

Ch3. Estimation by direct maximization of the likelihood

HyeRim Park

July 28, 2023

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

Introduction

- The likelihood of an HMM is given by

$$L_T = \Pr \left(\mathbf{X}^{(T)} = \mathbf{x}^{(T)} \right) = \boldsymbol{\delta} \mathbf{P}(x_1) \boldsymbol{\Gamma} \mathbf{P}(x_2) \boldsymbol{\Gamma} \mathbf{P}(x_3) \cdots \boldsymbol{\Gamma} \mathbf{P}(x_T) \mathbf{1},$$

where $\boldsymbol{\delta}$ is the initial distribution and $\mathbf{P}(x)$ the $m \times m$ diagonal matrix with i th diagonal element the state-dependent probability $p_i(x)$

- Then, in the likelihood formula:

$$\begin{aligned} \boldsymbol{\alpha}_1 &= \boldsymbol{\delta} \mathbf{P}(x_1); \\ \boldsymbol{\alpha}_t &= \boldsymbol{\alpha}_{t-1} \boldsymbol{\Gamma} \mathbf{P}(x_t) \quad \text{for } t = 2, 3, \dots, T; \\ L_T &= \boldsymbol{\alpha}_T \mathbf{1}. \end{aligned}$$

- Parameter estimation can be performed by numerical maximization of the likelihood with respect to the parameters.
- But there are several problems that need to be addressed when parameters are estimated.
 - 1 Numerical underflow
 - 2 Constraints on the parameters
 - 3 Multiple local maxima in the likelihood function

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

- In the case of discrete state-dependent distributions, the elements of α_t become smaller as t increases and are eventually rounded to zero.
- Then, likelihood approaches 0 exponentially fast with probability 1 (or possibly ∞ in the continuous case)
- We confine our attention to underflow.

- To solve the problem, we compute the logarithm of L_T by using a strategy of scaling the vector of forward probabilities α_t .
- Effectively we scale the vector α_t at each time t so that its elements add to 1.
- Define the vector

$$\phi_t = \alpha_t / w_t$$

where $w_t = \sum_i \alpha_t(i) = \alpha_t \mathbf{1}$

Scaling the likelihood computation

- Using the definitions of ϕ_t and w_t :

$$w_0 = \alpha_0 \mathbf{1} = \delta \mathbf{1} = 1$$

$$\phi_0 = \delta$$

$$\alpha_t = \alpha_{t-1} \Gamma \mathbf{P}(x_t)$$

$$w_t \phi_t = w_{t-1} \phi_{t-1} \Gamma \mathbf{P}(x_t) \tag{1}$$

$$L_T = \alpha_T \mathbf{1} = w_T (\phi_T \mathbf{1}) = w_T.$$

Scaling the likelihood computation

- Hence $L_T = w_T = \prod_{t=1}^T (w_t/w_{t-1})$. From (1) it follows that

$$w_t = \boldsymbol{\alpha}_t \mathbf{1} = w_{t-1} \phi_{t-1} \mathbf{\Gamma P}(x_t) \mathbf{1}$$

$$w_t = w_{t-1} (\phi_{t-1} \mathbf{\Gamma P}(x_t) \mathbf{1}),$$

and so we conclude that

$$\log L_T = \sum_{t=1}^T \log (w_t/w_{t-1}) = \sum_{t=1}^T \log (\phi_{t-1} \mathbf{\Gamma P}(x_t) \mathbf{1}).$$

Scaling the likelihood computation

- The computation of the log-likelihood is summarized below in the form of an algorithm.

δ : initial distribution $\mathbf{\Gamma}$ and $\mathbf{P}(x_t) : m \times m$ matrix

\mathbf{v} and ϕ_t : vectors of length m

l is the scalar in which the log-likelihood is accumulated.

$$w_1 \leftarrow \delta \mathbf{P}(x_1) \mathbf{1}; \phi_1 \leftarrow \delta \mathbf{P}(x_1) / w_1; l \leftarrow \log w_1$$

for $t = 2, 3, \dots, T$

$$\mathbf{v} \leftarrow \phi_{t-1} \mathbf{\Gamma} \mathbf{P}(x_t)$$

$$u \leftarrow \mathbf{v} \mathbf{1}$$

$$l \leftarrow l + \log u$$

$$\phi_t \leftarrow \mathbf{v} / u$$

return l

- The log-likelihood $\log L_T$ is given by the final value of l .

