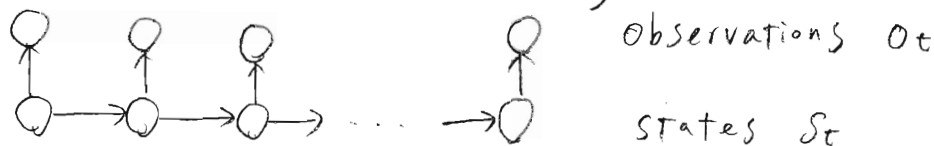


Review

* Hidden Markov Models (HMMs)



* Parameters

$$a_{ij} = P(S_{t+1}=j | S_t=i) \quad \text{transition matrix}$$

$$b_{ik} = P(O_t=k | S_t=i) \quad \text{emission matrix}$$

$$\pi_i = P(S_1=i) \quad \text{initial state distribution}$$

* EM algorithm

Choose parameters $\{a_{ij}, b_{ik}, \pi_i\}$ to maximize log-likelihood $\log P(O_1, O_2, \dots, O_T)$.

* Forward-backward procedures to compute:

$$\alpha_t = P(O_1, O_2, \dots, O_t, S_t=i)$$

$$\beta_t = P(O_{t+1}, \dots, O_T | S_t=i)$$

* E-step: Compute posterior probabilities

$$P(S_t=i, S_{t+1}=j | O_1, O_2, \dots, O_T) = \frac{P(S_t=i, S_{t+1}=j, O_1, \dots, O_T)}{P(O_1, O_2, \dots, O_T)}$$

$$= \frac{\alpha_t a_{ij} b_j(O_{t+1}) \beta_{j,t+1}}{\sum_i \alpha_{i,T}}$$

$$P(S_t=i | O_1, O_2, \dots, O_T) = \sum_j P(S_t=i, S_{t+1}=j | O_1, O_2, \dots, O_T)$$

* M-step: update parameters

$$\pi_i \leftarrow P(S_1=i | O_1, O_2, \dots, O_T)$$

$$a_{ij} \leftarrow \frac{\sum_t P(S_t=i, S_{t+1}=j | O_1, O_2, \dots, O_T)}{\sum_t P(S_t=i | O_1, O_2, \dots, O_T)}$$

$$\text{bit} \leftarrow \frac{\sum_i p(s_t = i | o_1, o_2, \dots, o_T) I(o_t, k)}{\sum_i p(s_t = i | o_1, o_2, \dots, o_T)}$$

Complexity of inference / learning in HMMs

1) compute $p(o_1, o_2, \dots, o_T)$

2) decode $s^* = \underset{s}{\operatorname{argmax}} p(s_1, s_2, \dots, s_T | o_1, \dots, o_T)$

3) performing iteration of EM

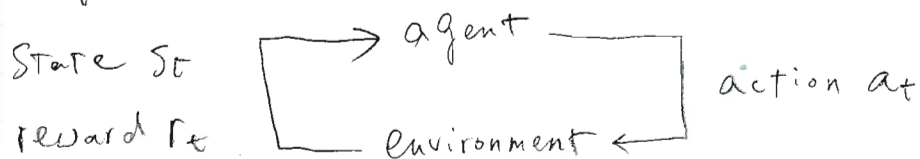
All these have same complexity: $O(N^2 T)$

↖ # states

↖ sequence length.

Reinforcement Learning

Q. How should embodied / embedded / situated / decision-making agents act and learn from experience in the world?



Ex: robot navigation

Ex: backgammon

* Challenges

- handling of uncertainty
- exploration vs. exploitation dilemma
- delayed vs. immediate rewards: "temporal credit assignment"
- evaluative vs. instructive feedback
- complex worlds, computational guarantees.

Markov decision processes (MDPs)

* Definition

- state space \mathcal{S} with states $s \in \mathcal{S}$
- action space \mathcal{A} with actions $a \in \mathcal{A}$
- Transition probabilities

for all state-action pairs (s, a) ,

$$P(s' | s, a) = P(s_{t+1} = s' | s_t = s, a_t = a)$$

probability moving from state s to state s' after taking action a .

