

1. INTRODUCTION

1.1 What is Statistics?

STATISTICS: a subject dealing with the **collection** and the **use** of data in order to (help to) answer the questions like:

- Do stock markets rise and fall randomly?
- Is global warming really happening?
- Does a certain new drug prolong life for AIDS sufferers?
- Are GCSE and A level examinations standards declining?
- Is the national lottery making us a nation of compulsive gamblers?
- Is the gap between rich and poor widening in Britain?
- Do Persil adverts really make us want to buy Persil?
-

Those questions are difficult to study in laboratory, and admit no self-evident axioms

Typically *Data* are subject to **uncertainty** — *random variables*

Statistics: **Let the Data Speak!**

- making guesses about the process generating the data – *Estimation*
- testing the guesses – *Testing hypotheses*
- forecasting future – *Prediction*

Note. Statistics also deals with *experiment design* which will be untouched in this course.

Statistics offers quantitative approaches to solve practical problems. Therefore, **a good knowledge** on the practical world behind data is essential for decent statistics practice.

No unique statistical solution for most practical problems! Therefore

Statistics is also an art

Some guidelines for learning/applying statistics:

- Understand what data say in each specific context. All the methods are just tools to help understand data
- Concentrate on what to do and why, rather than concrete calculation and graphing
- It may take a while to catch the essences of statistics – Keep thinking!

What is the difference between *probability* and *statistical inference*?

Probability — a mathematical subject

Statistics — an applied oriented subject: inference based on data plus ‘assumptions’

Note. The assumptions imposed in statistical inference are typically based on probability theory.

Consider a simple experiment: a coin is tossed n times.

Define a random variable X = the number of ‘heads’.

Probability Questions:

- What is $P(X = 1)$, $P(X \leq 1)$, $P(X = x)$ or $P(X \leq x)$?
- What is $E(X)$, $\text{Var}(X)$, $E(X^m)$ or $E(e^{tX})$?

Probability Answers

If we **assume** that the tosses are independent with a constant probability of success π , then

- $P(X = 1) = n\pi(1 - \pi)^{n-1}$,
 $P(X \leq 1) = (1 - \pi)^{n-1}(1 - \pi + n\pi)$,
 $P(X = x) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}$, etc
- $E(X) = n\pi$, $\text{Var}(X) = n\pi(1 - \pi)$, etc

Statistical Questions

- What is the value of π ?
- Is there any evidence that the coin is ‘biased’ i.e. $\pi \neq 0.5$?
- Are the tosses independent?

Statistical Answers

- Estimate π by the estimator x/n . Give an interval in which π may be expected to fall.
- Reject the Null Hypothesis that $\pi = 0.5$ if x is too far away from $n/2$
- Reject the model of independent trials if there is too much regularity in the observations - e.g. HTHTHTHTHTHTHT....

1.2 Terminology and Notation

The classic statistical inference is based on a **sample** (i.e. a set of **data**)

$$\mathbf{X} = (X_1, \dots, X_n)^T$$

and the **assumption** that $\{X_i\}$ i.i.d. with

$$X_i \sim F_X(x; \theta), \quad \theta \in \Theta,$$

where $F_X(\cdot; \theta)$ is a probability distribution and is known up to an *unknown* **parameter** θ , Θ is the **parameter space**.

Population: $\{F_X(\cdot; \theta), \theta \in \Theta\}$, which is assumed to contain the **true** distribution of X_i

Sample space: consisting of all possible values of \mathbf{X}

Statistic: any (known) function of \mathbf{X} only

Dual identity of a sample:

- a set of observed real numbers in practice, and
- a set of i.i.d. r.v.s in theoretical exploration.

Goal of statistical inference in this course: estimating the true value of θ or testing the hypotheses about θ

Note. A sample or a random sample refers to a set of i.i.d. observations.

Example 1. Let X_1, \dots, X_n be i.i.d. from $N(\mu, \sigma^2)$. Then (X_1, \dots, X_n) is a sample, R^n is the sample space, and

$$\theta = \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}, \quad \Theta = R \times R^+.$$

The population is

$$\{N(\mu, \sigma^2), (\mu, \sigma^2) \in \Theta\}.$$

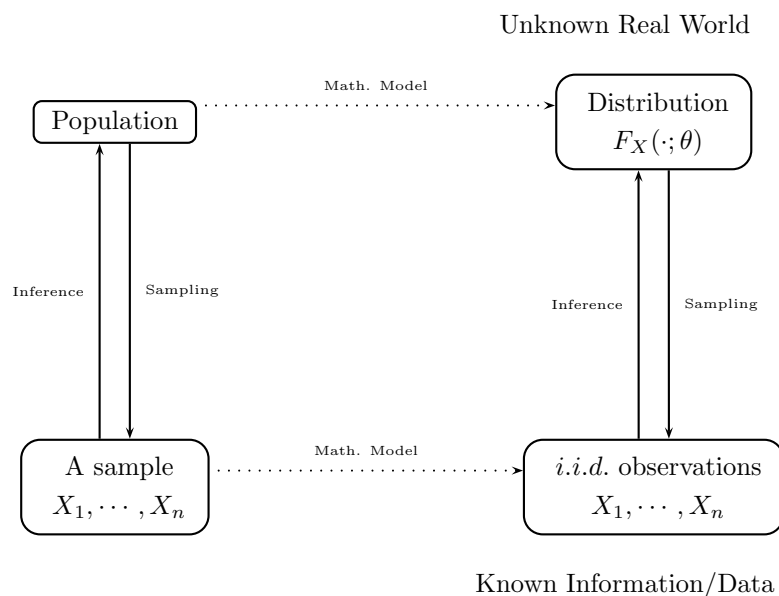
Both sample mean $\bar{X} = n^{-1} \sum_{1 \leq i \leq n} X_i$ and sample variance $S = \frac{1}{n-1} \sum_{1 \leq i \leq n} (X_i - \bar{X})^2$ are statistics. But $\bar{X} - \mu$ is not!

Note. The above setting is referred as a **parametric model** since the population is known up to unknown parameters.

In contrast, a **nonparametric model** specifies that the true distribution belongs to a class which cannot be specified by a finite number of parameters only. For example, a population may be

{all continuous density functions in one variable}.

We deal with parametric models first.



1.3 Statistics from Normal samples

In this section we always assume that $\mathbf{X} = (X_1, \dots, X_n)^\tau$ be a random sample from $N(0, 1)$.

Then the two summary statistics are sample mean and sample variance:

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n X_j, \quad S^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2.$$

Theorem 1. The following statements hold.

- (i) $\bar{X} \sim N(0, 1/n)$,
- (ii) $(n-1)S^2 \sim \chi_{n-1}^2$,
- (iii) \bar{X} and S^2 are independent.

Note. The distribution of $\sum_{1 \leq j \leq n} X_j^2$ is called **the χ^2 -square distribution with n degrees of freedom**, denoted as χ_n^2 or $\chi^2(n)$.

Proof. We make a linear transformation from X_1, \dots, X_n to $\mathbf{Z} = (Z_1, \dots, Z_n)^\tau$ in such a way that the Z_i 's are still i.i.d. $N(0, 1)$. This can be achieved by any transformation $\mathbf{Z} = \mathbf{H}\mathbf{X}$ with \mathbf{H} being a $n \times n$ orthogonal matrix (i.e. $\mathbf{H}^\tau \mathbf{H} = \mathbf{I}_n$). Let the first row of \mathbf{H} be

$$n^{-1/2}(1, 1, \dots, 1).$$

Then $Z_1 = \sqrt{n}\bar{X}$. Hence (i) holds. Since

$$\sum_{j=1}^n Z_j^2 = \mathbf{X}^\tau \mathbf{H}^\tau \mathbf{H} \mathbf{X} = \mathbf{X}^\tau \mathbf{X} = \sum_{j=1}^n X_j^2,$$

it holds that

$$\sum_{j=2}^n Z_j^2 = \sum_{j=1}^n X_j^2 - Z_1^2 = \sum_{j=1}^n X_j^2 - n\bar{X}^2 = (n-1)S^2.$$

Therefore (ii) and (iii) also hold.

Remark. The above proof indicates the following decomposition

$$\begin{array}{rcl} \sum_{i=1}^n X_i^2 & = & \sum_{i=1}^n (X_i - \bar{X})^2 + n\bar{X}^2 \\ \chi_n^2 & & \chi_{(n-1)}^2 \quad \chi_{(1)}^2, \end{array}$$

where the two χ^2 random variables on the RHS are independent.

Example 1. Let Y_1, \dots, Y_n iid from $N(\mu, \sigma^2)$. Define

$$X_i = (Y_i - \mu)/\sigma.$$

Then $\{X_i\}$ i.i.d. $N(0, 1)$, and

$$\bar{Y} = \sigma\bar{X} + \mu, \quad S_y^2 = \sigma^2 S_x^2.$$

Hence

- (i) $\bar{Y} \sim N(\mu, \sigma^2/n)$,
- (ii) $(n-1)S_y^2/\sigma^2 \sim \chi_{n-1}^2$,
- (iii) \bar{Y} and S_y^2 are independent.

Further,

$$\sum_{i=1}^n \frac{(Y_i - \mu)^2}{\sigma^2} = \sum_{i=1}^n \frac{(Y_i - \bar{Y})^2}{\sigma^2} + n \frac{(\bar{Y} - \mu)^2}{\sigma^2}.$$

F distribution

Definition. Let $Y \sim \chi_n^2$ and $X \sim \chi_m^2$, and X and Y are independent. Then the random variable

$$U \equiv \frac{X/m}{Y/n} = \frac{n}{m} \frac{X}{Y}$$

has an F -distribution with degrees of freedom (m, n) . We write $U \sim F(m, n)$ or $U \sim F_{m,n}$.

Obviously, $1/U \sim F_{n,m}$.

A formal notation: $F_{m,n} = \frac{\chi_m^2/m}{\chi_n^2/n}$.