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

Reparametrization to avoid constraints

- In Poisson–HMM, the elements of $\mathbf{\Gamma}$ and those of $\boldsymbol{\lambda}$, the vector of state-dependent means in \mathbf{a} are subject to non-negativity and the row sums of $\mathbf{\Gamma}$ equal 1.
- Estimates of parameters should satisfy such constraints.
- Thus, when maximizing the likelihood we need to solve a constrained optimization problem.

Reparametrization to avoid constraints

- Special-purpose software can be used to maximize a function of several variables which are subject to constraints.
- However, depending on the implementation and the nature of the data, constrained optimization can be slow.
- For example, the constrained optimizer is slow if the optimum lies on the boundary of the parameter space.
- We shall focus on the use of the unconstrained optimizer **nlm**.

Reparametrization to avoid constraints

- Constraints depends on which state-dependent distribution, and Constraints that apply to the parameters of the Markov chain
- **Example: Poisson-HMM**
 - The means λ_i of the state-dependent distributions must be non-negative.
 - The rows of the transition probability matrix $\mathbf{\Gamma}$ must add to 1, and all the parameters γ_{ij} must be non-negative.
- The constraints can be imposed by making transformations.

Reparametrization to avoid constraints

- Define $\eta_i = \log \lambda_i$, for $i = 1, \dots, m$. Then $\eta_i \in \mathbb{R}$
- After maximizing the likelihood with the unconstrained parameters, the constrained parameter estimates can be obtained by transforming back: $\hat{\lambda}_i = \exp \hat{\eta}_i$.
- Note that $\mathbf{\Gamma}$ has m^2 entries but only $m(m-1)$ free parameters, as there are m row-sum constraints

$$\sum_{j=1}^m \gamma_{ij} = 1 \quad (i = 1, \dots, m)$$

Reparametrization to avoid constraints

- Show transformation between the m^2 constrained probabilities γ_{ij} and $m(m-1)$ unconstrained real numbers τ_{ij} , $i \neq j$
- If $m = 3$, define the matrix

$$\mathbf{T} = \begin{pmatrix} - & \tau_{12} & \tau_{13} \\ \tau_{21} & - & \tau_{23} \\ \tau_{31} & \tau_{32} & - \end{pmatrix}$$

, a matrix with $m(m-1)$ entries $\tau_{ij} \in \mathbb{R}$. Now let $g : \mathbb{R} \rightarrow \mathbb{R}^+$ be a strictly increasing function, for example, $g(x) = e^x$

Reparametrization to avoid constraints

- Define

$$\nu_{ij} = \begin{cases} g(\tau_{ij}) & \text{for } i \neq j \\ 1 & \text{for } i = j. \end{cases}$$

- Set $\gamma_{ij} = \nu_{ij} / \sum_{k=1}^m \nu_{ik}$ (for $i, j = 1, 2, \dots, m$) and $\mathbf{\Gamma} = (\gamma_{ij})$ and the transformation in the opposite direction is

$$\tau_{ij} = \log(\gamma_{ij}/\gamma_{ii}), \quad \text{for } i \neq j$$

- We shall refer to the parameters η_i and τ_{ij} as **working parameters**, and to the parameters λ_i and γ_{ij} as **natural parameters**.

Reparametrization to avoid constraints

- Using the transformations of \mathbf{T} and $\boldsymbol{\lambda}$, the calculation of the likelihood-maximizing parameters in two steps.
 - 1 Maximize L_T with the working parameters $\mathbf{T} = \{\tau_{ij}\}$ and $\boldsymbol{\eta} = (\eta_1, \dots, \eta_m)$.
 - 2 Transform the estimates of the working parameters to estimates of the natural parameters:

$$\hat{\mathbf{T}} \rightarrow \hat{\mathbf{\Gamma}}, \quad \hat{\boldsymbol{\eta}} \rightarrow \hat{\boldsymbol{\lambda}}$$

Reparametrization to avoid constraints

- Considering initial distribution δ is not supplied.
- If δ is stationary distribution, compute by solving

$$\delta(\mathbf{I}_m - \mathbf{\Gamma} + \mathbf{J}_m) = \mathbf{1}'$$

- Otherwise, transform δ with the constraints $\delta_i \geq 0$ and $\sum_i \delta_i = 1$ as unconstrained parameters.

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

Multiple maxima in the likelihood

- The likelihood of an HMM frequently has several local maxima.
- The goal is to find the global maximum, but there is no simple method of determining in general whether a numerical maximization algorithm has reached the global maximum.
- Depending on the starting values, it can easily happen that the algorithm identifies a local, but not the global, maximum.
- A sensible strategy is therefore to use a range of starting values for the maximization, and to see whether the same maximum is identified in each case.

