# Reframing Control as an Inference Problem

CS 285: Deep Reinforcement Learning, Decision Making, and Control Sergey Levine

#### **Class Notes**

- 1. Homework 3 is out! Due Oct 21
  - Start early, this one will take a bit longer!

# Today's Lecture

- 1. Does reinforcement learning and optimal control provide a reasonable model of human behavior?
- 2. Is there a better explanation?
- 3. Can we derive optimal control, reinforcement learning, and planning as *probabilistic inference*?
- 4. How does this change our RL algorithms?
- 5. (next week) We'll see this is crucial for *inverse* reinforcement learning
- Goals:
  - Understand the connection between inference and control
  - Understand how specific RL algorithms can be instantiated in this framework
  - Understand why this might be a good idea

### Optimal Control as a Model of Human Behavior



# What if the data is **not** optimal?



some mistakes matter more than others!

behavior is **stochastic** 

but good behavior is still the most likely



# A probabilistic graphical model of decision making



# Why is this interesting?

![](_page_6_Figure_1.jpeg)

- Can model suboptimal behavior (important for inverse RL)
- Can apply inference algorithms to solve control and planning problems
- Provides an explanation for why stochastic behavior might be preferred (useful for exploration and transfer learning)

# Inference = planning

![](_page_7_Figure_1.jpeg)

#### how to do inference?

- 1. compute backward messages  $\beta_t(\mathbf{s}_t, \mathbf{a}_t) = p(\mathcal{O}_{t:T} | \mathbf{s}_t, \mathbf{a}_t)$
- 2. compute policy  $p(\mathbf{a}_t | \mathbf{s}_t, \mathcal{O}_{1:T})$
- 3. compute forward messages  $\alpha_t(\mathbf{s}_t) = p(\mathbf{s}_t | \mathcal{O}_{1:t-1})$

#### Backward messages

![](_page_8_Figure_1.jpeg)

#### A closer look at the backward pass

1. set  $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$ 2. set  $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ for t = T - 1 to 1:  $\beta_t(\mathbf{s}_t, \mathbf{a}_t) = p(\mathcal{O}_t | \mathbf{s}_t, \mathbf{a}_t) E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\beta_{t+1}(\mathbf{s}_{t+1})]$  $\beta_t(\mathbf{s}_t) = E_{\mathbf{a}_t \sim p(\mathbf{a}_t | \mathbf{s}_t)}[\beta_t(\mathbf{s}_t, \mathbf{a}_t)]$ "optimistic" transition (not a good idea!)  $Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E[\exp(V_{t+1}(\mathbf{s}_{t+1}))]$ let  $V_t(\mathbf{s}_t) = \log \beta_t(\mathbf{s}_t)$ deterministic transition:  $Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + V_{t+1}(\mathbf{s}_{t+1})$ let  $Q_t(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t)$ we'll come back to the stochastic case later!  $V_t(\mathbf{s}_t) = \log \int \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t)) d\mathbf{a}_t$ 

value iteration algorithm:

 $V_t(\mathbf{s}_t) \to \max_{\mathbf{a}_t} Q_t(\mathbf{s}_t, \mathbf{a}_t)$  as  $Q_t(\mathbf{s}_t, \mathbf{a}_t)$  gets bigger!

#### Backward pass summary

![](_page_10_Figure_1.jpeg)

 $\beta_t(\mathbf{s}_t, \mathbf{a}_t) = p(\mathcal{O}_{t:T} | \mathbf{s}_t, \mathbf{a}_t)$ 

probability that we can be optimal at steps t through T given that we take action  $\mathbf{a}_t$  in state  $\mathbf{s}_t$ 

for t = T - 1 to 1:  $\beta_t(\mathbf{s}_t, \mathbf{a}_t) = p(\mathcal{O}_t | \mathbf{s}_t, \mathbf{a}_t) E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\beta_{t+1}(\mathbf{s}_{t+1})]$  compute recursively from t = T to t = 1 $\beta_t(\mathbf{s}_t) = E_{\mathbf{a}_t \sim p(\mathbf{a}_t | \mathbf{s}_t)} [\beta_t(\mathbf{s}_t, \mathbf{a}_t)]$ 

let  $V_t(\mathbf{s}_t) = \log \beta_t(\mathbf{s}_t)$ let  $Q_t(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t)$  log of  $\beta_t$  is "Q-function-like"

#### The action prior

remember this?

