

CSE 250A. Principles of AI

Probabilistic Reasoning and Decision-Making

Lecture 14 – Hidden Markov models

Lawrence Saul
Department of Computer Science and Engineering
University of California, San Diego

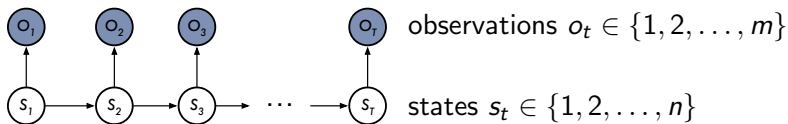
Fall 2021

Outline

- 1 Review & HW
- 2 Forward algorithm
- 3 Viterbi algorithm
- 4 Backward algorithm

Hidden Markov models

- Belief network**



- Parameters**

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

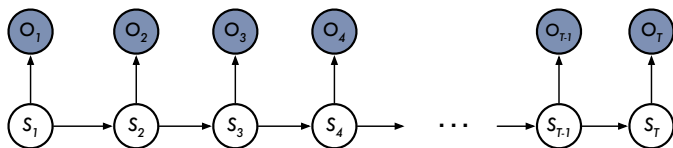
$$b_{ik} = P(O_t=k|S_t=i) \quad \text{emission matrix}$$

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

- Notation**

Sometimes we'll write $b_i(k) = b_{ik}$ to avoid double subscripts.

Key computations in HMMs



Inference

- 1 How to compute the likelihood $P(o_1, o_2, \dots, o_T)$?
- 2 How to compute the most likely state sequence $\arg\max_{\vec{s}} P(\vec{s}|\vec{o})$?
- 3 How to update beliefs by computing $P(s_t|o_1, o_2, \dots, o_t)$?

Learning

How to estimate parameters $\{\pi_i, a_{ij}, b_{ik}\}$ that maximize the log-likelihood of observed sequences?

HW 7.1 Decoding a hidden message

- What you are given

HMM parameters $\{a_{ij}, b_{ik}, \pi_i\}$ where $n=27$ and $m=2$

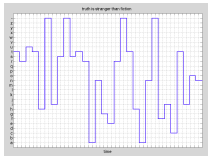
Binary observations $\{o_1, o_2, \dots, o_T\}$ where $T \approx 10^5$

- What you will compute

Most likely sequence of hidden states

$$\arg \max_{s_1, s_2, \dots, s_T} P(s_1, s_2, \dots, s_T | o_1, o_2, \dots, o_T)$$

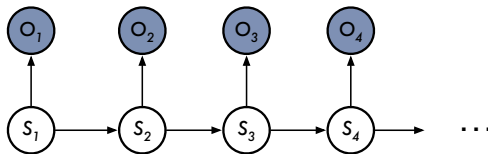
- What you will plot



Convert states 1–27 to letters A–Z, $\langle \text{SPC} \rangle$.
Ignore repeated states to reveal a message:

Truth is stranger than fiction.

HW 7.4 Belief updating



How to compute $P(s_t | o_1, o_2, \dots, o_t)$?

Important for real-time monitoring of dynamical systems.



Outline

- ① Review & HW
- ② **Forward algorithm**
- ③ Viterbi algorithm
- ④ Backward algorithm

Computing the likelihood $P(o_1, o_2, \dots, o_T)$

- One strategy that is not quite right

Compute $P(o_1)$.

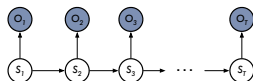
Compute $P(o_1, o_2)$.

Compute $P(o_1, o_2, o_3)$.

Etc.

- Why this is appealing

We expect an iterative approach of this kind in a BN with a chain-like DAG.



- Why this doesn't work

We need to build up a slightly different quantity.

This is what the **forward algorithm** does.

Computing $P(o_1, o_2, \dots, o_t, S_t = i)$

• Definition

For a particular sequence of observations $\{o_1, o_2, \dots, o_T\}$,
 define the matrix with elements:

$$\alpha_{it} = P(o_1, o_2, \dots, o_t, S_t = i) \quad n \text{ rows} \left\{ \begin{bmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1,T-1} & \alpha_{1T} \\ \alpha_{21} & \alpha_{22} & \cdots & \alpha_{2,T-1} & \alpha_{2T} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \alpha_{n1} & \alpha_{n2} & \cdots & \alpha_{n,T-1} & \alpha_{nT} \end{bmatrix} \right.$$

