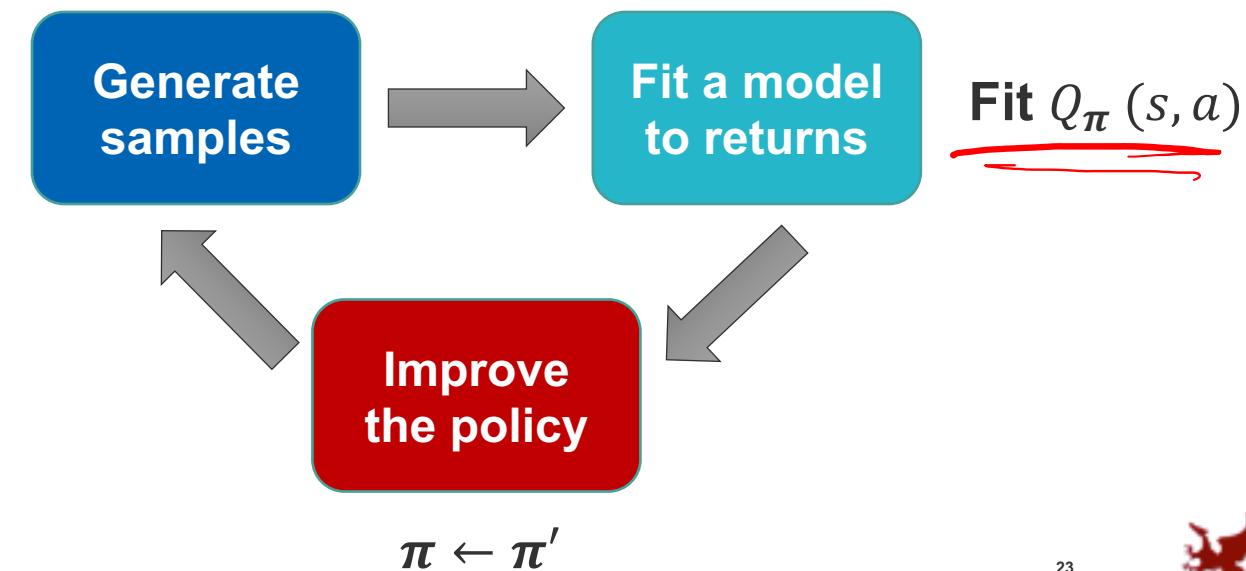
Policy iteration

- Can we learn a Q-function through interaction with the environment without a policy?

$$\pi'(a_t | s_t) = \delta \left(a_t = \arg \max_a [Q_{\underline{\pi}}(a, s_t)] \right) \geq \bar{\pi}$$

Policy iteration:

1. evaluate $\underline{Q_{\pi}(s, a)}$
2. update $\pi \leftarrow \pi'$

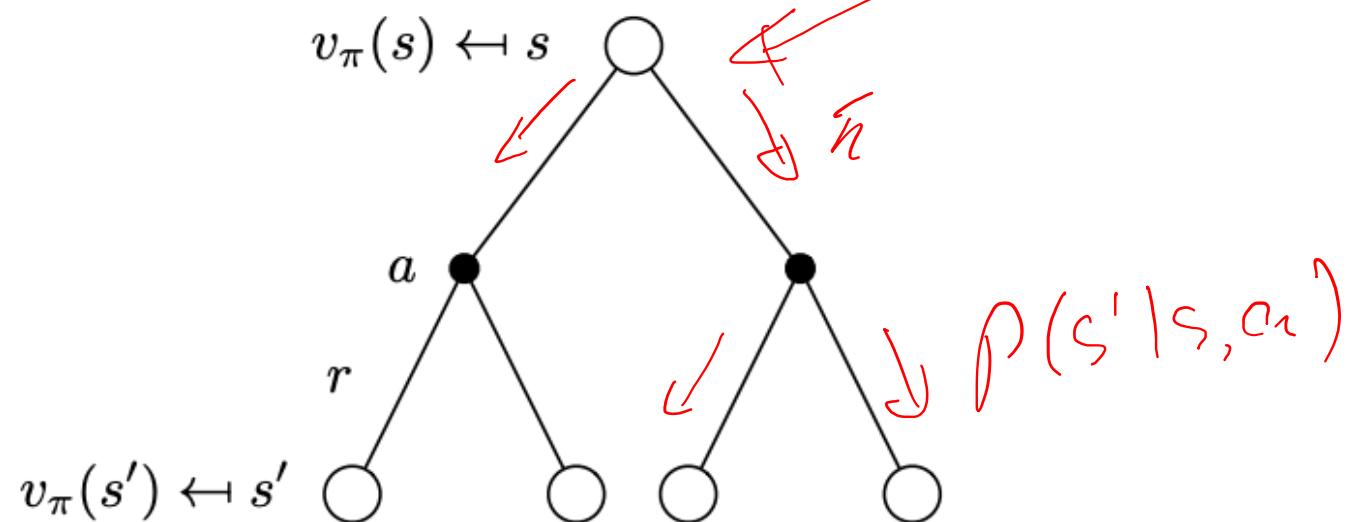




Policy iteration via dynamic programming

Policy iteration:

