

TEMPORAL PROBABILITY MODELS

CHAPTER 15, SECTIONS 1–5

Chapter 15, Sections 1–5 1

Outline

- ◇ Time and uncertainty
- ◇ Inference: filtering, prediction, smoothing
- ◇ Hidden Markov models
- ◇ Dynamic Bayesian networks

Time and uncertainty

The world changes; we need to track and predict it

Diabetes management vs vehicle diagnosis

Basic idea: copy state and evidence variables for each time step

X_t = set of unobservable state variables at time t
 e.g., *BloodSugar*, *StomachContents*, etc.

E_t = set of observable evidence variables at time t
 e.g., *MeasuredBloodSugar*, *PulseRate*, *FoodEaten*

This assumes **discrete time**: step size depends on problem

Notation: $X_{0:t} = X_0, X_{0+1}, \dots, X_{t-1}, X_t$

Chapter 15, Sections 1–5 2

Chapter 15, Sections 1–5 3

Markov processes (Markov chains)

Construct a Bayes net from these variables: parents? CPTs?

Chapter 15, Sections 1–5 4

Markov processes (Markov chains)

Construct a Bayes net from these variables: parents? CPTs?

Markov assumption: X_t depends on **bounded** subset of $X_{0:t-1}$

First-order Markov process: $P(X_t | X_{0:t-1}) = P(X_t | X_{t-1})$

Second-order Markov process: $P(X_t | X_{0:t-1}) = P(X_t | X_{t-2}, X_{t-1})$

First-order



Second-order



Stationary process: transition model $P(X_t | X_{t-1})$ fixed for all t

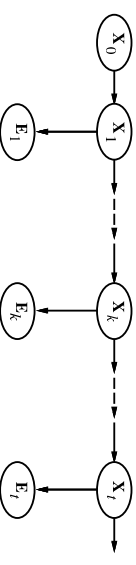
Chapter 15, Sections 1–5 5

Hidden Markov Model (HMM)

Sensor Markov assumption: $P(E_t | X_{0:t}, E_{1:t-1}) = P(E_t | X_t)$

Stationary process: transition model $P(X_t | X_{t-1})$ and sensor model $P(E_t | X_t)$ fixed for all t

HMM is a special type of Bayes net, X_t is single discrete random variable:



with joint probability distribution

$P(X_{0:t}, E_{1:t}) = ?$

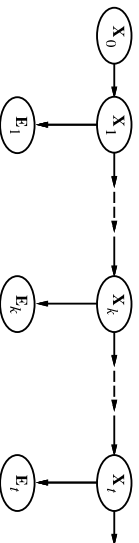
Chapter 15, Sections 1–5 6

Hidden Markov Model (HMM)

Sensor Markov assumption: $P(E_t | X_{0:t}, E_{1:t-1}) = P(E_t | X_t)$

Stationary process: transition model $P(X_t | X_{t-1})$ and sensor model $P(E_t | X_t)$ fixed for all t

HMM is a special type of Bayes net, X_t is single discrete random variable:



with joint probability distribution

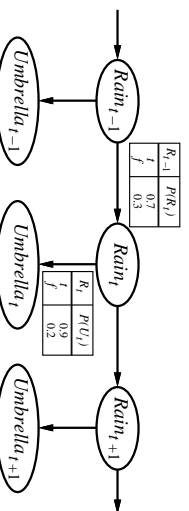
$$P(X_{0:t}, E_{1:t}) = P(X_0) \prod_{t=1}^t P(X_t | X_{t-1}) P(E_t | X_t)$$

Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(X_{t+1} | e_{1:t+1}) = f(e_{t+1}, P(X_t | e_{1:t}))$$

Example



First-order Markov assumption not exactly true in real world!

Possible fixes:

1. **Increase order** of Markov process
2. **Augment state**, e.g., add *Temp*, *Pressure*, etc.

Example: robot motion.

Augment position and velocity with *Battery*

Inference tasks

Filtering: $P(X_t | e_{1:t})$

belief state—input to the decision process of a rational agent

Prediction: $P(X_{t+k} | e_{1:t})$ for $k > 0$

evaluation of possible action sequences;

like filtering without the evidence

Smoothing: $P(X_k | e_{1:t})$ for $0 \leq k < t$

better estimate of past states, essential for learning

Most likely explanation: $\text{argmax}_{x_{1:t}} P(x_{1:t} | e_{1:t})$

speech recognition, decoding with a noisy channel

Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(X_{t+1} | e_{1:t+1}) = f(e_{t+1}, P(X_t | e_{1:t}))$$

$$P(X_{t+1} | e_{1:t+1}) = P(X_{t+1} | e_{1:t}, e_{t+1})$$

$$= \alpha P(e_{t+1} | X_{t+1}, e_{1:t}) P(X_{t+1} | e_{1:t})$$

$$= \alpha P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t})$$

Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(X_{t+1} | e_{1:t+1}) = f(e_{t+1}, P(X_t | e_{1:t}))$$