• First column ($t = 1$)

$$\begin{aligned} \alpha_{i1} &= P(o_1, S_1 = i) \\ &= P(S_1 = i) P(o_1 | S_1 = i) \\ &= \pi_i b_i(o_1) \end{aligned}$$

product rule

CPTs

Computing $\alpha_{it} = P(o_1, o_2, \dots, o_t, S_t = i)$

- Next columns ($t > 1$)

$$\begin{aligned}
 \alpha_{j,t+1} &= P(o_1, o_2, \dots, o_{t+1}, S_{t+1}=j) \\
 &= \sum_{i=1}^n P(o_1, o_2, \dots, o_{t+1}, S_t=i, S_{t+1}=j) \quad \boxed{\text{marginalization}} \\
 &= \sum_{i=1}^n \left[P(o_1, o_2, \dots, o_t, S_t=i) \cdot \right. \\
 &\quad P(S_{t+1}=j | o_1, o_2, \dots, o_t, S_t=i) \cdot \\
 &\quad \left. P(o_{t+1} | o_1, o_2, \dots, o_t, S_t=i, S_{t+1}=j) \right] \quad \boxed{\text{product rule}} \\
 &= \sum_{i=1}^n \left[P(o_1, o_2, \dots, o_t, S_t=i) P(S_{t+1}=j | S_t=i) P(o_{t+1} | S_{t+1}=j) \right] \quad \boxed{\text{CI}} \\
 &= \sum_{i=1}^n \alpha_{it} a_{ij} b_j(o_{t+1}) \quad \boxed{\text{CPTs}}
 \end{aligned}$$

Forward algorithm

The forward algorithm fills in the matrix of α_{it} elements one column at a time:

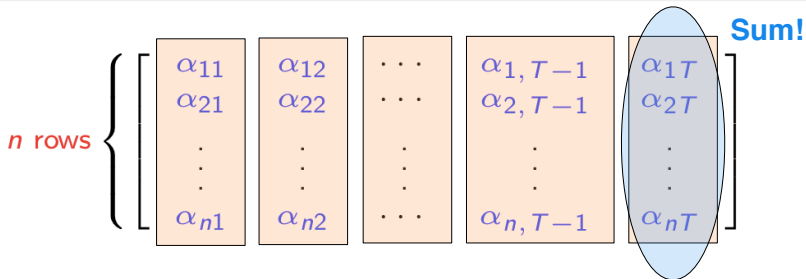
$$\begin{aligned}\alpha_{i1} &= \pi_i b_i(o_1) \\ \alpha_{j,t+1} &= \sum_{i=1}^n \alpha_{it} a_{ij} b_j(o_{t+1})\end{aligned}$$

n rows

$$\left\{ \begin{bmatrix} \alpha_{11} \\ \alpha_{21} \\ \vdots \\ \alpha_{n1} \end{bmatrix} \begin{bmatrix} \alpha_{12} \\ \alpha_{22} \\ \vdots \\ \alpha_{n2} \end{bmatrix} \begin{bmatrix} \cdots \\ \cdots \\ \vdots \\ \cdots \end{bmatrix} \begin{bmatrix} \alpha_{1,T-1} \\ \alpha_{2,T-1} \\ \vdots \\ \alpha_{n,T-1} \end{bmatrix} \begin{bmatrix} \alpha_{1T} \\ \alpha_{2T} \\ \vdots \\ \alpha_{nT} \end{bmatrix} \right\}$$

Warning: for long sequences, beware of numerical underflow ...

Computing the likelihood $P(o_1, o_2, \dots, o_T)$



$$P(o_1, o_2, \dots, o_T)$$

$$= \sum_{i=1}^n P(o_1, o_2, \dots, o_T, s_T = i)$$

marginalization

$$= \sum_{i=1}^n \alpha_{iT}$$

sum of last column

Outline

- ① Review & HW
- ② Forward algorithm
- ③ **Viterbi algorithm**
- ④ Backward algorithm

The most likely sequence of hidden states

$$\{s_1^*, s_2^*, \dots, s_T^*\}$$

$$= \operatorname{argmax}_{s_1, s_2, \dots, s_T} P(s_1, s_2, \dots, s_T | o_1, o_2, \dots, o_T)$$

$$= \operatorname{argmax}_{s_1, s_2, \dots, s_T} \left[\frac{P(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T)}{P(o_1, o_2, \dots, o_T)} \right]$$

product
rule

$$= \operatorname{argmax}_{s_1, s_2, \dots, s_T} P(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T)$$

no states in
denominator

$$= \operatorname{argmax}_{s_1, s_2, \dots, s_T} \log P(s_1, s_2, \dots, s_T, o_1, o_2, \dots, o_T)$$

log is
monotonic

Q: Why is this progress?