1. evaluate $Q_{\pi}(s, a) = \underline{r(s, a)} + \gamma \overline{\mathbb{E}_{s' \sim p(s'|s, a)} [V_{\pi}(s')]} \quad \text{with } \overbrace{\gamma}^{\leftarrow} \quad \overbrace{\mathbb{E}}^{\leftarrow} \quad \overbrace{p(s'|s, a)}^{\leftarrow}$
2. update $\pi \leftarrow \pi'$





Policy iteration via dynamic programming

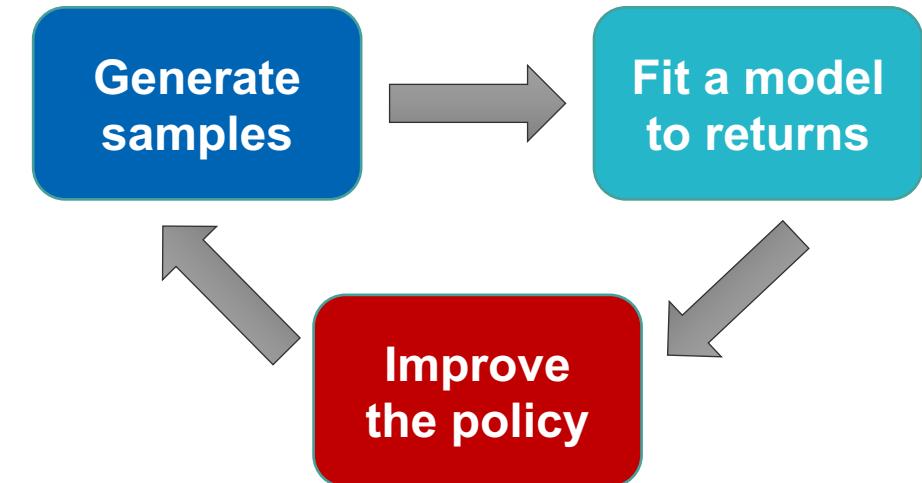
$$Q_{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [V_{\pi}(s')]$$

Policy iteration:

- ✓ 1. evaluate $Q_{\pi}(s, a)$
- 2. update $\underline{\pi \leftarrow \pi'}$

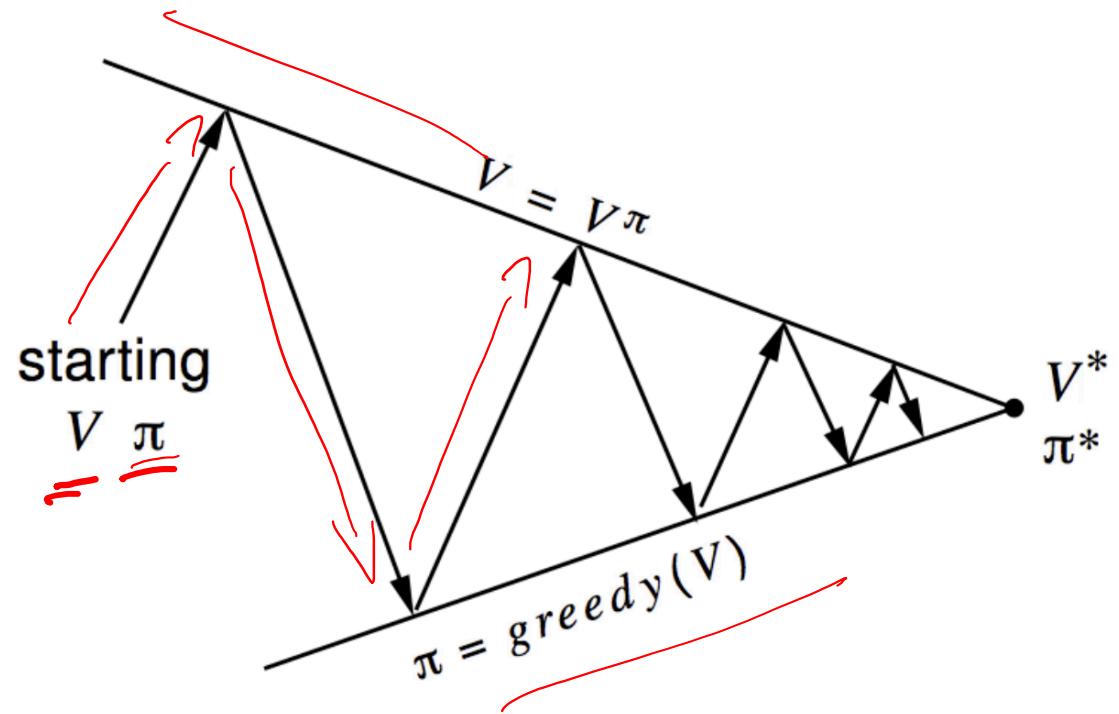
Policy evaluation:

