

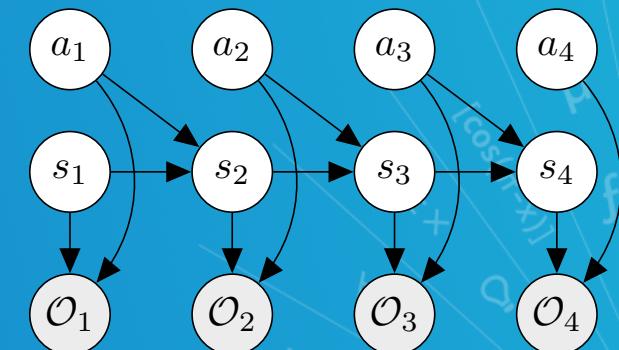
Probabilistic Graphical Models

Reinforcement Learning & Control Through Inference in GM

Maruan Al-Shedivat

Lecture 20, April 1, 2019

Reading: see class homepage





What we've covered so far

Module 1: Representation & Exact Inference

Module 2: Approximate Inference (VI, MC, SMC, MCMC)

Module 3 [introduced in 2015]: Deep Learning & Generative Models

Module 4 [new and experimental]: RL and Control as Inference in GM

Goal: gentle introduction to basic concepts of RL with a focus on connections between control and inference in a probabilistic GM





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).





RL has already come up previously (in text generation)

- The generalized ERPO objective [Tan, Hu, et al., 2018]:

$$\mathcal{L}(q, \theta) = \mathbb{E}_{\overbrace{q}}[R(\mathbf{y} | \mathbf{y}^*)] - \alpha \text{KL}[q(\mathbf{y} | \mathbf{x}) || p_\theta(\mathbf{y} | \mathbf{x})] + \beta H(q)$$

- The **reward term** is a non-differentiable function of the sample \mathbf{y} and the true sequence \mathbf{y}^* (e.g., represented by the BLEU metric).
- More generally, objectives of this form are called stochastic objectives:

$$\mathcal{L}(\theta) = \mathbb{E}_{x \sim p_\theta(x)}[f(x)] \quad \Bigg| \Bigg($$

- As we will see, such objectives often come up in RL.

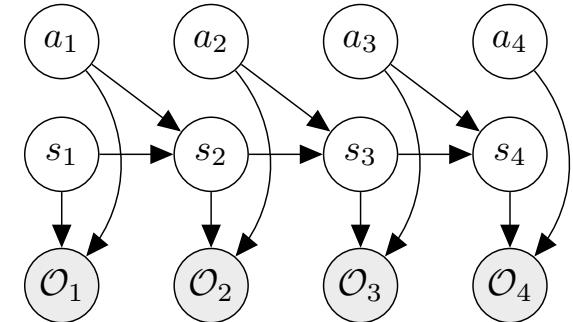




Plan

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

- ❑ Classical Q-learning and policy gradient methods
- ❑ Derivation of the soft Q-learning and soft policy gradients
- ❑ Algorithms and applications

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
```

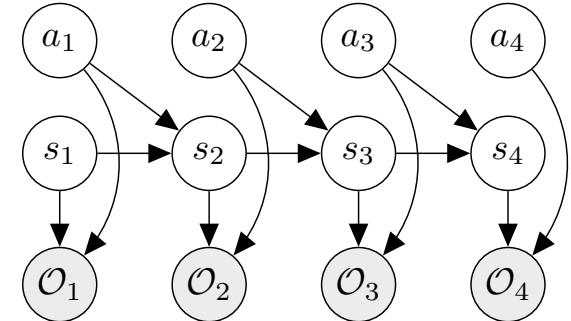




Plan

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

- ❑ Classical Q-learning and policy gradient methods
- ❑ Derivation of the soft Q-learning and soft policy gradients
- ❑ Algorithms and applications

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
```





Paradigms of machine learning

- Supervised learning

Given: a collection of data $D = \{(x_i, y_i)\}_{i=1}^N$

Goal: learn a model that approximates $P(y | x)$

- Unsupervised learning

Given: a collection of data $D = \{(x_1, x_2, \dots, x_d)_i\}_{i=1}^N$

Goal: learn a model that approximates $P(x_1, x_2, \dots, x_d)$

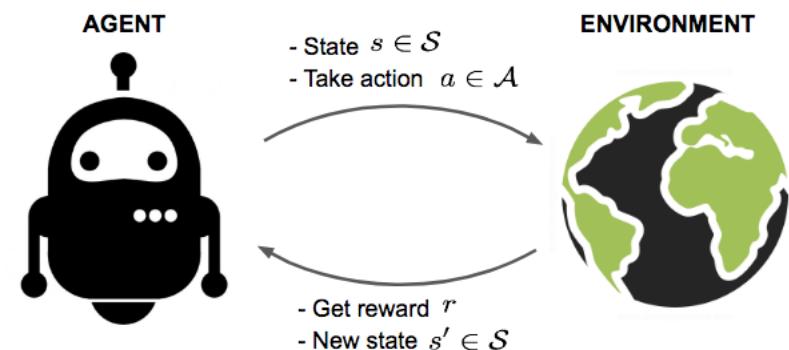
- Reinforcement learning

Given: an environment that an agent can perceive and interact with

Goal: learn a controller (policy) that can maximize the utility (reward) in the given environment

GMs allow us to efficiently represent, manipulate, and perform learning and inference on these probabilistic models.

DL gives the tools for learning expressive latent representations that lead to more accurate probabilistic models of the data.

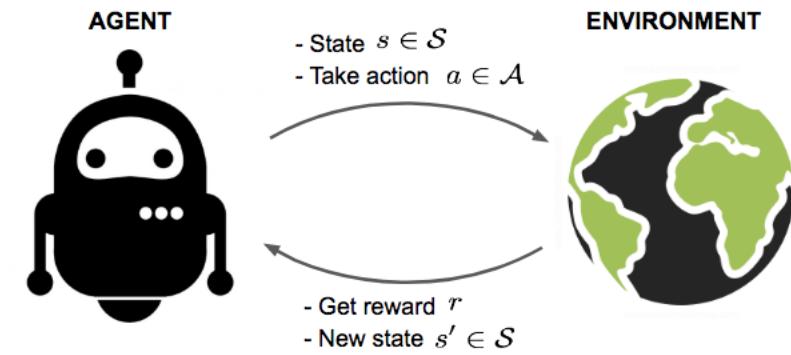




Why sequential decision making and RL?

