

Maximum Likelihood Estimation

1. INTRODUCTION

- Scientists use **models** to understand the phenomena they study.
- Any number computed with sample data → **statistics**.
- We can distinguish between **DETERMINISTIC MODELS & STATISTICAL MODELS**.
- **Probability:**
 - Known Population → **DEDUCTION** → Sample.
- **Statistics:**
 - Unknown Population ← **INDUCTION** ← Sample.
- **Parameters** → are fixed, unknown quantities that specify the population.
- **Estimators** → statistics that are used to estimate the unknown parameters.

2. MODEL ESTIMATION

- Statistical models have **unknown parameters** that need to be specified:
 - **Point Estimation** → estimating a population parameter with a **single value**.
 - MAXIMUM LIKELIHOOD ESTIMATOR.
 - METHOD OF MOMENTS.
 - **Interval Estimation** → estimating a population parameter with a **range of plausible values**.

3. MAXIMUM LIKELIHOOD ESTIMATOR

- Let X_1, \dots, X_n be a random sample from a distribution $f(x|\theta_1, \dots, \theta_k)$. *k = number of parameters*

- The **likelihood function** $L(\theta|\mathbf{x})$ is defined as

$$L(\theta|\mathbf{x}) = L(\theta_1, \dots, \theta_k | x_1, \dots, x_n) = \prod_{i=1}^n f(x_i | \theta_1, \dots, \theta_k)$$

product (pointing to the product symbol)
** must maximize*
→ d/dθ f(θ) = 0
→ d²/dθ² < 0 max

- This is in fact, **the joint density function** considering the data as given.
- We will often work with the **log-likelihood function** $\ell(\theta|\mathbf{x})$, defined correspondingly as

$$\ell(\theta|\mathbf{x}) = \ln(L(\theta|\mathbf{x})) = \ln(L(\theta_1, \dots, \theta_k | x_1, \dots, x_n)) = \sum_{i=1}^n \ln(f(x_i | \theta_1, \dots, \theta_k))$$

sum (pointing to the sum symbol)
Change Π to Σ

EXAMPLE BERNOULLI DISTRIBUTION:

Let X_1, \dots, X_n be a random sample with $X_i \sim \text{Bern}(p)$

$$P(X_1 = x_1 | p) = p^{x_1} (1 - p)^{1-x_1}$$

parameter (pointing to p)

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n | p) &= \prod_{i=1}^n p^{x_i} (1 - p)^{1-x_i} \\ &= p^{\sum_{i=1}^n x_i} (1 - p)^{n - \sum_{i=1}^n x_i} \end{aligned}$$

$$L(p | x_1, \dots, x_n) = p^{\sum_{i=1}^n x_i} (1 - p)^{n - \sum_{i=1}^n x_i}$$

- The maximum likelihood estimator $\hat{\theta}$ maximizes $L(\theta|\mathbf{x})$ as a function of θ .
- The method selects a value for θ such that the sample is most likely.
- Obtaining a maximum likelihood estimator is an optimization problem.
- In practice, it is often easier (and equivalent) to maximize the natural logarithm of the likelihood function, thus maximize $\ell(\theta|\mathbf{x})$.

- To find candidates for MLE:

1. Partial derivative of the parameter and equal to 0.
2. Do the second derivative and check if it's negative → **maximum**.

$$\frac{\partial}{\partial \theta} L(\theta|\mathbf{x}) = 0, \quad i = 1, \dots, k. \text{ and } \frac{\partial^2}{\partial \theta^2} L(\theta|\mathbf{x})|_{\theta=\hat{\theta}} < 0$$

- A point estimate obtained by ML is, by itself, not very informative.
- We need to specify its precision → with indicators
- The precision depends on the variance or the Fisher information of the ML estimator.

- FISCHER INFORMATION:

Let X_1, \dots, X_n be a random sample with

$$f(\mathbf{x} | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

theoretic formula

The Fisher information about θ contained in \mathbf{x} is defined by

$$I_{\mathbf{x}}(\theta) = E_{\theta} \left[\left(\frac{\partial}{\partial \theta} \ln(f(\mathbf{x} | \theta)) \right)^2 \right]$$

derivative of log-likelihood depending on θ
expected value
FISHER INFORMATION

- CRAMÉR-RAO LOWER BOUND:

tells which is the lowest variance

- For any unbiased estimator ($E(\hat{\theta}) = \theta$), there exists a lower bound on its variance.
- This bound equals the reciprocal of the Fisher information.

$$V(\hat{\theta}) \geq \frac{1}{I_{\mathbf{x}}(\theta)}$$

** choose the estimator with lowest variance*

Definition • An unbiased estimator that attains the Cramér-Rao lower bound is called **efficient**.

best estimator

- CONFIDENCE-INTERVALS:

- Having the variance and the distribution of the ML estimator, we can now say something about uncertainty.
- A **confidence interval** is an expression of the uncertainty of the estimate.
- A classical result, with $X_i \sim N(\mu, \sigma^2)$, is

$$CI(\mu)_{1-\alpha} = \bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

sample mean
confidence

where $\bar{X} = \hat{\mu}_{ML}$, and $\frac{\sigma}{\sqrt{n}} = \sqrt{\frac{\sigma^2}{n}} = \sqrt{V(\hat{\mu})}$.

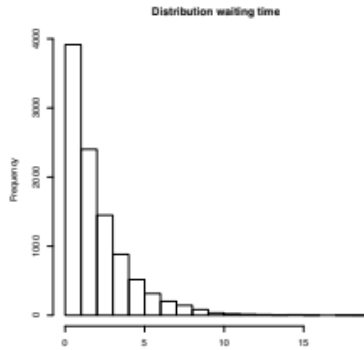
- Term $\frac{\sigma}{\sqrt{n}}$ (σ estimated by s) is called the **standard error of the mean**.
- Term $z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \approx 2 \frac{\sigma}{\sqrt{n}}$ when $\alpha = 0.05$ is the **error margin**.
- Equation (1) holds in general for ML estimators:

*α = error
 $1-\alpha$ = confidence*

$$CI(\theta)_{1-\alpha} = \hat{\theta} \pm z_{\alpha/2} \sqrt{V(\hat{\theta})}$$

1/2

- R EXAMPLE OF MLE:



- What is the rate of decay?
- What is the precision of a rate estimate?

```
> fitdistr(x,"exponential")
rate
0.498116487
(0.004981165)
```

Density and likelihood:

$$f(x|\lambda) = \lambda e^{-\lambda x} \quad L(\lambda|\mathbf{x}) = \lambda^n e^{-\lambda \sum x_i}$$

With some algebra, it follows that

$$\hat{\lambda} = 1/\bar{x},$$

$$I_n(\lambda) = n/\lambda^2$$

$$V(\hat{\lambda}) = \lambda^2/n$$

rate of decay

likelihood info

$$CI_{1-\alpha}(\lambda) = \hat{\lambda} \pm z_{\alpha/2} \sqrt{V(\hat{\lambda})} = \hat{\lambda} \pm z_{\alpha/2} \frac{\hat{\lambda}}{\sqrt{n}}$$

Descriptive statistics of a sample of $n = 10.000$ waiting times

	N	N*	Mean	Stdev	Med	Q1	Q3	Min	Max
X	10000	0	2.0075	2	1.397	0.579	2.768	0.001	18.163

$$\hat{\lambda} = 1/2.0075 = 0.49812$$

$$CI_{0.95}(\lambda) = 0.49812 \pm 1.96 \frac{0.49812}{\sqrt{10000}} = (0.4884; 0.5079)$$

WHAT we must answer?