A: Because **logs** of **joint probabilities** are easy to compute.

Computing the most likely state sequence

- **One strategy that is not quite right**

Compute $\operatorname{argmax}_{s_1} \log P(s_1, o_1)$.

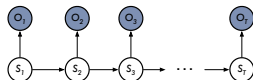
Compute $\operatorname{argmax}_{s_1, s_2} \log P(s_1, s_2, o_1, o_2)$.

Compute $\operatorname{argmax}_{s_1, s_2, s_3} \log P(s_1, s_2, s_3, o_1, o_2, o_3)$.

Etc.

- **Why this is appealing**

We expect an iterative approach of this kind in a BN with a chain-like DAG.



- **Why this doesn't work**

We need to build up a slightly different quantity.

This is what the **Viterbi algorithm** does.

The matrix ℓ^*

• Definition

For a particular sequence of observations $\{o_1, o_2, \dots, o_T\}$, we define the following matrix:

$$\ell_{it}^* = \max_{s_1, s_2, \dots, s_{t-1}} \log P(s_1, s_2, \dots, s_{t-1}, S_t = i, o_1, o_2, \dots, o_t)$$

$$n \text{ rows} \left\{ \begin{bmatrix} \ell_{11}^* & \ell_{12}^* & \cdots & \ell_{1,T-1}^* & \ell_{1T}^* \\ \ell_{21}^* & \ell_{22}^* & \cdots & \ell_{2,T-1}^* & \ell_{2T}^* \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \ell_{n1}^* & \ell_{n2}^* & \cdots & \ell_{n,T-1}^* & \ell_{nT}^* \end{bmatrix} \right.$$

• Intuition

ℓ_{it}^*	log-probability of the t -step path of hidden states s_1, s_2, \dots, s_t that best explains the observations o_1, o_2, \dots, o_t and ends at state $S_t = i$ at time t
---------------	--

Computing the matrix ℓ^*

$$\ell_{it}^* = \max_{s_1, s_2, \dots, s_{t-1}} \log P(s_1, s_2, \dots, s_{t-1}, S_t = i, o_1, o_2, \dots, o_t)$$

$$n \text{ rows} \left\{ \begin{bmatrix} \ell_{11}^* \\ \ell_{21}^* \\ \vdots \\ \ell_{n1}^* \end{bmatrix} \quad \ell_{12}^* \quad \cdots \quad \ell_{1,T-1}^* \quad \ell_{1T}^* \right. \\ \left. \begin{bmatrix} \ell_{22}^* \\ \vdots \\ \ell_{n2}^* \end{bmatrix} \quad \cdots \quad \ell_{2,T-1}^* \quad \ell_{2T}^* \right. \\ \left. \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad \cdots \quad \vdots \quad \vdots \right. \\ \left. \begin{bmatrix} \ell_{n,T-1}^* \quad \ell_{nT}^* \end{bmatrix} \right.$$

- First column ($t = 1$)

$$\ell_{i1}^* = \log P(S_1 = i, o_1)$$

$$= \log \left[P(S_1 = i) P(o_1 | S_1 = i) \right]$$

product rule

$$= \log \pi_i + \log b_i(o_1)$$

CPTs

Computing the matrix ℓ^* • Next columns ($t > 1$)

$$\begin{aligned}
\ell_{j,t+1}^* &= \max_{s_1, \dots, s_t} \log P(s_1, \dots, s_t, S_{t+1}=j, o_1, \dots, o_{t+1}) \\
&= \max_i \max_{s_1, \dots, s_{t-1}} \log P(s_1, \dots, s_{t-1}, S_t=i, S_{t+1}=j, o_1, \dots, o_{t+1}) \\
&= \max_i \max_{s_1, \dots, s_{t-1}} \log \left[P(s_1, \dots, s_{t-1}, S_t=i, o_1, \dots, o_t) \cdot \boxed{\text{product rule}} \right. \\
&\quad \left. P(S_{t+1}=j | s_1, \dots, s_{t-1}, S_t=i, o_1, \dots, o_t) \cdot \right. \\
&\quad \left. P(o_{t+1} | s_1, \dots, s_{t-1}, S_t=i, S_{t+1}=j, o_1, \dots, o_t) \right] \\
&= \max_i \max_{s_1, \dots, s_{t-1}} \log \left[P(s_1, \dots, s_{t-1}, S_t=i, o_1, \dots, o_t) \cdot \right. \\
&\quad \left. P(S_{t+1}=j | S_t=i) \cdot P(o_{t+1} | S_{t+1}=j) \right] \boxed{\text{CI}} \\
&= \max_i \max_{s_1, \dots, s_{t-1}} \left[\log P(s_1, \dots, s_{t-1}, S_t=i, o_1, \dots, o_t) + \log a_{ij} + \log b_j(o_{t+1}) \right] \\
&= \max_i \left[\ell_{it}^* + \log a_{ij} \right] + \log b_j(o_{t+1})
\end{aligned}$$