Ultimately, we want to build intelligent machines that:

- Can perceive and interact with the world
- Exhibit purposeful goal-directed behavior
- Learn from interactions, adapt to changes, plan and be able to maximize utility functions
(specified by humans or inferred from situations)





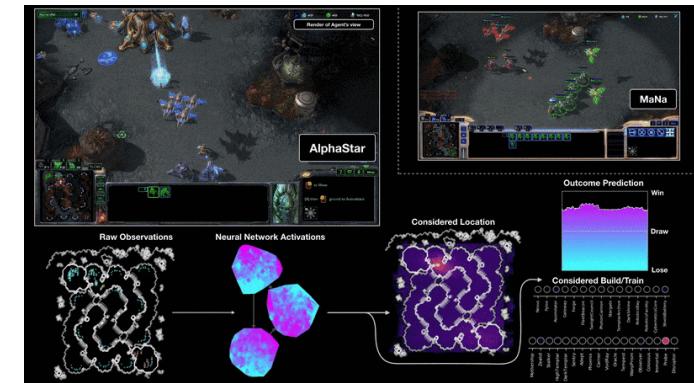
Some recent success stories of RL

Learning to play games

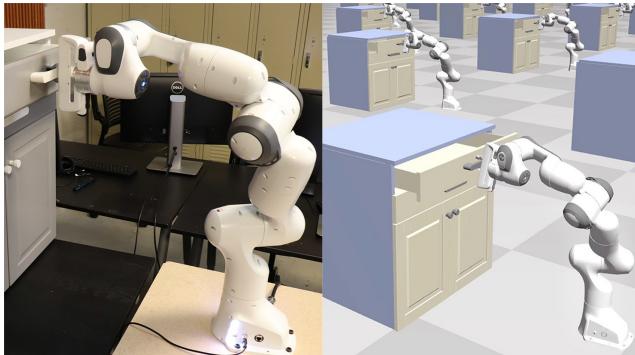
AlphaGo (DeepMind, 2016)



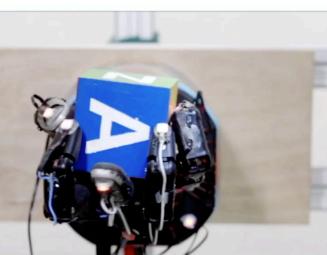
AlphaStar (DeepMind, 2019)



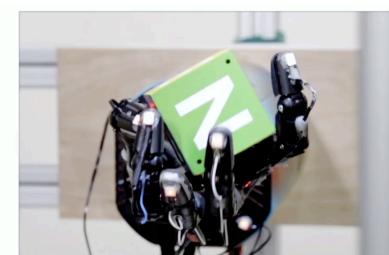
Robotics



Chebotar et al., 2018



FINGER PIVOTING



SLIDING



FINGER GAITING

OpenAI, 2018



Basic concepts of RL

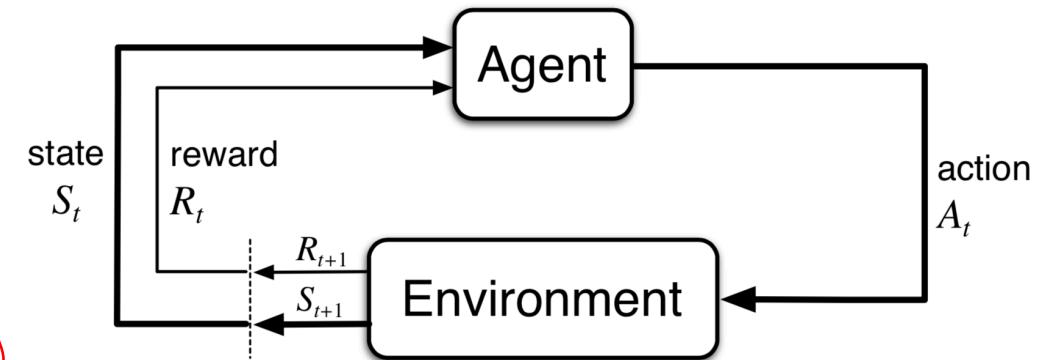


Markov Decision Processes (MDPs)

Markov Decision Process (MDP):

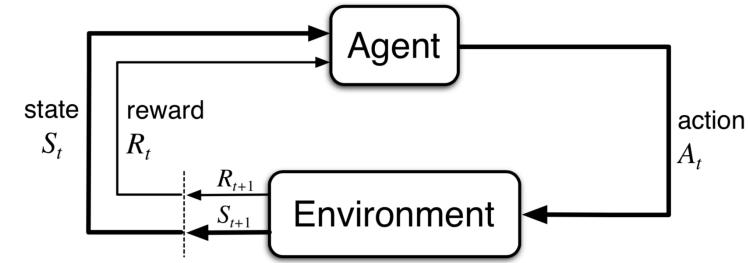
- Environment has a set of states \mathcal{S}
- Agent is given a set of possible actions \mathcal{A}
- Environment dynamics: transitions from state s_t into a new state s_{t+1} according to the transition probability $P(s_{t+1}|s_t, a_t)$ after agent takes action a_t
- Reward function: $r(s, a) = \mathbb{E}[r_{t+1}|s_t = s, a_t = a]$ provides scalar rewards to the agent at each time step
- “Life” of an agent (or trajectory):

$$\tau = (s_1, a_1, r_1, s_2, a_2, r_2, s_3, a_3, r_3, \dots)$$





Markov Decision Processes (MDPs)



What can we do with MDPs:

- (1) Policy search: Find a policy $\pi: \mathcal{S} \rightarrow \mathcal{A}$ that outputs actions for each given state such that the cumulative reward along the trajectory is maximized.
- (2) Inverse RL: Given a set of optimal trajectories (e.g., generated by a human expert), infer the corresponding MDP.





Returns and Episodes

Maximization of the return:

- Return (cumulative reward) starting step t : $G_t = r_{t+1} + r_{t+2} + \dots + r_T$
- If $T = \infty$, we can use the notion of discounted return:

$$\begin{aligned} G_t &= r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \\ &= r_{t+1} + \gamma G_{t+1} \end{aligned}$$

where $0 \leq \gamma \leq 1$ is called the discount rate





Policies and Value Functions

- Value function of a state s :

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

stochastic objective

- Value function of the state-action pair (s, a) :

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





Bellman Equation for $V_\pi(s)$

- Bellman equation for the value function of a state s :

$$\begin{aligned} \underline{\underline{V_\pi(s)}} &:= \mathbb{E}_\pi [G_t \mid s_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^T \gamma^k r_{t+k+1} \mid s_t = s \right] \\ &= \mathbb{E}_\pi [r_{t+1} + \gamma G_{t+1} \mid s_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s'} p(s' \mid s, a) [r(s, a) + \gamma \mathbb{E}_\pi [G_{t+1} \mid s_{t+1} = s']] \\ &= \sum_a \pi(a \mid s) \sum_{s'} p(s' \mid s, a) [r(s, a) + \gamma \underline{\underline{V_\pi(s')}}] \end{aligned}$$

Annotations: A red double underline is placed under $V_\pi(s)$. Red arrows point from the term r_{t+1} and G_{t+1} in the second line to their respective terms in the third line. Red arrows also point from the term $\pi(a \mid s)$ and $p(s' \mid s, a)$ in the fourth line to their respective terms in the third line.

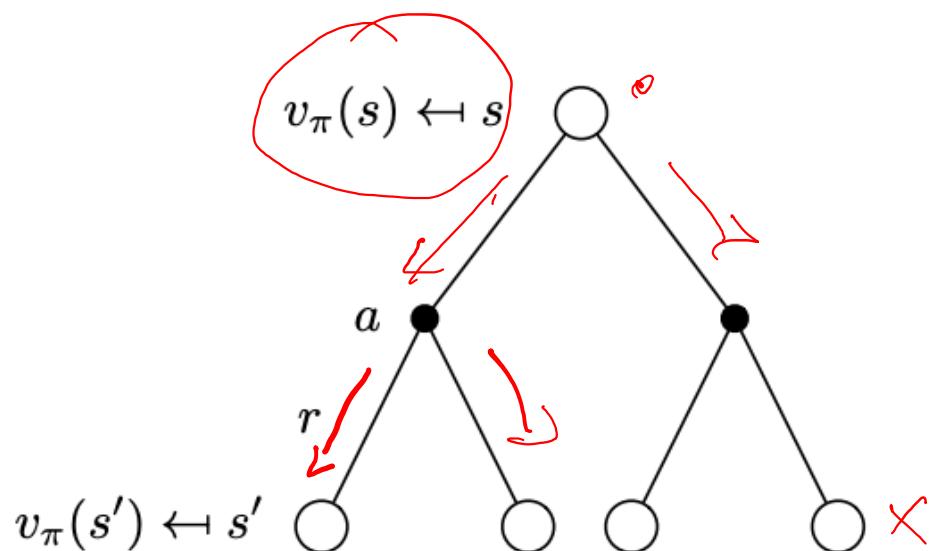




Bellman Equation for $V_\pi(s)$

- Bellman equation for the value function of a state s :

