

Introduction to statistical inference 2

Lauri Mehtätalo
University of Eastern Finland
School of Computing

May 28, 2018

Contents

1	Recap from “Introduction to statistical inference 1”	1
1.1	Random variable	1
1.2	Transformations of random variable	2
1.3	Expected values	3
1.4	Variance	3
1.5	Bivariate random variables	3
1.6	Independence	4
1.7	Covariance	4
1.8	Random vectors	5
1.9	Computing using expected values and variances	7
2	Random samples	9
2.1	Distribution of sum	14
2.2	Sampling from Normally distributed population	15
2.3	Convergence	18
3	Estimation	23
3.1	Two principles	24
3.2	The sufficient principle	24
3.3	The likelihood function	25
3.4	The likelihood principle	26
3.5	Point estimation	27
3.6	Estimating the population mean using least squares	27
3.7	The method of moments	28
3.8	Newton-Raphson	30
3.9	Maximum likelihood	31
3.10	Bayes estimators	36
3.11	Methods for evaluating estimators	39

Chapter 1

Recap from “Introduction to statistical inference 1”

1.1 Random variable

- Random variable is a function from sample space \mathcal{S} of an experiment to sample space of the random variable \mathcal{X} , which is set of real numbers.

$$X : \mathcal{S} \rightarrow \mathcal{X}$$

- The sample space is a set of all possible values random variable can get.
- \mathcal{X} can be
 - an interval of real axis (continuous random variable).

$$\mathcal{X} = [0, 10), \mathcal{X} = [0, 10], \mathcal{X} = (0, 10)$$

- An uncountable set of integers (discrete random variable)

$$\mathcal{X} = \{0, 1, 2, \dots\}$$

- A countable set of integers or real numbers (discrete random variable)

$$\mathcal{X} = \{0, 1\}, \mathcal{X} = \{0, 0.5, 1\}, \mathcal{X} = \{0, 1, \dots, 10\}$$

- Probabilities associated with each value of X are defined by the cumulative distribution function (cdf for short).

$$F_X(x) = P(X \leq x), \text{ where } -\infty < x < \infty$$

Note: $F_X(x)$ is a step function if X is discrete.

Note: $F_X(x)$ is a continuous function if X is continuous.

- $F_X(x)$ or $F(x)$ is cdf, if

- 1) $\lim_{x \rightarrow -\infty} F(x) = 0$ and $\lim_{x \rightarrow \infty} F(x) = 1$

- 2) $F(x)$ is non-decreasing

- 3) $F(x)$ is right-continuous

- Cdf is useful in calculation of any probabilities; for example

$$P(a < X \leq b) = F(b) - F(a)$$

Note: Be careful with $<$ and \leq when working with discrete random variables.

- The probability density function (pdf for short) is defined for continuous random variable as

$$f_X(x) = F'(x) = \frac{dF_X(x)}{dx}, \quad -\infty < x < \infty$$

and

$$\int_{-\infty}^x f_X(t) dt = F_X(x)$$

- The probability mass function (pmf for short) is defined for discrete random variables as

$$f_X(x) = P(X = x)$$

$$F_X(x) = \sum_{k=1}^x f_X(k)$$

1.2 Transformations of random variable

- Consider a monotonic function $g : \mathcal{X} \rightarrow \mathcal{Y}$
- $Y = g(X)$ is also a random variable; function g is called an transformation (muunnos).

- If $g(x)$ is a increasing function of x , then

$$F_Y(y) = F_X(g^{-1}(y))$$

- If $g(x)$ is a decreasing function of x , then

$$F_Y(y) = 1 - F_X(g^{-1}(y))$$

- The pdf of continuous Y is

$$f_Y(y) = F'_Y(y)$$

1.3 Expected values

$$E(X) = \mu_X = \begin{cases} \int_{-\infty}^{\infty} x f(x) \, dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} x f(x) & \text{if } X \text{ is discrete} \end{cases}$$

$$E(g(X)) = \begin{cases} \int_{-\infty}^{\infty} g(x) f(x) \, dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x) f(x) & \text{if } X \text{ is discrete} \end{cases}$$

1.4 Variance

$$\begin{aligned} \sigma_X^2 = \text{Var}(X) &= E(X - \mu_X)^2 \\ &= E(X^2 - 2X\mu_X + \mu_X^2) \\ &= E(X^2) - E(2X\mu_X) + E(\mu_X^2) \\ &= E(X^2) - 2\mu_X \underbrace{E(X)}_{\mu_X} + E(\mu_X^2) \\ &= E(X^2) - \mu_X^2 \end{aligned}$$

$$sd(X) = \sqrt{\text{Var}(X)} = \sigma_X$$

1.5 Bivariate random variables

- For two discrete random variables, the joint pmf is defined as

$$f_{X,Y}(x, y) = P(X = x, Y = y)$$

- For two continuous random variables, we define the joint pdf $f_{X,Y}(x, y)$ as

$$P((X, Y) \in A) = \iint_A f(x, y) \, dx \, dy$$

- The expected value for transformation $g(X, Y) : \mathbb{R}^2 \rightarrow \mathbb{R}$ (for example, $g(X, Y) = XY$ or $g(X, Y) = \frac{X}{Y}$) is

$$E(g(X, Y)) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) \, dx \, dy$$

if (X, Y) is continuous, and

$$E(g(X, Y)) = \sum_{x, y \in \mathbb{R}^2} g(x, y) f(x, y)$$

if (X, Y) is discrete.

- The marginal pmf / pdf for X are

$$f_X(x) = \sum_{y \in \mathbb{R}} f_{X,Y}(x, y) \quad (\text{pmf})$$

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy \quad (\text{pdf})$$

and correspondingly for Y .

- The conditional pmf / pdf are both defined as

$$f(y | x) = \frac{f_{X,Y}(x, y)}{f_X(x)} \quad (\text{for both discrete and continuous random variables})$$

and correspondingly for $x | y$.

1.6 Independence

- Random variables are said to be independent ($X \perp\!\!\!\perp Y$) if

$$f_{X,Y}(x, y) = f_X(x)f_Y(y)$$

- For independent random variables, conditional distribution $y | x$ is

$$f(y | x) = f_Y(y)$$

regardless of the value of x .

1.7 Covariance

- Covariance measures the linear association between two random variables X and Y

$$\text{cov}(X, Y) = E((X - \mu_X)(Y - \mu_Y))$$

$$= E(XY) - \mu_X\mu_Y$$

$$\text{cov}(X, Y) = \text{cov}(Y, X)$$

$$\text{corr}(X, Y) = \rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} \quad \text{where } -1 \leq \rho_{XY} \leq 1$$

Note: $\text{cov}(X, X) = \text{Var}(X)$

Note: If $X \perp\!\!\!\perp Y$, then $\text{cov}(X, Y) = 0$, but if $\text{cov}(X, Y) = 0$, it does not mean X and Y are necessarily independent.

1.8 Random vectors

- Random vectors generalize the bivariate random variables to a n -variate case.

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \quad \text{where } X_i, i = 1, 2, \dots, n \text{ are scalar random variables}$$

- If \mathbf{x} is a continuous random vector, then

$$P(\mathbf{X} \in A) = \int_A \cdots \int f(\mathbf{x}) \, dx_1 \cdots dx_n \quad f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

where $f(\mathbf{x})$ is a joint pdf.

- If \mathbf{x} is a discrete random vector, then

$$P(\mathbf{X} \in A) = \sum \cdots \sum f(\mathbf{x})$$

where $f(\mathbf{x})$ is a joint pmf.

- Let $g(\mathbf{x})$ be a transformation $g: \mathbb{R}^2 \rightarrow \mathbb{R}$. Then, the expected value is

$$E(g(\mathbf{X})) = \begin{cases} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} x f(\mathbf{x}) \, dx_1 \cdots dx_n & \text{if } \mathbf{X} \text{ is continuous} \\ \sum \cdots \sum_{\mathbf{X} \in \mathbb{R}^2} g(\mathbf{x}) f(\mathbf{x}) & \text{if } \mathbf{X} \text{ is discrete} \end{cases}$$

- Let us partition the n -variate random vector \mathbf{X} as follows:

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{pmatrix}$$

where \mathbf{X}_1 has length k and \mathbf{X}_2 has length $n - k$.

- The joint pdf of \mathbf{X} can be written as

$$f(\mathbf{x}) = f(\mathbf{x}_1, \mathbf{x}_2)$$

- The marginal density of \mathbf{X}_1 is

$$\begin{aligned} f(\mathbf{x}_1) &= \int_{\mathbb{R}^{n-k}} f(\mathbf{x}_1, \mathbf{x}_2) \, d\mathbf{x}_2 & \text{if } \mathbf{X} \text{ is continuous} \\ f(\mathbf{x}_1) &= \sum_{\mathbb{R}^{n-k}} f(\mathbf{x}_1, \mathbf{x}_2) & \text{if } \mathbf{X} \text{ is discrete} \end{aligned}$$

where $f(\mathbf{x}_1)$ is a k -variate pdf/pmf.

- The conditional pdf/pmf for $\mathbf{X}_2 \mid \mathbf{X}_1$ is

$$f(\mathbf{X}_2 \mid \mathbf{X}_1) = \frac{f(\mathbf{X}_1, \mathbf{X}_2)}{f(\mathbf{X}_1)}$$

- Expected value of random vectors is defined as vector

$$E(\mathbf{x}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_n) \end{bmatrix}_{n \times 1} = \begin{bmatrix} E(\mathbf{X}_1) \\ E(\mathbf{X}_2) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$$

- The variance of a random vector is $n \times n$ symmetric matrix called variance-covariance matrix.

$$\begin{aligned} Var(\mathbf{x})_{n \times n} &= \begin{bmatrix} Var(X_1) & cov(X_1, X_2) & cov(X_1, X_3) & \dots & cov(X_1, X_n) \\ cov(X_2, X_1) & Var(X_2) & cov(X_2, X_3) & \dots & cov(X_2, X_n) \\ cov(X_3, X_1) & cov(X_3, X_2) & Var(X_3) & \dots & cov(X_3, X_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ cov(X_n, X_1) & cov(X_n, X_2) & cov(X_n, X_3) & \dots & Var(X_n) \end{bmatrix} \\ &= \begin{bmatrix} Var(\mathbf{x}_1)_{k \times k} & cov(\mathbf{x}_1, \mathbf{x}'_2)_{k \times (n-k)} \\ cov(\mathbf{x}_2, \mathbf{x}'_1)_{(n-k) \times k} & Var(\mathbf{x}_2)_{(n-k) \times (n-k)} \end{bmatrix} \end{aligned}$$

- The correlation matrix is defined as

$$corr(\mathbf{x})_{n \times n} = \begin{bmatrix} 1 & corr(X_1, X_2) & corr(X_1, X_3) & \dots & corr(X_1, X_n) \\ corr(X_2, X_1) & 1 & corr(X_2, X_3) & \dots & corr(X_2, X_n) \\ corr(X_3, X_1) & corr(X_3, X_2) & 1 & \dots & corr(X_3, X_n) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ corr(X_n, X_1) & corr(X_n, X_2) & corr(X_n, X_3) & \dots & 1 \end{bmatrix}$$

- If \mathbf{X} has a n -variate normal distribution, then, with

$$E(\mathbf{x}) = \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \quad \text{and} \quad Var(\mathbf{x})_{n \times n} = \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_1 & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_2 \end{bmatrix}$$

then $\mathbf{X}_1 \mid \mathbf{X}_2$ has a k -variate normal distribution with

$$E(\mathbf{X}_1 \mid \mathbf{X}_2) = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_2^{-1}(\mathbf{X}_2 - \boldsymbol{\mu}_2)$$

$$\text{Var}(\mathbf{X}_1 \mid \mathbf{X}_2) = \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_{21}$$

Note: $\text{Var}(\mathbf{X}_1 \mid \mathbf{X}_2) \leq \text{Var}(\mathbf{X}_1)$

1.9 Computing using expected values and variances

Let a , b and c be constants, and let X , Y and Z be (scalar) random variables. The following rules hold regardless of the distribution of random variables X , Y and Z .

$$E(c) = c$$

$$E(cX) = cE(X)$$

$$E(X + Y) = E(X) + E(Y)$$

$$E(X + c) = E(X) + c$$

$$E(XY) = E(X)E(Y)$$

Only if $X \perp\!\!\!\perp Y$.

$$E(g(X)) = g(E(X))$$

Only in some special cases, like when $g(X)$ is a linear transformation.