$$p(\mathcal{O}_{t+1:T}|\mathbf{s}_{t+1}) = \int p(\mathcal{O}_{t+1:T}|\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) p(\mathbf{a}_{t+1}|\mathbf{s}_{t+1}) d\mathbf{a}_{t+1}$$
$$\beta_t(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$$

("soft max")

what if the action prior is not uniform?

$$V(\mathbf{s}_t) = \log \int \exp(Q(\mathbf{s}_t, \mathbf{a}_t) + \log p(\mathbf{a}_t | \mathbf{s}_t)) \mathbf{a}_t$$
$$Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log F[\exp(V(\mathbf{s}_{t+1}))]$$

$$Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E[\exp(V(\mathbf{s}_{t+1}))]$$

let 
$$\tilde{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log p(\mathbf{a}_t | \mathbf{s}_t) + \log E[\exp(V(\mathbf{s}_{t+1}))]$$

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

can always fold the action prior into the reward! uniform action prior can be assumed without loss of generality

# Policy computation

![](_page_12_Figure_1.jpeg)

2. compute policy  $p(\mathbf{a}_t | \mathbf{s}_t, \mathcal{O}_{1:T})$ 

 $\beta_t(\mathbf{s}_t, \mathbf{a}_t) = p(\mathcal{O}_{t:T} | \mathbf{s}_t, \mathbf{a}_t)$  $\beta_t(\mathbf{s}_t) = p(\mathcal{O}_{t:T} | \mathbf{s}_t)$ 

 $p(\mathbf{a}_{t}|\mathbf{s}_{t}, \mathcal{O}_{1:T}) = \pi(\mathbf{a}_{t}|\mathbf{s}_{t})$   $= p(\mathbf{a}_{t}|\mathbf{s}_{t}, \mathcal{O}_{t:T})$   $= \frac{p(\mathbf{a}_{t}, \mathbf{s}_{t}|\mathcal{O}_{t:T})}{p(\mathbf{s}_{t}|\mathcal{O}_{t:T})}$   $= \frac{p(\mathcal{O}_{t:T}|\mathbf{a}_{t}, \mathbf{s}_{t})p(\mathbf{a}_{t}, \mathbf{s}_{t})/p(\mathcal{O}_{t:T})}{p(\mathcal{O}_{t:T}|\mathbf{s}_{t})p(\mathbf{s}_{t})/p(\mathcal{O}_{t:T})}$   $= \frac{p(\mathcal{O}_{t:T}|\mathbf{a}_{t}, \mathbf{s}_{t})}{p(\mathcal{O}_{t:T}|\mathbf{s}_{t})} \frac{p(\mathbf{a}_{t}, \mathbf{s}_{t})}{p(\mathbf{s}_{t})} = \frac{\beta_{t}(\mathbf{s}_{t}, \mathbf{a}_{t})}{\beta_{t}(\mathbf{s}_{t})}$ 

#### Policy computation with value functions

for t = T - 1 to 1:

$$Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E[\exp(V_{t+1}(\mathbf{s}_{t+1}))]$$
$$V_t(\mathbf{s}_t) = \log \int \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t))\mathbf{a}_t$$

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \frac{\beta_t(\mathbf{s}_t, \mathbf{a}_t)}{\beta_t(\mathbf{s}_t)} \qquad V_t(\mathbf{s}_t) = \log \beta_t(\mathbf{s}_t)$$
$$Q_t(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t)$$

 $\pi(\mathbf{a}_t|\mathbf{s}_t) = \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t) - V_t(\mathbf{s}_t)) = \exp(A_t(\mathbf{s}_t, \mathbf{a}_t))$ 

# Policy computation summary

$$\pi(\mathbf{a}_t|\mathbf{s}_t) = \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t) - V_t(\mathbf{s}_t)) = \exp(A_t(\mathbf{s}_t, \mathbf{a}_t))$$

with temperature:  $\pi(\mathbf{a}_t | \mathbf{s}_t) = \exp(\frac{1}{\alpha} Q_t(\mathbf{s}_t, \mathbf{a}_t) - \frac{1}{\alpha} V_t(\mathbf{s}_t)) = \exp(\frac{1}{\alpha} A_t(\mathbf{s}_t, \mathbf{a}_t))$ 

- Natural interpretation: better actions are more probable
- Random tie-breaking
- Analogous to Boltzmann exploration
- Approaches greedy policy as temperature decreases

# Forward messages

![](_page_15_Figure_1.jpeg)