Typical use of F distribution

Let X_1, X_2, \dots, X_m and Y_1, Y_2, \dots, Y_n be two indep samples from $N(\mu_x, \sigma_x^2)$ and $N(\mu_y, \sigma_y^2)$ respectively. Then

$$\frac{\frac{1}{(m-1)\sigma_x^2} \sum_{i=1}^m (X_i - \bar{X})^2}{\frac{1}{(n-1)\sigma_y^2} \sum_{i=1}^n (Y_i - \bar{Y})^2} \sim F_{(m-1), (n-1)}.$$

Student's t -distribution

Definition. Let $Z \sim N(0, 1)$ and $U \sim \chi_{(k)}^2$, and Z and U be independent. Then the random variable

$$T = \frac{Z}{[\frac{1}{k}U]^{\frac{1}{2}}},$$

has a t -distribution with k degrees of freedom. We write $T \sim t_k$ or $T \sim t(k)$.

Obviously, $T^2 \sim F_{1,k}$.

Also T and $-T$ have the same distribution, i.e. the pdf is symmetric. In fact $ET = 0$ for $k > 1$.

A formal notation: $t_k = \frac{N(0,1)}{\{\chi_k^2/k\}^{1/2}}$.

Typical use of t distribution

If Y_1, Y_2, \dots, Y_n are iid $N(\mu, \sigma^2)$, then

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

and

$$(n-1)S_y^2/\sigma^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{\sigma^2} \sim \chi_{(n-1)}^2$$

and these two are independent. It follows that

$$\frac{\bar{Y} - \mu}{\sqrt{S_y^2/n}} \sim t_{n-1}.$$

2. MAXIMUM LIKELIHOOD ESTIMATION

2.1 Likelihood

Likelihood is one of the most fundamental concept in all types of statistical inference.

Definition 1 Suppose that \mathbf{X} has density function or probability function $f(\mathbf{x}; \boldsymbol{\theta})$. We have observed $\mathbf{X} = \mathbf{x}$. Then the likelihood function with observation \mathbf{x} is defined as

$$L(\boldsymbol{\theta}) \equiv L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}).$$

Density/probability function: a function of \mathbf{x} , specifying the distribution of random variable \mathbf{X}

Likelihood: a function of $\boldsymbol{\theta}$, reflecting information on $\boldsymbol{\theta}$ contained in observation \mathbf{x}

Note. A likelihood function represents the uncertainty on a unknown non-random constant $\boldsymbol{\theta}$, and it is **not** a density or probability function! It provides

- a rational degree of belief, or
- an order of preferences

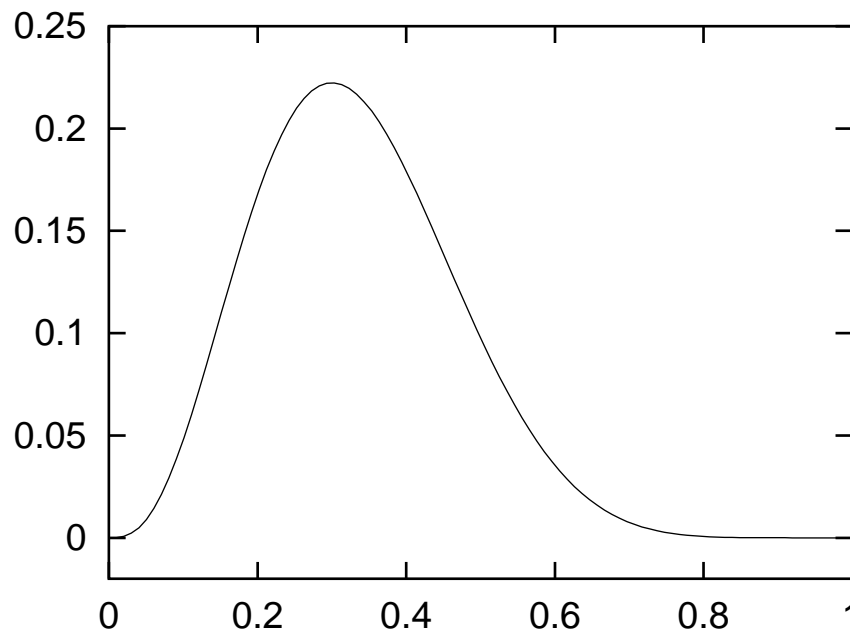
on possible values of the parameter $\boldsymbol{\theta}$. This can be seen more clearly in the simple example on next slide.

In fact, a likelihood function is often defined up to a positive constant — the constant here refers to a quantity independent of $\boldsymbol{\theta}$. But it may depending on \mathbf{x} . (Note \mathbf{x} is a given constant.)

Example 1. Suppose that x is the number of successes from a known number n of independent trials with unknown probability of success π . The probability function, and so the likelihood function is

$$L(\pi) = f(x; \pi) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}.$$

The likelihood function $L(\pi; x)$ can be graphed as a function of π . It changes shape for different values of x . A likelihood function for $x = 3$ when $n = 10$ is shown in the Figure below.



Notice that the likelihood function shown above is *not* a density function. It does not have an area of 1 below it.

We use the likelihood function to compare the plausibility of different possible parameter values. For instance, the likelihood is much larger for $\pi = 0.3$ than for $\pi = 0.8$, that is the data $x = 3$ have a greater probability of being observed if $\pi = 0.3$ than if $\pi = 0.8$. This makes $\pi = 0.3$ much more likely as the true value for π than 0.8.

Note. In the above argument, we do not need to calculate exact probabilities under different values of θ . Only the order of those quantities matters!

Let X_1, \dots, X_n be i.i.d. with pdf $f(\cdot, \theta)$. Write $\mathbf{X} = (X_1, \dots, X_n)^\tau$. Then the likelihood function is

$$L(\theta) = L(\theta; \mathbf{X}) = \prod_{i=1}^n f(X_i, \theta),$$

which is a product of n terms. Then the log-likelihood function is

$$l(\theta) = l(\theta; \mathbf{X}) \equiv \log\{L(\theta; \mathbf{X})\} = \sum_{i=1}^n \log\{f(X_i, \theta)\},$$

which is a sum of n terms.

This explains why log-likelihood functions are often used with independent observations.

Definition 2. *Likelihood Principle*

For two observed values \mathbf{x} and \mathbf{y} , if

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

the inferences for θ based on \mathbf{x} and \mathbf{y} should be the same.

2.2 Sufficiency and Data Compression

If we start with n observations and construct from them a k -dimensional statistic, if $k < n$ we expect that in general some information in the sample about θ would be lost. The statistics is less complicated than the original sample because it is in a smaller dimension. It reduces the data, and in general will reduce the information about θ . A simple example may help to make clear what is happening.

Example 2 Based on a sample $\{X_1, \dots, X_n\}$, consider a sequence of statistics of increasing dimensions:

$$\begin{aligned} T_1(\mathbf{X}) &= \left(\sum_{i=1}^n X_i \right) \\ T_2(\mathbf{X}) &= \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2 \right) \\ T_3(\mathbf{X}) &= \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i^3 \right) \\ &\dots \\ T_{n-1}(\mathbf{X}) &= \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i^3, \dots, \sum_{i=1}^n X_i^{n-1} \right) \\ T_n(\mathbf{X}) &= \left(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2, \sum_{i=1}^n X_i^3, \dots, \sum_{i=1}^n X_i^n \right) \end{aligned}$$

Here we see that T_1 is very simple and in 1 dimension, whereas T_n is in n dimensions. Actually, from T_n one can reconstruct the original observations X_1, X_2, \dots, X_n , except for their ordering. So T_n contains all

the information in the sample about θ . But except in special cases T_1 will not include all the information about θ in the sample.

If a statistic contains all the information in a sample about θ we shall say that it is a **sufficient statistic** for θ .

Example 3 If X_1, X_2 are independent observations from $N(\mu, 1)$ then the statistic $X_1 + X_2$ contains all the information in the sample about μ . This is intuitive, because knowing $X_1 + X_2$ and $X_1 - X_2$ is equivalent to having both observations X_1, X_2 , but $X_1 - X_2$ has a $N(0, 2)$ distribution and so contains no information about μ . This conclusion is made more secure by the independence of $X_1 + X_2$ and $X_1 - X_2$.

Now we make the idea of sufficiency more precise with one of the great ideas from the work of R.A.Fisher. We seek to define an idea of ‘sufficiency’ such that a statistic T is ‘sufficient’ for θ if it contains all the information about θ that was in the whole set of observations.

Definition 3. Sufficient Statistic

Suppose $\mathbf{X} \sim f(\cdot, \theta)$. $\mathbf{T}(\mathbf{X})$ is said to be a sufficient statistic for θ if the conditional distribution of the sample \mathbf{X} given $\mathbf{T}(\mathbf{X})$ does not depend on θ .

The likelihood principle implies the sufficiency principle which indicates that we only need to use sufficient statistics in inference.

Definition 4. Sufficiency Principle

All sufficient statistics based on \mathbf{X} should lead to the same inferences for θ .

Theorem 1. Factorisation Criterion

Let $\mathbf{X} \sim f(\mathbf{x}, \theta)$. Then $\mathbf{T} = \mathbf{T}(\mathbf{X})$ is a sufficient statistic for θ iff

$$f(\mathbf{x}; \theta) = g(\mathbf{T}(\mathbf{x}), \theta)h(\mathbf{x}).$$

We will not prove this result. To see how a proof might go, note that it is always possible, but not necessary, to take $h(\mathbf{x})$ as the conditional density of \mathbf{X} given $\mathbf{T}(\mathbf{X})$ and $g(\mathbf{T}(\mathbf{x}), \theta)$ as the density of $\mathbf{T}(\mathbf{X})$.

Example 4. Suppose that we have a random sample $\mathbf{X} = (X_1, \dots, X_n)^\tau$ from $N(\mu, 1)$. Then the joint density is

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{X}; \theta) &= \frac{1}{(\sqrt{2\pi})^n} \exp \left[-\frac{1}{2} \sum_{i=1}^n (X_i - \mu)^2 \right] \\ &= \exp \left[-\frac{1}{2} n(\bar{X} - \mu)^2 \right] \frac{1}{(\sqrt{2\pi})^n} \exp \left[-\frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})^2 \right]. \end{aligned}$$

So we can take $T(\mathbf{X}) = \bar{X}$ and

$$g(\bar{X}, \theta) = \exp \left[-\frac{1}{2} n(\bar{X} - \mu)^2 \right],$$

and

$$h(\mathbf{X}) = \frac{1}{(\sqrt{2\pi})^n} \exp \left[-\frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})^2 \right].$$

Two samples with the same sample mean should lead to exactly the same inferences about μ if the sufficiency principle holds.