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y)$$

$$\text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2\text{cov}(X, Y)$$

$$\text{Var}(aX) = a^2 \cdot \text{Var}(X)$$

$$\text{cov}(X, Y + Z) = \text{cov}(X, Y) + \text{cov}(X, Z)$$

$$\text{cov}(aX, bY) = ab \cdot \text{cov}(XY)$$

$$\text{Var}(X + a) = \text{Var}(X)$$

$$\text{cov}(X + a, Y + b) = \text{cov}(X, Y)$$

$$E(X) = E_Y \left[E_{X|Y}(X \mid Y) \right]$$

$$\text{Var}(X) = E_Y \left[\text{Var}_{X|Y}(X \mid Y) \right] + \text{Var}_Y \left[E_{X|Y}(X \mid Y) \right]$$

Let \mathbf{a} and \mathbf{b} be fixed vectors and \mathbf{X} and \mathbf{Y} random vectors so that the dimensions in the equations match.

$$E(\mathbf{a}'\mathbf{X}) = \mathbf{a}'E(\mathbf{X})$$

$$\text{Var}(\mathbf{a}'\mathbf{X}) = \mathbf{a}'\text{Var}(\mathbf{X})\mathbf{a} \quad \text{Compare to } \text{Var}(aX) = a^2 \cdot \text{Var}(X) = a \cdot \text{Var}(X) \cdot a$$

$$\text{cov}(\mathbf{a}'\mathbf{X}, \mathbf{b}'\mathbf{Y}) = \mathbf{a}'\text{cov}(\mathbf{X}, \mathbf{Y})\mathbf{b}$$

Note: These equations need to be remembered by heart!

Chapter 2

Random samples

Definition 2.1. Random variables X_1, \dots, X_n are called random sample of size n from population $f(x)$, if X_1, \dots, X_n are mutually independent random variables and the marginal pdf/pmf of each X_i is the same function $f(x)$. Alternatively, X_1, \dots, X_n are called independent and identically distributed random variables (i.i.d.) with pdf/pmf $f(x)$.

Note: Sample X_1, \dots, X_n can also be denoted by \mathbf{X} , where $\mathbf{X} = [X_1 \ \dots \ X_n]^T$

- If follows from the mutual independence, of X_1, \dots, X_n that the joint pdf or pmf of \mathbf{X} is

$$f(x_1, \dots, x_n) = f(x_1)f(x_2) \dots f(x_n) = \prod_{i=1}^n f(x_i)$$

- All univariate marginal distributions $f(x_i)$ are the same by definition 2.1.

Example 2.1. Let X_1, \dots, X_n be a random sample from *Exponential*(β) population. X_i specifies the time until failure for n identical cellphones. The exponential pdf is

$$f(x_i) = \frac{1}{\beta} e^{-x_i/\beta}$$

The joint pdf of the sample is

$$f(x_1, x_2, \dots, x_n \mid \beta) = \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} = \frac{1}{\beta^n} e^{-\frac{1}{\beta} \sum x_i} \quad \text{Recall: } a^b a^c = a^{b+c}$$

What is the probability that all n cellphones last more than 2 years?

$$\begin{aligned}
P(X_1 > 2, \dots, X_n > 2) &= \int_2^\infty \dots \int_2^\infty f(x_1, \dots, x_n) dx_1 \dots dx_n \\
&= \int_2^\infty \dots \int_2^\infty \prod_{i=1}^n \frac{1}{\beta} e^{-x_i/\beta} dx_1 \dots dx_n \\
&= \int_2^\infty \dots \int_2^\infty \prod_{i=2}^n \frac{1}{\beta} e^{-x_i/\beta} \underbrace{\int_2^\infty \frac{1}{\beta} e^{-x_1/\beta} dx_1}_{\substack{\text{Integral over the exponential pdf} \\ \int_2^\infty f(x_1) dx_1 = 1 - F(2) \\ = 1 - (1 - e^{-\frac{1}{\beta^2}}) = e^{-\frac{2}{\beta^2}}} } dx_2 \dots dx_n \\
&= e^{-\frac{2}{\beta}} \underbrace{\int_2^\infty \dots \int_2^\infty \prod_{i=2}^n \frac{1}{\beta} e^{-x_i/\beta} dx_2 \dots dx_n}_{\substack{\text{Identical to original integral except that} \\ \text{there are only } n-1 \text{ terms to be integrated}}} = \dots \\
&= e^{-2/\beta} \cdot e^{-2/\beta} \int_2^\infty \dots \int_2^\infty \prod_{i=3}^n \frac{1}{\beta} e^{-x_i/\beta} dx_3 \dots dx_n \\
&= (e^{-2/\beta})^n = e^{-2n/\beta}
\end{aligned}$$

A much more simpler solution: notice that $P(X_i > 2) = 1 - F(2) = e^{-1/\beta}$. Because X_i 's are independent, then also events $(X_i > 2)$ are independent event, and so

$$P(\text{All } X_i > 2) = \prod_{i=1}^n P(X_i > 2) = (e^{-2/\beta})^n = e^{-2n/\beta}$$

Illustration See `ExponentialSample.R` for the case where $n = 2$ and $\beta = 3$. Notice that R uses parametrization $\lambda = \frac{1}{\beta}$ for the exponential distribution, so in R, exponential pdf is $f(x | \lambda) = \lambda e^{-\lambda x}$.

Note: Definition 2.1 assumes that X_1, \dots, X_n are independent. In practice, we often do analysis on dependent data. For that use, we could define term “dependent random sample”. Mathematical treatment of dependent sample requires more specific description of dependence structure, e.g. through spatial or temporal autocorrelation models, or explicit models for grouped data.

Dependent data are modeled by assuming that random vectors $\mathbf{X}_1, \dots, \mathbf{X}_n$ are vectors of appropriate length, and there are n independent realizations of them in the data. Quite often $n = 1$ and each replicate to \mathbf{X}_1 , which is a rather long vector.

Example 2.2. Let $\mathbf{X} = [X(u_1) \ X(u_2) \ X(u_3)]^T$ include random variables X at locations u_1, u_2 and u_3 . Assume that \mathbf{X}_1 is normally distributed and the correlation $\rho_{ij} = \text{corr}(X(u_i), X(u_j))$ depends only on the spatial distance $s_{ij} = \|u_i - u_j\|$ between locations u_i and u_j .

The marginal means and variances of $X(u_i)$ are $E(X(u_i)) = \mu$ for all i and $\text{Var}(X(u_i)) = \sigma^2$ for all i . Consider case where $u_1 = (1, 0), u_2 = (0, 0), u_3 = (0, 2)$ and $\rho_{ij} = e^{-s_{ij}/2}$. Correlations and covariances related to different random variables can be seen on table 2.1.

Table 2.1: Correlations between the random variable pairs of example 2.2

Pair	s_{ij}	ρ_{ij}	$cov(X(u_i), X(u_j))$
1,2	1	0.61	$0.61\sigma^2$
1,3	$\sqrt{5}$	0.33	$0.33\sigma^2$
2,3	2	0.37	$0.37\sigma^2$

Let

$$\boldsymbol{\mu} = \begin{bmatrix} \mu \\ \mu \\ \mu \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Sigma} = \sigma^2 \begin{bmatrix} 1 & 0.61 & 0.33 \\ 0.61 & 1 & 0.37 \\ 0.33 & 0.37 & 1 \end{bmatrix}$$

The joint pdf of \mathbf{X} is 3 variate normal distribution with mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$

$$f(x_1 | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^3}} |\boldsymbol{\Sigma}|^{-1/2} e^{\frac{1}{2}(\mathbf{x}_1 - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu})}$$

Note: In random sample, we assume that X_1, \dots, X_n are identically distributed. In the case of random vectors, this is specified by saying that the marginal distributions of the elements of \mathbf{X}_i are identical.

Note: We can also have independent replicates of random vectors, e.g. if we have grouped data where the groups are independent replicates from the process that generates the groups, and the observations within the groups are dependent. This is related to mixed-effects models.

Note: Definition 2.1 specifies sampling from an infinite population (or population model): the sampling procedure does not change the population.

The population may also be finite, like numbers in hat. In that case, the sampling can be made with replacement: whenever a number has been drawn, the value is recorded but the number is put back to the hat. This procedure, which is useful e.g. in Bootstrapping, fulfils the conditions of definition 2.1.

If we do the sampling without replacement (which often makes much sense), then the conditions of definition 2.1 are not fulfilled: each draw changes the population by removing one unit from the finite population. This leads to so called design-based inference, which is covered in the literature of sampling theory.

The difference to definition 2.1 is important especially if the sample size n is large compared to the size of population, but may be unimportant if $N \gg n$. In this course, we are considering only sampling according to definition 2.1.

Random samples can be summarized to a well defined summary called statistic.

Definition 2.2. Let $\mathbf{X} = X_1, \dots, X_n$ be a random sample of size n from population $f(x)$ and let $T(\mathbf{X})$ be a real-valued or vector-valued function, whose domain includes

the sample space of \mathbf{X} . The random variable $Y = T(\mathbf{X})$ is called a statistic (statistikka, tunnusluku). The probability distribution of Y is called the sampling distribution of Y (otantajakauma).

Examples of statistics

Sample minimum $\min(\mathbf{X})$

Sample maximum $\max(\mathbf{X})$

Sample median $\text{median}(\mathbf{X})$

Also other like sample mean, variance, standard deviation, quartiles etc.

Example 2.3. Let \mathbf{X} be a random sample of size 5 from *Exponential*(2) population. R script `statistics.R` illustrates the distributions of $\min(\mathbf{X})$, $\max(\mathbf{X})$, $\text{median}(\mathbf{X})$ and $\text{mean}(\mathbf{X})$ through simulation. Script repeats the sampling from Exponential distribution $M = 10000$ times and illustrates the distribution of the above mentioned sample statistic when $n = 5$.