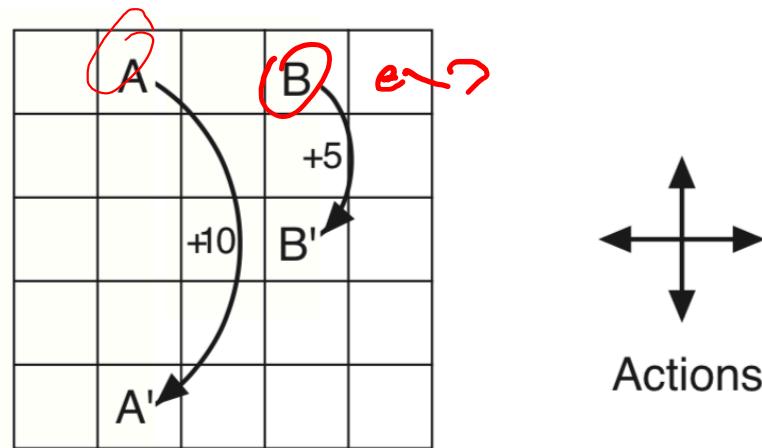
$$V_\pi(s) = \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a) + \gamma V_\pi(s')]$$





Example: Grid World and a Random Policy

- Getting off the grid results in -1 reward and no change in position.
- Any action in A and B result in +10 and +5 and move the agent to A' and B', respectively.



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

$\checkmark \pi_{\text{rand}}$

Figure 3.2: Gridworld example: exceptional reward dynamics (left) and state-value function for the equiprobable random policy (right).





Bellman Equation for $Q_\pi(s, a)$

- Bellman equation for the value function of the state-action pair (s, a) :

$$\begin{aligned} Q_\pi(s, a) &:= \mathbb{E}_\pi [G_t \mid s_t = s, a_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^T \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right] \\ &= r(s, a) + \gamma \mathbb{E}_\pi \left[\underbrace{G_{t+1}}_{\cdot} \mid s_t = s, a_t = a \right] \\ &= r(s, a) + \gamma \sum_{s'} p(\underbrace{s' \mid s, a}_{\cdot}) \sum_{a'} \underbrace{\pi(a' \mid s')}_{\cdot} \mathbb{E}_\pi [G_{t+1} \mid s_{t+1} = s', a_{t+1} = a'] \\ &= r(s, a) + \gamma \sum_{s'} p(s' \mid s, a) \sum_{a'} \pi(a' \mid s') Q_\pi(s', a') \end{aligned}$$

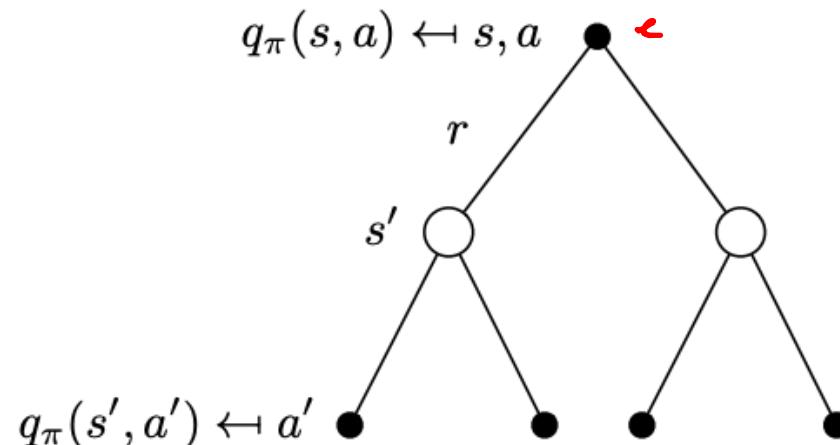




Bellman Equation for $Q_\pi(s, a)$

- Bellman equation for the value function of the state-action pair (s, a) :

$$Q_\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} \pi(a' | s') Q_\pi(s', a')$$





Optimal Policies and Value Functions

- Solving an RL task means finding an optimal policy that achieves high reward in the long run.
- Policy π is better or equal to π' ($\pi \geq \pi'$) if its expected return is greater or equal to that of π' for all states:

$$\pi \geq \pi' \Leftrightarrow V_\pi(s) \geq V_{\pi'}(s) \quad \forall s \in \mathcal{S}$$

- Optimal value functions (Bellman optimality):

$$\rightarrow V_\star(s) := \max_{\pi} V_{\pi}(s) = \max_a \sum_{s'} p(s' | s, a) [r(s, a) + \gamma V_\star(s')] \quad ||$$

$$\rightarrow Q_\star(s, a) := \max_{\pi} Q_{\pi}(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a) + \gamma \max_{a'} Q_\star(s', a') \right] \quad ||$$





Optimal Policies and Value Functions

- Optimal value functions (Bellman optimality):

$$V_*(s) := \max_{\pi} V_{\pi}(s) = \max_a \sum_{s'} p(s' | s, a) [r(s, a) + \gamma V_*(s')]$$

$$Q_*(s, a) := \max_{\pi} Q_{\pi}(s, a) = \sum_{s'} p(s' | s, a) \left[r(s, a) + \gamma \max_{a'} Q_*(s', a') \right]$$

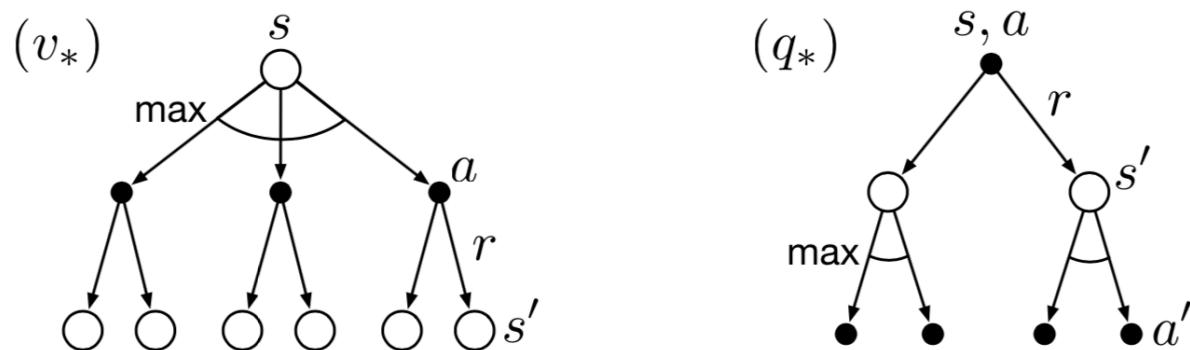


Figure 3.4: Backup diagrams for v_* and q_*





How to recover optimal policy and trajectories?