# Forward/backward message intersection

![](_page_16_Figure_1.jpeg)

# Forward/backward message intersection

![](_page_17_Figure_1.jpeg)

# Summary

1. Probabilistic graphical model for optimal control

![](_page_18_Figure_2.jpeg)

2. Control = inference (similar to HMM, EKF, etc.)

3. Very similar to dynamic programming, value iteration, etc. (but "soft")

#### Break

# The optimism problem

the inference problem:  $p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T} | \mathcal{O}_{1:T})$ 

marginalizing and conditioning, we get:  $p(\mathbf{a}_t | \mathbf{s}_t, \mathcal{O}_{1:T})$  (the policy)

"given that you obtained high reward, what was your action probability?"

marginalizing and conditioning, we get:  $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t, \mathcal{O}_{1:T}) \neq p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$ 

"given that you obtained high reward, what was your transition probability?"

# Addressing the optimism problem

marginalizing and conditioning, we get:  $p(\mathbf{a}_t | \mathbf{s}_t, \mathcal{O}_{1:T})$  (the policy)  $\leftarrow$  we want this

"given that you obtained high reward, what was your action probability?"

marginalizing and conditioning, we get:  $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t, \mathcal{O}_{1:T}) \neq p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t) \longleftarrow$  but not this!

"given that you obtained high reward, what was your transition probability?"

"given that you obtained high reward, what was your action probability, given that your transition probability did not change?"

can we find another distribution  $q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})$  that is close to  $p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T} | \mathcal{O}_{1:T})$  but has dynamics  $p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ 

where have we seen this before? let  $\mathbf{x} = \mathcal{O}_{1:T}$  and  $\mathbf{z} = (\mathbf{s}_{1:T}, \mathbf{a}_{1:T})$  find  $q(\mathbf{z})$  to approximate  $p(\mathbf{z}|\mathbf{x})$ 

let's try variational inference!

# Control via variational inference

![](_page_22_Figure_1.jpeg)

#### The variational lower bound

 $\log p(\mathbf{x}) \ge E_{\mathbf{z} \sim q(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z})]$ let  $\mathbf{x} = \mathcal{O}_{1:T}$  and  $\mathbf{z} = (\mathbf{s}_{1:T}, \mathbf{a}_{1:T})$ the entropy  $\mathcal{H}(q)$ let  $q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) = p(\mathbf{s}_1) \prod_t p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t | \mathbf{s}_t)$  $\log p(\mathcal{O}_{1:T}) \ge E_{(\mathbf{s}_{1:T},\mathbf{a}_{1:T})\sim q}[\log p(\mathbf{s}_1) + \sum_{t=1}^T \log p(\mathbf{s}_{t+1}|\mathbf{s}_t,\mathbf{a}_t) + \sum_{t=1}^T \log p(\mathcal{O}_t|\mathbf{s}_t,\mathbf{a}_t)$  $-\log p(\mathbf{s}_1) - \sum_{t=1}^T \log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) - \sum_{t=1}^T \log q(\mathbf{a}_t | \mathbf{s}_t)]$  $= E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q} \left| \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) - \log q(\mathbf{a}_{t} | \mathbf{s}_{t}) \right|$  $= \sum_{t} E_{(\mathbf{s}_{t},\mathbf{a}_{t})\sim q} \left[ r(\mathbf{s}_{t},\mathbf{a}_{t}) + \mathcal{H}(q(\mathbf{a}_{t}|\mathbf{s}_{t})) \right] \longleftarrow \text{maximize reward and maximize action entropy!}$ 

# Optimizing the variational lower bound

let  $q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) = p(\mathbf{s}_1) \prod_t p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t | \mathbf{s}_t)$ 

$$\log p(\mathcal{O}_{1:T}) \ge \sum_{t} E_{(\mathbf{s}_t, \mathbf{a}_t) \sim q} \left[ r(\mathbf{s}_t, \mathbf{a}_t) + \mathcal{H}(q(\mathbf{a}_t | \mathbf{s}_t)) \right]$$

base case: solve for  $q(\mathbf{a}_T | \mathbf{s}_T)$ :

$$q(\mathbf{a}_{T}|\mathbf{s}_{T}) = \arg\max E_{\mathbf{s}_{T}\sim q(\mathbf{s}_{T})} \left[ E_{\mathbf{a}_{T}\sim q(\mathbf{a}_{T}|\mathbf{s}_{T})} [r(\mathbf{s}_{T},\mathbf{a}_{T})] + \mathcal{H}(q(\mathbf{a}_{T}|\mathbf{s}_{T}))] \right]$$