$$P(X_{t+1} | e_{1:t+1}) = P(X_{t+1} | e_{1:t}, e_{t+1})$$

$$= \alpha P(e_{t+1} | X_{t+1}, e_{1:t}) P(X_{t+1} | e_{1:t})$$

$$= \alpha P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t})$$

i.e., **prediction + estimation**. Prediction by summing out X_t :

$$P(X_{t+1} | e_{1:t+1}) = \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1}, x_t | e_{1:t})$$

$$= \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1} | x_t, e_{1:t}) P(x_t | e_{1:t})$$

$$= \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1} | x_t) P(x_t | e_{1:t})$$

Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(\mathbf{X}_{t+1}|e_{1:t+1}) = f(e_{t+1}, P(\mathbf{X}_t|e_{1:t}))$$

$$\begin{aligned} P(\mathbf{X}_{t+1}|e_{1:t+1}) &= P(\mathbf{X}_{t+1}|e_{1:t}, e_{t+1}) \\ &= \alpha P(e_{t+1}|\mathbf{X}_{t+1}, e_{1:t})P(\mathbf{X}_{t+1}|e_{1:t}) \\ &= \alpha P(e_{t+1}|\mathbf{X}_{t+1})P(\mathbf{X}_{t+1}|e_{1:t}) \end{aligned}$$

I.e., **prediction + estimation**. Prediction by summing out \mathbf{X}_t :

$$\begin{aligned} P(\mathbf{X}_{t+1}|e_{1:t+1}) &= \alpha P(e_{t+1}|\mathbf{X}_{t+1})\sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}, \mathbf{x}_t|e_{1:t}) \\ &= \alpha P(e_{t+1}|\mathbf{X}_{t+1})\sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t, e_{1:t})P(\mathbf{x}_t|e_{1:t}) \\ &= \alpha P(e_{t+1}|\mathbf{X}_{t+1})\sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t)P(\mathbf{x}_t|e_{1:t}) \end{aligned}$$

$f_{t+1} = \text{FORWARD}(f_{1:t}, e_{t+1})$ where $f_{1:t} = P(\mathbf{X}_t|e_{1:t})$
Time and space **constant** (independent of t)

Most likely explanation

Most likely sequence \neq sequence of most likely states!!!

Most likely path to each \mathbf{X}_{t+1}
= most likely path to **some** \mathbf{x}_t plus one more step

$$\begin{aligned} &\max_{\mathbf{x}_t} P(\mathbf{x}_t, \dots, \mathbf{x}_t, \mathbf{X}_{t+1}|e_{1:t+1}) \\ &= P(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (P(\mathbf{X}_{t+1}|\mathbf{x}_t) \max_{\mathbf{x}_{1:t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{x}_t|e_{1:t})) \end{aligned}$$

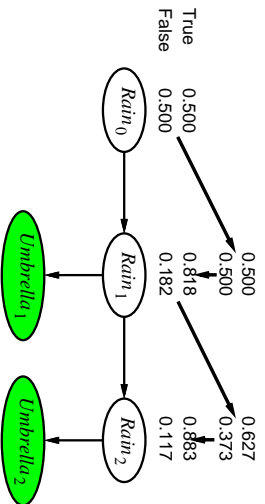
Identical to filtering, except $f_{1:t}$ replaced by

$$\mathbf{m}_{1:t} = \max_{\mathbf{x}_{1:t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{X}_t|e_{1:t}),$$

I.e., $\mathbf{m}_{1:t}(i)$ gives the probability of the most likely path to state i .
Update has sum replaced by max, giving the Viterbi algorithm:

$$\mathbf{m}_{1:t+1} = P(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (P(\mathbf{X}_{t+1}|\mathbf{x}_t)\mathbf{m}_{1:t})$$

Filtering example

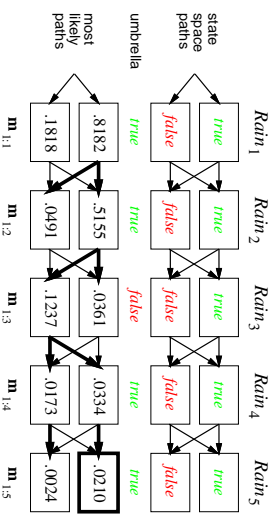


$$P(\mathbf{X}_{t+1}|e_{1:t+1}) = \alpha P(e_{t+1}|\mathbf{X}_{t+1})\sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t)P(\mathbf{x}_t|e_{1:t})$$

e_{t-1}	$P(e_{t-1})$	e_t	$P(e_t)$
f	0.7	f	0.9
f	0.3	f	0.2

Most likely explanation

Viterbi example



Implementation Issues

Viterbi message: $\mathbf{m}_{1:t+1} = P(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (P(\mathbf{X}_{t+1}|\mathbf{x}_t)\mathbf{m}_{1:t})$

or filtering update: $P(\mathbf{X}_{t+1}|e_{1:t+1}) = \alpha P(e_{t+1}|\mathbf{X}_{t+1})\sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1}|\mathbf{x}_t)P(\mathbf{x}_t|e_{1:t})$

What is $10^{-6} \cdot 10^{-6} \cdot 10^{-6}$?

