

Ch2. Hidden Markov models: definition and properties

HyeRim Park

July 21, 2023

Table of contents

1. The basics: Hidden Markov model
2. The likelihood

Table of contents

1. The basics: Hidden Markov model

2. The likelihood

Definition and notation

- A **hidden Markov model** $\{X_t : t \in \mathbb{N}\}$ is a particular kind of dependent mixture.
- The model consists of two parts:
 - Unobserved 'parameter process' $\{C_t : t = 1, 2, \dots\}$ satisfying the Markov property.
 - The 'state-dependent process' $\{X_t : t = 1, 2, \dots\}$
- With $\mathbf{X}^{(t)}$ and $\mathbf{C}^{(t)}$ representing the histories from time 1 to time t ,

$$\Pr(C_t \mid \mathbf{C}^{(t-1)}) = \Pr(C_t \mid C_{t-1}), \quad t = 2, 3, \dots \quad (1)$$

$$\Pr(X_t \mid \mathbf{X}^{(t-1)}, \mathbf{C}^{(t)}) = \Pr(X_t \mid C_t), \quad t \in \mathbb{N} \quad (2)$$

Definition and notation

- The distribution of X_t depends only on the current state C_t .
- If the Markov chain $\{C_t\}$ has m states, we call $\{X_t\}$ an m -state HMM.

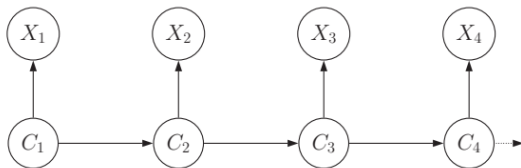


Figure 1: Directed graph of basic HMM.

Definition and notation

- Introduce some notation which cover both discrete- and continuous-valued observations.
- In the case of discrete observations, we define

$$p_i(x) = \Pr(X_t = x \mid C_t = i) \quad \text{for } i = 1, 2, \dots, m$$

- p_i is the probability mass function of X_t if the Markov chain is in state i at time t .
- The m distributions p_i as the **state-dependent distributions** of the model.

Homogeneous and Stationary Markov chain

- Define **transition probabilities**:

$$\gamma_{ij}(t) = \Pr(C_{s+t} = j | C_s = i)$$

- If these probabilities do not depend on s , the Markov chain is called **homogeneous**.
- A Markov chain with transition probability matrix Γ is said to have **stationary distribution** δ if $\delta\Gamma = \delta$ and $\delta\mathbf{1}' = 1$.

Homogeneous and Stationary Markov chain

- Denote these by the row vector

$$u(t) = (\Pr(C_t = 1), \dots, \Pr(C_t = m)), \quad t \in \mathbb{N}$$

- Homogeneity alone would not be sufficient to render the Markov chain a stationary process.
- Stationary for homogeneous Markov chains that have the additional property that the initial distribution $u(1)$ is the stationary distribution.

Marginal distributions

- We shall often need the marginal distribution of X_t and (X_t, X_{t+k}) .
- Assume that the Markov chain is homogeneous but not necessarily stationary.
- For convenience the derivation is given only for discrete state-dependent distributions.

Marginal distributions: Univariate distributions

- For discrete-valued observations X_t ,

$$\begin{aligned}\Pr(X_t = x) &= \sum_{i=1}^m \Pr(C_t = i) \Pr(X_t = x \mid C_t = i) \\ &= \sum_{i=1}^m u_i(t) p_i(x)\end{aligned}$$

where $u_i(t) = \Pr(C_t = i)$ for $t = 1, \dots, T$.

Marginal distributions: Univariate distributions

- This expression can conveniently be rewritten in matrix notation:

$$\begin{aligned}\Pr(X_t = x) &= (u_1(t), \dots, u_m(t)) \begin{pmatrix} p_1(x) & & 0 \\ & \ddots & \\ 0 & & p_m(x) \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ &= \mathbf{u}(t) \mathbf{P}(x) \mathbf{1},\end{aligned}$$

where $\mathbf{P}(x)$ is defined as the diagonal matrix with i th diagonal element $p_i(x)$.