$$\underline{\underline{V_{\pi}(s) \leftarrow r(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim p(s'|s,\pi(s))} [V_{\pi}(s')]}}$$





Policy iteration



Policy evaluation Estimate v_π

Iterative policy evaluation

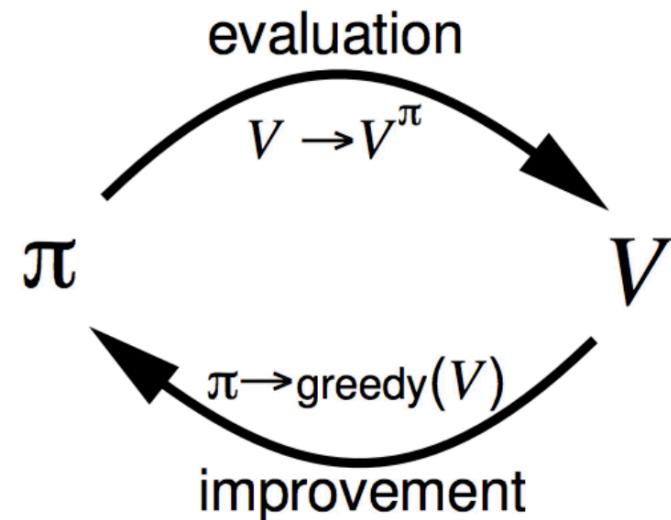
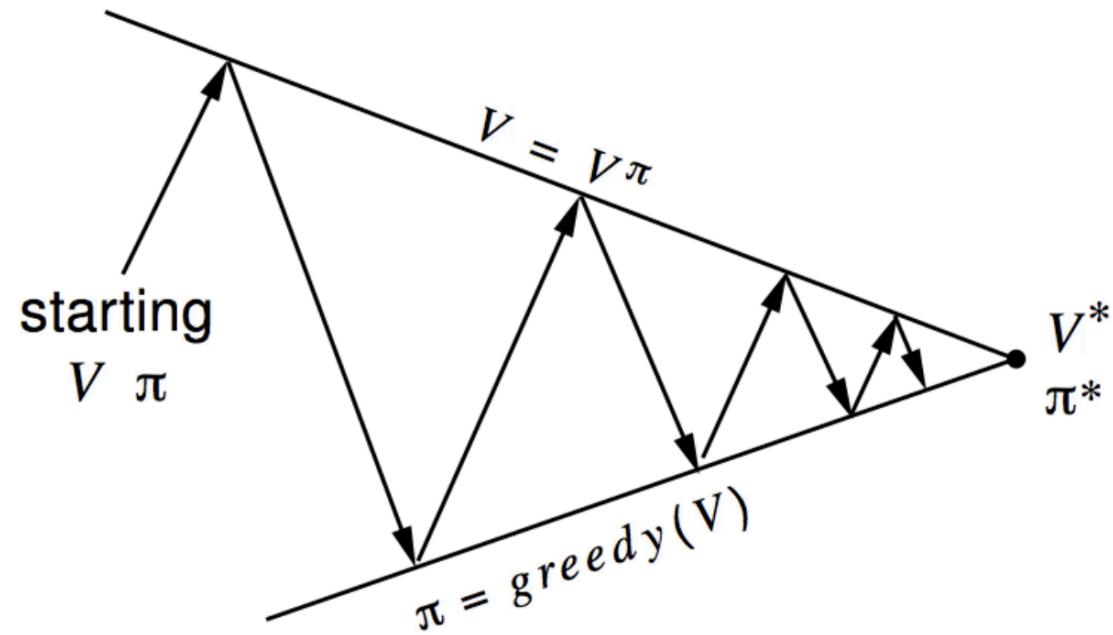
Policy improvement Generate $\pi' \geq \pi$