- We can recover an optimal policy from the optimal $Q_*(s, a)$:

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

\uparrow
Dirac delta

- To recover a set of optimal trajectories, just execute the optimal policy:

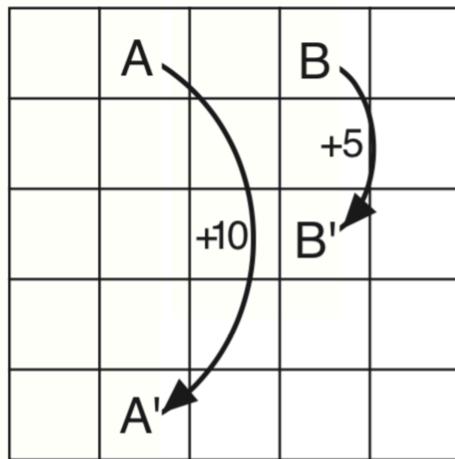
$$(\tau_*) = (s_1^*, a_1^*, r_1^*, s_2^*, a_2^*, r_2^*, \dots)$$

$$s_{t+1}^* \sim p(s_{t+1} \mid s_t, a_t^* = \arg \max_a Q_*(s, a))$$





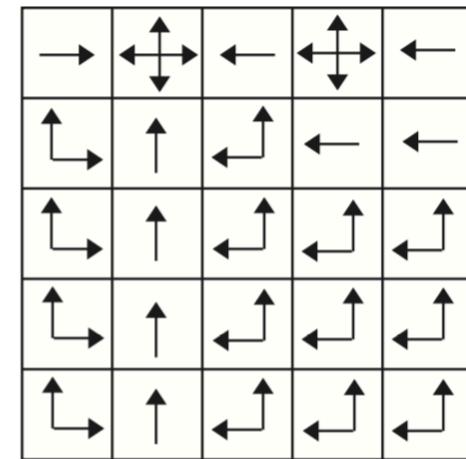
Example: Grid World and an Optimal Policy



Gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

v_*



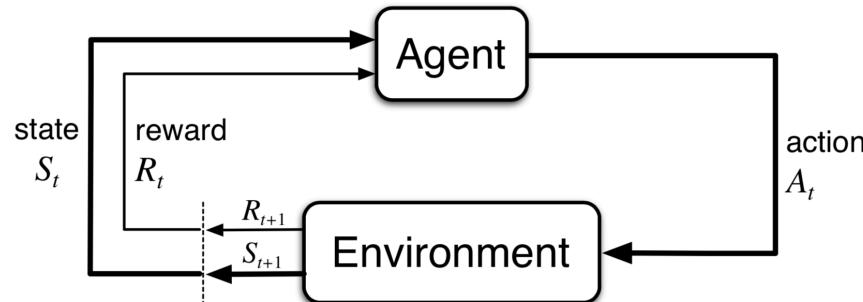
π_*

Figure 3.5: Optimal solutions to the gridworld example.





Recap



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

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

- Value functions:

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

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

- Recursive notion of optimality:

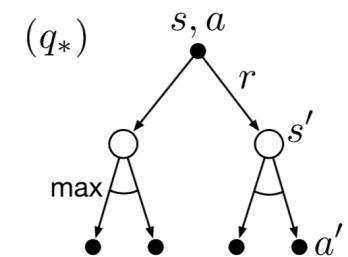
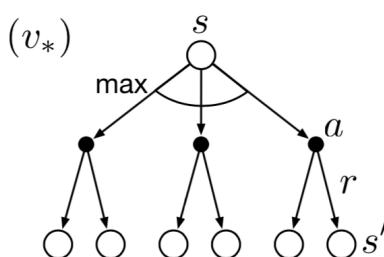


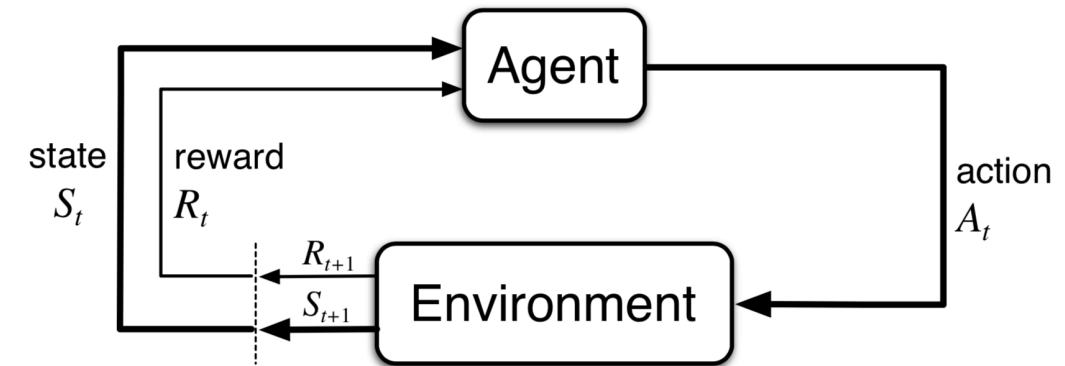
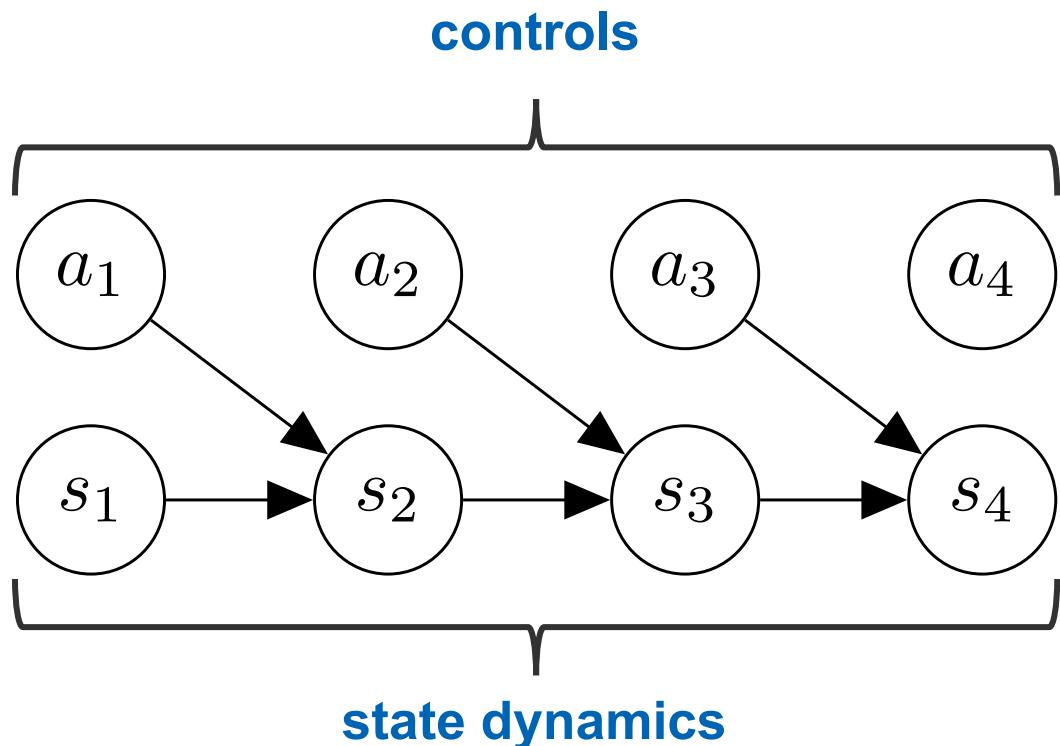
Figure 3.4: Backup diagrams for v_* and q_*



RL & Control as Inference in GM



MDP as a Graphical Model

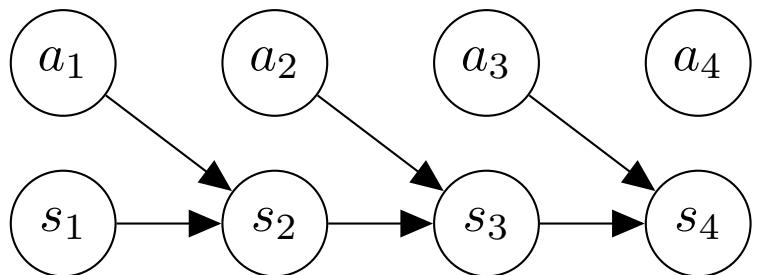


How do we define a distribution over rational/optimal trajectories?