Marginal distributions: Univariate distributions

- From equation $\mathbf{u}(t+1) = \mathbf{u}(t)\mathbf{\Gamma}$, $\mathbf{u}(t) = \mathbf{u}(1)\mathbf{\Gamma}^{t-1}$ and then

$$\Pr(X_t = x) = \mathbf{u}(1)\mathbf{\Gamma}^{t-1}\mathbf{P}(x)\mathbf{1} \quad (3)$$

- If the Markov chain is stationary, with stationary distribution $\boldsymbol{\delta}$, in that case $\boldsymbol{\delta}\mathbf{\Gamma}^{t-1} = \boldsymbol{\delta}$ for all $t \in \mathbb{N}$, and so

$$\Pr(X_t = x) = \boldsymbol{\delta}\mathbf{P}(x)\mathbf{1} \quad (4)$$

Marginal distributions: Bivariate distributions

- In any directed graphical model, the joint distribution of a set of random variables V_i is given by

$$\Pr(V_1, V_2, \dots, V_n) = \prod_{i=1}^n \Pr(V_i \mid \text{pa}(V_i)), \quad (5)$$

where $\text{pa}(V_i)$ denotes the set of all 'parents' of V_i in the set V_1, V_2, \dots, V_n .

- In the directed graph of the four random variables $X_t, X_{t+k}, C_t, C_{t+k}$ for positive integer k ,
 $\text{pa}(X_t) = \{C_t\}, \text{pa}(C_{t+k}) = \{C_t\}, \text{pa}(X_{t+k}) = \{C_{t+k}\}.$

Marginal distributions: Bivariate distributions

- $\Pr(X_t, X_{t+k}, C_t, C_{t+k})$
$$= \Pr(C_t) \Pr(X_t | C_t) \Pr(C_{t+k} | C_t) \Pr(X_{t+k} | C_{t+k})$$
- $\Pr(X_t = v, X_{t+k} = w)$
$$\begin{aligned} &= \sum_{i=1}^m \sum_{j=1}^m \Pr(X_t = v, X_{t+k} = w, C_t = i, C_{t+k} = j) \\ &= \sum_{i=1}^m \sum_{j=1}^m \underbrace{\Pr(C_t = i) p_i(v)}_{u_i(t)} \underbrace{\Pr(C_{t+k} = j | C_t = i) p_j(w)}_{\gamma_{ij}(k)} \\ &= \sum_{i=1}^m \sum_{j=1}^m u_i(t) p_i(v) \gamma_{ij}(k) p_j(w) \end{aligned}$$

where $\gamma_{ij}(k)$ denotes the (i, j) element of $\mathbf{\Gamma}^k$.

Marginal distributions: Bivariate distributions

- Writing the above double sum as a product of matrices

$$\Pr(X_t = v, X_{t+k} = w) = \mathbf{u}(t) \mathbf{P}(v) \mathbf{\Gamma}^k \mathbf{P}(w) \mathbf{1} \quad (6)$$

- If the Markov chain is stationary,

$$\Pr(X_t = v, X_{t+k} = w) = \boldsymbol{\delta} \mathbf{P}(v) \mathbf{\Gamma}^k \mathbf{P}(w) \mathbf{1} \quad (7)$$

- We note that

$$\begin{aligned} \mathbb{E}(X_t) &= \sum_{i=1}^m \mathbb{E}(X_t | C_t = i) \Pr(C_t = i) \\ &= \sum_{i=1}^m u_i(t) \mathbb{E}(X_t | C_t = i) \end{aligned}$$

- In the stationary case, $\mathbb{E}(X_t) = \sum_{i=1}^m \delta_i \mathbb{E}(X_t | C_t = i)$

