



Expectation-Maximization Algorithm

MODULE DES130: COMPUTATIONAL STATISTICS

Dr. Erick A. Chacón Montalván (echacon@uni.edu.pe)

Escuela de Profesional de Ingeniería Estadística
Facultad de Ingeniería Económica, Estadística y Ciencias Sociales
Universidad Nacional de Ingeniería (UNI)
Lima – Perú

Likelihood function

Likelihood function

Motivation

Concepts

Maximum likelihood estimation

Expectation-Maximization algorithm

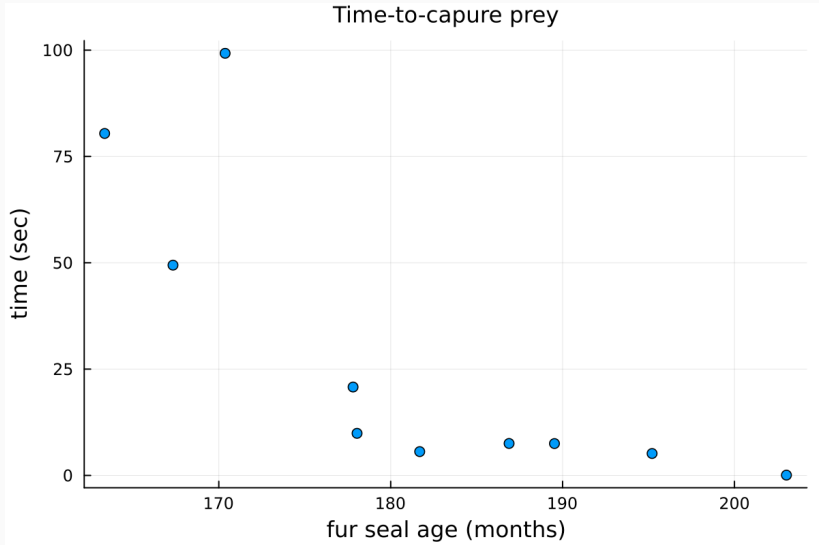
Theory

Modelling time to capture with age

A biologist is studying the time that fur seals need to capture their first prey. The observed times (seconds) for 10 fur seals and their age (weeks) are below.

fur seal	age	time
1	177.811	20.7883
2	195.207	5.15748
3	167.333	49.4457
4	178.041	9.90466
5	203.024	0.0770514
6	163.358	80.4105
7	186.886	7.51927
8	170.368	99.2618
9	181.688	5.60274
10	189.528	7.49648

Modelling time to capture with age



Model-based inference

Given a collection of observations y_i for $i = 1, \dots, n$. In a model-based inference framework

- We assume a known **stochastic model** for the data. For example,

$$Y_i \sim \text{Exponential}(\theta).$$

- The stochastic model have **unknown parameters**.

θ : Expected time-to-capture.

- Likelihood inference is one way to obtain **estimates** about the parameters from the data. For example,

$$\hat{\theta} = \frac{\sum_{i=1}^n y_i}{n}.$$

Some stochastic models

Poisson model

$$\Pr(Y_i = y; \lambda) = \frac{\exp(-\lambda)\lambda^y}{y!}$$

Exponential model

$$f_{Y_i}(y; \theta) = \frac{1}{\theta} \exp\left(-\frac{x}{\theta}\right)$$

Bernoulli model

$$\Pr(Y_i = y; \pi) = y^\pi (1 - y)^{1-\pi}$$

Likelihood function

Motivation

Concepts

Maximum likelihood estimation

Expectation-Maximization algorithm

Theory

Likelihood and log-likelihood functions

Likelihood function

Denoted as $L(\theta)$, is the chance (probability) of observing the data given a specific value of the parameter θ .

Notice that:

- $L(\theta) : \Theta \rightarrow \mathbb{R}^+$ is a function with respect to θ .
- $L(\theta)$ depends on the assumed stochastic model.