Definition 2.3. The sample mean is the statistic defined by

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

Definition 2.4. The sample variance is the statistic defined by

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

- The sample standard deviation is the statistic defined by $S = \sqrt{S^2}$.
- The observed values of these random variables are denoted by \bar{X} , s^2 and s .

Theorem 2.1. Let X_1, \dots, X_n be any numbers. Then

$$\begin{aligned} \text{a) } \min_a \sum_{i=1}^n (x_i - a)^2 &= \sum_{i=1}^n (x_i - \bar{x})^2 \\ \text{b) } (n-1)s^2 &= \sum_{i=1}^n (x_i - \bar{x})^2 = \underbrace{\sum_{i=1}^n x_i^2 - n\bar{x}^2}_{\text{Compare to } \text{Var}(\mathbf{X}) = E(\mathbf{X}^2) - (E(\mathbf{X}))^2} \end{aligned}$$

Lemma 2.1. Let X_1, \dots, X_n be a random sample from a population and the $g(X)$ be a function such that $E(g(X_i))$ and $\text{Var}(g(X_i))$ exist. Then

$$E \left[\sum_{i=1}^n g(X_i) \right] = n \cdot E(g(X_i))$$

and

$$Var\left(\sum_{i=1}^n g(X_i)\right) = n \cdot Var(g(X_i))$$

Proof

$$E\left(\sum g(X_i)\right) = \sum E(g(X_i)) = *$$

Because X_i 's are identically distributed

$$* = nE(g(X_i))$$

$$\begin{aligned} Var(\sum g(X_i)) &= E\left[\sum g(X_i) - E(\sum g(X_i))\right]^2 \\ &= E\left\{\sum \left[\underbrace{g(X_i) - E(g(X_i))}_{t_i}\right]\right\}^2 = * \end{aligned}$$

With t_i , we can get following

$$\begin{aligned} E(\sum t_i \sum t_i) \\ &= E(\sum t_i^2 + \sum_{i \neq j} t_i t_j) \\ &= \sum E(t_i^2) + \sum_{i \neq j} E(t_i t_j) \end{aligned}$$

which we can expand back to original equation

$$\begin{aligned} * &= \underbrace{\sum E([g(X_i) - E(g(X_i))]^2)}_{\sum E(t_i^2)} + \underbrace{\sum_{i \neq j} E[(g(X_i) - E(g(X_i)))(g(X_j) - E(g(X_j)))]}_{\sum_{i \neq j} E(t_i t_j)} \\ &= \sum Var(g(X_i)) + \sum_{i \neq j} cov(g(X_i), g(X_j)) \\ &= \underbrace{n \cdot Var(X_1)}_{\text{Because } X_i \text{'s are identically distributed}} + n(n-1) \cdot \underbrace{0}_{\text{Because } X_i \text{'s are independent}} = n \cdot Var[g(X_1)] \end{aligned}$$

Note: The proof of $E(\sum g(X_i))$ did not utilize the assumption of independence. Therefore it is valid also for dependent samples. The poof of $Var(\sum g(X_i))$ used the assumption of independence. Therefore it is not valid for dependent samples.

- For dependent samples, we get

$$Var(\sum g(X_i)) = \sum_{i=1}^n \sum_{j=1}^n cov[g(X_i), g(X_j)]$$

- Proof left as an exercise (recall, that $Var(X_i) = cov(X_i, X_i)$).

$$\begin{aligned}
\text{Var}(\mathbf{x}) &= \begin{bmatrix} \text{Var}(X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_n) \\ \text{cov}(X_2, X_1) & \text{Var}(X_2) & \dots & \text{cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(X_n, X_1) & \text{cov}(X_n, X_2) & \dots & \text{Var}(X_n) \end{bmatrix} \\
\text{Var}(\sum X_i) &= \sum \sum \text{cov}(X_i, X_j)
\end{aligned}$$

- That is, the sum of all elements of matrix $\text{Var}(\mathbf{X})$.

Note: This extends the rule

$$\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) + 2 \cdot \text{cov}(X_1, X_2)$$

to general sum $\sum X_i$.

Theorem 2.2. Let X_1, \dots, X_n be a random sample from a population with mean μ and variance $\sigma^2 < \infty$. Then

- $E(\bar{X}) = \mu$
- $\text{Var}(\bar{X}) = \frac{\sigma^2}{n} \quad E(\bar{X} - \mu)^2$
- $E(S^2) = \sigma^2$

- Proof for a and b are familiar from ISI1 course. Proof of c applies the theorem 2.1 and is left as an exercise.

Note: Because $E(\bar{X}) = \mu$, the statistic \bar{X} is said to be an unbiased estimator (harhaton estimaattori) of population mean μ . Also S^2 is an unbiased estimator of population variance σ^2 .

Example 2.4. Illustrate these results through simulation. Intuitively, mean of an observed value of a statistic over a large number of samples should be close to expected value of the statistic in question. Therefore, (a), (b) and (c) of theorem 2.2 can be demonstrated by simulating M samples of a fixed size n and exploring how the means of the sample values of \bar{x} , $(\bar{x} - \mu)^2$ and s^2 behave as $M \rightarrow \infty$.

R-script `demonstratebias.R` implements this by assuming that the population has the *Uniform*(0, 10) distribution and $n = 10$.

2.1 Distribution of sum

Theorem 2.3. (Convolution formula)

If X and Y are independent, continuous random variables with pdf's $f_X(x)$ and $f_Y(y)$, then the pdf of the sum $Z = x + Y$ is

$$f_Z(z) = \int_{-\infty}^{\infty} f_X(w)f_Y(z-w) dw$$

Proof: See Casella & Berger, p 215-216.

Illustration: See example 1.20 of `notes.pdf`

Note: To compute the complete distribution of a sum $Z = \sum_{i=1}^n X_i$, where X_i 's do not need to be identically distributed, but they are independent, we need to apply theorem 2.3 iteratively. E.g. to find the distribution of $X_1 + X_2 + X_3$, you may first find the distribution of $Z = X_1 + X_2$ using theorem 2.3 and thereafter use theorem 2.3 again to find the pdf of $Z + X_3$. This is not trivial in very general case, but there are easier ways to do this, e.g. when X_i 's are identically distributed and independent (i.i.d).

However, recall that in any case, the moments are easy

$$E(\sum X_i) = \sum E(X_i)$$

$$Var(\sum X_i) = \sum Var(X_i) \quad (\text{If } X_i\text{'s are independent})$$

2.2 Sampling from Normally distributed population

Theorem 2.4. Let $X_{n \times 1}$ be random sample of size n from a $N(\mu, \sigma^2)$ distribution and let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ and $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Then

- \bar{X} and S^2 are independent random variables.
- \bar{X} has $N(\mu, \frac{\sigma^2}{n})$
- $(n-1)S^2/\sigma^2$ has chi-squared distribution with parameter $n-1$. Parameter $n-1$ is commonly called the “degrees of freedom”.

Note: The chi-squared distribution with p degrees of freedom has the pdf

$$f_X(x) = \frac{1}{\underbrace{\Gamma(p/2)}_{\text{Gamma -function}}} X^{p/2-1} e^{-x/2}$$

- Gamma -function: $\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx$

Lemma 2.2. We use notation χ_p^2 for a chi-squared random variables with p degrees of freedom

- If Z is a $N(0, 1)$ random variable, then $Z^2 \sim \chi_1^2$, that is, the square of a standard normal random variable is a chi-squared random variable.

- If X_1, \dots, X_n are independent and $X_i \sim \chi_{p_i}^2$ then $\sum X_i \sim \chi_{\sum p_i}^2$. That is, independent chi-squared random variables add to a chi-squared random variable, and the degrees of freedom also add.

Example 2.5. Illustrate the result of theorem 2.4 using R-script. In file `normalchi.R`, we implement the following:

Repeat M times

1. Generate random variables $X_i \sim N(0, 1)$, where $n = 1, \dots, 9$
2. Compute $Z = \sum x_i^2$ for each sample
3. Plot the histogram of the M obtained values of Z and compare to the distribution χ_n^2 .

Note: If \bar{X}_n is a random sample from $N(\mu, \sigma^2)$ population, then the random variable

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$$

where $\mu = E(\bar{X})$ and $\sigma/\sqrt{n} = \text{Var}(\bar{X})$.

- This transformation could be used to make inference on μ using the observed \bar{x} , if σ^2 is known.
- However, that is not the case, and therefore we use $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ which can also be written as

$$\frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{\sqrt{S^2/\sigma^2}}$$

where $(\bar{X} - \mu)/(\sigma/\sqrt{n}) \sim N(0, 1)$ and $\sqrt{S^2/\sigma^2} \sim \chi_{n-1}^2$.

Theorem 2.5. Let \bar{X} be a random sample from $N(\mu, \sigma^2)$. The random variable

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

has Student's t-distribution with parameter (n-1) "degrees of freedom".

In general, a random variable T has Student's t-distribution with p degrees of freedom ($T \sim t_p$) if T has the pdf

$$f_T(t) = \frac{\Gamma(\frac{p+1}{2})}{\Gamma(\frac{p}{2})} \frac{1}{(p\pi)^{1/2}} \frac{1}{(1 + t^2/p)^{(p+1)/2}}$$

Proof: See Casella & Berger, p.223-224.

Illustration: File `normalchi.R`

- We use samples of sizes 3, 4, \dots , 100 and $M = 10000$.

- Plot empirical histogram of $\frac{\bar{X}-\mu}{s/\sqrt{n}}$
- Compare to t_{n-1} distribution and to $N(0, \sigma^2)$.
- $\mu = 5$ and $\sigma^2 = 3^2$.
- Because $S^2 \rightarrow \sigma^2$ as n increases, the difference between $N(0, 1)$ and t_{n-1} becomes meaningless as $n \rightarrow \infty$.

Note: t_p has only p moments. Especially $E(T_p) = 0$ if $p > 1$.

- A third important derived distribution under sampling from a normal population is the Snedecor's/Fisher's F-distribution.
- It is the theoretical distribution of the ratio of variances.

Example 2.6. Let \mathbf{X} be a random sample from $N(\mu_X, \sigma_X^2)$ and let \mathbf{Y} be a random sample from $N(\mu_Y, \sigma_Y^2)$. If we were interested in comparing the variability in these two populations, a quality of interest would be σ_X^2/σ_Y^2 . Information about σ_X^2/σ_Y^2 is contained in s_X^2/s_Y^2 ; ratio of the sample variances.

The F-distribution allows this comparison by giving the distribution for random variable

$$\frac{s_X^2/s_Y^2}{\sigma_X^2/\sigma_Y^2} = \frac{s_X^2/\sigma_X^2}{s_Y^2/\sigma_Y^2}$$

where ratios s_X^2/σ_X^2 and s_Y^2/σ_Y^2 are independent, scaled χ^2 -distributed random variables.