Greedy policy improvement





Policy iteration



Policy evaluation Estimate v_π

Iterative policy evaluation

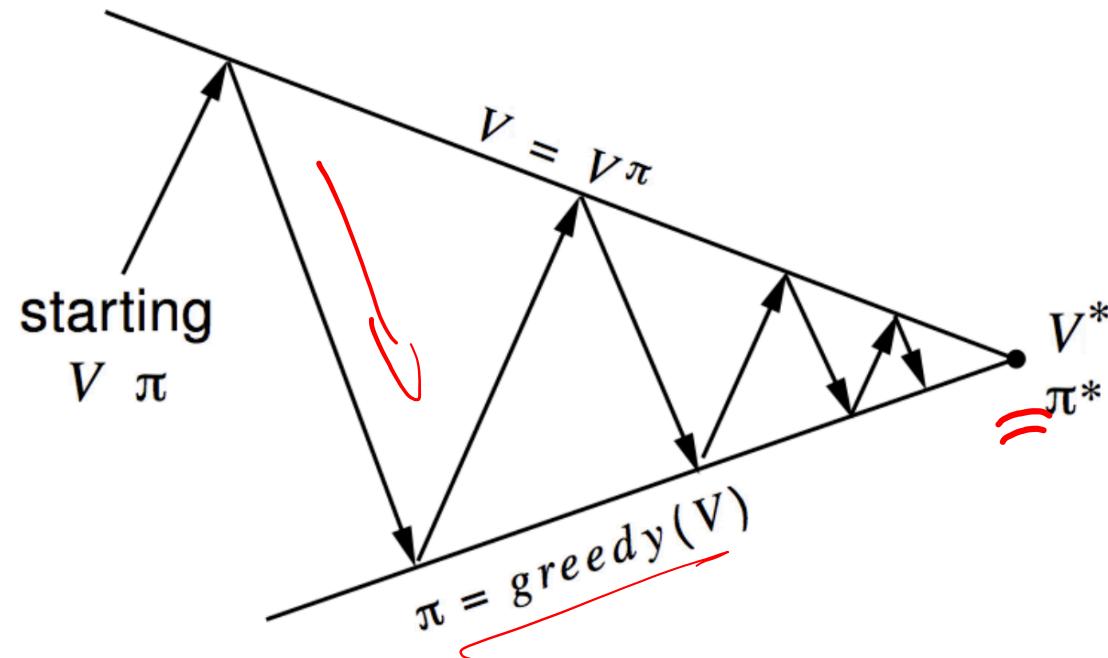
Policy improvement Generate $\pi' \geq \pi$