Log-likelihood function

Denoted as $l(\theta)$ is the logarithm of the likelihood function.

$$l(\theta) = \ln L(\theta).$$

Likelihood function

Discrete random variables

Given a set of observations y_1, \dots, y_n , then the likelihood function is

$$L(\theta) = \Pr(Y_1 = y_1, \dots, Y_n = y_n; \theta).$$

Continuous random variables

Given a set of observations y_1, \dots, y_n , then the likelihood function is

$$\begin{aligned} L(\theta) &= \Pr(y_1 - \epsilon < Y_1 < y_1 + \epsilon, \dots, y_n - \epsilon < Y_n < y_n + \epsilon; \theta) \\ &\approx (2 * \epsilon)^n f_{Y_1, \dots, Y_n}(y_1, \dots, y_n; \theta) \\ &\propto f_{Y_1, \dots, Y_n}(y_1, \dots, y_n; \theta). \end{aligned}$$

Likelihood function a vector parameter

In real-world applications, we usually work with models than have more than one parameter. We consider it a vector parameter

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} .$$

In those case, we will denote the likelihood function as

$$L(\boldsymbol{\theta})$$

and the log-likelihood function as

$$l(\boldsymbol{\theta}).$$

Time-to-capture model

Let consider that the time-to-capture observations y_i .

Assumed stochastic model

$$Y_i \stackrel{iid}{\sim} \text{Exponential}(\theta), \quad \text{for } \theta > 0 \text{ and } i = 1, \dots, 10.$$

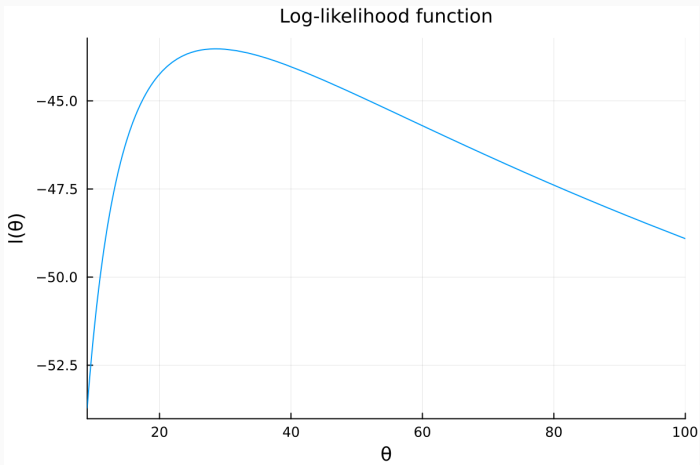
Likelihood function

$$\begin{aligned} L(\theta) &= f_{Y_1, \dots, Y_{10}}(y_1, \dots, y_n) \\ &= \prod_{i=1}^{10} f_{Y_i}(y_i) && \text{(due to independence)} \\ &= \prod_{i=1}^{10} \frac{1}{\theta} \exp\left(-\frac{y_i}{\theta}\right) \\ &= \theta^{-10} \exp\left(-\frac{\sum_{i=1}^n y_i}{\theta}\right). \end{aligned}$$