Definition 2.5. Let $\mathbf{X}_{n \times 1}$ be a random sample from a $N(\mu_X, \sigma_X^2)$ population and let $\mathbf{Y}_{m \times 1}$ be a random sample from an independent $N(\mu_Y, \sigma_Y^2)$ population. The random variable

$$F = \frac{s_X^2/\sigma_X^2}{s_Y^2/\sigma_Y^2}$$

has the F-distribution with parameters (numerator and denominator degrees of freedom) $n - 1$ and $m - 1$. Equivalently, random variable F has the F-distribution with p and q degrees of freedom ($F \sim F_{p,q}$) if F has the pdf

$$f_F(x) = \frac{\Gamma(\frac{p+q}{2})}{\Gamma(\frac{p}{2})\Gamma(\frac{q}{2})} \left(\frac{p}{q}\right)^{p/2} \frac{x^{p/2-1}}{\left[1 + \frac{p}{q}x\right]^{(p+q)/2}}$$

Note:

- If $X \sim F_{p,q}$, then $1/X \sim F_{q,p}$
- If $X \sim t_q$, then $X^2 \sim \chi_1^2$ (recall that if $X \sim N(0, 1)$, then $X^2 \sim \chi_1^2$)

2.3 Convergence

- What happens if $n \rightarrow \infty$.

Theorem 2.6. Markov inequality

Let X be a random variable such that $P(X \geq 0) = 1$ (i.e. X gets only positive values). Then, for every number $t > 0$

$$P(x > t) \leq \frac{E(X)}{t}$$

Proof: Consider only case where X is discrete random variable.

$$\begin{aligned} E(X) &= \sum_X x f(x) \\ &= \sum_{x < t} x f(x) + \sum_{x \geq t} x f(x) \end{aligned}$$

Because $X \geq 0$, both terms are positive

$$\begin{aligned} E(X) &\geq \sum_{X \geq t} x f(x) \geq \sum_{X \geq t} t f(x) = t \underbrace{\sum_{X \geq t} f(x)}_{P(X \geq t)} \\ &= t \cdot P(X \geq t) \end{aligned}$$

$$\begin{aligned} E(X) &\geq t \cdot P(X \geq t) \quad || : t \ (> 0) \\ \frac{E(X)}{t} &\geq P(X \geq t) \end{aligned}$$

Example 2.7. Let X be a non-negative random variable with $E(X) = 1$

$$P(X \geq 100) \leq \frac{E(X)}{100} = 0.01$$

Theorem 2.7. Chebyshev's inequality

Let X be a random variable such that $Var(X)$ exists. Then for every number $t > 0$

$$P(|X - E(X)| \geq t) \leq \frac{\sigma^2}{t^2}$$

Proof: Let $Y = (X - E(X))^2$. Now $P(Y \geq 0) = 1$ and $E(Y) = Var(X)$

$$P(|X - E(X)| \geq t) = P(Y \geq t^2) \leq \frac{E(Y)}{t^2}$$

$$P(|X - E(X)| \geq t) \leq \frac{Var(X)}{t^2}$$

Example 2.8. If $Var(X) = \sigma^2$ and we select $t = 3\sigma$

$$P(|X - E(X)| \geq 3\sigma) \leq \frac{\sigma^2}{(3\sigma)^2} = \frac{1}{9}$$

Definition 2.6. A sequence of random variables X_1, X_2, \dots converges in probability (konvergoi todennäköisyyksimielessä) to a random variable X if for every $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$$

or equivalently

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \epsilon) = 1$$

A commonly used notation for this kind of convergence is $X_n \xrightarrow{P} X$. The sequence X_1, \dots, X_n is not usually an i.i.d. sample but e.g. a statistic based on a sample of size n . Also, X is often a fixed constant, as it is in the following theorem.

Theorem 2.8. (The weak law of large numbers (Heikko suurten lukujen laki), WLLN)

Let X_1, X_2, \dots be i.i.d. random variables with expected value $E(X_i) = \mu$ and variance $Var(X_i) = \sigma^2 < \infty$. Define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. Then, for every $\epsilon > 0$

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \epsilon) = 1$$

meaning sample mean convergences in probability to population mean ($\bar{X}_n \xrightarrow{P} \mu$). We will have the sample means of an i.i.d. sample arbitrarily close to μ if just n is high enough.

Proof: Use Chebyshev's inequality for the complement event:

$$P(|\bar{X}_n - \mu| \geq \epsilon) \leq \frac{\overbrace{Var(\bar{X}_n)}^{=\sigma^2/n}}{\epsilon^2} = \frac{\sigma^2}{n\epsilon^2}$$

$$P(|\bar{X}_n - \mu| < \epsilon) = 1 - P(|\bar{X}_n - \mu| \geq \epsilon) \geq 1 - \frac{\sigma^2}{n\epsilon^2}$$

where $\lim_{n \rightarrow \infty} \frac{\sigma^2}{n\epsilon^2} = 0$ and $\lim_{n \rightarrow \infty} 1 - \frac{\sigma^2}{n\epsilon^2} = 1$.

Note The property that the same sample quantity approaches a fixed constant as $n \rightarrow \infty$ is called consistency (konsistenssi).

Note WLLN was already used in R-script `demonstrateBias.R`

Note WLLN also justifies the use of a histogram from a large number of replicates as an approximation of the true density as $n \rightarrow \infty$. This is because every class of the histogram gives $P(c_1 \leq x < c_2) = F(c_2) - F(c_1)$. Whether x belongs to a class $[c_1, c_2[$ is a Bernoulli(p) distributed random variable with

$$p = E(Y) = F(c_2) - F(c_1)$$

Note Another way of specifying the law of large numbers is the strong law of large numbers (SLLN).

$$P(\lim_{n \rightarrow \infty} \bar{X}_n = \mu) = 1.$$

\bar{X} converges to μ with probability 1, meaning almost sure convergence.

Definition 2.7. A sequence of random variables converges in distribution to random variable X if

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

at all points of x where $F_X(x)$ is a continuous cdf. Notation used for convergence in distribution is $X_n \xrightarrow{d} X$.

Example 2.9. (Maximum of uniforms)

If X_1, X_2, \dots are i.i.d. Uniform(0,1) distributed random variables and let $X_{(n)} = \max_{1 \leq i \leq n} X_i$. Let us explore if and to where $X_{(n)}$ converges in distribution.

When $n \rightarrow \infty$, $X_{(n)} \rightarrow 1$ and as $X_{(n)} < 1$, we have for any $\epsilon > 0$

$$P(|X_{(n)} - 1| \geq \epsilon) = P(1 - X_{(n)} \geq \epsilon) = P(X_{(n)} \leq 1 - \epsilon)$$

Because sample is independent

$$\begin{aligned} P(X_{(n)} < 1 - \epsilon) &= P(\text{All } X_i \leq 1 - \epsilon) \\ &= \left[\underbrace{P(X \leq 1 - \epsilon)}_{F_X(1-\epsilon)} \right]^n && X \sim \text{Unif}(0, 1) \\ &= (1 - \epsilon)^n && \lim_{n \rightarrow \infty} (1 - \epsilon)^n = 0 \\ &\rightarrow X_{(n)} \xrightarrow{P} 1 \end{aligned}$$

Let us take $\epsilon = \frac{t}{n}$ to rewrite $P(X_{(n)} \leq 1 - \epsilon)$ as

$$P(X_{(n)} \leq 1 - \frac{t}{n}) = (1 - \frac{t}{n})^n$$

for which $\lim_{n \rightarrow \infty} (1 - \frac{t}{n})^n = e^{-t}$.

$$\begin{aligned} X_{(n)} &\leq 1 - \frac{t}{n} \\ X_{(n)} - 1 &\leq -\frac{t}{n} && | \cdot (-1) \\ 1 - X_{(n)} &\leq \frac{t}{n} && | \cdot n \\ n(1 - X_{(n)}) &\leq t \\ P(n(1 - X_{(n)}) \leq t) &= e^{-t} \\ P(n(1 - X_{(n)}) \geq t) &= \underbrace{1 - e^{-t}}_{\text{Exponential}(1) \text{ cdf}} \end{aligned}$$

Illustration: See R-script `MaxOfUniforms.R`

How fast does the convergence occur, i.e. how large does the n need to be to have $n(1 - X_{(n)}) \sim \text{Exponential}(1)$?

We generate $M = 10000$ samples of each of the sample sizes $n = 2, 3, 4, 5, 10, 20$. For each sample, find $X_{(n)}$ and compute $y = n(1 - X_{(n)})$. Plot the histogram of y and compare to Exponential(1) pdf.

The approximation looks quite nice already, when $n \geq 10$. But notice the behaviour in the tails. Especially $n(1 - X_{(n)})$ is bounded to interval $[0, n]$.

Theorem 2.9. (The central limit theorem (CTL) (Keskeinen raja-arvolause))

Let X_1, X_2, \dots be a sequence of i.i.d. random variables with finite $E(X) = \mu$ and finite variance $Var(X) = \sigma^2 > 0$. Let $\bar{X}_n = \sum_{i=1}^n X_i$. Let $G_{Z_n}(z)$ denote the cdf of

$$Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} = \frac{\overbrace{\bar{X}_n - \mu}^{E(\bar{X}_n)}}{\underbrace{\sigma/\sqrt{n}}_{sd(\bar{X}_n)}}$$

Then $\lim_{n \rightarrow \infty} G_{Z_n}(z) = \underbrace{\int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-y^2/2} dy}_{\substack{N(0,1) \text{ pdf} \\ N(0,1) \text{ cdf}}}$ That is: Z_n has limited standard normal

distribution, meaning $Z_n \xrightarrow{d} N(0, 1)$. Equivalently, $\bar{X} \xrightarrow{d} N(\mu, \sigma^2/n)$. **Notes:**

- We assume only finite μ and σ^2 and end up to normality.
- The rate of convergence is affected by the original distribution of X_i
 - The closer the distribution of X is the Normal curve, the faster the convergence is.
 - The rate needs to be evaluated case by case.
- Provides an all purpose approximation of the distribution of sums of i.i.d. random variables.

Proof: See Casella & Berger pages 237-238.

Illustration: See R-script `unifmean.R`. Script takes $M = 10000$ samples using each of the following sample sizes $n = 1, 2, 3, 4, 5, 6$ from a Uniform(0,10) population, and demonstrates the distribution of \bar{X} using a histogram and compares it to $N(\mu, \sigma^2/n)$ distribution. Expected value and variance of Uniform(0,10) population are

$$\mu = \frac{b-a}{2} = \frac{10}{2} = 5$$

$$\sigma^2 = \frac{1}{12}(b-a)^2 = \frac{(10-0)^2}{12} = \frac{100}{12} = \frac{25}{3}$$

The visual evaluation of histogram gives an impression that the approximation is good already when $n > 4$. However, we are usually interested in the behaviour in the tails,

which needs more careful evaluation, e.g. compute the 0.95th quartile of the Normal distribution. Check which proportion of the simulated samples had the mean above this quartile. The proportion should be close to 0.05 if the approximation is good.

Theorem 2.10. (Slutsky) If $X_n \xrightarrow{d} X$ and if $Y_n \xrightarrow{P} a$, where a is a constant, then

- a) $Y_n X_n \xrightarrow{d} aX$
- b) $X_n + Y_n \xrightarrow{d} X + a$

Note Proof for this is not shown.