Example 5. Here is another example where a sufficient statistic is of the same dimension as the sample size, but is not the whole sample. Let $\mathbf{X} = (X_1, \dots, X_n)$ be a sample from $U(-\theta, \theta)$ and $\Theta = R^+$. Then the absolute values of the observations $|X_1|, |X_2|, \dots, |X_n|$ are from $U(0, \theta)$. One loses the information on the signs of the observations, but retains full information on θ . Here the statistic $\mathbf{T}(\mathbf{X}) = (|X_1|, |X_2|, \dots, |X_n|)$ is a sufficient statistic for θ .

In fact the conditional distribution of \mathbf{X} given $\mathbf{T}(\mathbf{X}) = (T_1, T_2, \dots, T_n) = (|X_1|, |X_2|, \dots, |X_n|)$ is the distribution where independently for all i , $X_i = T_i$ with probability 0.5 and $X_i = -T_i$ with probability 0.5. This conditional distribution does not depend on θ , so the definition is satisfied.

Example 6. (Continuing Example) Another sufficient statistic of θ is $\max_{1 \leq i \leq n} |X_i|$. To this end, note that

$$f_{X_i}(x_i) = \frac{1}{2} \theta I_{(-\theta, \theta)}(x_i) = \begin{cases} \frac{1}{2\theta} & |x_i| < \theta, \\ 0 & \text{otherwise.} \end{cases}$$

Hence

$$\begin{aligned} f_{\mathbf{X}}(x_1, \dots, x_n) &= \prod_{i=1}^n f_{X_i}(x_i) \\ &= \begin{cases} \frac{1}{(2\theta)^n} & \max_{1 \leq i \leq n} |x_i| < \theta, \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{1}{(2\theta)^n} I(\max_{1 \leq i \leq n} |x_i| < \theta). \end{aligned}$$

Now the conclusion follows from the factorisation criterion.

How do you estimate θ ?

Example 7. Suppose that the X_1, \dots, X_n are iid $U(0, \theta)$, where $\theta > 0$. Then

$$f(x_i; \theta) = \begin{cases} \frac{1}{\theta} & x_i \in [0, \theta] \\ 0 & \text{otherwise} \end{cases}$$

which can be written

$$f(x_i; \theta) = \frac{1}{\theta} \times I_{\{x_i \leq \theta\}} I_{\{x_i \geq 0\}}.$$

The joint density is

$$\begin{aligned} f(\mathbf{x}; \theta) &= \frac{1}{\theta^n} \prod_{i=1}^n I_{\{x_i \leq \theta\}} \prod_{i=1}^n I_{\{x_i \geq 0\}} \\ &= \frac{1}{\theta^n} I_{\{\max_i x_i \leq \theta\}} I_{\{\min_i x_i \geq 0\}} \\ &= g(T(\mathbf{x}); \theta) \times h(\mathbf{x}) \end{aligned}$$

where $T(\mathbf{x}) = \max_i x_i$. In this case the (at first sight) rather surprising result is that $\max_i X_i$ is a sufficient statistic for θ .

2.3 Maximum likelihood estimator (MLE)

The MLE is by far the most popular method for deriving estimators.

Definition 5 — MLE

A *Maximum Likelihood Estimator* (MLE), $\hat{\theta} = \hat{\theta}(\mathbf{X}) \in \Theta$, of parameter θ is an estimator satisfying

$$L(\hat{\theta}; \mathbf{X}) \geq L(\theta; \mathbf{X}) \text{ for all } \theta \in \Theta, \text{ or equivalently } l(\hat{\theta}; \mathbf{X}) \geq l(\theta; \mathbf{X}) \text{ for all } \theta \in \Theta.$$

Obviously, a maximum likelihood estimator is the most plausible value for θ as judged by the likelihood function. In many cases where Θ is continuous and the maximum does not occur at a boundary of Θ , $\hat{\theta}$ is often the solution of the equation

$$s(\theta; \mathbf{X}) = \frac{\partial}{\partial \theta} l(\theta; \mathbf{X}) = 0.$$

We call $s(\theta) \equiv s(\theta; \mathbf{X})$ a **score function**.

Example 8. Suppose that Y_1, Y_2, \dots, Y_n is a random sample from $N(\mu, \sigma^2)$ where neither μ or σ^2 is known. Then we can find the maximum likelihood estimator from the log-likelihood

$$\begin{aligned} l(\mu, \sigma^2) &= -n \log \sqrt{2\pi} - n/2 \log \sigma^2 - \sum_{i=1}^n (Y_i - \mu)^2 / (2\sigma^2) \\ &= -n \log \sqrt{2\pi} - n/2 \log \sigma^2 - \sum_{i=1}^n (Y_i - \bar{Y})^2 / (2\sigma^2) \\ &\quad - n(\bar{Y} - \mu)^2 / (2\sigma^2). \end{aligned}$$

This is maximised by choosing $\mu = \bar{Y}$, so $\hat{\mu} = \bar{Y}$ is the MLE for μ . It is easy to see

$$E(\hat{\mu}) = \frac{1}{n} \sum_{i=1}^n E(X_i) = \mu.$$

Such a estimator is called **unbiased**.

Note. The above MLE for μ is equivalent to minimise

$$\sum_{i=1}^n (Y_i - \mu)^2.$$

Therefore $\hat{\mu}$ is also the **least squares estimator** (LSE). Note that LSE is derived from a simple empirical/geometric rule, which makes **no use** of the underlying distribution.

The **profile log-likelihood** remaining is

$$l(\hat{\mu}, \sigma^2) = -n \log \sqrt{2\pi} + (n/2)(\log \sigma^{-2} - \hat{\sigma}^2 \sigma^{-2}),$$

where $\hat{\sigma}^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 / n$. By the lemma below, the MLE for σ^2 is $\hat{\sigma}^2$. Note that the MLE of σ^2 is *biased* since

$$E(\hat{\sigma}^2) = (1 - 1/n)\sigma^2 \neq \sigma^2.$$

Lemma. Define $L(x) = \log(x^{-1}) - b/x$, where $n \geq 1$ and $b > 0$ are constants. Then $L(b) \geq L(x)$ for all $x > 0$.

Example 9. Let X_1, \dots, X_n be i.i.d. Bernoulli(π). Then

$$L(\pi) = \prod_{i=1}^n \pi^{X_i} (1 - \pi)^{1-X_i} = \pi^{n\bar{X}} (1 - \pi)^{n(1-\bar{X})}.$$

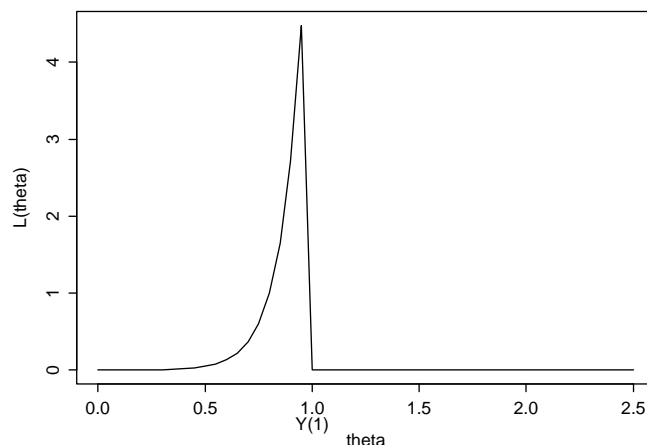
$$l(\pi) = n\bar{X} \log \pi + n(1 - \bar{X}) \log(1 - \pi).$$

Let $s(\pi) = \frac{\partial}{\partial \pi} l(\pi) = 0$, leading to $\hat{\pi} = \bar{X}$.

Example 10. Suppose that Y_1, Y_2, \dots, Y_n is a random sample from an exponential distribution with density function $e^{-(y-\theta)}$ for $y \geq \theta$. This is the usual exponential distribution shifted to start at θ . The Likelihood is

$$L(\theta; \mathbf{Y}) = e^{-n(\bar{Y}-\theta)} I_{\{(\theta, \infty)\}}(Y_{(1)}),$$

where $Y_{(1)}$ is the smallest observation. This likelihood is zero for $\theta > Y_{(1)}$ and increases in θ for $\theta \leq Y_{(1)}$. So the MLE $\hat{\theta} = Y_{(1)}$, which is a boundary maximum.



Invariance property of MLEs

Suppose $\mathbf{X} \sim f(\mathbf{x}, \boldsymbol{\theta})$, and $\boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\theta})$ is a 1-1 transform. Let $\hat{\boldsymbol{\theta}}$ be the MLE for $\boldsymbol{\theta}$, i.e.

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} f(\mathbf{X}, \boldsymbol{\theta}).$$

It is obvious to see that the MLE for $\boldsymbol{\mu}$ is

$$\hat{\boldsymbol{\mu}} = \boldsymbol{\mu}(\hat{\boldsymbol{\theta}}).$$

2.4 Numerical computation of MLEs

In modern statistical applications, it is typically difficult to find explicit analytic forms for the maximum likelihood estimators. These estimators are found more often by iterative procedures built into computer software. An iterative scheme starts with some guess at the MLE and then steadily improves it with each iteration. The estimator is considered found when it has become numerically stable. Sometimes the iterative procedures become trapped at a local maximum which is not a global maximum. There may be a very large number of parameters in a model, which makes such local entrapment much more common.