Moments

- More generally, for any functions g in the stationary case

$$\mathbb{E}(g(X_t)) = \sum_{i=1}^m \delta_i \mathbb{E}(g(X_t) \mid C_t = i) \quad (8)$$

$$\mathbb{E}(g(X_t, X_{t+k})) = \sum_{i,j=1}^m \mathbb{E}(g(X_t, X_{t+k}) \mid C_t = i, C_{t+k} = j) \delta_i \gamma_{ij}(k) \quad (9)$$

- If a function g factorizes as

$$g(X_t, X_{t+k}) = g_1(X_t)g_2(X_{t+k}),$$

$$\mathbb{E}(g(X_t, X_{t+k})) = \sum_{i,j=1}^m \mathbb{E}(g_1(X_t) \mid C_t = i) \mathbb{E}(g_2(X_{t+k}) \mid C_{t+k} = j) \delta_i \gamma_{ij}(k) \quad (10)$$

Table of contents

1. The basics: Hidden Markov model
2. The likelihood

The likelihood

- The aim of this section is to develop a convenient formula for the likelihood L_T of T consecutive observations x_1, x_2, \dots, x_T assumed to be generated by an m -state HMM.
- The computation of the likelihood appears to require $O(Tm^T)$ operations.

$$\begin{aligned} \Pr(X_1 = x_1, \dots, X_T = x_T) \\ = \sum_{c_1, \dots, c_T=1}^m \Pr(X_1 = x_1, \dots, X_T = x_T, C_1 = c_1, \dots, C_T = c_T) \end{aligned}$$

$$\begin{aligned} \Pr(X_1, \dots, X_T, C_1, \dots, C_T) \\ = \underbrace{\Pr(X_1|C_1) \cdots \Pr(X_T|C_T)}_{\text{product of } T \text{ factors}} \underbrace{\Pr(C_1) \Pr(C_2|C_1) \cdots \Pr(C_T|C_{T-1})}_{\text{product of } T \text{ factors}} \end{aligned}$$

The likelihood

- It is our purpose here to demonstrate that L_T can in general be computed relatively simply in $O(Tm^2)$ operations.
- First the likelihood of a two-state model will be explored.
- Then the general formula will be presented.

The likelihood of a two-state Bernoulli-HMM

- **Example (Bernoulli-HMM)**

Consider the stationary two-state HMM with t.p.m.

$$\mathbf{\Gamma} = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{4} & \frac{3}{4} \end{pmatrix}$$

and state-dependent distributions given by

$$\Pr(X_t = x \mid C_t = 1) = \frac{1}{2} \quad (\text{for } x = 0, 1),$$

$$\Pr(X_t = 1 \mid C_t = 2) = 1$$

- The stationary distribution of the Markov chain is $\boldsymbol{\delta} = \frac{1}{3}(1, 2)$

The likelihood of a two-state Bernoulli-HMM

- Note that

$$\begin{aligned}\Pr(X_1, X_2, X_3, C_1, C_2, C_3) \\ = \Pr(C_1) \Pr(X_1|C_1) \Pr(C_2|C_1) \Pr(X_2|C_2) \Pr(C_3|C_2) \Pr(X_3|C_3)\end{aligned}$$

$$\begin{aligned}\Pr(X_1 = 1, X_2 = 1, X_3 = 1) \\ = \sum_{i,j,k=1}^2 \Pr(X_1 = 1, X_2 = 1, X_3 = 1, C_1 = i, C_2 = j, C_3 = k) \\ = \sum_{i,j,k=1}^2 \delta_i p_i(1) \gamma_{ij} p_j(1) \gamma_{jk} p_k(1)\end{aligned}\tag{11}$$

- Notice that the triple sum (11) has $m^T = 2^3$ terms, each of which is a product of $2T = 2 \times 3$ factors.