Example 2.10. (Normal approximation with estimated variance)

The central limit theorem said that under mild conditions

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

However, σ^2 is often unknown, but can be unbiasedly and consistently estimated by $s_n^2 = \frac{1}{n-1} \sum (X_i - \bar{X}_n)^2$. It can also be shown that (see Casella & Berger example 5.5.3 and exercise 3.32)

$$\frac{\sigma}{s_n} \xrightarrow{P} 1$$

Note

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{s_n} = \underbrace{\frac{\sigma^2}{s_n}}_{\frac{\sigma}{s_n} \xrightarrow{P} 1} \underbrace{\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma^2}}_{\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0,1)} \xrightarrow{d} N(0, 1)$$

So even though the variance is estimated, $\frac{\bar{X} - \mu}{s/\sqrt{n}}$ converges to a normal distribution. However, intuitively, the convergence should be more slow than when σ is known.

Chapter 3

Estimation

In estimation, we want to use sample X_1, \dots, X_n to make inference on an unknown population parameter θ . It means that we want to find a statistic $T(\mathbf{X})$ that somewhat optimally captures the information of the sample on the parameter of interest, θ .

Most commonly, θ is the population mean, but it can also be related to population variance, covariance, median, maximum etc. Very often the population parameters, especially the mean, is not a scalar number, but it's function of some known characteristics of sampled units i , and therefore θ includes parameters of that function.

Example 3.1. Simple linear regression

Consider case of simple linear regression, where X is a fixed known predictor, and Y is a random variable we want to predict with values of X . The model can be of following for

$$Y_i = \mu(X_i) + e_i$$

Where $\mu(X_i)$ is a fixed mean of Y

$$\mu_i = \mu(X_i) = \beta_0 + \beta_1 X_i$$

and e_i is the random error related to value of Y_i

$$e_i \sim N(0, \sigma^2)$$

This would mean that distribution of Y_i would be

$$Y_i \sim N(\mu(x_i), \sigma^2)$$

In this case, we are interested in estimating the the regression coefficients

$$\boldsymbol{\theta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

Example 3.2. Lets consider a random variable

$$Y_i = \begin{cases} 1 & \text{if tree } i \text{ is dead} \\ 0 & \text{if tree } i \text{ is alive} \end{cases}$$

Parameter of interest might be the probability of a tree being dead.

$$Y_i \sim \text{Bernoulli}(p), \theta = p$$

If we also know the age of the tree, and it is justified to assume that older trees are more commonly dead than younger trees, then

$$Y_i \sim \text{Bernoulli}(p_i)$$

$$p_i = \beta_0 + \beta_1 \text{Age}_i, \rightarrow \boldsymbol{\theta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

Because $\beta_0 + \beta_1 \text{Age}_i$ is not restricted to $[0, 1]$, a better model uses some function for $p(\text{Age}_i)$. Most commonly we use

$$\log \left(\frac{p_i}{1 - p_i} \right) = \beta_0 + \beta_1 \text{Age}_i$$

Which is a logit transformation $[0, 1] \rightarrow \mathbb{R}$. This provides us with logistic regression.

Note: See also example 1.2 and 1.3 in `notes.pdf`

3.1 Two principles

A sufficient statistic (tyhjentävä statistiikka) for a population parameter θ is such statistic $T(\mathbf{X})$, that in certain sense it captures all the information about θ contained in the sample.

3.2 The sufficient principle

If $T(\mathbf{X})$ is a sufficient statistic for θ , then any inference about θ should depend on \mathbf{X} only through the value of $T(\mathbf{X})$.

Definition 3.1. A statistic $T(\mathbf{X})$ is sufficient statistic for θ if the conditional distribution of the sample given $T(\mathbf{X})$ [$\mathbf{X}|T(\mathbf{X})$] does not depend on θ .

Example 3.3. Let X_1, \dots, X_n be i.i.d. from $N(\mu, \sigma^2)$ population, where σ^2 is known. Now the sample mean \bar{X} is a sufficient statistic for μ .

That is, if we have the complete data and you know only the sample mean, we both still have all information on μ the sample includes. However, if a third person knows

only $\min(\mathbf{X})$, $\max(\mathbf{X})$ and $\text{median}(\mathbf{X})$, he still does not have much information on μ as you have when you know \bar{X} .

If also σ^2 is unknown, then $T(\mathbf{X}) = \begin{bmatrix} \bar{X} \\ s^2 \end{bmatrix}$ is sufficient for $\theta = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}$

For more formal discussion on the topic, see Casella & Berger section 6.2.

3.3 The likelihood function

Definition 3.2. Let $f(\mathbf{x}, \theta)$ denote the joint pdf or pmf of the sample $\mathbf{X} = X_1, \dots, X_n$. Then Given that $\mathbf{X} = \mathbf{x}$ has been observed, the function of θ defined by

$$L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$$

is called the likelihood function (uskottavuusfunktio).

Note We do not assume X_1, \dots, X_n to be independent.

- If \mathbf{X} is a discrete random vector, then

$$L(\theta|\mathbf{X}) = P_\theta(\mathbf{X} = \mathbf{x})$$

- If we compare the likelihood at two parameter points $\theta = \theta_1$ and $\theta = \theta_2$, and find out that

$$P_{\theta_1}(\mathbf{X} = \mathbf{x}) = L(\theta_1|\mathbf{x}) > L(\theta_2|\mathbf{x}) = P_{\theta_2}(\mathbf{X} = \mathbf{x})$$

then the sample we observed is more likely to have occurred when $\theta = \theta_1$ than when $\theta = \theta_2$.

→ θ_1 is more plausible (uskottava) value for θ than θ_2 .

- If X is a continuous, real valued random variable, then for a small value ϵ

$$P_\theta(X - \epsilon < X < X + \epsilon) \approx 2\epsilon f(X, \theta) = 2\epsilon L(\theta|X)$$

therefore

$$LR = \frac{P_{\theta_1}(X - \epsilon < X < X + \epsilon)}{P_{\theta_2}(X - \epsilon < X < X + \epsilon)} \approx \frac{L(\theta_1|X)}{L(\theta_2|X)}$$

- If the $LR > 1$, then the sample was more likely to occur when $\theta = \theta_1$, and if $LR < 1$, then the sample was more likely to occur when $\theta = \theta_2$.

Example 3.4. Let X have negative binomial distribution. The Negative binomial distribution specifies the pmf for the number of successes before a fixed number of a failures (r) in a series of Bernoulli trials with success probability p . The pmf

$$P(X = x|r, p) = \binom{r+x-1}{x} p^r (1-p)^x, \text{ where } x = 0, 1, \dots \text{ and } 0 \leq p \leq 1$$

Consider an experiment where $r = 3$ and it was observed that $X = 2$. The likelihood as a function of $\theta = p$

$$L(p|2) = P_p(X = 2) = \binom{4}{2} p^3 (1-p)^2$$

→ The likelihood is a 5th order polynomial with respect to p .

In general, if $X = x$ has been observed,

$$L(p|x) = P_p(X = x) = \binom{3+x-1}{2} p^3 (1-p)^x$$

in which case the likelihood is a polynomial of order $x + 3$.

Example 3.5. Let $\mathbf{X}_{3 \times 1}$ have 3-variate normal distribution so that all components have a common, unknown mean of μ and a known, common variance $\sigma^2 = 3.69$. In addition, $\text{cov}(X_1, X_2) = 2.25$ and $\text{cov}(X_1, X_3) = \text{cov}(X_2, X_3) = 0$.

This kind of situation could stem from a grouped structure of the data: X_1 and X_2 might belong to the same group (e.g. school class) and X_3 may originate from another group. That is

$$\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ where } \boldsymbol{\mu} = \begin{bmatrix} \mu \\ \mu \\ \mu \end{bmatrix} \text{ and } \boldsymbol{\Sigma} = \begin{bmatrix} 3.69 & 2.25 & 0 \\ 2.25 & 3.69 & 0 \\ 0 & 0 & 3.69 \end{bmatrix}$$

The likelihood of this model for observation $\mathbf{x} = \begin{bmatrix} 6 \\ 5 \\ 4 \end{bmatrix}$ has been illustrated in R-script `MVNlikelihood.R`.

Note The likelihood function is not a pdf as a function of θ . Therefore we say, that θ_1 is more plausible than θ_2 (not more probable). Note also, that θ is thought to be fixed, therefore it makes no sense to specify a pdf for it. (In Bayesian statistics, θ is thought to be random.)

3.4 The likelihood principle

If \mathbf{x} and \mathbf{y} are two observed sample points such that $L(\theta|\mathbf{x})$ is proportional to $L(\theta|\mathbf{y})$, that is there exists a constant $C(\mathbf{x}, \mathbf{y})$ such that

$$L(\theta|\mathbf{x}) = C(\mathbf{x}, \mathbf{y}) L(\theta|\mathbf{y}) \text{ for all } \theta$$

then the conclusions drawn from \mathbf{x} and \mathbf{y} should be identical. **Note** $C(\mathbf{x}, \mathbf{y})$ may be different for different (\mathbf{x}, \mathbf{y}) pairs, but does not depend on θ .

3.5 Point estimation

Definition 3.3. A point estimator is any function $W(X_1, \dots, X_n)$. That is, any statistic is a point estimator.

- We now first consider the methods to find estimators and there after consider criteria to evaluate, whether the estimators are good or bad.

Note Estimator is a function of the sample X_1, \dots, X_n and estimate is the numerical value the estimator got after the sample was actually taken.

→ Estimator is a random variable, and estimate is its realized value in the sample.

- Often, there is a natural, intuitive candidate for a point estimator, e.g. \bar{X} is a intuitive choice for a point estimator of μ .

3.6 Estimating the population mean using least squares

Based on theorem 2.1, the ordinary least squares estimator of parameter μ using observed sample values y_1, \dots, y_n is found by

$$\min_{\theta} \sum_{i=1}^n (y_i - \theta)^2$$

The solution is $\theta = \bar{y}$. The approach also generalized to cases where $E(Y_i)$ is a function of a fixed variable x_i . Consider a model $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, where Y_i is a random variable and x_i is another fixed observable variable. We can β_0 such that $E(\epsilon_i) = 0$ and

$$E(Y_i) = E(\underbrace{\beta_0 + \beta_1 x_i}_{\text{fixed}} + \underbrace{\epsilon_i}_{E(\epsilon_i)=0}) = \beta_0 + \beta_1 x_i = E(Y_i | x_i) = \mu(x_i, \beta)$$

Parameters β are “second order” parameters. They quantify the (assumed) linear relationship between x_i and $E(Y_i)$. Using the idea of theorem 2.1, β can be estimated using

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$

The solutions are

$$\widehat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \text{ and } \widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x}$$

- A more general formulation results by defining

$$\mathbf{X} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}, \mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}, \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix},$$

and then the model is $\mathbf{Y} = \mathbf{X}'\boldsymbol{\beta} + \boldsymbol{\epsilon}$

- The least squares estimator of $\boldsymbol{\beta}$ is found by minimizing

$$\boldsymbol{\epsilon}'\boldsymbol{\epsilon} = (\mathbf{y} - \mathbf{X}'\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}'\boldsymbol{\beta})$$

- The solution is $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$
- This is a general solution for the multiple regression case as well.

The least squares method is optimal and intuitive way to find the estimators for population mean. However, it does not provide tools to estimate other parameters of distribution. For that purpose, we will next present three other approaches for parameter estimation in general case:

- The method of moments
- Maximum likelihood
- Bayesian approach

3.7 The method of moments

- This method dates back to late 1800's.