* Assumptions

- time-independent

$$P(s_{t+1} = s' | s_t = s, a_t = a) = P(s_t = s' | s_{t-1} = s, a_{t-1} = a)$$

- Markov condition

$$P(s_{t+1} | s_t, a_t) = P(s_{t+1} | s_t, a_t, s_{t-1}, a_{t-1}, s_{t-2}, a_{t-2}, \dots)$$

(conditional independence)

* Def (Cont)

- reward function

$R(s, s', a)$ = real-valued reward after taking action a in state s and moving to state s' .

- Simplifications for CSE 150

- reward function $R(s, s', a) = R(s) = R_s$
only depends on current state.

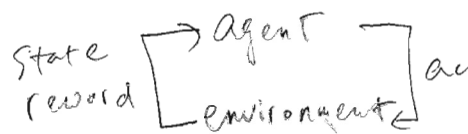
- bounded, deterministic rewards $\max_s |R_s| < \infty$

- more Simplifications

* discrete, finite state space } vs. continuous, infinite
* discrete, finite action space }

Example: backgammon

\mathcal{S} = board position and roll of dice
 \mathcal{A} = set of possible moves



$P(S'|S, a)$ = how state changes due to Agent's move,
Opponent's roll of dice, Opponent's move,
agent's roll of dice.

$$R(S) = \begin{cases} +1 & \text{win} \\ -1 & \text{lose} \\ 0 & \text{otherwise} \end{cases}$$

* Decision making.

- policy: deterministic mapping from states to actions

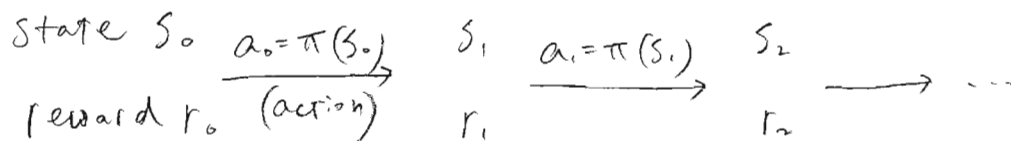
$$\pi: \mathcal{S} \rightarrow \mathcal{A}$$

- # policies: $|\mathcal{A}|^{|\mathcal{S}|}$

- dynamics

$P(S'|S, \pi(S))$ ← action in state S under policy π .

- experience under policy π



* How to measure accumulated rewards over time?

long term discounted return

discount factor $0 \leq \gamma \leq 1$

$$\text{return} = \sum_{t=0}^{\infty} \gamma^t r_t$$

Possibilities:

$\gamma = 0 \rightarrow$ Only immediate reward at $t=0$ matters

$\gamma \ll 1 \rightarrow$ near-sighted agent

$\gamma \approx 1 \rightarrow$ far-sighted agents

Intuitively: near future weighted more heavily than distant future.

Mathematically Convenient: leads to recursive algorithms.

State Value function

$V^\pi(s)$ = expected discounted return following policy π from initial state s .

$$V^\pi(s) = E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s \right]$$

\uparrow expectation operation (computes mean)

* Relating Value function in different states:

$$V^\pi(s) = E^\pi [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s]$$

$$= R(s) + \gamma E^\pi [R(s_1) + \gamma R(s_2) + \gamma^2 R(s_3) + \dots \mid s_0 = s]$$

$$= R(s) + \gamma \sum_{s'} P(s' \mid s, \pi(s)) E^\pi [R(s_1) + \gamma R(s_2) + \dots \mid s_1 = s']$$

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P(s' \mid s, \pi(s)) V^\pi(s')$$

Bellman equation

Action-Value function

$Q^\pi(s, a)$ = expected return from initial state s ,
taking action a , then following policy π .

$$Q^\pi(s, a) = E^\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_0 = a \right]$$

↖ it might be that $a \neq \pi(s)$

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s'} P(s' | s, a) V^\pi(s')$$

Optimality in MDPs

Thm: there is always at least one policy π^*

for which $V^{\pi^*}(s) \geq V^\pi(s)$ for all states and policies.

↖ optimal policy

Goal: how to compute π^* ?