The likelihood of a two-state Bernoulli–HMM

Table 1: Example of a likelihood computation.

i	j	k	$p_i(1)$	$p_j(1)$	$p_k(1)$	δ_i	γ_{ij}	γ_{jk}	Product
1	1	1	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{2}{4}$	$\frac{2}{4}$	$\frac{1}{96}$
1	1	2	$\frac{1}{2}$	$\frac{1}{2}$	1	$\frac{1}{3}$	$\frac{2}{4}$	$\frac{2}{4}$	$\frac{1}{48}$
1	2	1	$\frac{1}{2}$	1	$\frac{1}{2}$	$\frac{1}{3}$	$\frac{2}{4}$	$\frac{1}{4}$	$\frac{1}{96}$
1	2	2	$\frac{1}{2}$	1	1	$\frac{1}{3}$	$\frac{2}{4}$	$\frac{3}{4}$	$\frac{1}{16}$
2	1	1	1	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{48}$
2	1	2	1	$\frac{1}{2}$	1	$\frac{2}{3}$	$\frac{1}{4}$	$\frac{2}{4}$	$\frac{1}{24}$
2	2	1	1	1	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{3}{4}$	$\frac{1}{4}$	$\frac{1}{16}$
2	2	2	1	1	1	$\frac{2}{3}$	$\frac{3}{4}$	$\frac{3}{4}$	$\frac{3}{8}$
									$\frac{29}{48}$

The likelihood of a two-state Bernoulli-HMM

- The state sequence that maximizes the joint probability

$$\Pr(X_1 = 1, X_2 = 1, X_3 = 1, C_1 = i, C_2 = j, C_3 = k)$$

is therefore the sequence $i = 2, j = 2, k = 2$.

- More convenient way to present the sum is to use matrix notation. Let $\mathbf{P}(u)$ be defined (as before) as $\text{diag}(p_1(u), p_2(u))$. Then

$$\mathbf{P}(0) = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 0 \end{pmatrix} \quad \text{and} \quad \mathbf{P}(1) = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{pmatrix},$$

and the triple sum (11) can be written as a matrix product:

$$\Pr(X_1 = 1, X_2 = 1, X_3 = 1) = \delta \mathbf{P}(1) \mathbf{\Gamma} \mathbf{P}(1) \mathbf{\Gamma} \mathbf{P}(1) \mathbf{1}'$$

The likelihood in general

- Consider the likelihood of an HMM in general.
- Suppose there is an observation sequence x_1, x_2, \dots, x_T generated by such a model.
- Seek the probability L_T of observing that sequence under an m -state HMM that has initial distribution δ and t.p.m. Γ for the Markov chain, and state-dependent probability (density) functions p_i .

The likelihood in general

Proposition 1

The likelihood is given by

$$L_T = \delta \mathbf{P}(x_1) \mathbf{\Gamma P}(x_2) \mathbf{\Gamma P}(x_3) \cdots \mathbf{\Gamma P}(x_T) \mathbf{1}. \quad (12)$$

If δ , the distribution of C_1 , is the stationary distribution of the Markov chain, then in addition

$$L_T = \delta \mathbf{\Gamma P}(x_1) \mathbf{\Gamma P}(x_2) \mathbf{\Gamma P}(x_3) \cdots \mathbf{\Gamma P}(x_T) \mathbf{1}. \quad (13)$$

Proof

Only the case of discrete observations. First, note that

$$L_T = \Pr \left(\mathbf{X}^{(T)} = \mathbf{x}^{(T)} \right) = \sum_{c_1, c_2, \dots, c_T=1}^m \Pr \left(\mathbf{X}^{(T)} = \mathbf{x}^{(T)}, \mathbf{C}^{(T)} = \mathbf{c}^{(T)} \right),$$

The likelihood in general

Proof (Cont.)

and that, by equation (5),

$$\Pr(\mathbf{X}^{(T)}, \mathbf{C}^{(T)}) = \Pr(C_1) \prod_{k=2}^T \Pr(C_k | C_{k-1}) \prod_{k=1}^T \Pr(X_k | C_k). \quad (14)$$