Let X_1, \dots, X_n be a sample from population with pmf or pdf $f(x|\theta_1, \dots, \theta_k)$. The method of moments (momenttimenetelmä), moment matching, finds estimates of $\theta_1, \dots, \theta_k$ by equating the first k sample moments to the corresponding population moments.

Define

$$\begin{aligned} m_1 &= \frac{1}{n} \sum_{i=1}^n x_i & \mu'_1 &= E(X^1) \\ m_2 &= \frac{1}{n} \sum_{i=1}^n x_i^2 & \mu'_2 &= E(X^2) \\ \vdots & & \vdots & \\ m_k &= \frac{1}{n} \sum_{i=1}^n x_i^k & \mu'_k &= E(X^k) \end{aligned}$$

- The population moments will typically be functions of $\boldsymbol{\theta}$, $\mu'_j(\theta_1, \dots, \theta_k)$

- The moment estimator $(\tilde{\theta}_1, \dots, \tilde{\theta}_k)$ of $(\theta_1, \dots, \theta_k)$ is found by solving the following system of equations for $(\theta_1, \dots, \theta_k)$

$$\begin{cases} m_1 = \mu'_1(\boldsymbol{\theta}) \\ m_2 = \mu'_2(\boldsymbol{\theta}) \\ \vdots \\ m_k = \mu'_k(\boldsymbol{\theta}) \end{cases}$$

Example 3.6. Let X_1, \dots, X_n be i.i.d. with $X_i \sim N(\mu, \sigma^2)$, and both μ and σ^2 are unknown.

$$\begin{aligned} m_1 &= \bar{X} & m_2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 \\ \mu'_1 &= \mu & \mu'_2 &= E(X^2) = \mu^2 + \sigma^2 \end{aligned}$$

Recall $Var(X) = E(X^2) - (E(X))^2 \rightarrow E(X^2) = Var(X) + (E(X))^2 = \sigma^2 + \mu^2$
Following system of equation needs to be solved for μ and σ^2

$$\begin{cases} \bar{X} = \mu \\ \frac{1}{n} \sum X^2 = \mu^2 + \sigma^2 \end{cases}$$

$$\begin{cases} \bar{X} - \mu = 0 \\ \frac{1}{n} \sum X^2 - \mu^2 - \sigma^2 = 0 \end{cases}$$

We get

$$\begin{aligned} \tilde{\mu} &= \bar{X} \\ \tilde{\sigma}^2 &= \underbrace{\frac{1}{n} \sum (X_i - \bar{X})}_{\text{Biased estimator}} \neq \underbrace{\frac{1}{n-1} \sum (X_i - \bar{X})^2}_{\text{Unbiased estimator}} \end{aligned}$$

But the bias vanishes as $n \rightarrow \infty$.

This example was easy, since the system of equations we needed to solve was linear.

Example 3.7. A numerical solution Let X_1, \dots, X_n be i.i.d. from the Weibull(a, b) distribution, which has the pdf

$$f(x) = \frac{a}{b} \left(\frac{x}{b} \right)^{a-1} e^{-(x/b)^a} \quad \text{where } 0 < x < \infty \text{ and } a, b > 0$$

where a is called the shape parameter and b is called the scale parameter.

The first two moments are

$$\begin{aligned} E(X) &= \int_0^\infty f(x)x \, dx = \dots = b\Gamma\left(1 + \frac{1}{a}\right) \\ E(X^2) &= \int_0^\infty f(x)x^2 \, dx = \dots = b^2\Gamma\left(1 + \frac{2}{a}\right) \end{aligned}$$

where $\Gamma(z) = \int_0^\infty z^{u-1} e^{-z} dz$ cannot be written in closed form, but can be well approximated using function gamma in R.

The moment-matching system of equations is

$$\begin{cases} \frac{1}{n} \sum x_i - b\Gamma\left(1 + \frac{1}{a}\right) = 0 \\ \frac{1}{n} \sum X^2 - b^2\Gamma\left(1 + \frac{2}{a}\right) = 0 \end{cases}$$

This is nonlinear with respect to a !

We solve it for a and b using the Gauss-Newton method, which is a multivariate generalization of the Newton-Raphson method. Example in R-script `WeibullMoment.R` generates first a random sample of size 30 from `Weibull(2,10)` population. Using the generalized sample, it finds the moment estimates of a , b numerically using R-function `NRnum`.

The observed sample moments are

$$m_1 = \bar{X} = 9.604, \quad m_2 = \frac{1}{n} \sum x^2 = 113.22$$

We solve

$$\begin{cases} 9.604 - b\Gamma\left(1 + \frac{1}{a}\right) = 0 \\ 113.22 - b^2\Gamma\left(1 + \frac{2}{a}\right) = 0 \end{cases}$$

for (a, b) ($\theta = [a \ b]^T$). We use $a' = 1$ and $b' = 9.60$ as starting values, and we get

$$\tilde{a} = 2.21 \quad \tilde{b} = 10.84$$

3.8 Newton-Raphson

Consider function $f(a) = 0$, which needs to be solved for a .

1. Start with some guess a' , compute $f(a')$ and $f'(a')$ to approximate $f(a)$ by a linear function. Find a such that the approximator = 0.
2. Set the solution as a' and repeat step 1.

If the initial guess a' was good enough and f is continuous and differentiable, we will find the approximate solution within few iteration.

Gauss-Newton generalized this to vector valued \mathbf{a} and \mathbf{x} . The algorithm has been implemented in R-function `NRnum` of package `lmfor`. The function evaluates also the required derivatives $\frac{df(\mathbf{x})}{d\theta_i}$ numerically.

3.9 Maximum likelihood

Definition 3.2 defined the likelihood as

$$L(\boldsymbol{\theta}|\mathbf{x}) = f(\mathbf{x}|\boldsymbol{\theta})$$

where

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_k \end{bmatrix}$$

Note For i.i.d. data X_1, \dots, X_n from population $f(x|\theta)$, the likelihood simplifies to

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

Definition 3.4. For each sample point \mathbf{x} , let $\hat{\boldsymbol{\theta}}(\mathbf{x})$ be a parameter value at which $L(\boldsymbol{\theta}|\mathbf{x})$ attains its maximum value as a function of $\boldsymbol{\theta}$, with \mathbf{x} fixed. A maximum likelihood estimator of $\boldsymbol{\theta}$ based on sample \mathbf{x} is this function $\hat{\boldsymbol{\theta}}(\mathbf{x})$. The maximum likelihood estimate (ML estimate) is the numerical value of that function in the data point we have.

Note We use abbreviation MLE for both ML estimator and ML estimate.

- MLE is the parameter point for which the observed sample is most likely.
- The MLE also fulfills some optimality criteria, which will be discussed later.
- How to find MLE?
 - If $L(\boldsymbol{\theta}|\mathbf{x})$ is differentiable with respect to all θ_i , then the candidates for MLE are the values that solve

$$\frac{d}{d\theta_i} L(\boldsymbol{\theta}|\mathbf{x}) = 0, \quad i = 1, \dots, k$$

- These together with the boundaries of the parameter space provide candidates of MLE.
- Of the candidates the one that gives the maximum value of $L(\boldsymbol{\theta}|\mathbf{x})$ is the MLE.

Example 3.8. Example ML1

Let X_1, \dots, X_n be i.i.d. from $N(\theta, 1)$, meaning only the population mean θ is unknown.

$$\begin{aligned} L(\theta, \mathbf{x}) &= \prod_{i=1}^n \frac{1}{(2\pi)^{1/2}} e^{-\frac{1}{2}(x_i - \theta)^2} \\ &= \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2} \\ \frac{d}{d\theta} L(\theta|\mathbf{x}) &= \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \theta)^2} \left(-\frac{1}{2}\right) \sum_{i=1}^n (-2(x_i - \theta)) \\ &= L(\theta|\mathbf{x}) \sum_{i=1}^n (x_i - \theta) \end{aligned}$$

$$\frac{d}{d\theta} L(\theta|\mathbf{x}) = 0 \text{ if}$$

$$\text{either } L(\theta|\mathbf{x}) = 0 \quad \text{A}$$

$$\text{or } \sum_{i=1}^n (x_i - \theta) = 0 \quad \text{B}$$

But if $L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta) > 0$

→ There is no solution for A.

B gives

$$\sum (x_i - \theta) = 0$$

$$\sum x_i - \sum \theta = 0$$

$$\sum x_i - n\theta = 0 \iff \hat{\theta} = \frac{\sum x_i}{n} = \bar{x}$$

So the sample mean, which is the least squares and moment estimator for the population mean is also the ML estimator for the population mean.

To confirm that \bar{X} is the MLE, we still need to make sure that

- 1) \bar{x} is the maximum, not minimum
- 2) The boundaries of the support of θ do not provide higher value of θ

We can conclude that

- 1) This is maximum because

$$\frac{d^2}{d\theta^2} L(\theta|\mathbf{x})_{\theta=\bar{x}} < 0$$

- 2) We can notice that

$$\lim_{\theta \rightarrow \infty} L(\theta|\mathbf{x}) = \lim_{\theta \rightarrow -\infty} L(\theta|\mathbf{x}) = 0$$

and $L(\bar{x}|\mathbf{x}) > 0$.

- In most cases, the solution of the estimation problem becomes easier if the natural logarithm of the likelihood is used instead of likelihood.

- The maximum of $L(\theta|\mathbf{x})$ and $\ln(L(\theta|\mathbf{x}))$ coincide (are at the same point θ) because \ln -function is increasing on $(0, \infty)$.

Example 3.9. Example ML2 (Bernoulli MLE)

Let X_1, \dots, X_n Bernoulli(p), i.i.d.

$$L(p|\mathbf{x}) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = (*)$$

Recall that

$$\text{PMF} : f(x_i|p) = \begin{cases} p & \text{if } x_i = 1 \text{ (success)} \\ 1-p & \text{if } x_i = 0 \text{ (failure)} \end{cases}$$

$$(*) = p^{\sum x_i} (1-p)^{\sum (1-x_i)} = p^{\sum x_i} (1-p)^{n-\sum x_i}$$

The log-likelihood is

$$\ln(L(p|\mathbf{x})) = \ln(p^{\sum x_i} (1-p)^{n-\sum x_i}) = \sum x_i \ln(p) + (n - \sum x_i) \ln(1-p)$$

$$\frac{d}{dp} \ln(L) = \sum x_i \frac{1}{p} + (n - \sum x_i) \frac{-1}{1-p} = \dots = \frac{\sum x_i - np}{p(1-p)}$$

$$\frac{d}{dp} \ln(L) = 0$$

$$\frac{\sum x_i - np}{p(1-p)} = 0 \quad || \cdot p(1-p) \neq 0$$

$$\sum x_i - np = 0$$

$$p = \frac{\sum x_i}{n} = \bar{x}$$

Check that this is a global maximum and $\ln(L(\bar{x})) > \ln(L(0))$ and $\ln(L(\bar{x})) > \ln(L(1))$ to confirm that the solution is also the MLE.