$$= \arg\max E_{\mathbf{s}_{T}\sim q(\mathbf{s}_{T})} \left[ E_{\mathbf{a}_{T}\sim q(\mathbf{a}_{T}|\mathbf{s}_{T})} [r(\mathbf{s}_{T},\mathbf{a}_{T}) - \log q(\mathbf{a}_{T}|\mathbf{s}_{T})] \right]$$
optimized when  $q(\mathbf{a}_{T}|\mathbf{s}_{T}) \propto \exp(r(\mathbf{s}_{T},\mathbf{a}_{T}))$ 

$$q(\mathbf{a}_{T}|\mathbf{s}_{T}) = \frac{\exp(r(\mathbf{s}_{T},\mathbf{a}_{T}))}{\int \exp(r(\mathbf{s}_{T},\mathbf{a}_{T}))d\mathbf{a}} = \exp(Q(\mathbf{s}_{T},\mathbf{a}_{T}) - V(\mathbf{s}_{T})) \quad V(\mathbf{s}_{T}) = \log \int \exp(Q(\mathbf{s}_{T},\mathbf{a}_{T}))d\mathbf{a}_{T}$$

 $E_{\mathbf{s}_T \sim q(\mathbf{s}_T)} \left[ E_{\mathbf{a}_T \sim q(\mathbf{a}_T | \mathbf{s}_T)} [r(\mathbf{s}_T, \mathbf{a}_T) - \log q(\mathbf{a}_T | \mathbf{s}_T)] \right] = E_{\mathbf{s}_T \sim q(\mathbf{s}_T)} \left[ E_{\mathbf{a}_T \sim q(\mathbf{a}_T | \mathbf{s}_T)} [V(\mathbf{s}_T)] \right]$ 

# Optimizing the variational lower bound

$$\begin{split} \log p(\mathcal{O}_{1:T}) &\geq \sum_{t} E_{(\mathbf{s}_{t},\mathbf{a}_{t})\sim q} \left[ r(\mathbf{s}_{t},\mathbf{a}_{t}) + \mathcal{H}(q(\mathbf{a}_{t}|\mathbf{s}_{t})) \right] \\ q(\mathbf{a}_{T}|\mathbf{s}_{T}) &= \frac{\exp(r(\mathbf{s}_{T},\mathbf{a}_{T}))}{\int \exp(r(\mathbf{s}_{T},\mathbf{a}_{T}))d\mathbf{a}} = \exp(Q(\mathbf{s}_{T},\mathbf{a}_{T}) - V(\mathbf{s}_{T})) \\ E_{\mathbf{s}_{T}\sim q(\mathbf{s}_{T})} \left[ E_{\mathbf{a}_{T}\sim q(\mathbf{a}_{T}|\mathbf{s}_{T})}[r(\mathbf{s}_{T},\mathbf{a}_{T}) - \log q(\mathbf{a}_{T}|\mathbf{s}_{T})] \right] = E_{\mathbf{s}_{T}\sim q(\mathbf{s}_{T})} \left[ E_{\mathbf{a}_{T}\sim q(\mathbf{a}_{T}|\mathbf{s}_{T})}[V(\mathbf{s}_{T})] \right] \\ q(\mathbf{a}_{t}|\mathbf{s}_{t}) &= \arg \max E_{\mathbf{s}_{t}\sim q(\mathbf{s}_{t})} \left[ E_{\mathbf{a}_{t}\sim q(\mathbf{a}_{t}|\mathbf{s}_{t})}[r(\mathbf{s}_{t},\mathbf{a}_{t}) + E_{\mathbf{s}_{t+1}\sim p(\mathbf{s}_{t+1}|\mathbf{s}_{t},\mathbf{a}_{t})}[V(\mathbf{s}_{t+1})]] + \mathcal{H}(q(\mathbf{a}_{t}|\mathbf{s}_{t})) \right] \\ &= \arg \max E_{\mathbf{s}_{t}\sim q(\mathbf{s}_{t})} \left[ E_{\mathbf{a}_{t}\sim q(\mathbf{a}_{t}|\mathbf{s}_{t})}[Q(\mathbf{s}_{t},\mathbf{a}_{t})] + \mathcal{H}(q(\mathbf{a}_{t}|\mathbf{s}_{t}))] \right] \\ &= \arg \max E_{\mathbf{s}_{t}\sim q(\mathbf{s}_{t})} \left[ E_{\mathbf{a}_{t}\sim q(\mathbf{a}_{t}|\mathbf{s}_{t})}[Q(\mathbf{s}_{t},\mathbf{a}_{t}) - \log q(\mathbf{a}_{t}|\mathbf{s}_{t})] \right] \\ &\text{optimized when } q(\mathbf{a}_{t}|\mathbf{s}_{t}) \propto \exp(Q(\mathbf{s}_{t},\mathbf{a}_{t}) - \log q(\mathbf{a}_{t}|\mathbf{s}_{t})] \right] \\ &V_{t}(\mathbf{s}_{t}) = \log \int \exp(Q_{t}(\mathbf{s}_{t},\mathbf{a}_{t}))d\mathbf{a}_{t} \qquad \int regular \text{ Bellman backup} \\ &q(\mathbf{a}_{t}|\mathbf{s}_{t}) = \exp(Q(\mathbf{s}_{t},\mathbf{a}_{t}) - V(\mathbf{s}_{t})) & not \text{ optimistic} \end{split}$$