Summary

We have shown how to efficiently compute the matrix ℓ^* one column at a time, from left to right:

$$\begin{aligned}\ell_{i1}^* &= \log \pi_i + \log b_i(o_1) \\ \ell_{j,t+1}^* &= \max_i \left[\ell_{it}^* + \log a_{ij} \right] + \log b_j(o_{t+1})\end{aligned}$$

$$n \text{ rows} \left\{ \begin{bmatrix} \ell_{11}^* \\ \ell_{21}^* \\ \vdots \\ \ell_{n1}^* \end{bmatrix} \begin{bmatrix} \ell_{12}^* \\ \ell_{22}^* \\ \vdots \\ \ell_{n2}^* \end{bmatrix} \begin{bmatrix} \cdots \\ \cdots \\ \vdots \\ \cdots \end{bmatrix} \begin{bmatrix} \ell_{1,T-1}^* \\ \ell_{2,T-1}^* \\ \vdots \\ \ell_{n,T-1}^* \end{bmatrix} \begin{bmatrix} \ell_{1T}^* \\ \ell_{2T}^* \\ \vdots \\ \ell_{nT}^* \end{bmatrix} \right\}$$

But how do we derive $\{s_1^*, s_2^*, \dots, s_T^*\}$ from this matrix?

Computing $\{s_1^*, s_2^*, \dots, s_T^*\}$

- Form one more matrix:

$$\Phi_{t+1}(j) = \arg \max_i \left[\ell_{it}^* + \log a_{ij} \right]$$

Intuitively, given the observations o_1, o_2, \dots, o_{t+1} , record the most likely state at time t given that $S_{t+1}=j$.

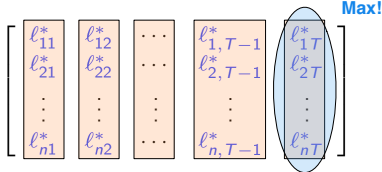
- Compute the most likely states by **backtracking**:

$$s_T^* = \arg \max_i \left[\ell_{iT}^* \right]$$

for $t = T-1$ **to** **1**

$$s_t^* = \Phi_{t+1}(s_{t+1}^*)$$

end



Summary of Viterbi algorithm

- **Fill ℓ^* matrix from left to right:**

$$\boxed{t = 1} \quad \ell_{i1}^* = \log \pi_i + \log b_i(o_1)$$

$$\boxed{t > 1} \quad \ell_{j,t+1}^* = \max_i \left[\ell_{it}^* + \log a_{ij} \right] + \log b_j(o_{t+1})$$

- **Backtrack through ℓ^* from right to left:**

$$\boxed{t = T} \quad s_T^* = \operatorname{argmax}_i \left[\ell_{iT}^* \right]$$

$$\boxed{t < T} \quad s_t^* = \operatorname{argmax}_i \left[\ell_{it}^* + \log a_{is_{t+1}^*} \right]$$

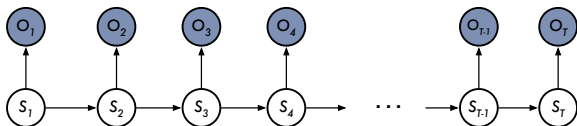
- **Where you've seen this before:**

This algorithm is an instance of **dynamic programming**.
Sometimes $\{s_1^*, s_2^*, \dots, s_T^*\}$ is called the **Viterbi path**.

Outline

- ① Review & HW
- ② Forward algorithm
- ③ Viterbi algorithm
- ④ **Backward algorithm**

Learning in HMMs



Given: one or more sequences of observations $\{o_1, o_2, \dots, o_T\}$.
 For simplicity, we'll assume just one.

Goal: estimate $\{\pi_i, a_{ij}, b_{ik}\}$ to maximize $P(o_1, o_2, \dots, o_T)$,
 the likelihood of the observed data.