Example 3.10. Example ML3 (MLE in dependent sample)

Recall the grouped data example 3.5, where $Y \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Assume that Y has been observed to be $\mathbf{y} = \begin{bmatrix} 6 \\ 5 \\ 4 \end{bmatrix}$ ($\bar{y} = 5$). Assume (unrealistically) that $\boldsymbol{\Sigma}$ is known to be

$$\boldsymbol{\Sigma} = \begin{bmatrix} 3.69 & 2.25 & 0 \\ 2.25 & 3.69 & 0 \\ 0 & 0 & 3.69 \end{bmatrix}$$

We want to estimate parameter θ of the unknown vector $\boldsymbol{\mu} = \begin{bmatrix} \theta \\ \theta \\ \theta \end{bmatrix}$ (i.e. we assume that all three components of $\boldsymbol{\mu}$ have the common mean of θ). The likelihood is directly

the pdf of the multivariate normal distribution and the log-likelihood is the natural logarithm of that.

R-script `MVNMLE.R` does the estimation. Define the likelihood as we did previously in the R-script `MVNLikelihood.R`. Find the value of θ that maximizes the likelihood.

1. Simple solution (for demonstration). Compute $L(\theta)$ for values $\theta = (0, 0.001, 0.002, \dots, 10)$ and select the one that provides the highest value. This direct search or grid search algorithm have the estimate $\hat{\theta} = 4.831 < \bar{y} = 5$, because the independent $Y_3 = 4$ gets more weight in estimation than the two mutually dependent elements $\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 5 \end{bmatrix}$. Therefore $\hat{\theta}$ is shrunken towards $Y_1 = 4$ compared to $\bar{y} = 5$.

2. We solved this also using function `mle` of R-package `stats4`.

1. Define a R-function that evaluates the negative log-likelihood $-\ln(L)$; call it e.g. `nll`.
2. Specify some good starting values of the parameters in a named list.
3. Make a call to `mle`: `mle(nll, start=list(theta=5))` → gives $\hat{\theta} = 4.831081$.

Note To get more stable and faster solution, you could also provide analytical derivative of the function you used

Note `mle` gives only the candidate of MLE, so you are responsible for checking that it really is the MLE.

- If $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ is multidimensional, then we need to maximize the $\ln(L)$ with respect to all components θ_i simultaneously.
- If the log-likelihood is differentiable with respect to all θ_i , we must need to find the practical derivatives of $\ln(L)$ (or L) with respect to all θ_i and set them to 0
→ We get system of k equations.

Example 3.11. Example ML4 Let X_1, X_2, \dots, X_n be i.i.d. form $N(\theta, \sigma^2)$ population, and let both θ and σ^2 be unknown. Now

$$L(\theta, \sigma^2 | \mathbf{x}) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{-\frac{1}{2\sigma^2} \sum (x_i - \theta)^2}$$

$$\ln(L(\theta, \sigma^2 | \mathbf{x})) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2} \frac{\sum (x_i - \theta)^2}{\sigma^2}$$

The ML equations becomes

$$\begin{cases} \frac{d}{d\theta} \ln(L(\theta, \sigma^2 | \mathbf{x})) = \frac{1}{\sigma^2} \sum (x_i - \theta) \\ \frac{d}{d\sigma^2} \ln(L(\theta, \sigma^2 | \mathbf{x})) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (x_i - \theta)^2 \end{cases}$$

Setting these to 0 and solving for $\theta = \begin{bmatrix} \theta \\ \sigma^2 \end{bmatrix}$ gives

$$\hat{\theta} = \bar{x} \text{ and } \sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

These are the candidates for the MLE, but it can be shown that they are also the MLE's. See Casella & Berger p.322 for more formal argumentation.

The so called invariance property of MLE's is highly useful in making restrictions to the parameters. Suppose that the distribution is indexed using θ , but we want to estimated a function of it, $\tau(\theta)$.

Theorem 3.1. If $\hat{\theta}$ is the MLE of θ , and $\tau(\theta)$ is any function of θ , then $\tau(\hat{\theta})$ is the MLE of $\tau(\theta)$.

Proof: See Casella & Berger, p.320.

Example 3.12. Example ML5

Recall the Weibull population from the example 3.7 for which we estimated the parameters a and b using the method of moments. Because a and b are both positive, it is convenient to reparametrize the pdf in terms of $\ln(a)$ and $\ln(b)$.

Weibull MLE: $x = (14.7, \dots, 17.56)$

The MLE's of $\theta = (a, b)$ is

Table 3.1: The MLE's of $\theta = (a, b)$.

	In the original scale	$\exp(\widehat{\ln(\theta)})$
\hat{a}	2.24	2.24
\hat{b}	10.88	10.88

In R-script `WeibullMLEInvariance.R` $A = \ln(a)$ and $B = \ln(b)$.

$$\begin{aligned} \widehat{\ln(a)} = \hat{A} &= 0.81 & \exp(\hat{A}) &= 2.243 \\ \widehat{\ln(b)} = \hat{B} &= 2.89 & \exp(\hat{B}) &= 10.876 \end{aligned}$$

Example 3.13. Example ML6 (MLE in dependent data)

Recall exercise 2.2 where we had spatial, multivariate normal data. Consider observations of Y_i in the following 2-dimensional grid, $i = 1, \dots, 25$. $\mathbf{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_{25} \end{bmatrix}$ is generated by following model

$$Y \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \text{ where } \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_{25} \end{bmatrix}, \mu_i = 0 \text{ for all } i, \text{ and } Var(Y_i) = \sigma_i^2 = 2^2 \text{ for all } i$$

$$\text{corr}(Y_i, Y_j) = \exp^{-\phi s_{ij}}, \text{ where } \phi = \frac{1}{2}$$

and s_{ij} is the euclidian distance between the location i and j , $s_{ij} = ||u_i - u_j||$.

Let us first generate the data and illustrate is graphically; see `R-script MLESpatial.R`. Next, given the simulated data, let us estimate the parameters μ (common mean for all i), σ^2 (common variance for all i) and ϕ (the parameter of correlation function).

Because $\sigma^2 > 0$ and $\phi > 0$, we parametrize the likelihood function in terms of $\ln(\sigma^2)$ and $\ln(\phi)$, and define the negative log-likelihood in R-function `nll`. The script provides following resulting estimates.

$$\begin{array}{ll} \hat{\mu} = -1.87 & \text{True: } 0 \\ \widehat{\ln(\phi)} = -2.03 \rightarrow \hat{\phi} = 0.13 & \text{True: } 0.5 \\ \widehat{\ln(\sigma^2)} = 2.35 \rightarrow \hat{\sigma}^2 = 10.47 & \text{True: } 4 \end{array}$$

3.10 Bayes estimators

In the Bayesian approach, the fundamental difference to the previous methods is to think the parameters θ as random variable (vector), not as fixed quantity.

- The variability of θ in the population is described by some distribution called the prior distribution.
 - It is the belief of the experimenter about θ , and is formulated before the data are seen.
 - This prior is updated with the data x to get the posterior distribution of θ
- There are two main reasons to use Bayesian approach:
 1. We may be willing to combine the information we already have on θ with the new information we get from the data.
 2. We do not actually believe on the prior, but we use the Bayesian approach because it provides tools for estimation in the situation where e.g. ML is very tricky.
 - Especially, the Bayesian computational methods provide estimation tools for such cases.
- Let $\pi(\theta)$ be the prior pdf and let $f(x|\theta)$ be the sampling pdf of the data x .

- The posterior distribution is the conditional distribution of θ given the sample \mathbf{x}

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})}$$

where $m(\mathbf{x}) = \int f(\mathbf{x}|\theta)\pi(\theta) d\theta$

Note θ means both the random variable and its realized value here.

Example 3.14. Bayes 1 (Binomial Bayes)

Let X_1, \dots, X_n be i.i.d. Bernoulli(p). Then

$$\sum X_i \sim \text{Binomial}(n, p)$$

Assume that prior of the p is the $\text{Beta}(\alpha, \beta)$ distribution with some values of α and β , which are set before the data are seen. The joint pdf of (y, p) is

$$\begin{aligned} f(y, p) &= \underbrace{f(y|p)}_{\text{Binomial}} \underbrace{\pi(p)}_{\text{Beta}} = \left[\binom{n}{y} p^y (1-p)^{n-y} \right] \left[\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \right] \\ &= \binom{n}{y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{y+\alpha-1} (1-p)^{n-y+\beta-1} \end{aligned}$$

The marginal pdf of Y is

$$f(y) = \int_0^1 f(y, p) dp = \binom{n}{y} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(y+\alpha)\Gamma(n-y+\beta)}{\Gamma(n+\alpha+\beta)}$$

which is called the beta-binomial pdf. The posterior distribution is

$$f(p|y) = \frac{f(y, p)}{f(y)} = \frac{\Gamma(n+\alpha+\beta)}{\Gamma(y+\alpha)\Gamma(n-y+\beta)} p^{y+\alpha-1} (1-p)^{n-y+\beta-1}$$

which is the $\text{Beta}(y+\alpha, n-y+\beta)$.

A natural point estimator of p is the expected value $E(p|\mathbf{x})$ of posterior distribution:

$$\hat{p}_\beta = \frac{y + \alpha}{\alpha + \beta + n}$$

It is a weighted mean of the prior mean $\frac{\alpha}{\alpha+\beta}$ and the ML-estimator $\frac{y}{n}$:

$$\frac{n}{\alpha + \beta + n} \frac{y}{n} + \frac{\alpha + \beta}{\alpha + \beta + n} \frac{\alpha}{\alpha + \beta} = w \frac{y}{n} + (1-w) \frac{\alpha}{\alpha + \beta}, \text{ where } w = \frac{n}{\alpha + \beta + n} \in [0, 1]$$

Note The Beta distribution is the conjugate prior for the Binomial distribution, because the posterior distribution belongs to the same family. The use of conjugate priors leads to mathematically beautiful results, but are not otherwise justified by any theory.

- In the conditional distribution

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})}$$

the denominator $m(\mathbf{x}) = \int f(\mathbf{x}|\theta)\pi(\theta) d\theta$ is hard to compute in general case.

- The modern computational methods, so called Markov chain Monte Carlo (MCMC) methods allow simulation of observations from $\pi(\theta|\mathbf{x})$ without the need to evaluate $m(\mathbf{x})$.
- The most simple MCMC method is the Metropolis algorithm. (Denote $f(\mathbf{x}|\theta)\pi(\theta)$ by $a(\theta)$)
 1. Choose an arbitrary starting point θ_0 from Θ (the sample space of θ). Choose an arbitrary symmetric pdf (the proposal density) $g(\theta_{t+1}|\theta_t)$, which suggests the new value of θ (θ_{t+1}) given the previous value θ_t . Usually we use $N(\theta_t|\sigma^2)$ as the proposal density.
 2. At each iteration t
 - 2.1 Generate a candidate of θ : θ' using $g(\theta_{t+1}|\theta_t)$.
 - 2.2 Compute the acceptance ratio

$$\alpha = \frac{a(\theta')}{a(\theta_t)}$$
 - 2.3 If $\alpha < 1$, then accept θ' as θ_{t+1} with probability *alpha*.

Repeat until a desired sample size has been obtained.

Example 3.15. Bayes 2

Consider a sample $\mathbf{x}' = [9.76 \ 9.08 \ 16.23 \ 10.28 \ 10.51]$. The aim is to estimate the population mean θ . We use restrictive prior for it:

$$\pi(\theta) \sim Uniform(7, 13)$$

In R-script `Bayes.R` there is function `postProfile` that implements function $a(\theta)$. $N(\theta_t, 1)$ is used as $g(\theta_{t+1}|\theta_t)$. We generate a sample size of $n = 10000$ and use the posterior sample mean as the point estimate of θ .