not optimistic

#### Backward pass summary - variational

for t = T - 1 to 1:

$$Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + E[(V_{t+1}(\mathbf{s}_{t+1}))]$$
$$V_t(\mathbf{s}_t) = \log \int \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t)) d\mathbf{a}_t$$

value iteration algorithm: 1. set  $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$ 2. set  $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ 

soft value iteration algorithm: 1. set  $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$ 2. set  $V(\mathbf{s}) \leftarrow \text{soft} \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$ 

# Summary

![](_page_27_Figure_1.jpeg)

#### variants:

discounted SOC:  $Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma E[V_{t+1}(\mathbf{s}_{t+1})]$ 

explicit temperature:  $V_t(\mathbf{s}_t) = \alpha \log \int \exp\left(\frac{1}{\alpha}Q_t(\mathbf{s}_t, \mathbf{a}_t)\right) d\mathbf{a}_t$ 

For more details, see: Levine. (2018). Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review.

# Q-learning with soft optimality

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

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

soft Q-learning: 
$$\phi \leftarrow \phi + \alpha \nabla_{\phi} Q_{\phi}(\mathbf{s}, \mathbf{a}) (r(\mathbf{s}, \mathbf{a}) + \gamma V(\mathbf{s}') - Q_{\phi}(\mathbf{s}, \mathbf{a}))$$
  
target value:  $V(\mathbf{s}') = \operatorname{soft} \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}', \mathbf{a}') = \log \int \exp(Q_{\phi}(\mathbf{s}', \mathbf{a}')) d\mathbf{a}'$   
 $\pi(\mathbf{a}|\mathbf{s}) = \exp(Q_{\phi}(\mathbf{s}, \mathbf{a}) - V(\mathbf{s})) = \exp(A(\mathbf{s}, \mathbf{a}))$ 

1. take some action a<sub>i</sub> and observe (s<sub>i</sub>, a<sub>i</sub>, s'<sub>i</sub>, r<sub>i</sub>), add it to R
2. sample mini-batch {s<sub>j</sub>, a<sub>j</sub>, s'<sub>j</sub>, r<sub>j</sub>} from R uniformly
3. compute y<sub>j</sub> = r<sub>j</sub> + γsoft max<sub>a'<sub>j</sub></sub> Q<sub>φ'</sub>(s'<sub>j</sub>, a'<sub>j</sub>) using target network Q<sub>φ'</sub>
4. φ ← φ − α ∑<sub>j</sub> dQ<sub>φ</sub>/dφ (s<sub>j</sub>, a<sub>j</sub>)(Q<sub>φ</sub>(s<sub>j</sub>, a<sub>j</sub>) − y<sub>j</sub>)
5. update φ': copy φ every N steps, or Polyak average φ' ← τφ' + (1 − τ)φ