Time-to-capture model

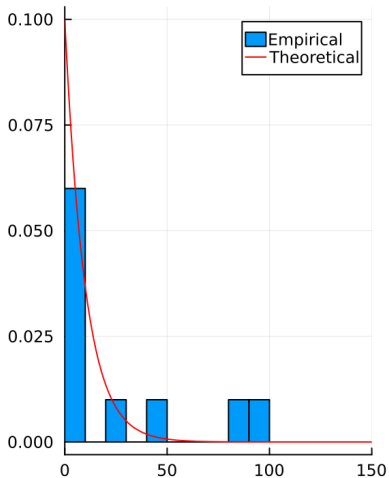
Log-likelihood function

$$l(\theta) = -10 \log(\theta) - \frac{\sum_{i=1}^n y_i}{\theta}$$

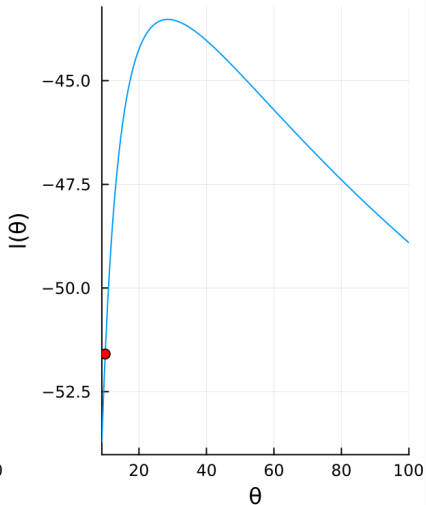


Time-to-capture model: Interpret $\theta = 10$

Comparison

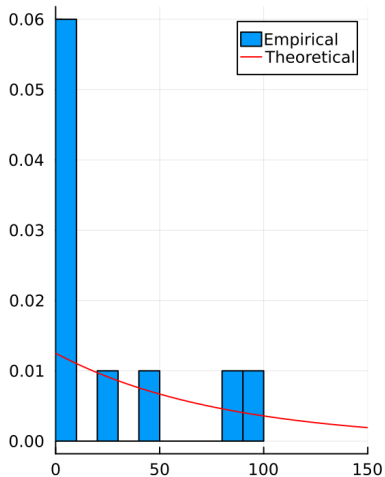


Log-likelihood function

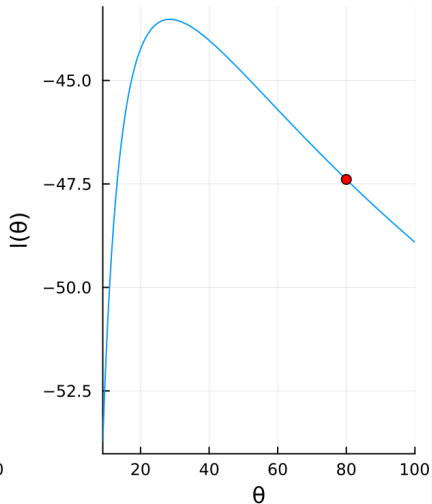


Time-to-capture model: Interpret $\theta = 80$

Comparison

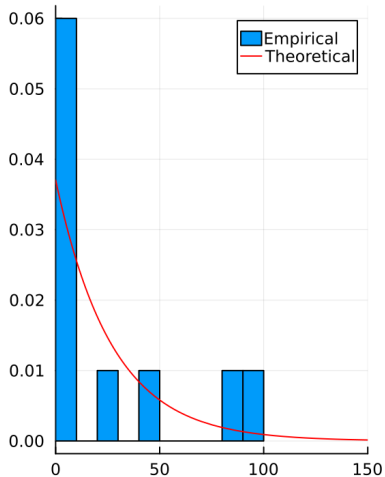


Log-likelihood function

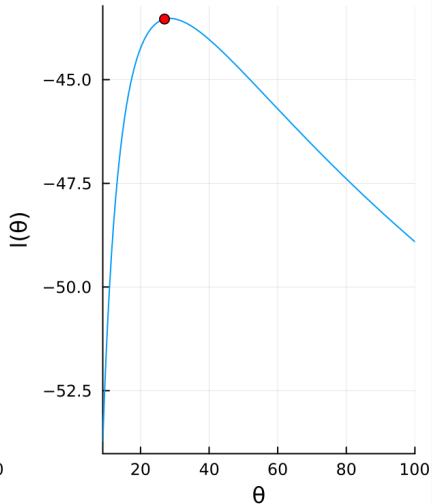


Time-to-capture model: Interpret $\theta = 27$

Comparison



Log-likelihood function



Time-to-capture model

Let consider that the time-to-capture observations y_i .

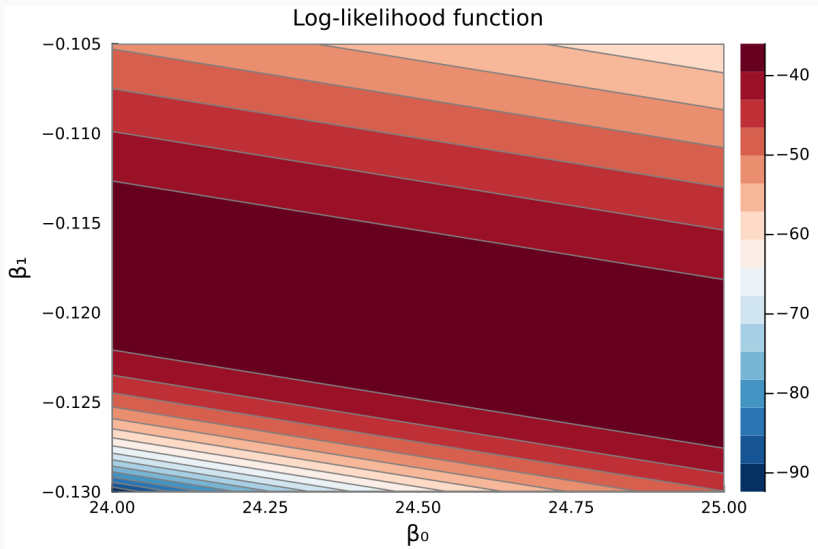
Assumed stochastic model

$$Y_i \stackrel{iid}{\sim} \text{Exponential}(\mu_i), \quad \text{for } \mu_i > 0 \text{ and } i = 1, \dots, 10.$$
$$\log(\mu_i) = \beta_0 + \beta_1 \text{age}_i$$