Implementation Issues

Viterbi message: $\mathbf{m}_{t:t+1} = \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) \mathbf{m}_{t:t})$
 or filtering update: $\mathbf{P}(\mathbf{X}_{t+1}|e_{t:t+1}) = \alpha \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) \mathcal{P}(\mathbf{x}_t|e_{1:t})$

What is $10^{-6} \cdot 10^{-6} \cdot 10^{-6}$?

What is floating point arithmetic precision?

Hidden Markov models

\mathbf{X}_t is a single, discrete variable (usually \mathbf{E}_t is too)
 Domain of X_t is $\{1, \dots, S\}$

Transition matrix $\mathbf{T}_{ij} = P(X_t = j | X_{t-1} = i)$, e.g., $\begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}$

Sensor matrix \mathbf{O}_t for each time step, diagonal elements $P(e_t | X_t = i)$
 e.g., with $U_1 = true$, $\mathbf{O}_1 = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}$

Forward messages as column vectors:

$$\mathbf{f}_{t:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^T \mathbf{f}_{t:t}$$

Implementation Issues

Viterbi message: $\mathbf{m}_{t:t+1} = \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) \mathbf{m}_{t:t})$
 or filtering update: $\mathbf{P}(\mathbf{X}_{t+1}|e_{t:t+1}) = \alpha \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) \mathcal{P}(\mathbf{x}_t|e_{1:t})$

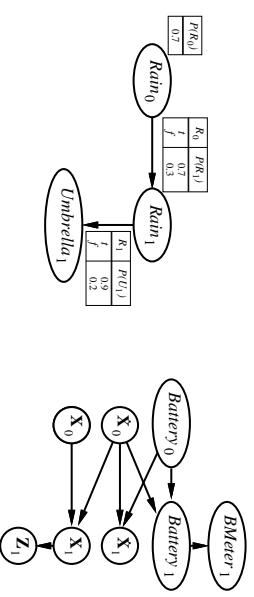
What is $10^{-6} \cdot 10^{-6} \cdot 10^{-6}$?

What is floating point arithmetic precision?

$10^{-6} \cdot 10^{-6} \cdot 10^{-6} = 0$

Dynamic Bayesian networks

$\mathbf{X}_t, \mathbf{E}_t$ contain arbitrarily many variables in a replicated Bayes net



Answer?

Use either:

- Rescaling, multiply values by a (large) constant
- Logsum trick (Assignment 5)

log is monotone increasing, so:

$$\arg \max f(x) = \arg \max \log f(x)$$

Also,

$$\log(a \cdot b) = \log a + \log b$$

Therefore, work with sums of logarithms of probabilities, rather than products of probabilities:

$$\begin{aligned} \mathbf{m}_{t:t+1} &= \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) \mathbf{m}_{t:t}) \\ \rightarrow \log \mathbf{m}_{t:t+1} &= \log \mathbf{P}(e_{t+1}|\mathbf{X}_{t+1}) + \max_{\mathbf{x}_t} (\log \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) + \log \mathbf{m}_{t:t}) \end{aligned}$$

Summary

Temporal models use state and sensor variables replicated over time

Markov assumptions and stationarity assumption, so we need

- transition model $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$
- sensor model $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$

Tasks are filtering, prediction, smoothing, most likely sequence:

all done recursively with constant cost per time step

Hidden Markov models have a single discrete state variable: used for speech recognition

Dynamic Bayes nets subsume HMMs: exact update intractable

Example Umbrella Problems

Filtering:

$$\mathbf{f}_{t+1} := \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) = \alpha \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \Sigma_{\mathbf{x}} \mathbf{P}(\mathbf{X}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t})$$

Viterbi: $\mathbf{m}_{1:t+1} = \mathbf{P}(\mathbf{e}_{1:t+1} | \mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{x}_t) \mathbf{m}_{1:t})$

R_{t+1}	$P(R_t=0)$	$P(R_t=1)$
t	0.7	0.3
f	0.3	0.7

R_t	$P(U_t=0)$	$P(U_t=1)$
t	0.9	0.1
f	0.2	0.8

$$\mathbf{P}(R_3 | \neg u_1, u_2, \neg u_3) = ?$$

$$\arg \max_{R_{1:3}} \mathbf{P}(R_{1:3} | \neg u_1, u_2, \neg u_3) = ?$$