Assume: the cardinality n of the hidden state space is fixed.
 $s_t \in \{1, 2, \dots, n\}$

EM algorithm for HMMs

- CPTs to re-estimate:

$$\pi_i = P(S_1 = i)$$

$$a_{ij} = P(S_{t+1} = j | S_t = i)$$

$$b_{ik} = P(O_t = k | S_t = i)$$

- How EM works in general:

To re-estimate $P(X_i = x | \text{pa}_i = \pi)$ in the M-step,
we must compute $P(X_i = x, \text{pa}_i = \pi | V)$ in the E-step.

- E-step in HMMs must compute:

$$P(S_1 = i | o_1, o_2, \dots, o_T)$$

$$P(S_{t+1} = j, S_t = i | o_1, o_2, \dots, o_T)$$

$$P(O_t = k, S_t = i | o_1, o_2, \dots, o_T)$$

How to efficiently compute these posteriors?

We need one more matrix ...

Analogous to $\alpha_{it} = P(o_1, o_2, \dots, o_t, S_t = i)$,
define $\beta_{it} = P(o_{t+1}, o_{t+2}, \dots, o_T | S_t = i)$.

$$n \text{ rows } \left\{ \begin{bmatrix} \beta_{11} & \beta_{12} & \cdots & \beta_{1,T-1} & \beta_{1T} \\ \beta_{21} & \beta_{22} & \cdots & \beta_{2,T-1} & \beta_{2T} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \beta_{n1} & \beta_{n2} & \cdots & \beta_{n,T-1} & \beta_{nT} \end{bmatrix} \right.$$

Understand the **differences** between these matrices:

- α_{it} predicts observations up to and including time t .
- β_{it} predicts observations from **time $t + 1$** to time T .

Computing $\beta_{it} = P(o_{t+1}, o_{t+2}, \dots, o_T | S_t = i)$ • Last column ($t = T$)

$$\beta_{iT} = P(______ | S_T = i) \quad \text{What does this mean?}$$

Note: β_{it} computes the probability of the future given $S_t = i$.

But we don't see *any* observations beyond time T .

Put another way, the future after time T is unspecified.

What is the probability of some unspecified future occurring?

By definition, we set:

$$\beta_{iT} = 1 \quad \text{for all } i \in \{1, 2, \dots, n\}$$

Computing $\beta_{it} = P(o_{t+1}, o_{t+2}, \dots, o_T | S_t = i)$ • Previous columns ($t < T$)

$$\begin{aligned}\beta_{it} &= P(o_{t+1}, o_{t+2}, \dots, o_T | S_t = i) \\&= \sum_{j=1}^n P(S_{t+1} = j, o_{t+1}, o_{t+2}, \dots, o_T | S_t = i) \quad \boxed{\text{marginalization}} \\&= \sum_{j=1}^n \left[P(S_{t+1} = j | S_t = i) \cdot \right. \\&\quad \left. P(o_{t+1} | S_t = i, S_{t+1} = j) \cdot \right. \\&\quad \left. P(o_{t+2}, \dots, o_T | S_t = i, S_{t+1} = j, o_{t+1}) \right] \quad \boxed{\text{product rule}} \\&= \sum_{j=1}^n \left[P(S_{t+1} = j | S_t = i) P(o_{t+1} | S_{t+1} = j) P(o_{t+2}, \dots, o_T | S_{t+1} = j) \right] \quad \boxed{\text{CI}} \\&= \sum_{j=1}^n a_{ij} b_j(o_{t+1}) \beta_{j,t+1} \quad \boxed{\text{CPTs}}\end{aligned}$$

Backward algorithm

The backward algorithm fills in the matrix of β_{it} elements one column at a time:

$$\beta_{iT} = 1 \quad \text{for } i \in \{1, 2, \dots, n\}$$

$$\beta_{it} = \sum_{j=1}^n a_{ij} b_j(o_{t+1}) \beta_{j,t+1}$$

n rows

$$\left\{ \begin{bmatrix} \beta_{11} \\ \beta_{21} \\ \vdots \\ \beta_{n1} \end{bmatrix} \begin{bmatrix} \beta_{12} \\ \beta_{22} \\ \vdots \\ \beta_{n2} \end{bmatrix} \begin{bmatrix} \cdots \\ \cdots \\ \vdots \\ \cdots \end{bmatrix} \begin{bmatrix} \beta_{1,T-1} \\ \beta_{2,T-1} \\ \vdots \\ \beta_{n,T-1} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \right\}$$

Next lecture: **EM with the forward-backward algorithm.**