Likelihood function

$$L(\beta_0, \beta_1) = f_{Y_1, \dots, Y_{10}}(y_1, \dots, y_n)$$
$$= \prod_{i=1}^{10} f_{Y_i}(y_i) \quad (\text{due to independence})$$
$$= \prod_{i=1}^{10} \frac{1}{\exp(\beta_0 + \beta_1 \text{age}_i)} \exp\left(-\frac{y_i}{\exp(\beta_0 + \beta_1 \text{age}_i)}\right).$$

Time-to-capture model with predictor



Score function and information

Score function

$$S_i(\theta) = \frac{\partial}{\partial \theta_i} \log L(\theta)$$

Observed information

$$I_{Oij}(\theta) = -\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log L(\theta)$$

Fisher information

$$I_{Eij}(\theta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log L(\theta) \right]$$

Maximum likelihood estimation

Likelihood function

Maximum likelihood estimation

Inference

Expectation-Maximization algorithm

Theory

Maximum likelihood estimate

Definition

Given a set of observations y_1, \dots, y_n , then the maximum likelihood estimate is obtained by getting the value of θ that maximizes the likelihood function,

$$\hat{\theta}_{ml} = \arg \max_{\theta} L(\theta).$$

It is equivalent to the value that maximizes the log-likelihood function,

$$\hat{\theta}_{ml} = \arg \max_{\theta} l(\theta).$$

The idea is to obtain the value of θ under which it is more likely to have observed the data y_1, \dots, y_n .

Maximum likelihood estimation

Score function

It is the gradient vector of the log-likelihood function with respect to the parameter,

$$U(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_1} l(\theta) \\ \frac{\partial}{\partial \theta_2} l(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} l(\theta). \end{bmatrix}$$

Due to the definition of the MLE, then

$$U(\hat{\theta}_{ml}) = \mathbf{0}.$$

Solving the previous equation system allow us to obtain the maximum and minimum local points of $l(\theta)$.

Maximum likelihood estimator

Maximum likelihood estimator (MLE)

It is a random variable used to make inference about a parameter. It is asymptotically normally distributed under regular conditions:

$$\hat{\theta}_{ml} \sim MVN(\theta_0, I_E(\hat{\theta})^{-1});$$

where θ_0 is the true parameter value.

In case $I_E(\theta_0)$ is not available, it is usually replaced by another asymptotically equivalent term. For example, $I_E(\hat{\theta})$, $I_O(\theta_0)$, $I_O(\hat{\theta})$. This result is used to obtain confidence intervals and perform hypothesis testing.

- Generalized linear models
- Generalized additive models
- Time series models
- Survival analysis
- Point processes
- Spatial models

Expectation-Maximization algorithm

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Introduction

Algorithm

Example

Theory

Species modelling

In an study of 4 animal species, it is know that the probability to belong to each category es $\frac{1}{2} + \frac{\theta}{4}$, $\frac{1-\theta}{4}$, $\frac{1-\theta}{4}$ and $\theta/4$. Estimate the value of θ given that, from a total of 197, we observed 125, 18, 20, 34 animales respectively.

Likelihood function

$$L(\theta; \mathbf{x}) = \frac{197!}{125!18!20!34!} \left(\frac{1}{2} + \frac{\theta}{4}\right)^{125} \left(\frac{1-\theta}{4}\right)^{18+20} \left(\frac{\theta}{4}\right)^{34}$$
$$\propto (2 + \theta)^{125} (1 - \theta)^{38} \theta^{34},$$

$$L(\theta; \mathbf{x}, y) = \frac{197!}{y!(125 - y)!18!20!34!} \left(\frac{\theta}{4}\right)^y \left(\frac{1}{2}\right)^{125-y} \left(\frac{1-\theta}{4}\right)^{38} \left(\frac{\theta}{4}\right)^{34}$$
$$\propto (1 - \theta)^{38} \theta^{y+34}.$$

Definition

The EM algorithm is an **iterative** algorithm to obtain **maximum likelihood estimates** under the presence of **missing data**.

Notes

- It is iterative.
- It converges to a local maximum, being a maximum likelihood estimate when the likelihood function is unimodal.
- It is used in incomplete data problems; for instance, mixture models.

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Introduction

Algorithm

Example

Theory

The algorithm

The name EM is used because it performs two steps at each iteration.