Newton-Raphson Scheme

Suppose that the log-likelihood function $l(\boldsymbol{\theta})$ is sufficiently smooth. Then

$$s(\hat{\boldsymbol{\theta}}) = 0,$$

where $\hat{\boldsymbol{\theta}}$ is the MLE and $s(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} l(\boldsymbol{\theta})$ is the score function. Let

$$\dot{s}(\boldsymbol{\theta}) = \ddot{l}(\boldsymbol{\theta}) = \frac{\partial^2}{\partial \theta \partial \theta^T} l(\boldsymbol{\theta}).$$

Suppose $\hat{\theta}$ is close to the true value θ_0 . By a simple Taylor expansion,

$$s(\theta_0) = \dot{s}(\theta_0)(\theta_0 - \hat{\theta}) + o_p(\|\hat{\theta} - \theta_0\|).$$

This leads to the approximation

$$\hat{\theta} \approx \theta_0 - \{\dot{s}(\theta_0)\}^{-1}s(\theta_0).$$

Since θ_0 is unknown, we use iterative estimators

$$\hat{\theta}_{k+1} = \hat{\theta}_k - \{\dot{s}(\hat{\theta}_k)\}^{-1}s(\hat{\theta}_k) \quad (1)$$

for $k = 1, 2, \dots$, where $\hat{\theta}_0$ is a prescribed initial value. We define $\hat{\theta} = \hat{\theta}_j$ if $\hat{\theta}_j$ and $\hat{\theta}_{j-1}$ differ by a small amount.

Fisher Score method: replace $\dot{s}(\hat{\theta}_k)$ in (1) by $E_{\theta}\{\dot{s}(\theta)\}$ under $\theta = \hat{\theta}_k$.

Like most iterative algorithms, the choice of appropriate initial values is important to ensure the convergence to right limits. In practice multiple initial values are often used.

The differences between the Newton-Raphson and score methods are subtle. We make observations below

- The convergence of the Newton-Raphson algorithm is often faster when both algorithms converge
- The radius of convergence for the score method is often larger, making the choice of initial values less important for the score method.

2.5 EM algorithms

Goal: to find the MLE $\hat{\theta} = \hat{\theta}(\mathbf{Y})$ for θ from the likelihood based on data \mathbf{Y} :

$$L(\theta; \mathbf{Y}) = f_{\mathbf{Y}}(\mathbf{Y}, \theta),$$

while the ‘complete’ data $\mathbf{X}^{\tau} = (\mathbf{Y}^{\tau}, \mathbf{Z}^{\tau})$ contain a ‘missing’ component \mathbf{Z} . The likelihood based on the complete data is

$$L(\theta; \mathbf{X}) = f_{\mathbf{X}}(\mathbf{X}, \theta).$$

EM (Expectation and Maximisation) algorithms

- *E-step:* compute the conditional expectation

$$Q(\theta) = Q(\theta|\mathbf{Y}, \theta_0) \equiv E\{\log L(\theta; \mathbf{X})|\mathbf{Y}, \theta_0\}$$

- *M-step:* maximise $Q(\theta)$ to give an updated value θ_1

then go to the E-step using $\theta_0 = \theta_1$, and keep iterating until convergence. The limit of θ_0 is taken as $\hat{\theta}(\mathbf{Y})$.

Example 11. The genetic example from (Rao 1973, p.396) assumes that the phenotype data

$$\begin{aligned} \mathbf{Y} &= (Y_1, Y_2, Y_3, Y_4)^{\tau} \sim \\ &\text{Multinomial}\left(4; \frac{1}{2} + \frac{\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4}\right), \end{aligned}$$

where $\theta \in (0, 1)$. Then log-likelihood is

$$l(\theta, \mathbf{Y}) = Y_1 \log(2 + \theta) + (Y_2 + Y_3) \log(1 - \theta) + Y_4 \log \theta + C,$$

which does not yield a closed form $\hat{\theta}$.

Now we treat \mathbf{Y} as incomplete data from $\mathbf{X} = (X_1, \dots, X_5)^\tau$ with multinomial probabilities

$$\left(\frac{1}{2}, \frac{\theta}{4}, \frac{1-\theta}{4}, \frac{1-\theta}{4}, \frac{\theta}{4}\right).$$

Then

$$Y_1 = X_1 + X_2, \quad Y_i = X_{i+1} \quad \text{for } i = 2, 3, 4.$$

The log-likelihood of based on \mathbf{X} is

$$l(\theta, \mathbf{X}) = (X_2 + X_5) \log \theta + (X_3 + X_4) \log(1 - \theta) + C,$$

which readily yields

$$\hat{\theta}(\mathbf{X}) = \frac{X_2 + X_5}{X_2 + X_3 + X_4 + X_5}.$$

Now the E-step is to find

$$\begin{aligned} Q(\theta) &= E\{l(\theta, \mathbf{X}) | \mathbf{Y}, \theta_0\} \\ &= \log \theta E(X_2 + X_5 | \mathbf{Y}, \theta_0) + \log(1 - \theta) E(X_3 + X_4 | \mathbf{Y}, \theta_0) \\ &= (\hat{X}_2 + Y_4) \log \theta + (Y_2 + Y_3) \log(1 - \theta), \end{aligned}$$

where $\hat{X}_2 = E(X_2 | \mathbf{Y}, \theta_0)$. Since the conditional distribution of X_2 given $Y_1 (= X_1 + X_2)$ is a binomial distribution with $n = Y_1$ and

$$p = \frac{\theta_0/4}{1/2 + \theta_0/4} = \frac{\theta_0}{2 + \theta_0}.$$

Hence

$$\hat{X}_2 = np = \frac{Y_1 \theta_0}{2 + \theta_0}. \quad (2)$$

The M-step leads to

$$\theta_1 = \frac{\hat{X}_2 + Y_4}{\hat{X}_2 + Y_4 + Y_2 + Y_3}. \quad (3)$$

For $Y = (125, 18, 20, 34)$ and the **initial value**

$$\theta_0 = 4 \times 34 / (125 + 18 + 20 + 34)$$

which is a relative frequency estimate, the first 5 iterations between (2) and (3) are 0.690, 0.635, 0.628, 0.627 and 0.627, giving the MLE $\hat{\theta} = 0.627$.

What can be said about *convergence properties* of the EM algorithm?

Let θ_0 be an arbitrary initial value and θ_1 be the updated value obtained from applying the iteration once. Then it can be shown that

$$L(\theta_1; \mathbf{Y}) \geq L(\theta_0; \mathbf{Y}).$$

Unfortunately, it does not imply that the iterations will always lead to the MLE eventually.

It is important to choose appropriate initial values to ensure the algorithm converges to the MLE. (This is also true for both Gaussian-Raphson and score methods!) In practice, it is a good idea to use a variety of initial values.

Further discussion on the convergence of EM algorithms, see Wu (1983) *Annals of Statistics, Vol.11, pp.95-103* and §12.4 of Pawitan (2001).

General comments on EM algorithms:

- The EM algorithm is a general procedure for computing MLEs. It is not really a numerical algorithm. The calculation of M-Step typically involves other numerical algorithms such as Newton-Raphson and the score methods.
- It can be applied when some data are *genuinely* missing. It can also be applied when the missing information is merely a concept based on which we transform a difficult optimisation problem into a sequence of easier problems; see Example 11 above. This is particularly relevant when $\hat{\boldsymbol{\theta}}(\mathbf{Y})$ is difficult to calculate while $\hat{\boldsymbol{\theta}}(\mathbf{X})$ is easier to obtain.
- The convergence of EM algorithm may be very slow, depending on the amount of missing information. It is often slower than the Newton-Raphson algorithm.

Example 12. Mixture distributions

Let Y_1, \dots, Y_n be i.i.d. with a mixture pdf

$$f(y) = \sum_{j=1}^k \alpha_j f_j(y, \boldsymbol{\lambda}_j),$$

where f_j are pdfs, $\alpha_j \geq 0$ and $\sum_j \alpha_j = 1$. The parameter $\boldsymbol{\theta}$ contains all α_j and $\boldsymbol{\lambda}_j$. This model becomes important because

1. it represents heterogeneous data well since each f_j represents one heterogeneous component, and
2. it provides very good approximations to a large class of distributions.

The likelihood based on $\mathbf{Y} = (Y_1, \dots, Y_n)^\tau$ is

$$L(\boldsymbol{\theta}; \mathbf{Y}) = \prod_{i=1}^n \left\{ \sum_{j=1}^k \alpha_j f_j(Y_i, \boldsymbol{\lambda}_j) \right\}.$$

Maximising this likelihood is difficult, due to the presence of the summations, which reflect the fact that we are typically lacking in knowledge of which component any particular sample value comes from. This is the missing information!

Let $\mathbf{X}_i^\tau = (Y_i, \mathbf{Z}_i^\tau)$, where $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{ik})^\tau$ is a $k \times 1$ vector with a 1 in the position corresponding to the component of the mixture that Y_i comes from, and 0 elsewhere.

Let \mathbf{e}_j be the $k \times 1$ vector with the j -th component 1 and all the other components 0. The joint probability-density function for (\mathbf{Z}_1, Y_1) is

$$\begin{aligned} & P\{Z_1 = \mathbf{e}_\ell, Y_1 \in [y_1, y_1 + dy_1]\} / dy_1 \\ &= P(Z_1 = \mathbf{e}_\ell) P\{Y_1 \in [y_1, y_1 + dy_1] | Z_1 = \mathbf{e}_\ell\} / dy_1 \\ &= \alpha_\ell f_\ell(y_1, \boldsymbol{\lambda}_\ell) = \prod_{j=1}^k \alpha_j^{e_{\ell j}} f_j(y_1, \boldsymbol{\lambda}_j)^{e_{\ell j}}, \end{aligned}$$

where $\mathbf{e}_\ell = (e_{\ell 1}, \dots, e_{\ell k})^\tau$.

