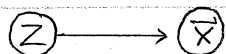


11/9

Review

Gaussian mixture model

* Belief network



$$P(Z=i) = \pi_i \quad P(\vec{X} | Z=i) = \frac{1}{(2\pi)^{n/2} |\Sigma_i|^{1/2}} e^{-\frac{1}{2} (\vec{X} - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{X} - \vec{\mu}_i)}$$

* ML estimation for complete data $\{(Z_t, \vec{X}_t)\}_{t=1}^T$

$$\pi_i = T_i / T \quad \text{where } T_i = \sum_{t=1}^T I(Z_t, i)$$

$$\vec{\mu}_i = \frac{1}{T_i} \sum_{t=1}^T \vec{X}_t I(Z_t, i)$$

$$\Sigma_i = \frac{1}{T_i} \sum_{t=1}^T (\vec{X}_t - \vec{\mu}_i)(\vec{X}_t - \vec{\mu}_i)^T I(Z_t, i)$$

* EM algorithm for incomplete data.

- E-step: Compute posterior probability: Gaussian

$$P(Z=i | \vec{X}_t) = \frac{P(\vec{X}_t | Z=i) P(Z=i)}{\sum_{j=1}^K P(\vec{X}_t | Z=j) P(Z=j)} \quad \text{Gaussian } \pi_i$$

- M-step: update CPTs (by analogy to complete data case)

$$\pi_i \leftarrow \frac{1}{T} \sum_t P(Z=i | \vec{X}_t)$$

$$\vec{\mu}_i \leftarrow \frac{\sum_t \vec{X}_t P(Z=i | \vec{X}_t)}{\sum_t P(Z=i | \vec{X}_t)}$$

effective # points assigned to i th cluster.

$$\Sigma_i \leftarrow \frac{\sum_t (\vec{X}_t - \vec{\mu}_i)(\vec{X}_t - \vec{\mu}_i)^T P(Z=i | \vec{X}_t)}{\sum_t P(Z=i | \vec{X}_t)}$$

from previous update

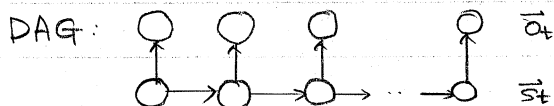
Linear dynamical systems

HMM = discrete dynamical system

How to model continuous states and observations?

Ex: missile locations and radar measurements

* Belief network

variables: $\vec{s}_t \in \mathbb{R}^n$ (hidden) $\vec{o}_t \in \mathbb{R}^m$ (observed)

$$\text{CPTs: } P(\vec{s}_1) = \mathcal{N}(\vec{s}_1; \vec{\mu}_1, \Sigma_1)$$

$$P(\vec{s}_{t+1} | \vec{s}_t) = \mathcal{N}(\vec{s}_{t+1}; A \vec{s}_t, \Sigma_H)$$

$$P(\vec{o}_t | \vec{s}_t) = \mathcal{N}(\vec{o}_t; B \vec{s}_t, \Sigma_O)$$

 $n \times n$ matrix parameter n -dim vector. $n \times n$ matrix $m \times n$ matrix $m \times m$ matrix

* Belief updating

What is $P(\vec{S}_t | \vec{O}_1, \vec{O}_2, \dots, \vec{O}_t)$? "Kalman filtering".

Recall key property: in a BN with all Gaussian CPTs, every marginal and posterior probability is also Gaussian.

Hence $P(\vec{S}_t | \vec{O}_1, \dots, \vec{O}_t)$ must also be Gaussian.

Enough to track $\vec{\mu}_t$ and Σ_t .

Recursion (stated w/o proof):

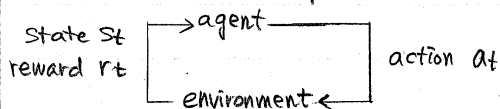
"Kalman gain" matrix corrects for errors

$$\vec{\mu}_{t+1} = \underbrace{A \vec{\mu}_t}_{\substack{\text{predicted state} \\ \text{from evolution of } \vec{\mu}_t \text{ at time } t \\ \text{mean of } P(\vec{S}_{t+1} | \vec{S}_t = \vec{\mu}_t)}} + K_{t+1} \underbrace{(O_{t+1} - BA \vec{\mu}_t)}_{\text{error of predicted observation.}}$$

Quiz #2 ↓

Reinforcement learning

Q: How should embodied / embedded / interactive decision-making agents learn from experience in the world?



Ex: robot navigation,
chess / backgammon,
elevator scheduling

* Challenges

- handling uncertainty
 - exploration vs. exploitation dilemma.
 - temporal credit assignment: delayed vs. immediate rewards.
 - evaluative vs. instructive feedback.
 - complex worlds: balance representational power vs. tractability
- computational guarantees: convergence, optimality, efficiency, etc...

Markov decision process (MDP)

* 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)
- $P(s' | s, a) \triangleq P(\vec{S}_{t+1} = s' | \vec{S}_t = s, a_t = a)$: probability of moving from state s to state s' given action a

* Assumptions

- time independent : $P(S_{t+1} = s' | S_t = s, a_t = a) = P(S_t = s' | S_{t-1} = s, a_{t-1} = a)$
- conditional independence : $P(S_{t+1} | S_t, a_t) = P(S_{t+1} | S_t, a_t, S_{t-1}, a_{t-1}, \dots, S_0, a_0)$

* Reward function

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

$$\text{MDP} = \{ \mathcal{S}, \mathcal{A}, P(s' | s, a), R(s, s', a) \}$$

* Simplifications (for CSE 250A)

- Reward function $R(s, s', a) = R(s) = R_s$ (only depends on current state)
- Bounded rewards $\max_s |R_s| < \infty$.
- Deterministic rewards
- Discrete & finite state space \mathcal{S}
- " " " action space \mathcal{A}

* Ex: backgammon

\mathcal{S} = board position & roll of dice

\mathcal{A} = set of possible moves

$P(s' | s, a)$ = agent's move, opponent rolls dice, opponent's move, agent rolls dice.

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

* Decision - making

- policy: deterministic mapping $\pi: \mathcal{S} \rightarrow \mathcal{A}$ from states to actions
- # policies = $|\mathcal{A}|^{|\mathcal{S}|}$
- dynamics $P(s' | s, \pi(s))$

- experience state $s_0 \xrightarrow{\pi(s_0)} s_1 \xrightarrow{\pi(s_1)} s_2 \rightarrow \dots$
 reward $r_0 \quad \quad \quad r_1 \quad \quad \quad r_2$

* How to measure long-term return (accumulated rewards)?

- return = $\frac{1}{T} [r_0 + r_1 + r_2 + \dots + r_{T-1}]$
 undiscounted return w/ finite horizon T

- return = $\lim_{T \rightarrow \infty} \frac{1}{T} [r_0 + r_1 + \dots + r_{T-1}]$
 undiscounted infinite horizon

- discounted infinite horizon return.

discount factor $0 \leq \gamma < 1$ return = $r_0 + \gamma r_1 + \gamma^2 r_2 + \dots = \sum_{t=0}^{\infty} \gamma^t r_t$

possibilities: $\gamma = 0 \rightarrow$ only immediate reward matters

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

$(1-\gamma) \ll 1 \rightarrow$ far-sighted agent

* Justification

intuitive: near future weighted more than distant future

mathematically convenient: leads to recursive algorithms

* State value function (over discounted infinite horizon)

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

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

- Maximizing expected return different than:

- maximizing worst-case return

- maximizing best-case return.