MDP as a Graphical Model



Initial state

$$s_0 \sim p_0(s)$$

Transition

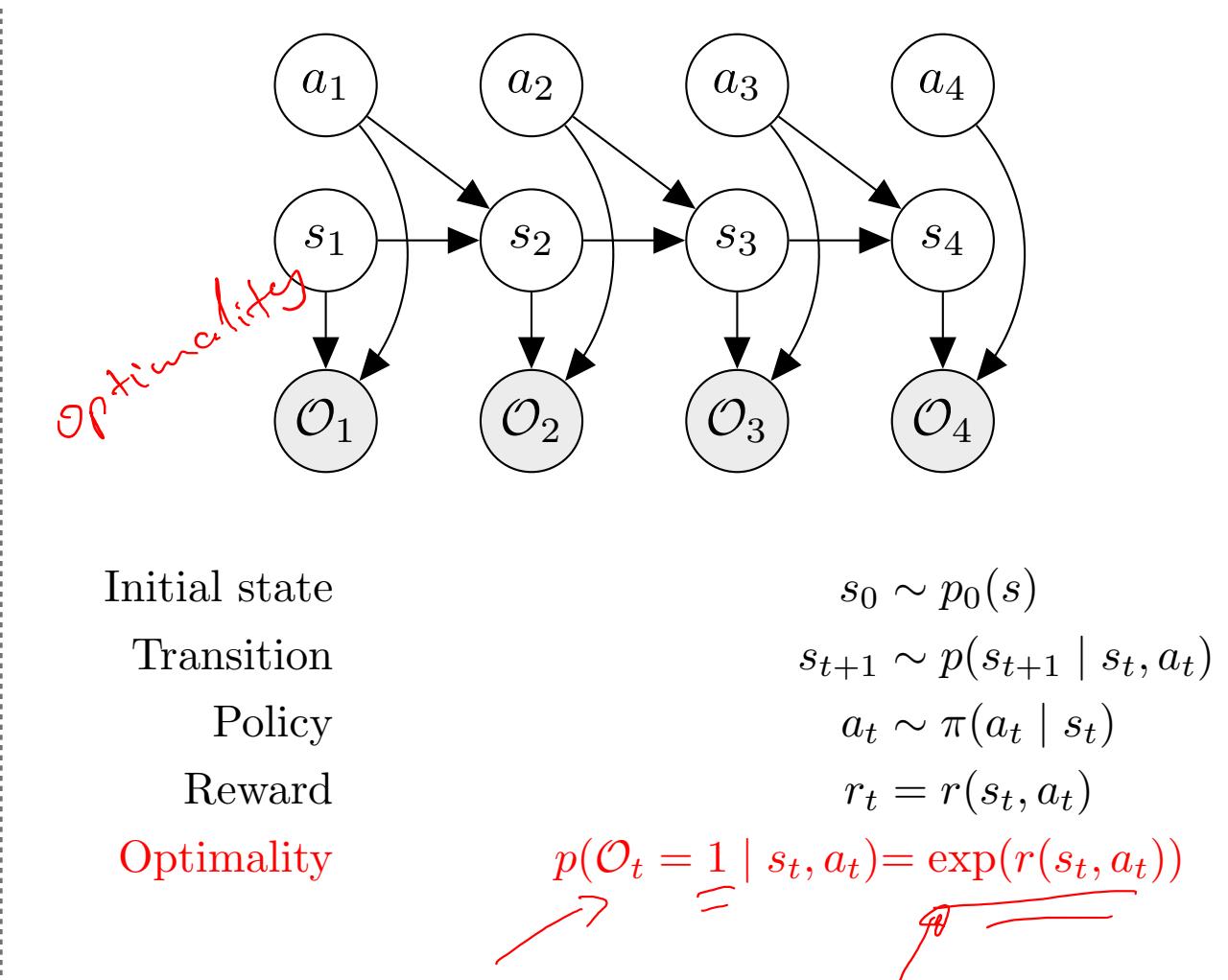
$$s_{t+1} \sim p(s_{t+1} | s_t, a_t)$$

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

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

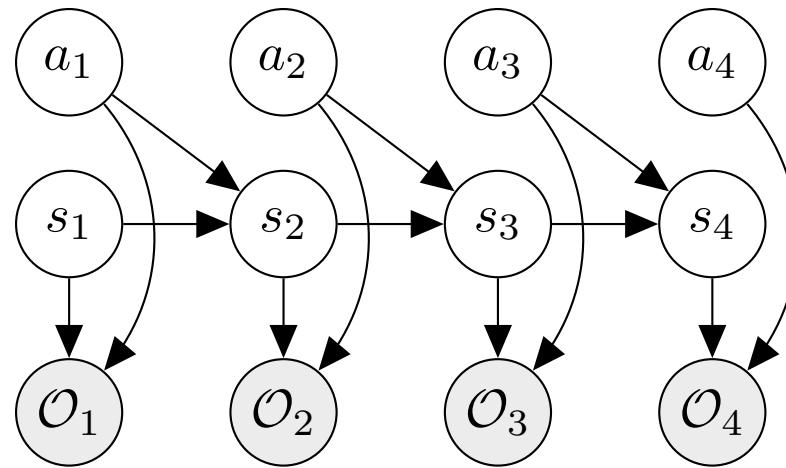
Optimality

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





Why is this model interesting?

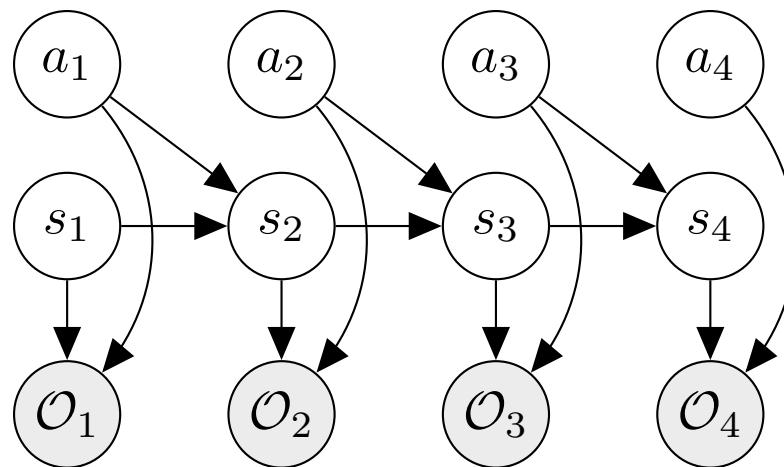


- Can solve control and planning problems using inference algorithms
- Allows to model suboptimal behavior (important for inverse RL)
- Provides explanation for why stochastic behavior might be preferred (from the exploration and transfer learning point of view)





What can we do with this graphical model?