Greedy policy improvement





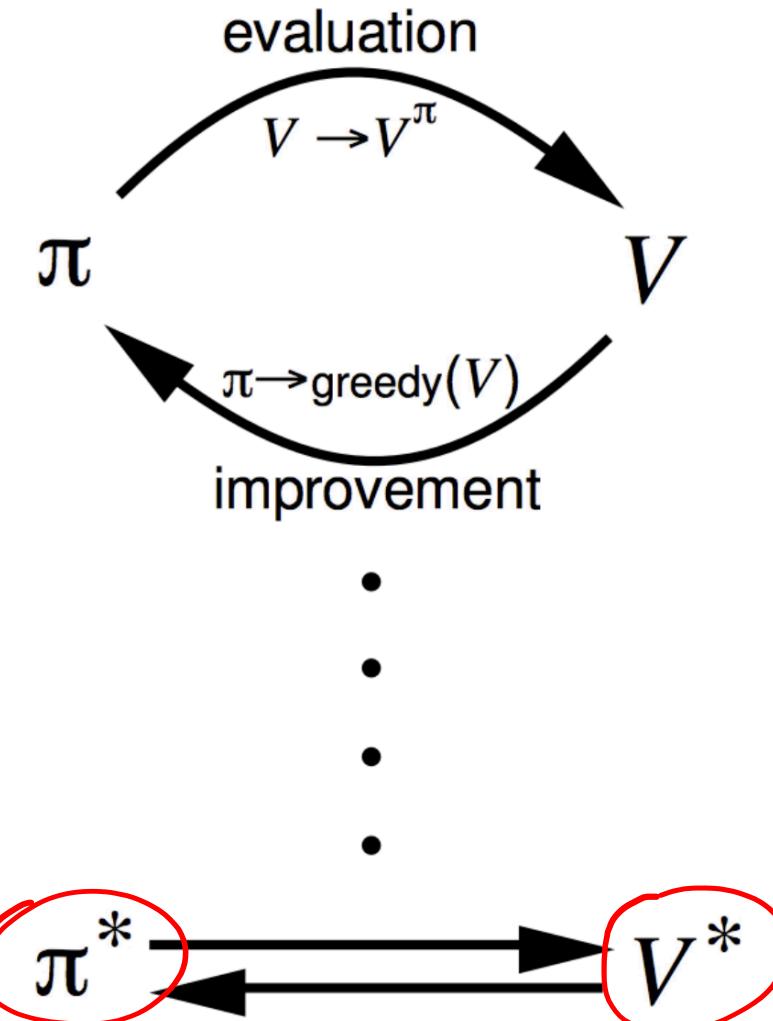
Policy iteration



Policy evaluation Estimate v_π

Iterative policy evaluation

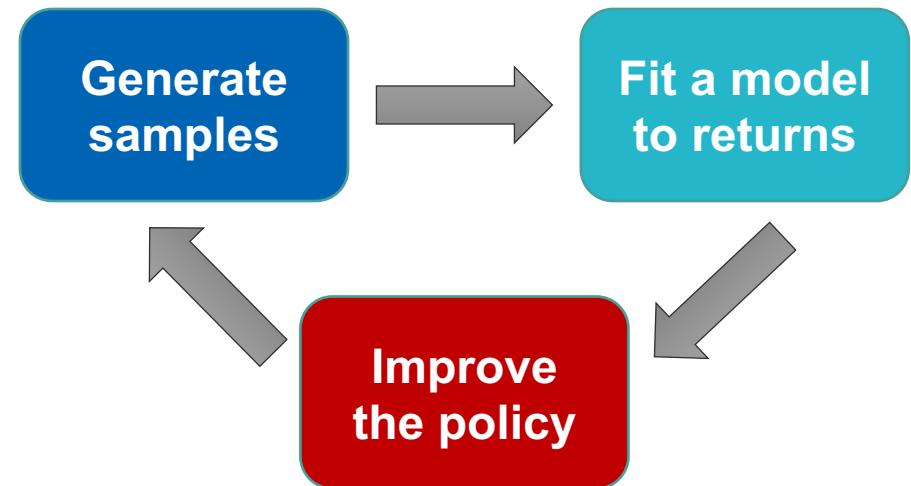
Policy improvement Generate $\pi' \geq \pi$
Greedy policy improvement





Value iteration

- Can we get rid of the policy?



$$\pi'(a_t | s_t) = \delta \left(a_t = \arg \max_a [Q_{\pi}(a, s_t)] \right)$$

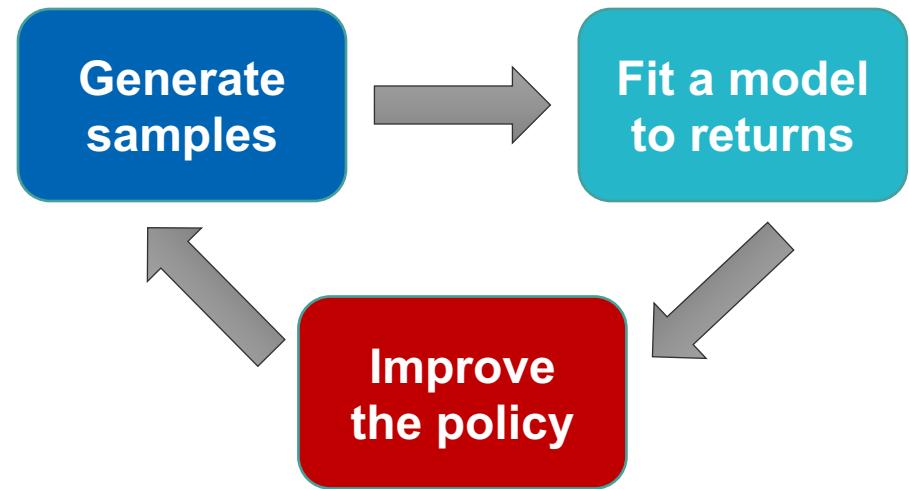
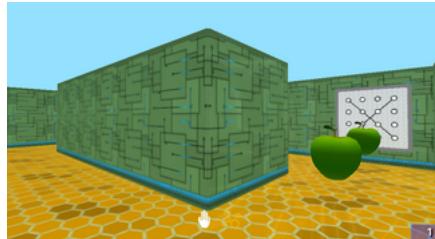
Value iteration:

1. set $\overline{Q(s, a)} \leftarrow r(\underline{s}, \underline{a}) + \gamma \mathbb{E}_{s' \sim p(s'|s, a)}[V(s')]$
2. set $\overline{V(s)} \leftarrow \max_a Q(s, a)$ ← \checkmark , \check{Q}_α





Fitted Q-iteration



- If the state space is high-dimensional, let's represent $Q_\phi(s, a)$ with a parametric function instead of a tabular representation.

Fitted Q-iteration:

1. set $y_i \leftarrow r(\underline{s}_i, \underline{a}_i) + \gamma \mathbb{E}_{s' \sim p(s'|s, a)} [V_\phi(\underline{s}'_i)]$
2. set $\underline{\phi} \leftarrow \arg \min_{\phi} \sum_i \|Q_\phi(\underline{s}_i, \underline{a}_i) - \underline{y}_i\|^2$

$\max_a Q_\phi(s', a)$

Handwritten red annotations: A red circle highlights the \max_a term. A red oval encloses the entire expression $Q_\phi(s', a)$. A red arrow points from the handwritten \max_a to the \max_a in the equation.