Let $\mathbf{X}^\tau = (\mathbf{X}_1^\tau, \dots, \mathbf{X}_n^\tau)$.

$$L(\boldsymbol{\theta}; \mathbf{X}) = \prod_{i=1}^n \prod_{j=1}^k \alpha_j^{Z_{ij}} f_j(Y_i, \boldsymbol{\lambda}_j)^{Z_{ij}},$$

$$l(\boldsymbol{\theta}; \mathbf{X}) = \sum_{i=1}^n \mathbf{Z}_i^\tau \left\{ \begin{pmatrix} \log \alpha_1 \\ \vdots \\ \log \alpha_k \end{pmatrix} + \begin{pmatrix} \log f_1(Y_i, \boldsymbol{\lambda}_1) \\ \vdots \\ \log f_k(Y_i, \boldsymbol{\lambda}_k) \end{pmatrix} \right\}.$$

Note that

$$E(\mathbf{Z}_i | \mathbf{Y}, \boldsymbol{\theta}_0) = \left(\frac{\alpha_1^0 f_1(Y_i | \boldsymbol{\lambda}_1^0)}{\sum_{\ell=1}^k \alpha_\ell^0 f_\ell(Y_i | \boldsymbol{\lambda}_\ell^0)}, \dots, \frac{\alpha_k^0 f_k(Y_i | \boldsymbol{\lambda}_k^0)}{\sum_{\ell=1}^k \alpha_\ell^0 f_\ell(Y_i | \boldsymbol{\lambda}_\ell^0)} \right)^\tau,$$

which are constants as far as the M-step is concerned.

Now the E-step implies that

$$Q(\boldsymbol{\theta}) = \sum_{i=1}^n E(\mathbf{Z}_i^\tau | \mathbf{Y}, \boldsymbol{\theta}_0) \left\{ \begin{pmatrix} \log \alpha_1 \\ \vdots \\ \log \alpha_k \end{pmatrix} + \begin{pmatrix} \log f_1(Y_i, \boldsymbol{\lambda}_1) \\ \vdots \\ \log f_k(Y_i, \boldsymbol{\lambda}_k) \end{pmatrix} \right\},$$

and the M-step requires to maximise this function, which is much easier than minimizing $l(\boldsymbol{\theta}, \mathbf{Y})$ directly for, for example, normal pdf f_j 's.

2.6 Evaluating estimation

To measure the accuracy of an MLE or, more general, any estimation procedure, we need to define some measures for the goodness (or badness) of an estimator.

Let $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{X})$ be an estimator of $\boldsymbol{\theta}$, and $\boldsymbol{\theta}_o$ be the (unknown) *true value* of $\boldsymbol{\theta}$. Note that

- (i) exact estimation error $\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_o$ is unknown, and
- (ii) $\hat{\boldsymbol{\theta}}$ is a random variable

we have to gauge the error

- (i) in terms of a probability average, and
- (ii) for all possible values of $\boldsymbol{\theta}_o \in \Theta$.

Let $P_{\boldsymbol{\theta}}$, $E_{\boldsymbol{\theta}}$ and $\text{Var}_{\boldsymbol{\theta}}$ denote the probability distribution, expectation and variance under $\boldsymbol{\theta}_o = \boldsymbol{\theta}$.

Bias: $\text{Bias}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = E_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) - \boldsymbol{\theta}$

Variance: $\text{Var}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})$

Standard error: $\{\text{Var}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})\}^{1/2}$

Mean square error (MSE): $E_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^2$

Mean absolute error (MAE): $E_{\boldsymbol{\theta}}|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}|$

Note that

- standard error is a meaningful measure of accuracy for unbiased estimators only, and
- MSE (it its squared-root) should be used in general as

$$\text{MSE}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = \{\text{Bias}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})\}^2 + \text{Var}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}).$$

Ideally we would seek for the estimator which minimises MSE or MAE for all $\boldsymbol{\theta} \in \Theta$ over all possible candidate estimators. Unfortunately such a global optimum rarely exists. However if we confine to some subclass of estimators, the MLE is often optimal or asymptotically optimal.

The MSE is most frequently used largely due to its technical tractability while the MAE leads to estimators which is more robust against outliers in observations.

Fisher Information

Suppose $\mathbf{X} \sim f(\mathbf{x}, \boldsymbol{\theta})$. The score function is

$$s(\boldsymbol{\theta}) = \dot{l}(\boldsymbol{\theta}; \mathbf{X}) = \frac{\partial}{\partial \boldsymbol{\theta}} \log\{f(\mathbf{X}, \boldsymbol{\theta})\}.$$

We assume certain regularity conditions so that we can take derivatives under the integral sign.

Mean of $s(\boldsymbol{\theta})$:

$$\begin{aligned} E_{\boldsymbol{\theta}}\{s(\boldsymbol{\theta})\} &= \int s(\boldsymbol{\theta}) f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= \int \frac{\partial}{\partial \boldsymbol{\theta}} \log\{f(\mathbf{x}, \boldsymbol{\theta})\} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= \int \frac{\partial}{\partial \boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= \frac{\partial}{\partial \boldsymbol{\theta}} \int f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} = 0. \end{aligned}$$

Variance of $s(\boldsymbol{\theta})$ — Fisher information matrix:

$$\begin{aligned} \mathcal{I}(\boldsymbol{\theta}) &= \mathcal{I}_{\mathbf{X}}(\boldsymbol{\theta}) = \text{Var}_{\boldsymbol{\theta}}\{s(\boldsymbol{\theta})\} \\ &= E_{\boldsymbol{\theta}}\{s(\boldsymbol{\theta})s(\boldsymbol{\theta})^{\tau}\} \\ &= -E_{\boldsymbol{\theta}}\left[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\tau}} \log\{f(\mathbf{X}, \boldsymbol{\theta})\}\right], \end{aligned}$$

because

$$\begin{aligned} &E_{\boldsymbol{\theta}}\left\{\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^{\tau}} l(\boldsymbol{\theta})\right\} \\ &= \int \frac{\ddot{L}(\boldsymbol{\theta})L(\boldsymbol{\theta}) - \dot{L}(\boldsymbol{\theta})\dot{L}(\boldsymbol{\theta})^{\tau}}{L(\boldsymbol{\theta})} d\mathbf{x} \\ &= \int \ddot{L}(\boldsymbol{\theta}) d\mathbf{x} - \int \frac{\dot{L}(\boldsymbol{\theta})\dot{L}(\boldsymbol{\theta})^{\tau}}{L(\boldsymbol{\theta})} d\mathbf{x} \\ &= - \int \frac{\dot{L}(\boldsymbol{\theta})\dot{L}(\boldsymbol{\theta})^{\tau}}{L(\boldsymbol{\theta})} d\mathbf{x} \\ &= - \int s(\boldsymbol{\theta})s(\boldsymbol{\theta})^{\tau} f(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} \\ &= -E_{\boldsymbol{\theta}}\{s(\boldsymbol{\theta})s(\boldsymbol{\theta})^{\tau}\}. \end{aligned}$$

Fisher information $\mathcal{I}(\boldsymbol{\theta})$ measures the information on $\boldsymbol{\theta}$ contained in data \mathbf{X} . Further if $\mathbf{X} = (X_1, \dots, X_n)^{\tau}$, and X_1, \dots, X_n are independent,

$$\mathcal{I}(\boldsymbol{\theta}) = \mathcal{I}_{\mathbf{X}}(\boldsymbol{\theta}) = \sum_{j=1}^n \mathcal{I}_{X_j}(\boldsymbol{\theta}).$$

For $\theta = \theta$ is a scalar, the Fisher information is

$$\mathcal{I}(\theta) = E_{\theta}\{s(\theta)^2\} = -E_{\theta}\{\ddot{l}(\theta)\}.$$

Theorem 2. (Cramér-Rao inequality)

Let $\mathbf{X} \sim f(\cdot, \theta)$ which satisfying some regularity conditions. Let $T = T(\mathbf{X})$ be a statistic with $g(\theta) = E_{\theta}(T)$. Then for any $\theta \in \Theta$,

$$\text{Var}_{\theta}(T) \geq \{\dot{g}(\theta)\}^2 / \mathcal{I}(\theta).$$

The Cramér-Rao inequality specifies a lower bound for any *unbiased estimator* for the parameter $g(\theta)$. When the equality holds, T is the **minimum variance unbiased estimator** (MVUE) of $g(\theta)$.

An MVUE is a member in the class of all unbiased estimators with the minimum variance. *Not all MVUEs can attain the Cramér-Rao lower bounds.*

Example 13. Let X_1, \dots, X_n be a sample from $N(\mu, \sigma^2)$. We consider estimators for μ , treating σ^2 as known. The score function (for one observation) is

$$\begin{aligned} s(\mu; X_1) &= \frac{\partial}{\partial \mu} \log[e^{-\frac{1}{2\sigma^2}(X_1 - \mu)^2} / \sqrt{2\pi\sigma^2}] \\ &= \frac{\partial}{\partial \mu} [-\frac{1}{2\sigma^2}(X_1 - \mu)^2] = (X_1 - \mu)/\sigma^2. \end{aligned}$$

Note $\ddot{l}(\mu) = \dot{s}(\mu) = -\sigma^{-2}$. Hence the Fisher information based on a single observation is $\mathcal{I}_{X_1}(\mu) = \sigma^{-2}$. Therefore

$$\mathcal{I}(\mu) = \mathcal{I}_{X_1, \dots, X_n}(\mu) = n/\sigma^2.$$

For any unbiased estimator $\hat{\mu}$ for μ , it holds that

$$\text{Var}_{\mu}(\hat{\mu}) \geq \sigma^2/n,$$

which is the variance of \bar{X} . Hence \bar{X} is the MVUE for μ .

Asymptotic properties of MLEs

Let X_1, \dots, X_n be i.i.d. with pdf $f(\cdot, \theta)$. Write

$$l(\theta) = l(\theta; \mathbf{X}) = \sum_{j=1}^n \log f(X_j, \theta).$$

Let $\hat{\theta}$ be the MLE which maximises $l(\theta)$. Suppose f fulfils certain regularity conditions.

(a) **Consistency.**

The MLE is consistent in the sense that as $n \rightarrow \infty$,

$$P_{\theta}\{||\hat{\theta} - \theta|| > \varepsilon\} \rightarrow 0$$

for any $\varepsilon > 0$.

Consistency requires that an estimator converges to the parameter to be estimated. It is a very mild and modest condition that any reasonable estimator should fulfil. The consistency condition is often used to *rule out bad estimators*.