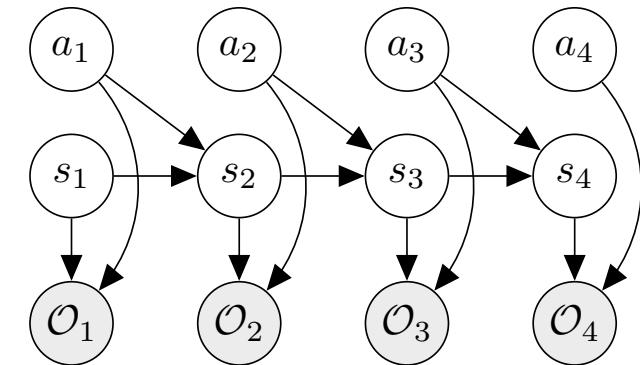
Here is what we can do:

- Given a reward, determine a likely optimal trajectory $p(\tau | \mathcal{O}_{1:T})$
- Given a collection of optimal trajectories, infer the reward and priors
- Given a reward, infer the optimal policy $p(a_t | s_t, \mathcal{O}_{t:T})$





Distribution over the optimal trajectories



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

$$p(\tau \mid \mathcal{O}_{1:T}) \propto \underbrace{p(s_1)}_{\text{action}} \prod_{t=1}^T \underbrace{p(a_t \mid s_t)}_{\text{prior}} p(\underbrace{s_{t+1} \mid s_t, a_t}_{\text{action prior}}) p(\mathcal{O}_t \mid s_t, a_t)$$

~~X~~

R

?





Inferring the reward & prior that generate trajectories

$$p(\tau \mid \mathcal{O}_{1:T}, \theta, \phi) \propto \left[p(s_1) \prod_{t=1}^T p(s_{t+1} \mid s_t, a_t) \right] \exp \left(\sum_{t=1}^T r_\phi(s_t, a_t) + \log p_\theta(a_t \mid s_t) \right)$$

parametrized

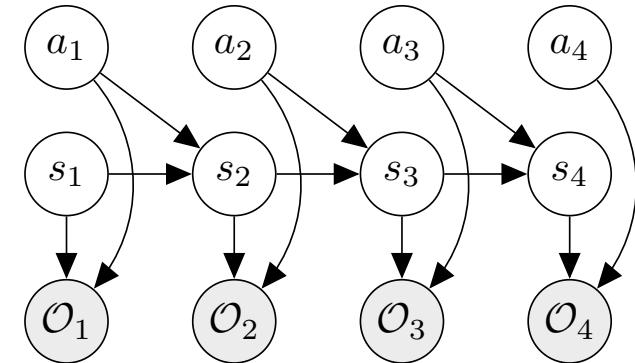


The model reminds a featurized CRF. (*Note: CRF is undirected and does not preserve the causal structure; this model is more restrictive and known as MEMM.)





Optimal policy and planning via inference



- Unroll the dynamics and compute backward messages:

$$\beta_t(s_t, a_t) := p(O_{t:T} \mid s_t, a_t)$$

- Compute optimal policy:

$$p(a_t \mid s_t, O_{t:T})$$

- Compute forward messages (state filtering under optimality constraint):

$$\underline{\alpha}_t(s_t) = p(s_t \mid O_{1:t-1})$$



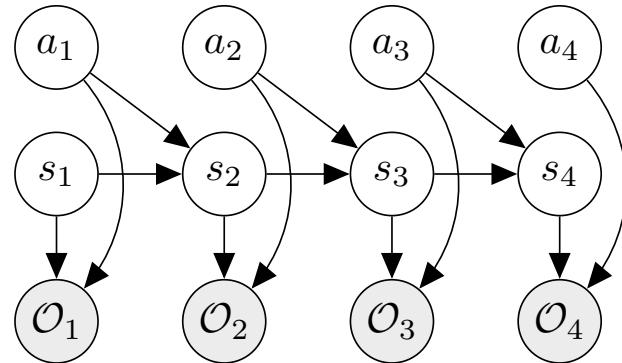


Backward messages

$$\begin{aligned}
 \beta_t(s_t, a_t) &= p(\mathcal{O}_{t:T} \mid s_t, a_t) \\
 &= \int_S p(\mathcal{O}_{t:T}, \underline{s_{t+1}} \mid s_t, a_t) ds_{t+1} \\
 &= \int_S \underbrace{p(\mathcal{O}_{t+1:T} \mid s_{t+1})}_{\beta_{t+1}(s_{t+1})} \underbrace{p(s_{t+1} \mid s_t, a_t)}_{\pi(s_{t+1} \mid s_t, a_t)} \underbrace{p(\mathcal{O}_t \mid s_t, a_t)}_{\mathbb{E}_{s_{t+1} \sim p(s_{t+1} \mid s_t, a_t)} [\beta_{t+1}(s_{t+1})]} ds_{t+1} \\
 p(\mathcal{O}_{t+1:T} \mid s_{t+1}) &= \int_A p(\mathcal{O}_{t+1:T} \mid s_{t+1}, a_{t+1}) p(a_{t+1} \mid s_{t+1}) da_{t+1} \quad || \\
 &\quad \beta(s_{t+1}, a_{t+1})
 \end{aligned}$$

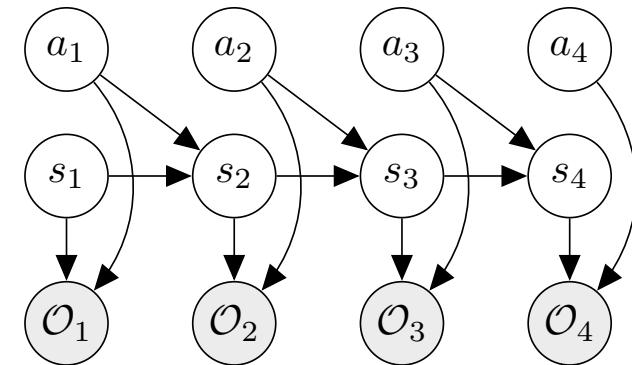
for $t = T - 1$ to 1 :

$$\begin{aligned}
 \beta_t(s_t, a_t) &= p(\mathcal{O}_t \mid s_t, a_t) \mathbb{E}_{s_{t+1} \sim p(s_{t+1} \mid s_t, a_t)} [\beta_{t+1}(s_{t+1})] \\
 \beta_t(s_t) &= \mathbb{E}_{a_t \sim p(a_t \mid s_t)} [\beta_t(s_t, a_t)]
 \end{aligned}$$





How are these messages related to RL?