Fitted Q-iteration

- Here's our policy-independent algorithm:
 1. collect a dataset $\{(s_i, a_i, s'_i, r_i)\}$ under some policy

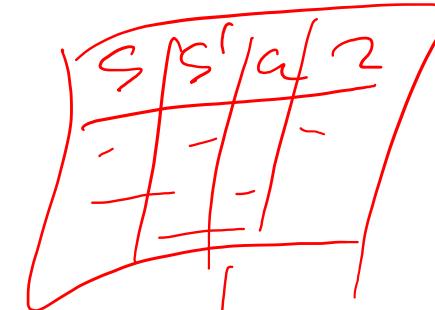




Fitted Q-iteration

- Here's our policy-independent algorithm:

- collect a dataset $\{(s_i, a_i, s'_i, r_i)\}$ under some policy
 - set $y_i \leftarrow r_i + \gamma \max_a Q_\phi(s'_i, a)$
 - set $\underline{\phi} \leftarrow \arg \min_{\phi} \sum_i \|Q_\phi(s_i, a_i) - y_i\|^2$
- $\times \times$



$$\mathcal{E} = \frac{1}{2} E_{(\mathbf{s}, \mathbf{a}) \sim \beta} \left[Q_\phi(\mathbf{s}, \mathbf{a}) - [r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q_\phi(\mathbf{s}', \mathbf{a}')] \right]$$

if $\mathcal{E} = 0$, then $Q_\phi(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q_\phi(\mathbf{s}', \mathbf{a}') \leftarrow$

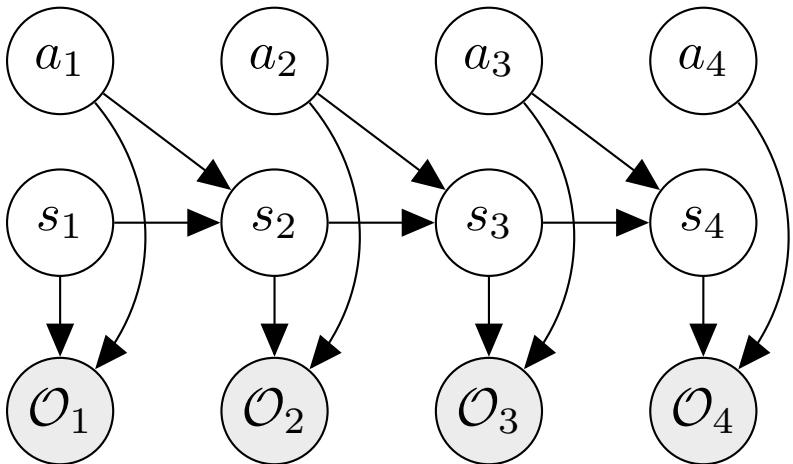
this is an *optimal* Q-function, corresponding to optimal policy π'



Soft Policy Gradients and Q-learning



Recap: Control as Inference



Initial state
Transition
Policy
Reward
Optimality

$$\begin{aligned}s_0 &\sim p_0(s) \\ s_{t+1} &\sim p(s_{t+1} \mid s_t, a_t) \\ a_t &\sim \pi(a_t \mid s_t) \\ r_t &= r(s_t, a_t) \\ p(\mathcal{O}_t = 1 \mid s_t, a_t) &= \exp(r(s_t, a_t))\end{aligned}$$

Which objective does inference optimize?

$$- D_{\text{KL}} (\hat{p}(\tau) \| p(\tau)) = \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim \hat{p}(s_t, a_t)} [r(s_t, a_t)] + \mathbb{E}_{(s_t) \sim \hat{p}(s_t)} [\mathcal{H}(\pi(a_t \mid s_t))]$$

- For deterministic dynamics, get it directly
- For stochastic dynamics, obtain it from the ELBO on the evidence