(b) **Asymptotic normality**

As $n \rightarrow \infty$,

$$n^{1/2}(\hat{\theta} - \theta) \xrightarrow{D} N(0, \{\mathcal{I}_{X_1}(\theta)\}^{-1}).$$

For large n , it holds approximately that

$$\hat{\boldsymbol{\theta}} \sim N(\boldsymbol{\theta}, \{\mathcal{I}_{X_1}(\boldsymbol{\theta})\}^{-1}/n).$$

Therefore *asymptotically* the MLE is unbiased and attains the Cramér-Rao lower bound. Any estimator fulfilling this condition is called **efficient**.

An approximate standard error of the j -th component of $\hat{\boldsymbol{\theta}}$ is the square-root of the (j, j) -th element of $\{\mathcal{I}_{X_1}(\boldsymbol{\theta})\}^{-1}$ divided by $n^{1/2}$.

2.7 Bootstrapping MSEs

An MSE provides a measure for the accuracy of the estimator. The standard errors of an MLE derived from the asymptotic normality typically depend on unknown parameters, which are replaced by their estimators in practice. Alternatively we may estimate the MSE by a simulation method called Bootstrap.

Let X_1, \dots, X_n be i.i.d. with pdf $f(\cdot, \theta)$. Let

$$\hat{\boldsymbol{\theta}} = T(X_1, \dots, X_n)$$

be an estimator. The *goal* here is to estimate

$$\nu \equiv \{\text{MSE}_{\theta_o}(\hat{\boldsymbol{\theta}})\}^{1/2},$$

where θ_o is the true value.

If we *knew* $f(\cdot, \theta_o)$ completely, ν is known in principle, and may be estimated easily via a repeated sampling as follows. We draw B independent samples of size n from $f(\cdot, \theta_o)$. For each sample, we calculate $\hat{\boldsymbol{\theta}}$, obtaining $\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_B$. Then the sample root-MSE

$$\left\{ \frac{1}{B} \sum_{b=1}^B (\hat{\boldsymbol{\theta}}_b - \theta_o)^2 \right\}^{1/2}$$

is a reasonable estimator for ν . By the LLN, this estimator converges to ν as $B \rightarrow \infty$.

The **basic** idea of bootstrap is to adopt the above sampling procedure in the so-called *bootstrap world*: now the population is $f(\cdot, \hat{\boldsymbol{\theta}})$ which is **known**. We draw a sample denoted as (X_1^*, \dots, X_n^*) from this distribution. Define the *bootstrap version* of the estimator

$$\hat{\boldsymbol{\theta}}^* = T(X_1^*, \dots, X_n^*).$$

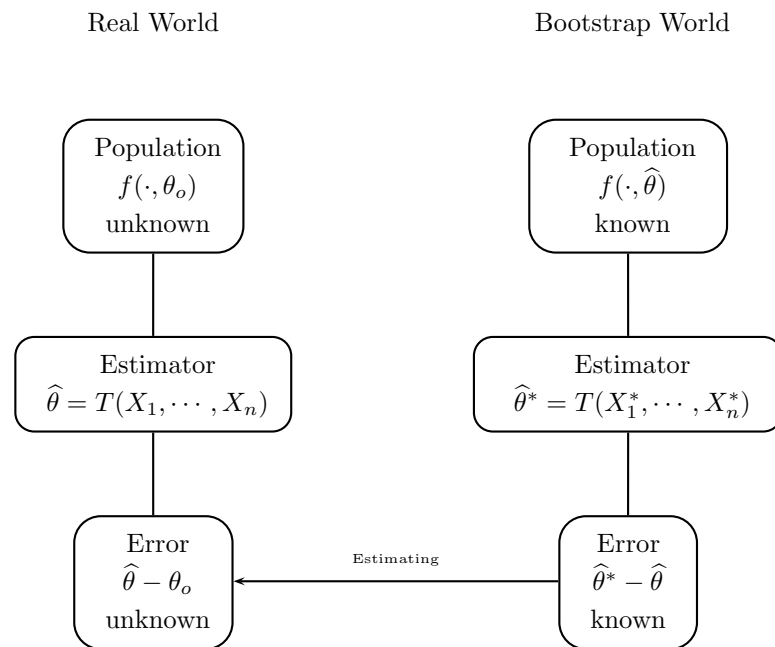
Then the quantity

$$\nu^* = \{\text{MSE}_{\hat{\boldsymbol{\theta}}}(\hat{\boldsymbol{\theta}}^*)\}^{1/2}$$

is known in principle since the distribution $f(\cdot, \hat{\boldsymbol{\theta}})$ is completely known. Define ν^* as a **bootstrap** estimator for ν .

In practice, we draw B sets samples from $f(\cdot, \hat{\boldsymbol{\theta}})$, forming B bootstrap versions of estimator $\hat{\boldsymbol{\theta}}_1^*, \dots, \hat{\boldsymbol{\theta}}_B^*$. The ν^* is calculated as

$$\nu^* = \left\{ \frac{1}{B} \sum_{j=1}^B (\hat{\boldsymbol{\theta}}_j^* - \hat{\boldsymbol{\theta}})^2 \right\}^{1/2}.$$



Bootstrap is a powerful tool for statistical inference. It has different forms for different applications. The one introduced here is in the form of *parametric bootstrap*. The special monographs on this topic include:

- Hall, P. (1992). The bootstrap and Edgeworth expansion. Springer, New York.
- Efron, B. and Tibshirani, R.J. (1993). An introduction to the bootstrap. Chapman and Hall, New York.
- Shao, J. and Tu, D. (1995). The jackknife and bootstrap. Springer, New York.
- Davison, A.C. and Hinkley, D.V. (1997). Bootstrap methods and their application. Cambridge University Press, Cambridge.

3. TESTING HYPOTHESES

3.1 Hypothesis testing: basic setting

Statistical estimation and tests address *different kind* of practical problems.

We work with a family of distributions $f(\mathbf{x}; \boldsymbol{\theta})$, where $\boldsymbol{\theta} \in \Theta$, and $f(\mathbf{x}; \boldsymbol{\theta})$ is completely known apart from $\boldsymbol{\theta}$.

Basic setting:

A statistical test make a decision between two sets of hypotheses on unknown parameter $\boldsymbol{\theta}$, namely

$$H_0 : \boldsymbol{\theta} \in \Theta_0 \quad \text{vs} \quad H_1 : \boldsymbol{\theta} \in \Theta_1,$$

where Θ_0, Θ_1 are subsets of Θ . We always assume

$$\Theta_0 \cap \Theta_1 = \emptyset \text{ (empty set),}$$

i.e. H_0 and H_1 are distinct, mutually exclusive, although we do not always require $\Theta_0 \cup \Theta_1 = \Theta$.

Based on a sample $\mathbf{X} = (X_1, \dots, X_n)^\tau$, a statistical test is to make a decision either

to reject H_0 , or
not to reject H_0 .

In latter case, we claim that there is **no** significant difference between the null hypothesis and the underlying distribution according to the observed data.

Remark. The setting is **not** symmetric in the sense that the two hypotheses H_0 and H_1 will be treated differently.

Some specific terms used statistical tests are now in order:

- A *Hypothesis* is an assumption made about the value of $\boldsymbol{\theta}$.
- A *Simple Hypothesis* completely specifies $\boldsymbol{\theta}$ to a known constant such as $H: \boldsymbol{\theta} = \boldsymbol{\theta}_0$.
- A *Composite Hypothesis* is any hypothesis that is not simple. For instance the hypothesis $H: \boldsymbol{\theta} \in [1, \infty)$.
- A *Test Statistic* is a statistic based on which the decision is made.
- A *Null Hypothesis*, usually denoted as H_0 , is the basic or default hypothesis.
- The *Alternative Hypothesis*, usually denoted as H_1 , is a hypothesis that is to be compared with the Null Hypothesis.

Mathematically speaking, a statistical test is equivalent to define a binary function

$$\phi(\mathbf{X}) = \begin{cases} 1 & \text{when reject } H_0, \\ 0 & \text{when not reject } H_0. \end{cases}$$

$\phi(\cdot)$ is called a **decision rule**, which practically divides the sample space of \mathbf{X} into two subspaces: *rejection region* (which is also called ‘critical region’) and its compliment which is often **imprecisely** labelled as ‘acceptance’ region.

There are two types of errors that can be made in a statistical test, which are displayed in the table below.

		Decision Made	
		H_0 not rejected Θ_0	H_0 rejected Θ_1
True State of Nature	H_0 Θ_0	Correct Decision	Type I Error
	H_1 Θ_1	Type II Error	Correct Decision

Ideally we would like to have a test that minimises the probabilities of making both types of errors, which unfortunately is not feasible.

Construction of a statistical test:

we control the probability of Type I error under a given level, say, α , and then we minimize the probability of Type II error.

Definition. A test is said to have *significance level* α if

$$\sup_{\theta \in \Theta_0} P_{\theta}(H_0 \text{ is rejected}) \leq \alpha.$$

Definition. A test T is said to be of *size* α if

$$\sup_{\theta \in \Theta_0} P_{\theta}(H_0 \text{ is rejected}) = \alpha.$$

Remark. (i) In practice, we usually take the significance level $\alpha = 0.1, 0.05$ or 0.01 .

(ii) The size of the test is no greater than the significance level.

Definition. The *power function* of a test is defined as

$$\beta(\theta) = P_{\theta}(H_0 \text{ is rejected}), \quad \theta \in \Theta_1.$$

In practice, it is useful to extend $\beta(\theta)$ for all $\theta \in \Theta$. For $\theta \in \Theta_0$, $\beta(\theta)$ must be smaller than the prescribed significance level, and is the probability of making a Type I error at the parameter value θ . For $\theta \in \Theta_1$, $1 - \beta(\theta)$ is the probability of making a Type II error.