# Policy gradient with soft optimality

 $\pi(\mathbf{a}|\mathbf{s}) = \exp(Q_{\phi}(\mathbf{s}, \mathbf{a}) - V(\mathbf{s})) \text{ optimizes } \sum_{t} E_{\pi(\mathbf{s}_{t}, \mathbf{a}_{t})}[r(\mathbf{s}_{t}, \mathbf{a}_{t})] + E_{\pi(\mathbf{s}_{t})}[\mathcal{H}(\pi(\mathbf{a}_{t}|\mathbf{s}_{t}))]$ policy entropy

intuition:  $\pi(\mathbf{a}|\mathbf{s}) \propto \exp(Q_{\phi}(\mathbf{s}, \mathbf{a}))$  when  $\pi$  minimizes  $D_{\mathrm{KL}}(\pi(\mathbf{a}|\mathbf{s}) \| \frac{1}{Z} \exp(Q(\mathbf{s}, \mathbf{a})))$  $D_{\mathrm{KL}}(\pi(\mathbf{a}|\mathbf{s}) \| \frac{1}{Z} \exp(Q(\mathbf{s}, \mathbf{a}))) = E_{\pi(\mathbf{a}|\mathbf{s})}[Q(\mathbf{s}, \mathbf{a})] - \mathcal{H}(\pi)$ 

often referred to as "entropy regularized" policy gradient

combats premature entropy collapse

turns out to be closely related to soft Q-learning: see Haarnoja et al. '17 and Schulman et al. '17

Ziebart et al. '10 "Modeling Interaction via the Principle of Maximum Causal Entropy"

#### Policy gradient vs Q-learning

policy gradient derivation:

Q-

$$\begin{split} J(\theta) &= \sum_{t} E_{\pi(\mathbf{s}_{t},\mathbf{a}_{t})}[r(\mathbf{s}_{t},\mathbf{a}_{t})] + E_{\pi(\mathbf{s}_{t})}[\mathcal{H}(\pi(\mathbf{a}|\mathbf{s}_{t}))] = \sum_{t} E_{\pi(\mathbf{s}_{t},\mathbf{a}_{t})}[r(\mathbf{s}_{t},\mathbf{a}_{t}) - \log \pi(\mathbf{a}_{t}|\mathbf{s}_{t})] \\ &= \sum_{t} E_{\pi(\mathbf{s}_{t},\mathbf{a}_{t})}[r(\mathbf{s}_{t},\mathbf{a}_{t}) - \log \pi(\mathbf{a}_{t}|\mathbf{s}_{t})] \\ &= \sum_{t} E_{\pi(\mathbf{s}_{t},\mathbf{a}_{t})}[r(\mathbf{s}_{t},\mathbf{a}_{t}) - \frac{\log \pi(\mathbf{a}_{t}|\mathbf{s}_{t})]}{\prod}] \\ &\approx \frac{1}{N} \sum_{i} \sum_{t} \nabla_{\theta} \log \pi(\mathbf{a}_{t}|\mathbf{s}_{t}) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + \left(\sum_{t'=t+1}^{T} r(\mathbf{s}_{t'},\mathbf{a}_{t'}) - \log \pi(\mathbf{a}_{t'}|\mathbf{s}_{t'})\right)\right) - \log \pi(\mathbf{a}_{t}|\mathbf{s}_{t}) - \frac{1}{N} \right) \\ &= \operatorname{recall:} \log \pi(\mathbf{a}_{t}|\mathbf{s}_{t}) = Q(\mathbf{s}_{t},\mathbf{a}_{t}) - V(\mathbf{s}_{t}) \\ &\approx \frac{1}{N} \sum_{i} \sum_{t} \sum_{t} \left( \nabla_{\theta}Q(\mathbf{a}_{t}|\mathbf{s}_{t}) - \nabla_{\theta}V(\mathbf{s}_{t}) \right) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + Q(\mathbf{s}_{t+1},\mathbf{a}_{t+1}) - Q(\mathbf{s}_{t},\mathbf{a}_{t}) + V(\mathbf{s}_{t}) \right) \\ &= \operatorname{Q-learning} \bigoplus_{i} \sum_{t} \sum_{t} \sum_{t} \nabla_{\theta}Q(\mathbf{a}_{t}|\mathbf{s}_{t}) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + \operatorname{soft} \max_{\mathbf{a}_{t+1}}Q(\mathbf{s}_{t+1},\mathbf{a}_{t+1}) - Q(\mathbf{s}_{t},\mathbf{a}_{t}) \right) \\ &= \operatorname{Q-learning} \left( \sum_{i} \sum_{t} \sum_{t} \sum_{t} \nabla_{\theta}Q(\mathbf{a}_{t}|\mathbf{s}_{t}) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + \operatorname{soft} \max_{\mathbf{a}_{t+1}}Q(\mathbf{s}_{t+1},\mathbf{a}_{t+1}) - Q(\mathbf{s}_{t},\mathbf{a}_{t}) \right) \\ &= \operatorname{Q-learning} \left( \sum_{i} \sum_{t} \sum_{t} \sum_{t} \nabla_{\theta}Q(\mathbf{a}_{t}|\mathbf{s}_{t}) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + \operatorname{soft} \max_{\mathbf{a}_{t+1}}Q(\mathbf{s}_{t+1},\mathbf{a}_{t+1}) - Q(\mathbf{s}_{t},\mathbf{a}_{t}) \right) \right) \\ &= \operatorname{Q-learning} \left( \sum_{i} \sum_{t} \sum_{t} \sum_{t} \sum_{t} \nabla_{\theta}Q(\mathbf{a}_{t}|\mathbf{s}_{t}) \left(r(\mathbf{s}_{t},\mathbf{a}_{t}) + \operatorname{soft} \max_{\mathbf{a}_{t+1}}Q(\mathbf{s}_{t+1},\mathbf{a}_{t+1}) - Q(\mathbf{s}_{t},\mathbf{a}_{t}) \right) \right) \\ &= \operatorname{Q-learning} \left( \sum_{i} \sum_{t} \sum_{t}$$