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

let $Q_t(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t)$

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

Softmax

Deterministic dynamics: $Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + V(\mathbf{s}_{t+1})$

Stochastic dynamics: $Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\exp(V(\mathbf{s}_{t+1}))]$

“optimistic” transition (not good)

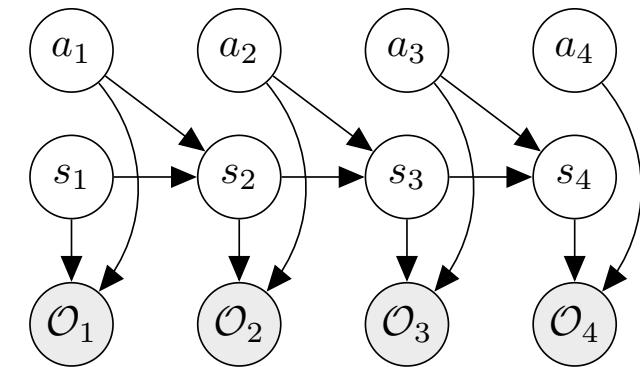




Optimal policy

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

$$| \rho(a_t | s_t, \theta_{t:T})$$



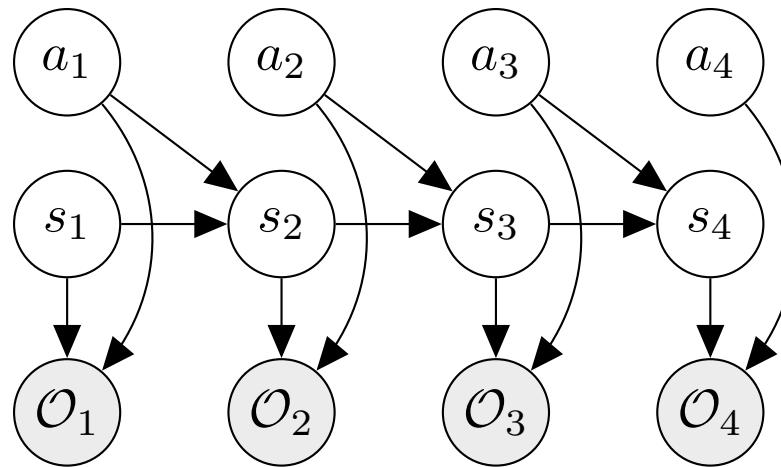
$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t) - \underbrace{V_t(\mathbf{s}_t)}_{\text{---}}) = \exp(\underbrace{A_t(\mathbf{s}_t, \mathbf{a}_t)}_{\text{---}}) \frac{1}{\alpha}$$

- (Derivation: HW4!)
- Natural interpretation: better actions are more probable + random tie breaking
- Approaches greedy policy as temperature decreases





Summary



- Using auxiliary potentials and/or optimality variables, we reduced optimal control to inference in a graphical model.
- “Solving MDP” becomes very similar to inference in HMM / MEMM / CRF.
- The approach is quite similar to dynamic programming, value iterations, etc.



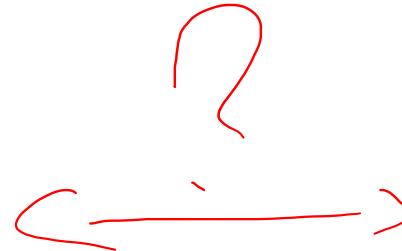
Control via Variational Inference



Which objective does inference optimize?

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

$$\pi_\star = \arg \max_{\pi} \mathbb{E} [V_\pi(s)]$$



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

- Q: Is there a way to find an objective function by optimizing which we can recover our inference-based policy?
- A: Yes! Let's take a look at the KL divergence between trajectory distributions.





Which objective does inference optimize?

[the case of deterministic dynamics]

Optimal

$$p(\tau) \cancel{=} \left[p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right] \exp \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right)$$

=

$$-D_{\text{KL}}(\hat{p}(\tau) \| p(\tau)) = E_{\tau \sim \hat{p}(\tau)} \left[\log p(\mathbf{s}_1) + \sum_{t=1}^T (\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + r(\mathbf{s}_t, \mathbf{a}_t)) - \log \cancel{p}(\mathbf{s}_1) - \sum_{t=1}^T (\log p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) + \log \pi(\mathbf{a}_t | \mathbf{s}_t)) \right]$$





Which objective does inference optimize?

[the case of deterministic dynamics]

$$-D_{\text{KL}}(\hat{p}(\tau) \| p(\tau)) = E_{\tau \sim \hat{p}(\tau)} \left[\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t) \right]$$

under policy $\hat{p}(\mathbf{a}^{\tau})$

$$= \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t | \mathbf{s}_t)]$$

$$= \sum_{t=1}^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim \hat{p}(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)] + E_{\mathbf{s}_t \sim \hat{p}(\mathbf{s}_t)} [\mathcal{H}(\pi(\mathbf{a}_t | \mathbf{s}_t))]$$

entropy of the policy

expected returns *standard deviation*



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



The problem of optimism in stochastic dynamics

Deterministic dynamics: $Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + V(\mathbf{s}_{t+1})$

Stochastic dynamics: $Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [\exp(V(\mathbf{s}_{t+1}))]$



“optimistic” transition (not good)

Why did this happen?

- Dichotomy between what resulted in the high reward:
was it a good policy or we just got lucky with the stochastic dynamics?
- The optimal policy: $p(\underline{a_t} \mid \mathbf{s}_t, \mathcal{O}_{t:T})$ 
“given that you obtained high reward, what was your action probability?”
- The “optimal” transition probability: $p(\underline{s_{t+1}} \mid \mathbf{s}_t, \mathbf{a}_t, \mathcal{O}_{t:T}) \neq p(\underline{s_{t+1}} \mid \mathbf{s}_t, \mathbf{a}_t)$ 
“given that you obtained ~~high reward~~, what was your transition probability?” 



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



Control via variational inference

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

Let's find $q(s_{1:T}, a_{1:T}) \leftarrow q(\tau)$

such that it approximates $p(s_{1:T}, a_{1:T} | \mathcal{O}_{1:T}) \leftarrow p(\tau | \mathcal{O}_{1:T})$

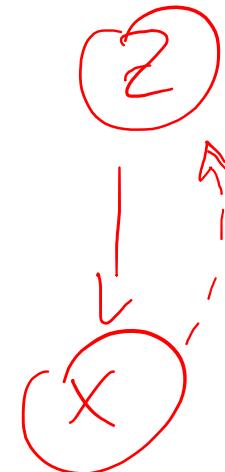
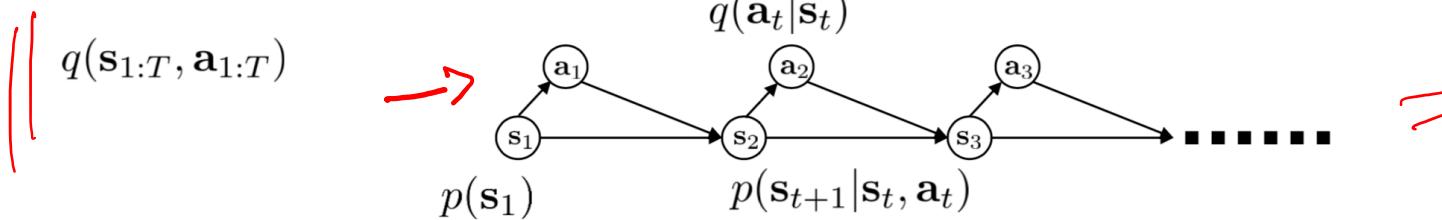
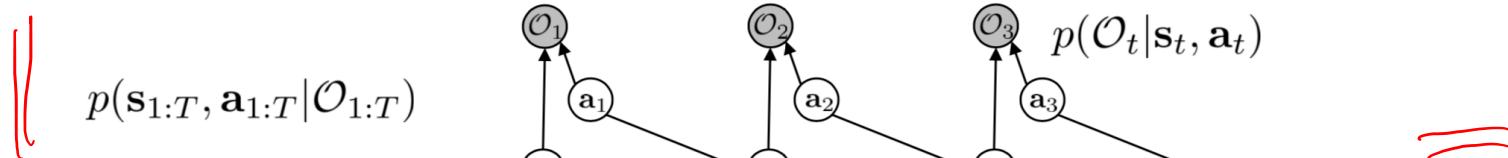
while the dynamics stays fixed to $p(s_{t+1} | s_t, a_t)$