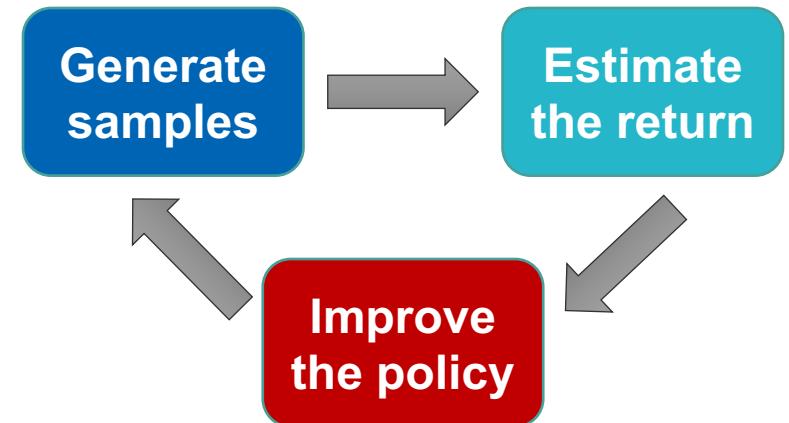


Soft policy gradients

$$\theta^* = \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_{\theta}(s_t, a_t)} [r(s_t, a_t)]$$

$J(\theta)$

$$\nabla_{\theta} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_{\theta}(s_t, a_t)} [r(s_t, a_t)] + \mathbb{E}_{(s_t) \sim \hat{p}(s_t)} [\mathcal{H}(\pi(a_t | s_t))] =$$



~~Play & Set $O_{\pi^{\theta}(T)}$~~

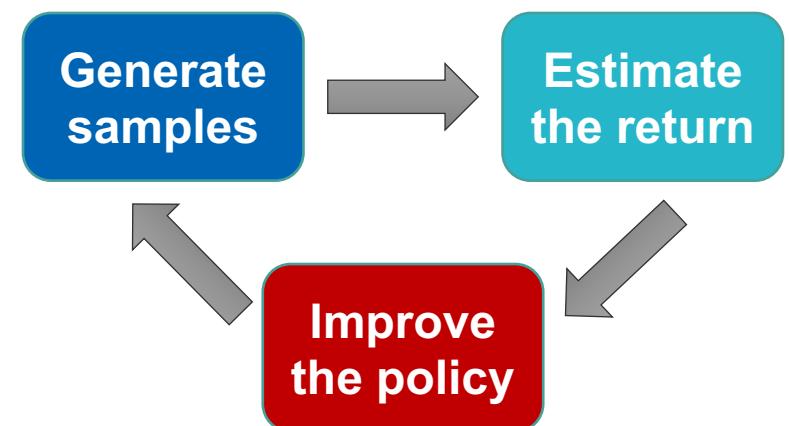




Soft policy gradients

$$\theta^* = \arg \max_{\theta} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_\theta(s_t, a_t)} [r(s_t, a_t)]$$

$J(\theta)$



$$\sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_\theta(s_t, a_t)} [r(s_t, a_t)] + \mathbb{E}_{(s_t) \sim \hat{p}(s_t)} [\mathcal{H}(\pi(a_t | s_t))] =$$
$$\sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_\theta(s_t, a_t)} [r(s_t, a_t) - \log \pi(a_t | s_t)]$$





$$p(a_t | s_t, \phi_{1:T}) = \exp(Q - V)$$

Relationship to Q-learning

$$\nabla_{\theta} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim p_{\theta}(s_t, a_t)} [r(s_t, a_t) - \log \pi(a_t | s_t)]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(r(s_t, a_t) + \underbrace{\left(\sum_{t'=t+1}^T r(s_{t'}, a_{t'}) - \log \pi_{\theta}(a_{t'} | s_{t'}) \right)}_{\approx Q(s_{t+1}, a_{t+1})} - \log \pi_{\theta}(a_t | s_t) - 1 \right)$$

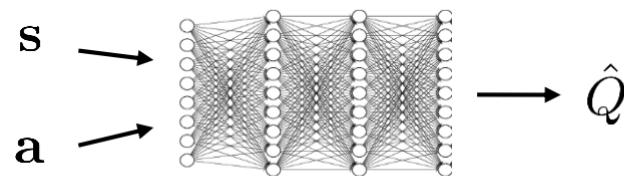
$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} Q_{\theta}(s_t, a_t) \left(r(s_t, a_t) + \text{soft max}_{a_{t+1}} Q_{\theta}(s_{t+1}, a_{t+1}) - Q_{\theta}(s_t, a_t) \right)$$
$$\hat{V}(s_t) = \text{soft sum } \exp(Q)$$





Soft Q-learning

Learned (neural network) Q-function: $Q_\theta(\mathbf{s}, \mathbf{a})$



1. collect a dataset $\{(s_i, a_i, s'_i, r_i)\}$ under some policy
2. set $y_i \leftarrow r_i + \gamma \max_a Q_\phi(s'_i, a)$
3. set $\phi \leftarrow \arg \min_\phi \sum_i \|Q_\phi(s_i, a_i) - y_i\|^2$