Starting values for the iterations

- It is often easy to find plausible starting values for some of the parameters of an HMM.
- For instance, in a Poisson-HMM with two states to fit the sample mean for the values of the two state-dependent means.
- To assign a common starting value to all the off-diagonal transition probabilities γ_{ij} .

Table of contents

1. Introduction
2. Scaling the likelihood computation
3. Maximization of the likelihood subject to constraint
4. Other problems
5. Example: earthquakes

Example: earthquakes

- Fitting stationary Poisson–hidden Markov models to the earthquakes series by means of the unconstrained optimizer **nlm**.
 - ① nature parameter \rightarrow working parameter
 - ② iteration of **nlm**(mllk,wp) algorithm;
mllk(working param \rightarrow nature param, log-likelihood)
 - ③ estimate of working param \rightarrow estimate of nature param

Example: earthquakes

- The three-state model is

$$\mathbf{\Gamma} = \begin{pmatrix} 0.955 & 0.024 & 0.021 \\ 0.050 & 0.899 & 0.051 \\ 0.000 & 0.197 & 0.803 \end{pmatrix}$$

with $\boldsymbol{\delta} = (0.4436, 0.4045, 0.1519)$, $\boldsymbol{\lambda} = (13.146, 19.721, 29.714)$, and log-likelihood given by $l = -329.4603$.

- The four-state is

$$\mathbf{\Gamma} = \begin{pmatrix} 0.805 & 0.102 & 0.093 & 0.000 \\ 0.000 & 0.976 & 0.000 & 0.024 \\ 0.050 & 0.000 & 0.902 & 0.048 \\ 0.000 & 0.000 & 0.188 & 0.812 \end{pmatrix}$$

with $\boldsymbol{\delta} = (0.0936, 0.3983, 0.3643, 0.1439)$, $\boldsymbol{\lambda} = (11.283, 13.853, 19.695, 29.700)$, and log-likelihood given by $l = -327.8316$.

Example: earthquakes

- The means and variances of the marginal distributions of the four models compare as follows with those of the observations.

	mean	variance
observations:	19.364	51.573
'one-state HMM':	19.364	19.364
two-state HMM:	19.086	44.523
three-state HMM:	18.322	50.709
four-state HMM:	18.021	49.837

Example: earthquakes

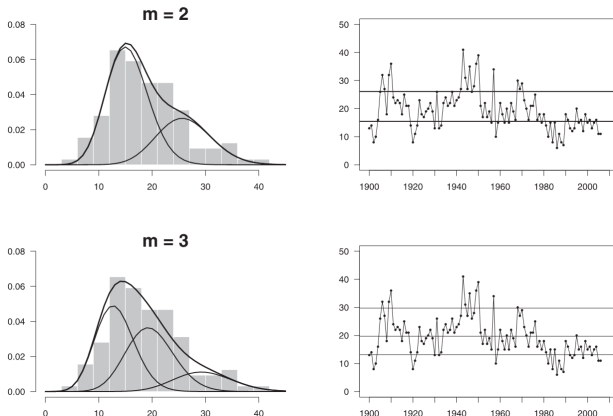


Figure 1: Earthquakes series. Left: marginal distributions of Poisson–HMMs with two and three states, and their components, compared with a histogram of the observations. Right: the state-dependent means (horizontal lines) compared to the observations.

Example: earthquakes

- The three-state model is

$$\Gamma = \begin{pmatrix} 0.955 & 0.024 & 0.021 \\ 0.050 & 0.899 & 0.051 \\ \mathbf{0.000} & 0.197 & 0.803 \end{pmatrix}$$

- The four-state is

$$\Gamma = \begin{pmatrix} 0.805 & 0.102 & 0.093 & \mathbf{0.000} \\ \mathbf{0.000} & 0.976 & \mathbf{0.000} & 0.024 \\ 0.050 & \mathbf{0.000} & 0.902 & 0.048 \\ \mathbf{0.000} & \mathbf{0.000} & 0.188 & 0.812 \end{pmatrix}$$

- When one fits models with three or more states is that the estimates of one or more of the transition probabilities turn out to be very close to zero.

Example: earthquakes

- In a stationary Markov chain, the expected number of transitions from state i to state j in a series of T observations is $(T - 1)\delta_i\gamma_{ij}$.
- For $\delta_3 = 0.152$ and $T = 107$ in our three-state model, this expectation will be less than 1 if $\gamma_{31} < 0.062$.
- Therefore, it is likely that if γ_{31} is fairly small there will be no transitions from state 3 to state 1, and so when seeking to estimate γ_{31} in an HMM the estimate is likely to be effectively zero.
- As m increases, the probabilities δ_i and γ_{ij} get smaller on average.