Notes

- MCMC methods require only that $\pi(\theta)f(\mathbf{x}|\theta)$ can be evaluated.
- The sample generated from the posterior is not independent: the sequence $\theta_1, \dots, \theta_t, \dots, \theta_N$ has autocorrelation so that successive values are dependent.
- Commonly used MCMC methods are the Metropolis-Hastings-algorithm and the Gibbs sampler.

3.11 Methods for evaluating estimators

Definition 3.5. The mean squared error (MSE) of an estimator W of parameter θ is the function of θ defined by

$$E_{\theta}(W - \theta)^2$$

- In general, any increasing function of the absolute distance $|W - \theta|$ is a reasonable measure for the quality of W (e.g. mean absolute error $E_{\theta}(W - \theta)$), but MSE is 1) mathematically factorable and 2) has the interpretation

$$E_{\theta}(W - \theta)^2 = \text{Var}_{\theta}(W) + (E_{\theta}(W - \theta))^2$$

$$MSE = \text{variance} + \text{bias}^2$$

Definition 3.6. The bias of the point estimator W is the difference

$$\text{Bias}_{\theta}(W) = E_{\theta}(W) - \theta$$

- Estimator W is said to be unbiased if $\text{Bias}_{\theta}(W) = 0$ for all θ
- MSE as a criterion for an estimator leads to estimators that have small bias and variance.

Example 3.16. (Normal MSE)

Let X_1, \dots, X_n be i.i.d. $N(\mu, \sigma^2)$. The statistics \bar{x} and s^2 are unbiased estimators of μ and σ^2 since

$$E(\bar{x}) - \mu = \mu - \mu = 0$$

$$E(s^2) = \sigma^2 \rightarrow E(s^2) - \sigma^2 = 0$$

Note These results are true even if the population is not normally distributed.

The MSE of these estimators are

$$(*) \text{MSE}(\bar{X}) = E(\bar{X} - \mu)^2 = \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

$$(**) \text{MSE}(s^2) = E(s^2 - \sigma^2)^2 = \text{Var}(s^2) = \dots = \frac{2\sigma^4}{n-1}$$

Note that (*) did not require normality, but (**) did. An alternative estimator for σ^2 is the MLE/moment estimator

$$\hat{\sigma}^2 = \frac{1}{n}(X_i - \bar{X})^2 = \frac{n-1}{n}s^2$$

$$E(\hat{\sigma}^2) = E\left(\frac{n-1}{n}s^2\right) = \frac{n-1}{n}E(s^2) = \frac{n-1}{n}\sigma^2$$

Now $E(\hat{\sigma}^2) - \sigma^2 = \frac{n-1}{n}\sigma^2 - \sigma^2 \neq 0$, meaning $\hat{\sigma}^2$ is a biased estimator of σ^2 .

The variance of $\hat{\sigma}^2$ is

$$Var(\hat{\sigma}^2) = Var\left(\frac{n-1}{n}s^2\right) = \left(\frac{n-1}{n}\right)^2 Var(s^2) = \frac{2(n-1)}{n^2}\sigma^4$$

The MSE of $\hat{\sigma}^2$ is

$$E(\hat{\sigma}^2 - \sigma^2)^2 = variance + bias^2 = \frac{2(n-1)}{n^2}\sigma^4 + \left(\frac{n-1}{n}\sigma^2 - \sigma^2\right)^2 = \frac{2n-1}{n^2}\sigma^2$$

$$\underbrace{E(\hat{\sigma}^2 - \sigma^2)^2}_{MSE(\hat{\sigma}^2)} = \frac{2n-1}{n^2}\sigma^2 < \frac{2\sigma^4}{n-1} = \underbrace{E(s^2 - \sigma^2)^2}_{MSE(s^2)}$$

$$\frac{2n-1}{n^2} = \frac{2 - \frac{1}{n}}{n} < \frac{2}{n-1}$$

→ Allowing bias in the estimator of σ^2 leads to gain in terms of MSE.

- The bias is negative → ML underestimates the variance on average.
- In general, MSE is justified choice to evaluate estimators of location parameters, but not the scale parameters.

It is often impossible to find a single estimator that is best in terms of MSE. In addition, the bias of an estimator is often seen as a severe problem, or inconvenience. Therefore, we often restrict consideration in the class of unbiased estimators. In that class

$$E_\theta(E - \theta)^2 = Var_\theta(W)$$

Therefore, form among unbiased estimators, we should select the one with the lowest variance.

Definition 3.7. An estimator W^* is the best unbiased estimator (BUE) of $\tau(\theta)$, if it satisfies $E_\theta(W^*) = \tau(\theta)$ for all θ and for any other estimator W with $E_\theta(W) = \tau(\theta)$ we have

$$Var_\theta(\tau^*) \leq Var_\theta(W) \text{ for all } \theta$$

W^* is also called uniform minimum variance unbiased estimator of $\tau(\theta)$ (UMVUE).

Example 3.17. Let X_1, \dots, X_n be i.i.d. $Poisson(\lambda)$ and let \bar{X} and s^2 be the sample mean and sample variance respectively. Recall that in $Poisson(\lambda)$ population $E(X) = \lambda$. Therefore, both \bar{X} and s^2 are unbiased estimators of λ . Which one to select?

$$Var_\lambda(\bar{X}) = \frac{\lambda}{n}$$

$$Var_\lambda(s^2) = \dots \text{ long calculation}$$

Conclusion however will be, that $Var_\lambda(\bar{X}) \leq Var_\lambda(s^2)$ for all λ . So, from among these two candidates, \bar{X} should be selected. But is \bar{X} really the best?

We could construct infinitely many unbiased estimators as weighted mean of \bar{X} and s^2 :

$$W_W = w\bar{X} + (1-w)s^2 \text{ where } w \in [0, 1]$$

Could some selection of w lead to smaller variance than $Var(\bar{X}) = \frac{\lambda}{n}$? Or is there some other estimator that has $Var(w) = \frac{\lambda}{n}$?

Theorem 3.2. (Cramer-Rao Inequality)

Let X_1, \dots, X_n be a sample with pdf $f(\mathbf{X}|\theta)$ and let $W(\mathbf{X}) = W(X_1, \dots, X_n)$ be any estimator satisfying

$$(*) \quad \frac{d}{d\theta} E_\theta(W(\mathbf{X})) = \int_{\mathcal{X}} \frac{d}{d\theta} W(\mathbf{x}) f(\mathbf{x}|\theta) d\mathbf{x}$$

and

$$Var_\theta(\mathbf{X}) < \infty$$

, then

$$Var_\theta(W(\mathbf{X})) \geq \frac{\left[\frac{d}{d\theta} E_\theta(W(\mathbf{X})) \right]^2}{E \left\{ \left[\frac{d}{d\theta} \ln(f(\mathbf{x}|\theta)) \right]^2 \right\}}$$

Note If $W(\mathbf{X})$ is unbiased for θ , then

$$\left[\frac{d}{d\theta} E_\theta(W(\mathbf{X})) \right]^2 = 1 \text{ and } Var_\theta(W(\mathbf{X})) \geq \frac{1}{E \left\{ \left[\frac{d}{d\theta} \ln(f(\mathbf{x}|\theta)) \right]^2 \right\}}$$

Note Integral in condition (*) can be replaced by sum to apply the theorem to discrete random variables.

Proof: Based on Cauchy-Schwarz inequality

$$cov(X, Y)^2 \leq Var(X)Var(Y) \quad (1)$$

where $X = W(\mathbf{X})$ and $Y = \frac{d}{d\theta} \ln[f(\mathbf{X}|\theta)]$.

First note that

$$(2) \quad \frac{d}{d\theta} E_\theta(W(\mathbf{X})) = \int_{\mathcal{X}} \frac{d}{d\theta} W(\mathbf{x}) f(\mathbf{x}|\theta) d\mathbf{x} \quad \text{condition } (*)$$

We multiply this by $1 = \frac{f(\mathbf{X}|\theta)}{f(\mathbf{X}|\theta)}$ and utilize $\frac{d \ln(f(x))}{dx} = \frac{1}{f(x)} f'(x) = \frac{\frac{d}{dx} f(x)}{f(x)}$ to get

$$= \int_{\mathcal{X}} W(\mathbf{X}) \frac{\frac{d}{d\theta} f(\mathbf{X}|\theta)}{f(\mathbf{X}|\theta)} f(\mathbf{X}|\theta) d\mathbf{x}$$

Which we can write as an expected value of transformation

$$= E_{\theta} \left[W_{\theta}(\mathbf{X}) \frac{d \ln(f(\mathbf{X}|\theta))}{d\mathbf{X}} \right] \quad (3)$$

If we consider here $X = W_{\theta}(\mathbf{X})$ and $Y = \frac{d \ln(f(\mathbf{X}|\theta))}{d\mathbf{X}}$, we have

$$E(Y) = E \left[\frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right]$$

and if we assume that $X = W_{\theta}(\mathbf{X}) = 1$, and utilize the equation (3), we get

$$= E \left[1 \cdot \frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right]$$

we can write this equation then in the form of equation (2)

$$= \frac{d}{d\theta} E_{\theta}(X) = \frac{d}{d\theta} 1 = 0 \quad (4)$$

Now, $E_{\theta} \left[W_{\theta}(\mathbf{X}) \frac{d \ln(f(\mathbf{X}|\theta))}{d\mathbf{X}} \right] = \text{cov}(W(\mathbf{X}), \frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)))$. Also, because of (4), we also have

$$(5) \text{Var}_{\theta} \left[\frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right] = \underbrace{E \left[\left(\frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right)^2 \right]}_{\text{Var}(Y) = E(Y^2) - [E(Y)]^2 = E(Y^2)} \\ = 0 \text{ because of (4)}$$

Using (5) and (2) in (1) gives

$$\text{Var}_{\theta}(W(\mathbf{X})) \geq \frac{\left[\text{cov} \left(W(\mathbf{X}), \frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right) \right]^2}{\text{Var} \left[\frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right]} \\ = \frac{\left[\frac{d}{d\theta} E(W(\mathbf{X})) \right]^2}{E \left\{ \left[\frac{d}{d\theta} \ln(f(\mathbf{X}|\theta)) \right]^2 \right\}}$$

Corollary 3.1. (Cramer-Rao i.i.d. case)

$$\text{Var}(W(\mathbf{X})) \geq \frac{\left[\frac{d}{d\theta} E(W(\mathbf{X})) \right]^2}{n \cdot E_{\theta} \left\{ \left[\frac{d}{d\theta} \ln(f(X_i|\theta)) \right]^2 \right\}}$$

Note The quantity $E_{\theta} \left\{ \left[\frac{d}{d\theta} \ln(f(X_i|\theta)) \right]^2 \right\}$ is called the information number or Fisher's information.

Note Any candidate estimator satisfying $E_{\theta}[W(\mathbf{X})] = \tau(\theta)$ and attaining the Cramer-Rao bound is the best unbiased estimator of $\tau(\theta)$.