

# Computational Statistics

## Project 2



The code for this project is available under  
<https://github.com/max607/computational-statistics-em>.

By  
Maximilian Schneider

31.01.2023

# Contents

1	Maximum Likelihood estimation of $\theta$	2
2	Estimation of standard error	3
3	EM	4

# 1 Maximum Likelihood estimation of $\theta$

The pdf is

$$f(y_i; \theta) = \frac{\theta^2}{\theta + 1} (1 + y_i) \exp(-\theta y_i), i = 1 \dots n. \quad (1.1)$$

The log likelihood is

$$\ell(\theta) = 2n \log(\theta) - n \log(\theta + 1) - \theta \sum_{i=1}^n y_i + c \quad (1.2)$$

$$\propto 2 \log(\theta) - \log(\theta + 1) - \theta \bar{y} + c. \quad (1.3)$$

The first derivative is

$$\ell'(\theta) = \frac{2}{\theta} - \frac{1}{\theta + 1} - \bar{y}, \quad (1.4)$$

where  $\bar{y}$  is the sample mean of  $\mathbf{y}$ . Note, we can drop the n.

Setting the derivative to zero leads to equation (1.5) which has to be solved for  $\theta$ .

$$\frac{\theta + 2}{\theta(\theta + 1)} = \bar{y}. \quad (1.5)$$

One approach is to use Newton-Raphson, which requires the second order derivative.

$$\ell''(\theta) = -\frac{2}{\theta^2} + \frac{1}{(\theta + 1)^2}, \quad (1.6)$$

## 2 Estimation of standard error

The pdf can be restated as

$$f(y_i; \theta) = \frac{\theta}{\theta + 1} \exp(-\theta y_i) + \frac{1}{\theta + 1} \theta^2 y_i \exp(-\theta y_i), \quad (2.1)$$

i.e., a mixture of two gamma distributions in shape and rate parameterization.  $\theta$  is the rate and the shapes are equal to 1 and 2.

It is straight forward to simulate from this, but starting from  $U \stackrel{iid}{\sim} U(0, 1)$  exponentially distributed variables can be obtained via inversion

$$f^{-1}(u; \theta) = -\frac{\log(u)}{\theta}, \quad (2.2)$$

which is the same as a gamma with shape one and a gamma with shape 2 is obtained via the sum of 2 exponentials. For optimizing computation time  $n$  observations are generated in the following way:

- 1) Draw the number of shape 2 gammas ( $n_2$ ) by counting the number of  $u < \frac{1}{\theta+1}$
- 2) Sample  $n$  and  $n_2$  uniforms
- 3) Transform the uniforms to exponentials using  $f^{-1}$
- 4) Add to the first  $n_2$  of the  $n$  exponentials the other exponentials
- 5) Return

The result is a sample with  $n_2$  observations of a gamma distribution with shape 2 and  $n - n_2$  observation of a gamma distribution with shape 1.

For the purpose of estimating  $\theta$  with the estimator of section 1 it is of no importance that the sample is sorted by shape.

### 3 EM

Consider the more complex pdf

$$f(y_i; \theta, \lambda, \pi) = \pi \frac{\theta^2}{\theta + 1} (1 + y_i) \exp(-\theta y_i) + (1 - \pi) \lambda \exp(\lambda - y_i), \quad (3.1)$$

$$y_i, \theta, \lambda \in \mathbb{R}^+, \pi \in [0, 1]. \quad (3.2)$$

The goal is to estimate  $\theta, \lambda$  and  $\pi$  applying EM. As all observations are independent we formulate the complete likelihood as

$$\mathcal{L}(\theta, \lambda, \pi | \mathbf{x}, \mathbf{y}) = \prod_{i=1}^n \left( \pi \frac{\theta^2}{\theta + 1} (1 + y_i) \exp(-\theta y_i) \right)^{x_i} + \left( (1 - \pi) \lambda \exp(\lambda - y_i) \right)^{1-x_i}, \quad (3.3)$$

with missing data  $\mathbf{x}$ . As parameters are independent the relevant loglikelihoods parts are

- $\ell(\theta | \mathbf{x}, \mathbf{y}) = 2 \log(\theta) \sum_{i=1}^n x_i - \log(\theta + 1) \sum_{i=1}^n x_i - \theta \sum_{i=1}^n x_i y_i + c$
- $\ell(\lambda | \mathbf{x}, \mathbf{y}) = \log(\lambda) \sum_{i=1}^n (1 - x_i) - \lambda \sum_{i=1}^n (1 - x_i) y_i + c$

$\theta$  can be estimated as in section 1 replacing  $\bar{y}$  with  $(\sum_{i=1}^n x_i y_i) / \sum_{i=1}^n x_i$ . The derivative of  $\ell(\lambda | \mathbf{x}, \mathbf{y})$  is

$$\ell'(\lambda | \mathbf{x}, \mathbf{y}) = \frac{1}{\lambda} \left( n - \sum_{i=1}^n x_i \right) - \sum_{i=1}^n y_i + \sum_{i=1}^n x_i y_i. \quad (3.4)$$

This can be solved analytically

$$\hat{\lambda} = \frac{n - \sum_{i=1}^n x_i}{\sum_{i=1}^n y_i - \sum_{i=1}^n x_i y_i} \quad (3.5)$$