# Benefits of soft optimality

- Improve exploration and prevent entropy collapse
- Easier to specialize (finetune) policies for more specific tasks
- Principled approach to break ties
- Better robustness (due to wider coverage of states)
- Can reduce to hard optimality as reward magnitude increases
- Good model for modeling human behavior (more on this later)

#### Review

- Reinforcement learning can be viewed as inference in a graphical model
  - Value function is a backward message
  - Maximize reward and entropy (the bigger the rewards, the less entropy matters)
  - Variational inference to remove optimism
- Soft Q-learning
- Entropy-regularized policy gradient

![](_page_32_Figure_7.jpeg)

# Stochastic models for learning control

![](_page_33_Picture_1.jpeg)

![](_page_33_Picture_2.jpeg)

![](_page_33_Picture_3.jpeg)

![](_page_33_Picture_4.jpeg)

 How can we track both hypotheses?

#### Stochastic energy-based policies

Q-function:  $Q(\mathbf{s}, \mathbf{a}) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ 

![](_page_34_Figure_2.jpeg)

![](_page_34_Picture_3.jpeg)

 $\pi(\mathbf{a}|\mathbf{s}) \propto \exp(Q(\mathbf{s},\mathbf{a}))$ 

$$\pi(\mathbf{a}_t|\mathbf{s}_t) = \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t) - V_t(\mathbf{s}_t)) = \exp(A_t(\mathbf{s}_t, \mathbf{a}_t))$$

$$Q_t(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + E[V_{t+1}(\mathbf{s}_{t+1})]$$
$$V_t(\mathbf{s}_t) = \log \int \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t))\mathbf{a}_t$$

Haarnoja\*, Tang\*, Abbeel, L., Reinforcement Learning with Deep Energy-Based Policies. ICML 2017

#### Stochastic energy-based policies provide pretraining

![](_page_35_Figure_1.jpeg)

![](_page_35_Figure_2.jpeg)

![](_page_35_Figure_3.jpeg)

# Soft optimality suggested readings

- Todorov. (2006). Linearly solvable Markov decision problems: one framework for reasoning about soft optimality.
- Todorov. (2008). General duality between optimal control and estimation: primer on the equivalence between inference and control.
- Kappen. (2009). Optimal control as a graphical model inference problem: frames control as an inference problem in a graphical model.
- Ziebart. (2010). Modeling interaction via the principle of maximal causal entropy: connection between soft optimality and maximum entropy modeling.
- Rawlik, Toussaint, Vijaykumar. (2013). On stochastic optimal control and reinforcement learning by approximate inference: temporal difference style algorithm with soft optimality.
- Haarnoja\*, Tang\*, Abbeel, L. (2017). Reinforcement learning with deep energy based models: soft Q-learning algorithm, deep RL with continuous actions and soft optimality
- Nachum, Norouzi, Xu, Schuurmans. (2017). Bridging the gap between value and policy based reinforcement learning.
- Schulman, Abbeel, Chen. (2017). Equivalence between policy gradients and soft Q-learning.
- Haarnoja, Zhou, Abbeel, L. (2018). Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor.
- Levine. (2018). Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review