Q-learning: $\theta \leftarrow \theta + \alpha \nabla_\theta Q_\theta(\mathbf{s}, \mathbf{a})(r(\mathbf{s}, \mathbf{a}) + \gamma V(\mathbf{s}') - Q_\theta(\mathbf{s}, \mathbf{a}))$

target value: $V(\mathbf{s}') = \max_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}')$

soft Q-learning: $\theta \leftarrow \theta + \alpha \nabla_\theta Q_\theta(\mathbf{s}, \mathbf{a})(r(\mathbf{s}, \mathbf{a}) + \gamma V(\mathbf{s}') - Q_\theta(\mathbf{s}, \mathbf{a}))$ ↗

target value: $V(\mathbf{s}') = \underline{\text{soft max}}_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}') = \log \int \exp(Q_\theta(\mathbf{s}', \mathbf{a}')) d\mathbf{a}'$

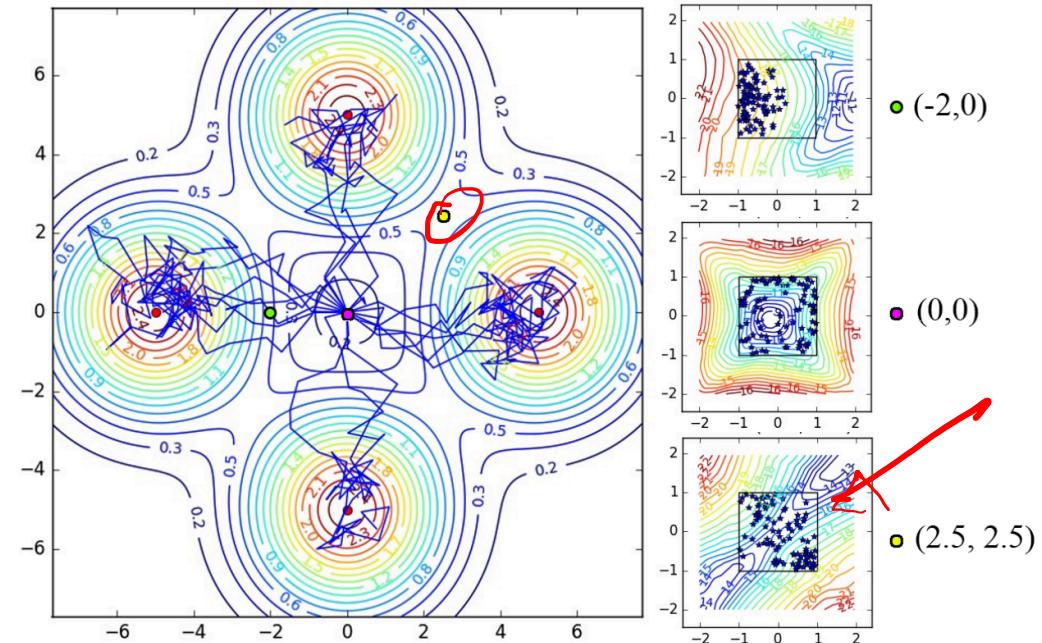
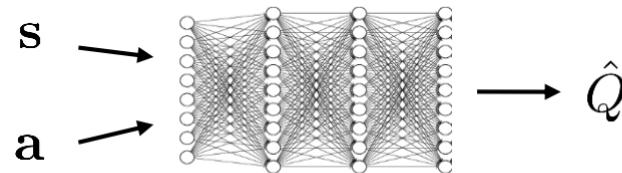
$$\exp\left(\frac{1}{\alpha} Q\right)$$





Soft Q-learning

Learned (neural network) Q-function: $Q_\theta(\mathbf{s}, \mathbf{a})$



$$\text{Q-learning: } \theta \leftarrow \theta + \alpha \nabla_\theta Q_\theta(\mathbf{s}, \mathbf{a})(r(\mathbf{s}, \mathbf{a}) + \gamma V(\mathbf{s}') - Q_\theta(\mathbf{s}, \mathbf{a}))$$

Haarnoja et al. (2017)

$$\text{target value: } V(\mathbf{s}') = \max_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}')$$

$$\text{soft Q-learning: } \theta \leftarrow \theta + \alpha \nabla_\theta Q_\theta(\mathbf{s}, \mathbf{a})(r(\mathbf{s}, \mathbf{a}) + \gamma V(\mathbf{s}') - Q_\theta(\mathbf{s}, \mathbf{a}))$$

$$\text{target value: } V(\mathbf{s}') = \text{soft max}_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}') = \log \int \exp(Q_\theta(\mathbf{s}', \mathbf{a}')) d\mathbf{a}'$$





Benefits of soft optimality

- Improves exploration and prevents entropy collapse
- Empirically, policies are easier to finetune for more specific tasks
- Better robustness (due to wider coverage of states)
- Reduces to hard optimality (by increasing the magnitude of the rewards)
- Good model for human behavior