If a test is defined in terms of decision rule ϕ , then

$$\beta(\theta) = E_{\theta}[\phi(\mathbf{X})],$$

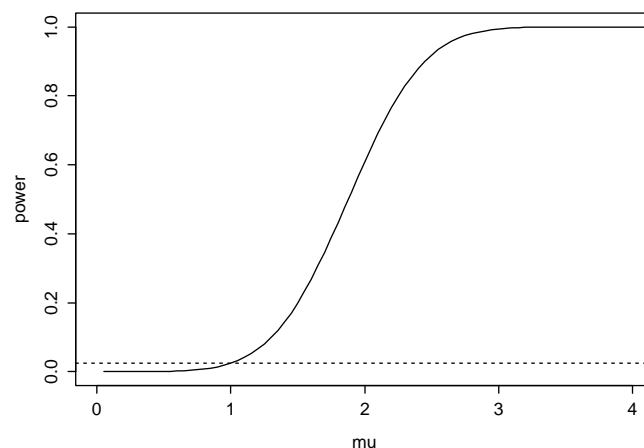
because

$$\begin{aligned} E[\phi(\mathbf{X})] &= \int \phi(\mathbf{x})f(\mathbf{x}; \theta)d\mathbf{x} \\ &= P_{\theta}\{\phi(\mathbf{X}) = 1\} = P_{\theta}\{H_0 \text{ is rejected}\}. \end{aligned}$$

Example 1. Suppose that we have a random sample of size 5 from a normal distribution with mean μ and variance 1. Let us consider:

$$\begin{aligned} H_0: & \quad \mu = 1 \\ H_1: & \quad \mu > 1 \end{aligned}$$

One possible test statistic is \bar{X} . Intuitively we would reject H_0 for large values of \bar{X} , say, $\bar{X} > K$. The constant K is determined by the prescribed significance level α . Let $\alpha = 0.025$. Since $\sqrt{5}(\bar{X} - 1) \sim N(0, 1)$, $K = 1 + 1.96/\sqrt{5}$. We reject H_0 if $\bar{X} > 1 + 1.96/\sqrt{5}$. The power function is shown below.



The power is small for $\mu = 1$, and rises quite fast as μ increases past 1.

Question: What should we do if we change the alternative hypothesis to $H_1 : \mu < 1$ or $H_1 : \mu \neq 1$.

The above null hypothesis is simple, because $\Theta_0 = \{1\}$, but the alternative hypothesis is composite because $\Theta_1 = (1, \infty)$. If instead we used the Null Hypothesis $H_0 : \mu \leq 1$, that would be a composite null hypothesis. It is convenient to have one value to summarise the probability of Type I error, even when the null hypothesis is composite.

Question: Is the test constructed above still reasonable for testing

$$H_0 : \mu \leq 1 \quad \text{vs} \quad H_1 : \mu > 1 ?$$

If so, at which value of μ under H_0 does the probability of type I error obtain the maximum?

3.2 Neyman-Pearson Lemma for testing simple hypotheses

Most powerful test (MPT): among all the tests of size α , the one with maximum power $\beta(\theta)$ for all $\theta \in \Theta_1$ is called the MPT. When Θ_1 consists of more than one point, it is often called the uniformly MPT (UMPT).

Note. UMPT may not exist.

Let $L(\theta; \mathbf{X})$ be the likelihood function with observation \mathbf{X} .

Theorem 1. (Neyman-Pearson Lemma)

The MPT for hypotheses

$$H_0 : \theta = \theta_0 \quad \text{against} \quad H_1 : \theta = \theta_1$$

rejects H_0 when

$$L(\theta_1; \mathbf{X}) > K L(\theta_0; \mathbf{X}),$$

and does not reject H_0 when

$$L(\theta_1; \mathbf{X}) < K L(\theta_0; \mathbf{X}),$$

where constant K is determined by the size of the test.

Remark. In general the MPT rejects H_0 if the likelihood ratio

$$LR(\mathbf{X}) = \frac{L(\boldsymbol{\theta}_1; \mathbf{X})}{L(\boldsymbol{\theta}_0; \mathbf{X})} > K,$$

which indicates that $\boldsymbol{\theta}_1$ is relatively favoured over $\boldsymbol{\theta}_0$ according to the likelihood function.

Proof. For each test T , we define the decision rule $\phi_T(\mathbf{x})$ as follows

$$\phi_T(\mathbf{X}) = \begin{cases} 1 & \text{Reject } H_0 \\ 0 & \text{Do not Reject } H_0. \end{cases}$$

The power $\beta(\boldsymbol{\theta})$ is equal to $E\phi_T(\mathbf{X})$.

Suppose that T is a test of size α that satisfies the conditions of the Neyman-Pearson Lemma, and S is some other test of size α . Then

$$\begin{aligned} \beta_S(\boldsymbol{\theta}_1) &= E\phi_S(\mathbf{x}) = \int \phi_S(\mathbf{x})f(\mathbf{x}; \boldsymbol{\theta}_1)d\mathbf{x} \\ &= \int \phi_S(\mathbf{x})f(\mathbf{x}; \boldsymbol{\theta}_1)d\mathbf{x} - K \int \phi_S(\mathbf{x})f(\mathbf{x}; \boldsymbol{\theta}_0)d\mathbf{x} + K\alpha \\ &= \int \phi_S(\mathbf{x})[f(\mathbf{x}; \boldsymbol{\theta}_1) - Kf(\mathbf{x}; \boldsymbol{\theta}_0)]d\mathbf{x} + K\alpha \\ &= \int \phi_S(\mathbf{x})[L(\boldsymbol{\theta}_1; \mathbf{x}) - KL(\boldsymbol{\theta}_0; \mathbf{x})]d\mathbf{x} + K\alpha \\ &\leq \int_{\{L(\boldsymbol{\theta}_1; \mathbf{x}) > KL(\boldsymbol{\theta}_0; \mathbf{x})\}} \phi_S(\mathbf{x})[L(\boldsymbol{\theta}_1; \mathbf{x}) - KL(\boldsymbol{\theta}_0; \mathbf{x})]d\mathbf{x} + K\alpha \\ &\leq \int \phi_T(\mathbf{x})[L(\boldsymbol{\theta}_1; \mathbf{x}) - KL(\boldsymbol{\theta}_0; \mathbf{x})]d\mathbf{x} + K\alpha \\ &= \int \phi_T(\mathbf{x})[f(\mathbf{x}; \boldsymbol{\theta}_1) - Kf(\mathbf{x}; \boldsymbol{\theta}_0)]d\mathbf{x} + K\alpha \\ &= \int \phi_T(\mathbf{x})f(\mathbf{x}; \boldsymbol{\theta}_1)d\mathbf{x} - K \int \phi_T(\mathbf{x})f(\mathbf{x}; \boldsymbol{\theta}_0)d\mathbf{x} + K\alpha \\ &= \beta_T(\boldsymbol{\theta}_1) - K\alpha + K\alpha \\ &= \beta_T(\boldsymbol{\theta}_1). \end{aligned}$$

Example 2. Suppose that X_1, X_2, \dots, X_n are a random sample from $N(\mu, \sigma^2)$, where we treat σ^2 as a known quantity. Suppose we want to test:

$$H_0 : \mu = 0 \quad \text{against} \quad H_1 : \mu = 5.$$

The likelihood ratio is

$$\begin{aligned} LR &= \frac{\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n (X_i - 5)^2/(2\sigma^2)\right)}{\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\sum_{i=1}^n X_i^2/(2\sigma^2)\right)} \\ &= \exp\left(5n(2\bar{X} - 5)/(2\sigma^2)\right). \end{aligned}$$

We can look for a most powerful test by noting that $LR > K$ is equivalent to $\bar{X} > K_1$, and the critical value K_1 is determined by the size of the test. So, if we find a test of this form which has size α , it will

be a most powerful test of size α . It is easy to see that the test that rejects H_0 iff $\bar{X} > z_\alpha \sigma / \sqrt{n}$ is of this form and is therefore a most powerful test of size α .

Question: If we change the alternative hypothesis to $H_1 : \mu = 50$, what is the MPT then?

Example 3. Let (X_1, \dots, X_n) be a sample from Poisson distribution with mean μ . To test

$$H_0 : \mu = 2 \quad \text{against} \quad H_1 : \mu = 6,$$

the likelihood ratio is

$$LR = \frac{6^{\sum X_i} e^{-6n} / (X_1! \dots X_n!)}{2^{\sum X_i} e^{-2n} / (X_1! \dots X_n!)} = 3^{\sum X_i} e^{-4n}.$$

According to the N-P lemma, the MPT will reject H_0 if $LR > K$, or equivalently, $\sum_i X_i > K_1$, where the critical value K_1 is determined by the significance level of the test.

Suppose $n = 4$, then $\sum_{i=1}^4 X_i \sim \text{Poisson}(8)$ under H_0 . From the table for Poisson CDFs, the size of test $\alpha = 0.064$ with $K_1 = 12$, and $\alpha = 0.034$ with $K_1 = 13$.

When $n = 8$, $\sum_{i=1}^8 X_i \sim \text{Poisson}(16)$ under H_0 . The size of test $\alpha = 0.058$ with $K_1 = 22$, and $\alpha = 0.037$ with $K_1 = 23$.

3.3 Uniformly Most Powerful Tests

Although the Neyman-Pearson Lemma was designed for testing simple hypotheses, it is possible to use it to construct the UMPT for, for example, *one-sided null hypothesis against one-sided alternative*.

A general setting: For hypotheses

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_1 : \theta \in \Theta_1, \tag{4}$$

$\phi(\mathbf{X})$ is the decision rule of a test at significant level α , i.e.

$$E_\theta\{\phi(\mathbf{X})\} \leq \alpha \quad \text{for all } \theta \in \Theta_0.$$

If $\phi(\mathbf{X})$ is also the MPT of size α for *simple hypotheses*

$$H_0 : \theta = \theta_0 \quad \text{vs} \quad H_1 : \theta = \theta_1$$

for **some** $\theta_0 \in \Theta_0$ and **all** $\theta_1 \in \Theta_1$. Then $\phi(\mathbf{X})$ is the UMPT for hypotheses (4).

We look at a simple scenario first — *UMPTs for simple H_0 and one-sided H_1* .

Suppose the MPT for simple hypotheses

$$H_0 : \theta = \theta_0 \quad \text{against} \quad H_1 : \theta = \theta_1$$

does not change its form for all $\theta_1 \in \Theta_1$. Then it is also the UMPT for

$$H_0 : \theta = \theta_0 \quad \text{against} \quad H_1 : \theta \in \Theta_1.$$

Such an UMPT often exists for $\Theta_1 = (\theta_0, \infty)$ or $\Theta_1 = (-\infty, \theta_0)$.