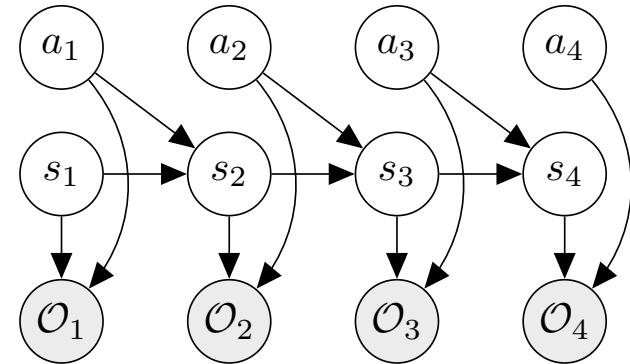
Control via variational inference

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





Control via variational inference



Optimal

$$p(\tau) = \left[p(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \right] \exp \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right)$$

Produced by a policy

$$q(\tau) = q(\mathbf{s}_1) \prod_{t=1}^T p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) \underbrace{q(\mathbf{a}_t | \mathbf{s}_t)}_{\text{policy}}$$

$$\begin{aligned} \log p(\mathcal{O}_{1:T}) &= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \\ &= \log \int \int p(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \underbrace{\frac{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})}{q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})}}_{\text{ELBO}} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T} \end{aligned}$$

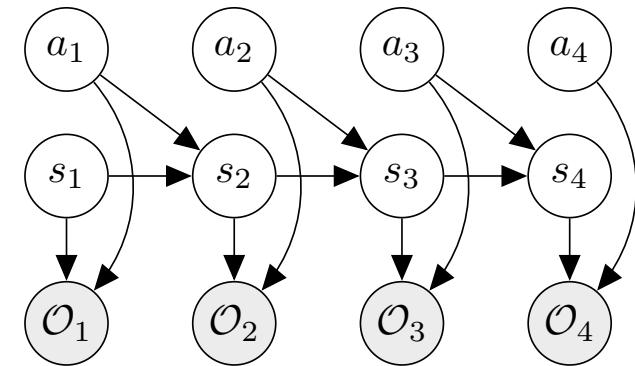
Jensen's

ELBO





Control via variational inference



$$\begin{aligned}
 \log p(\mathcal{O}_{1:T}) &\geq \mathbb{E}_{\tau \sim q} \left[\cancel{\log p(s_1)} + \sum_{t=1}^T \cancel{\log p(s_{t+1} \mid s_t, a_t)} + \sum_{t=1}^T \cancel{\log p(\mathcal{O}_t \mid s_t, a_t)} \right. \\
 &\quad \left. - \cancel{\log p(s_1)} - \sum_{t=1}^T \cancel{\log p(s_{t+1} \mid s_t, a_t)} - \sum_{t=1}^T \cancel{\log q(a_t \mid s_t)} \right] \\
 &\geq \mathbb{E}_{\tau \sim q} \left[\sum_{t=1}^T r(s_t, a_t) - \log q(a_t \mid s_t) \right] \\
 &= \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim q} [r(s_t, a_t)] + \underline{\mathcal{H}(q(a_t \mid s_t))} \\
 &\quad \text{expected return} \quad \text{entropy of the actions.}
 \end{aligned}$$





Optimal policy with respect to ELBO

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

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

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

$$\begin{aligned} q(\mathbf{a}_T | \mathbf{s}_T) &= \arg \max E_{\mathbf{s}_T \sim q(\mathbf{s}_T)} [E_{\mathbf{a}_T \sim q(\mathbf{a}_T | \mathbf{s}_T)} [r(\mathbf{s}_T, \mathbf{a}_T)] + \mathcal{H}(q(\mathbf{a}_T | \mathbf{s}_T))] \\ &= \arg \max E_{\mathbf{s}_T \sim q(\mathbf{s}_T)} [E_{\mathbf{a}_T \sim q(\mathbf{a}_T | \mathbf{s}_T)} [r(\mathbf{s}_T, \mathbf{a}_T) - \log q(\mathbf{a}_T | \mathbf{s}_T)]] \end{aligned}$$

minimized when $q(\mathbf{a}_T | \mathbf{s}_T) \propto \exp(r(\mathbf{s}_T, \mathbf{a}_T))$ \times

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

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



Optimal policy with respect to ELBO

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

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

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

$$\begin{aligned} \underline{q(\mathbf{a}_t | \mathbf{s}_t)} &= \arg \max E_{\mathbf{s}_t \sim q(\mathbf{s}_t)} [E_{\mathbf{a}_t \sim q(\mathbf{a}_t | \mathbf{s}_t)} [r(\mathbf{s}_t, \mathbf{a}_t) + E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} [V(\mathbf{s}_{t+1})]] + \mathcal{H}(q(\mathbf{a}_t | \mathbf{s}_t))] \\ &= \arg \max E_{\mathbf{s}_t \sim q(\mathbf{s}_t)} [E_{\mathbf{a}_t \sim q(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)] + \mathcal{H}(q(\mathbf{a}_t | \mathbf{s}_t))] \\ &= \arg \max E_{\mathbf{s}_t \sim q(\mathbf{s}_t)} [E_{\mathbf{a}_t \sim q(\mathbf{a}_t | \mathbf{s}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t) - \log q(\mathbf{a}_t | \mathbf{s}_t)]] \end{aligned}$$

minimized when $q(\mathbf{a}_t | \mathbf{s}_t) \propto \exp(Q(\mathbf{s}_t, \mathbf{a}_t))$

$$V_t(\mathbf{s}_t) = \log \int \exp(Q_t(\mathbf{s}_t, \mathbf{a}_t)) d\mathbf{a}_t$$

$$q(\mathbf{a}_t | \mathbf{s}_t) = \exp(Q(\mathbf{s}_t, \mathbf{a}_t) - V(\mathbf{s}_t))$$

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

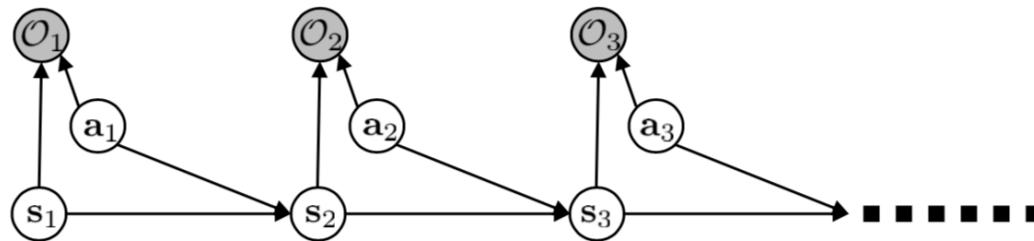
regular Bellman backup

not optimistic

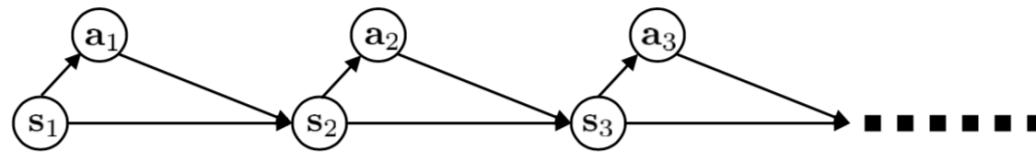


Summary

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



$$q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})$$



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