It follows that

$$\begin{aligned} L_T &= \sum_{c_1, \dots, c_T=1}^m (\delta_{c_1} \gamma_{c_1, c_2} \gamma_{c_2, c_3} \cdots \gamma_{c_{T-1}, c_T}) (p_{c_1}(x_1) p_{c_2}(x_2) \cdots p_{c_T}(x_T)) \\ &= \sum_{c_1, \dots, c_T=1}^m \delta_{c_1} p_{c_1}(x_1) \gamma_{c_1, c_2} p_{c_2}(x_2) \gamma_{c_2, c_3} \cdots \gamma_{c_{T-1}, c_T} p_{c_T}(x_T) \\ &= \boldsymbol{\delta} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \boldsymbol{\Gamma} \mathbf{P}(x_3) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1}, \end{aligned}$$

If $\boldsymbol{\delta}$ is the stationary distribution, $\boldsymbol{\delta} \mathbf{P}(x_1) = \boldsymbol{\delta} \boldsymbol{\Gamma} \mathbf{P}(x_1)$



The likelihood in general

- A consequence of the matrix expression for the likelihood is the ‘forward algorithm’ for recursive computation of the likelihood.
- The recursive nature of likelihood evaluation via either (12) is computationally much more efficient than brute-force summation over all possible state sequences.
- To state the forward algorithm, define the vector α_t , for $t = 1, 2, \dots, T$, by

$$\alpha_t = \delta \mathbf{P}(x_1) \mathbf{\Gamma P}(x_2) \mathbf{\Gamma P}(x_3) \cdots \mathbf{\Gamma P}(x_t) = \delta \mathbf{P}(x_1) \prod_{s=2}^t \mathbf{\Gamma P}(x_s), \quad (15)$$

The likelihood in general

- Then, in the likelihood formula (12) :

$$\begin{aligned}\alpha_1 &= \delta \mathbf{P}(x_1); \\ \alpha_t &= \alpha_{t-1} \mathbf{\Gamma P}(x_t) \quad \text{for } t = 2, 3, \dots, T; \\ L_T &= \alpha_T \mathbf{1}.\end{aligned}$$

- That the number of operations involved is of order Tm^2 can be deduced thus.
- The elements of the vector α_t are usually referred to as **forward probabilities**.
- The multiple-sum expression for the likelihood \implies the matrix expression \implies the forward recursion.

The likelihood when data are missing

- If some of the data are missing, the likelihood computation turns out to be a simple one.

- **Example**

Suppose that one has available the observations $x_1, x_2, x_4, x_7, x_8, \dots, x_T$ of an HMM, but x_3, x_5 and x_6 are missing. Then the likelihood of the observations is given by

$$\begin{aligned} & \Pr(X_1 = x_1, X_2 = x_2, X_4 = x_4, X_7 = x_7, \dots, X_T = x_T) \\ &= \sum \delta_{c_1} \gamma_{c_1, c_2} \gamma_{c_2, c_4} (2) \gamma_{c_4, c_7} (3) \gamma_{c_7, c_8} \cdots \gamma_{c_{T-1}, c_T} \\ & \quad \times p_{c_1}(x_1) p_{c_2}(x_2) p_{c_4}(x_4) p_{c_7}(x_7) \cdots p_{c_T}(x_T), \end{aligned}$$

where the sum is taken over all indices c_t other than c_3, c_5 and c_6 .

The likelihood when data are missing

- This is just

$$\begin{aligned} & \sum \delta_{c_1} p_{c_1}(x_1) \gamma_{c_1, c_2} p_{c_2}(x_2) \gamma_{c_2, c_4} p_{c_4}(x_4) \gamma_{c_4, c_7} p_{c_7}(x_7) \\ & \quad \times \cdots \times \gamma_{c_{T-1}, c_T} p_{c_T}(x_T) \\ & = \boldsymbol{\delta} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \boldsymbol{\Gamma}^2 \mathbf{P}(x_4) \boldsymbol{\Gamma}^3 \mathbf{P}(x_7) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1} \end{aligned}$$

- In general, in the expression for the likelihood the diagonal matrices $\mathbf{P}(x_t)$ corresponding to missing observations x_t are replaced by the identity matrix.
- Thus, even if some observations are missing, the likelihood of an HMM can be computed easily.