- E-step: Expectation
- M-step: Maximization

Instead of using $l(\theta; \mathbf{x})$ such as

$$\hat{\theta}_{ml} = \arg \max l(\theta; \mathbf{x}).$$

We can think of an additional data that has not been observed to obtain

$$l(\theta; \mathbf{x}, \mathbf{y})$$

to perform the estimation. These augmented log-likelihood should be easier to work with but unfortunately \mathbf{y} is not available, so we will take the expected value.

EM algorithm

Let θ^* be the current estimate of θ , and consider

$$Q(\theta, \theta^*) = \int l(\theta; x, y) f(y | x; \theta^*) dy, \quad (1)$$

which is the expected value of $l(\theta; x, y)$ with respect to $f(y | x; \theta^*)$.

EM algorithm

1. Set-up an initial estimate θ^*
2. Estimate $l(\theta; x, y)$ by taking the expectation with respect Y and conditioning in $X = x, \theta = \theta^*$. $Q(\theta, \theta^*)$.
3. Maximize the estimated log-likelihood $Q(\theta, \theta^*)$ with respect to θ to update it.
4. Repeats steps until converge.

Then the steps at each iteration are

- **E-step:** Compute $Q(\theta, \theta^*)$ as a function of θ only.
- **M-step:** Maximize $Q(\theta, \theta^*)$ with respect to θ to get θ^{**} .

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Introduction

Algorithm

Example

Theory

Species modelling

In an study of 4 animal species, it is know that the probability to belong to each category es $\frac{1}{2} + \frac{\theta}{4}$, $\frac{1-\theta}{4}$, $\frac{1-\theta}{4}$ and $\theta/4$. Estimate the value of θ given that, from a total of 197, we observed 125, 18, 20, 34 animales respectively.

$$L(\theta; x) = \frac{197!}{125!18!20!34!} \left(\frac{1}{2} + \frac{\theta}{4}\right)^{125} \left(\frac{1-\theta}{4}\right)^{18+20} \left(\frac{\theta}{4}\right)^{34} \\ \propto (2 + \theta)^{125} (1 - \theta)^{38} \theta^{34}.$$

$$L(\theta; x, y) = \frac{197!}{y!(125 - y)!18!20!34!} \left(\frac{\theta}{4}\right)^y \left(\frac{1}{2}\right)^{125-y} \left(\frac{1-\theta}{4}\right)^{38} \left(\frac{\theta}{4}\right)^{34} \\ \propto (1 - \theta)^{38} \theta^{y+34}.$$

Example

$$l(\theta; \mathbf{x}, y) = c + 83 \log(1 - \theta) + (y + 34) \log(\theta).$$

Expectation-Step:

$$\mathbb{E} [l(\theta; \mathbf{x}, y) \mid \mathbf{x}, \theta^*] = c + 83 \log(1 - \theta) + (\mathbb{E} [y \mid \mathbf{x}, \theta^*] + 34) \log(\theta)$$

$$\Pr(z \in A \mid z \in G1, \theta^*) = \frac{\theta^*/4}{1/2 + \theta^*/4} = \frac{\theta^*}{2 + \theta^*}$$

$$y \mid \mathbf{x}, \theta^* = \text{Binomial} \left(125, \frac{\theta^*}{2 + \theta^*} \right)$$

$$\mathbb{E} [l(\theta; \mathbf{x}, y) \mid \mathbf{x}, \theta^*] = c + 83 \log(1 - \theta) + \left(\frac{125\theta^*}{2 + \theta^*} + 34 \right) \log(\theta)$$

$$Q(\theta, \theta^*) = c + 83 \log(1 - \theta) + \left(\frac{125\theta^*}{2 + \theta^*} + 34 \right) \log(\theta).$$

Maximization-Step:

$$\begin{aligned}\frac{dQ(\theta, \theta^*)}{d\theta} &= -\frac{83}{1-\theta} + \left(\frac{125\theta^*}{2+\theta^*} + 34\right) \frac{1}{\theta} \\ 0 &= -\frac{83}{1-\theta^{**}} + \left(\frac{125\theta^*}{2+\theta^*} + 34\right) \frac{1}{\theta^{**}} \\ \theta^{**} &= \frac{125\theta^*/(2+\theta^*) + 34}{125\theta^*/(2+\theta^*) + 34 + 38}.\end{aligned}$$

Iterating the previous equation until converge perform the Expectation and Maximization steps implicitly.

Theory

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Theory

Proof

Estandard Error

Examples

Why does EM work?

We can prove that at each iteration the log-likelihood function is increasing.

$$\begin{aligned} Q(\theta^{**}, \theta^*) - Q(\theta^*, \theta^*) &\geq 0, \\ \int \log(f(x, y; \theta^{**}))f(y | x; \theta^*)dy - \int \log(f(x, y; \theta^*))f(y | x; \theta^*)dy &\geq 0, \\ \int \log\left(\frac{f(x, y; \theta^{**})}{f(x, y; \theta^*)}\right) f(y | x; \theta^*)dy &\geq 0. \end{aligned}$$

Considering that $f(x, y; \theta) = f(y | x; \theta)f(x; \theta)$,

$$\int \log\left(\frac{f(y | x; \theta^{**})}{f(y | x; \theta^*)}\right) f(y | x; \theta^*)dy + \int \log\left(\frac{f(x; \theta^{**})}{f(x; \theta^*)}\right) f(y | x; \theta^*)dy \geq 0.$$

Given that $\log(x) \leq x - 1$,

$$\begin{aligned} \int \log\left(\frac{f(x; \theta^{**})}{f(x; \theta^*)}\right) f(y | x; \theta^*)dy &\geq 0 \\ \log(f(x; \theta^{**})) - \log(f(x; \theta^*)) &\geq 0. \end{aligned}$$

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Theory

Proof

Estandard Error

Examples

The covariance matrix of the estimator is approximated by the inverse of the observed information evaluated at $\hat{\theta}$:

$$\left[-\frac{\partial^2 \log f(x; \theta)}{\partial \theta_i \partial \theta_j} \right]^{-1}.$$

$$f(x; \theta) = f(x, y; \theta) / f(y | x; \theta)$$

$$-\log f(x; \theta) = -\log f(x, y; \theta) - \{-\log f(y | x; \theta)\}$$

$$-\frac{d^2}{d\theta^2} \log f(x; \theta) = -\frac{d^2}{d\theta^2} \log f(x, y; \theta) - \left\{ -\frac{d^2}{d\theta^2} \log f(y | x; \theta) \right\}$$

Multiplying by $f(y | x, \phi)$ and integrate with respect to y .

$$\begin{aligned} -\frac{d^2}{d\theta^2} \log f(x; \theta) &= -\int \frac{d^2}{d\theta^2} f(y | x, \phi) \log f(x, y; \theta) dy - \int \left\{ -\frac{d^2}{d\theta^2} \log f(y | x; \theta) \right\} f(y | x, \phi) dy \\ -\frac{d^2}{d\theta^2} \log f(x; \theta) &= -\frac{d^2}{d\theta^2} Q(\theta, \phi) - -\frac{d^2}{d\theta^2} H(\theta, \phi). \end{aligned}$$

Which can be computed considering that

$$\begin{aligned} -\frac{d^2}{d\theta^2} H(\theta, \phi) |_{\hat{\theta}} &= \mathbb{V} \left[\frac{d \log f(y | x, \theta)}{d\theta} \mid x, \hat{\theta} \right] \\ &= \mathbb{V} \left[\frac{d \log f(y, x | \theta)}{d\theta} - \frac{d \log f(x; \theta)}{d\theta} \mid x, \hat{\theta} \right] \\ &= \mathbb{V} \left[\frac{d \log f(y, x | \theta)}{d\theta} \mid x, \hat{\theta} \right]. \end{aligned}$$

Example

For the species modelling example:

$$-\frac{d^2}{d\theta^2}Q(\theta, \phi)|_{\hat{\theta}} = \frac{\mathbb{E}[Y | x_1; \hat{\theta}] + x_4}{\hat{\theta}^2} + \frac{x_2 + x_3}{(1 - \hat{\theta}^2)}.$$

$$\begin{aligned}\mathbb{V}\left[\frac{d \log f(y | x, \theta)}{d\theta} \mid x, \hat{\theta}\right] &= \mathbb{V}\left[(Y + x_4)/\hat{\theta} - (x_2 + x_3)/(1 - \theta) \mid x_1, \hat{\theta}\right] \\ &= \frac{\mathbb{V}[Y | x_1, \hat{\theta}]}{\hat{\theta}^2}, \\ &= \frac{125}{\hat{\theta}^2} \left(\frac{\hat{\theta}}{2 + \hat{\theta}}\right) \left(\frac{2}{2 + \hat{\theta}}\right).\end{aligned}$$

Likelihood function

Maximum likelihood estimation

Expectation-Maximization algorithm

Theory

Proof

Estandard Error

Examples

Examples

- Augmented data.
- Censored data.
- Mixtures.
- Random effect models.
- Missing data.