Example 4. If X_1, X_2, \dots, X_n is a random sample from $N(\mu, 1)$, then the most powerful test for testing

$$H_0 : \mu = 1 \quad \text{vs} \quad H_1 : \mu = 2$$

is obtained by rejecting H_0 if and only if $\bar{X} > 1 + z_\alpha / \sqrt{n}$. The same test is obtained for

$$H_0 : \mu = 1 \quad \text{vs} \quad H_1 : \mu = \mu_1$$

for every $\mu_1 > 1$. So this test is uniformly most powerful for

$$H_0 : \mu = 1 \quad \text{vs} \quad H_1 : \mu > 1.$$

Although the UMPT does not depend on the value μ_1 specified under H_1 , its power varies over $\{\mu | \mu > 1\}$. In fact

$$\begin{aligned} \beta(\mu) &= P_\mu\{\bar{X} > 1 + z_\alpha/\sqrt{n}\} \\ &= P_\mu\{\sqrt{n}(\bar{X} - \mu) > \sqrt{n}(1 - \mu) + z_\alpha\} \\ &= 1 - \Phi\{z_\alpha - \sqrt{n}(\mu - 1)\} = \Phi\{\sqrt{n}(\mu - 1) - z_\alpha\}, \end{aligned}$$

where $\Phi(\cdot)$ is the CDF of $N(0,1)$. Thus the power increases as μ increases.

Note that the test is not the UMPT for $H_1 : \mu \neq 1$.

Note. The UMPT usually does not exist for two-sided (composite) alternative hypothesis.

Example 5. Let Y have the binomial distribution $\text{Bin}(n, p)$. Find the UMPT for testing

$$H_0 : p = p_0 \quad \text{vs} \quad H_1 : p > p_0.$$

Let $p_1 > p_0$. The LR for testing the H_0 above against $H_1 : p = p_1$ is

$$\begin{aligned} LR &= \frac{\binom{n}{Y} p_1^Y (1 - p_1)^{n-Y}}{\binom{n}{Y} p_0^Y (1 - p_0)^{n-Y}} \\ &= \left(\frac{p_1(1 - p_0)}{p_0(1 - p_1)} \right)^Y \left(\frac{1 - p_1}{1 - p_0} \right)^{n-Y} \end{aligned}$$

Note that $p_1(1 - p_0) > p_0(1 - p_1)$ since $p_1 > p_0$. Thus $LR > K$ is equivalent to $Y > K_1$. Thus the UMPT rejects H_0 iff $Y > K_1$ with the size

$$P(Y > K_1 | p = p_0).$$

Example 6. Let (X_1, \dots, X_n) be a random sample from an exponential distribution with mean $1/\lambda$. We are interested in testing

$$H_0 : \lambda \leq \lambda_0 \quad \text{vs} \quad H_1 : \lambda > \lambda_0.$$

For

$$H_0 : \lambda = \lambda_0 \quad \text{vs} \quad H_1 : \lambda = \lambda_1,$$

the MPT rejects H_0 iff $\sum_{i=1}^n X_i \leq K$ for any $\lambda_1 > \lambda_0$, where K is determined by

$$P_{\lambda_0}\left\{\sum_{i=1}^n X_i < K\right\} = \alpha.$$

It is easy to verify that for $\lambda < \lambda_0$,

$$P_\lambda\left\{\sum_{i=1}^n X_i < K\right\} < \alpha.$$

Hence the MPT for the simple null hypothesis against simple alternative is also the UMPT for the composite hypotheses.

3.4 Likelihood Ratio Tests

We now deal with the most popular ways of constructing tests when both null and alternative hypotheses are composite. There are no guaranteed optimality properties in small samples for these tests, but in regular cases they usually have good power for large sample sizes.

Let $\mathbf{X} \sim f(\cdot, \boldsymbol{\theta})$. Consider hypotheses

$$H_0 : \boldsymbol{\theta} \in \Theta_0 \quad \text{vs} \quad H_1 : \boldsymbol{\theta} \in \Theta - \Theta_0.$$

The likelihood ratio test will reject H_0 for the large values of the statistic

$$\begin{aligned} LR = LR(\mathbf{X}) &\equiv \frac{\sup_{\boldsymbol{\theta} \in \Theta} f(\mathbf{X}, \boldsymbol{\theta})}{\sup_{\boldsymbol{\theta} \in \Theta_0} f(\mathbf{X}, \boldsymbol{\theta})} \\ &= f(\mathbf{X}, \hat{\boldsymbol{\theta}}) / f(\mathbf{X}, \tilde{\boldsymbol{\theta}}), \end{aligned}$$

where $\hat{\boldsymbol{\theta}}$ the (unconstrained) MLE, and $\tilde{\boldsymbol{\theta}}$ is the constrained MLE under hypothesis H_0 .

Remark. (i) It is easy to see that $LR \geq 1$.

(ii) The exact sampling distributions of LR are usually unknown, except in a few special cases.

Example 7. (One-sample t -test)

Let $\mathbf{X} = (X_1, \dots, X_n)^\tau$ be a random sample from $N(\mu, \sigma^2)$. We are interested in testing hypotheses

$$H_0 : \mu = \mu_0 \quad \text{against} \quad H_1 : \mu \neq \mu_0,$$

where μ_0 is given, and σ^2 is unknown and is a nuisance parameter. Now both H_0 and H_1 are composite. The likelihood function is

$$L(\mu, \sigma^2) = C\sigma^{-n} \exp \left\{ -\frac{1}{2}\sigma^2 \sum_{j=1}^n (X_j - \mu)^2 \right\}.$$

The unconstrained MLEs are

$$\hat{\mu} = \bar{X}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \bar{X})^2,$$

and the constrained MLE is

$$\tilde{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (X_j - \mu_0)^2.$$

The LR-ratio statistic is then

$$LR = \frac{L(\hat{\mu}, \hat{\sigma}^2)}{L(\mu_0, \tilde{\sigma}^2)} = (\tilde{\sigma}^2 / \hat{\sigma}^2)^{n/2}.$$

Since

$$n\tilde{\sigma}^2 = n\hat{\sigma}^2 + n(\bar{X} - \mu_0)^2,$$

it holds that $\tilde{\sigma}^2 / \hat{\sigma}^2 = 1 + T^2 / (n - 1)$, where

$$T = \sqrt{n}(\bar{X} - \mu_0) / \left\{ \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X})^2 \right\}^{1/2}.$$

Note that $T \sim t_{n-1}$ under H_0 . The LRT will reject H_0 iff $|T| > t_{n-1, \alpha/2}$, where $t_{k, \alpha}$ is the upper α -point of the t -distribution with k degrees of freedom.

Asymptotic Distribution of Likelihood ratio test statistic

Let $\mathbf{X} = (X_1, \dots, X_n)^\tau$, and assume certain regularity conditions. Then as $n \rightarrow \infty$, the distribution

of $2\log(LR)$ under H_0 converges to the χ^2 -distribution with $d - d_0$ degrees of freedom, where d is the ‘dimension’ of Θ and d_0 is the ‘dimension’ of Θ_0 .

To make the computation of ‘dimension’ easy, reparametrisation is often adopted. Suppose that the parameter θ may be written in two parts

$$\theta = (\psi, \lambda)$$

where ψ is $k \times 1$ parameter of interest, and λ is of little interest and is called *nuisance parameters*. The hypotheses to be tested may be expressed as

$$H_0 : \psi = \psi_0 \quad \text{vs} \quad H_1 : \psi \neq \psi_0.$$

Now the LR-statistic is of the form

$$LR = \frac{L(\hat{\psi}, \hat{\lambda}; \mathbf{X})}{L(\psi_0, \tilde{\lambda}; \mathbf{X})},$$

where $(\hat{\psi}, \hat{\lambda})$ is unconstrained MLE while $\tilde{\lambda}$ is the constrained MLE of λ subject to $\psi = \psi_0$. Then as $n \rightarrow \infty$,

$$2\log(LR) \xrightarrow{D} \chi_k^2 \quad \text{under } H_0.$$

Example 8. Let X_1, \dots, X_n be independent, and $X_j \sim N(\mu_j, 1)$. Consider the null hypothesis

$$H_0 : \mu_1 = \dots = \mu_n.$$

The likelihood function is

$$L(\mu_1, \dots, \mu_n) = C \exp \left\{ -\frac{1}{2} \sum_{j=1}^n (X_j - \mu_j)^2 \right\},$$

where $C > 0$ is a constant independent of μ_j . Then the unconstrained MLE are $\hat{\mu}_j = X_j$ and the constrained MLE is $\tilde{\mu} = \bar{X}$. Hence

$$\begin{aligned} LR &= \frac{L(\hat{\mu}_1, \dots, \hat{\mu}_n)}{L(\tilde{\mu}, \dots, \tilde{\mu})} \\ &= \exp \left\{ \frac{1}{2} \sum_{j=1}^n (X_j - \bar{X})^2 \right\}. \end{aligned}$$

Hence

$$2\log(LR) = \sum_{j=1}^n (X_j - \bar{X})^2 \sim \chi_{n-1}^2 \quad \text{under } H_0,$$

which is true for any finite n as well.

How to *calculate the degree of freedom*?

Since

$$d = n, \quad d_0 = 1,$$

the d.f. is $d - d_0 = n - 1$.

Alternatively we may adopt the following reparametrisation:

$$\mu_j = \mu_1 + \psi_j \quad \text{for } 2 \leq j \leq n.$$

Then the null hypothesis can be expressed as

$$H_0 : \psi_2 = \cdots = \psi_n = 0.$$

Therefore $\boldsymbol{\psi} = (\psi_2, \cdots, \psi_n)^\tau$ has $n - 1$ component, i.e. $k = n - 1$.

4. INTERVAL ESTIMATION

Interval estimation is more informative than point estimation, and is very important in practice.

Confidence Interval is the most commonly used interval estimation. More general type is *Confidence Set* which may consist of several intervals.

4.1 What is an interval estimator

Example 1. Let us start with a simple example. A random sample X_1, \dots, X_n are drawn from $N(\mu, 1)$. Then

$$\sqrt{n}(\bar{X} - \mu) \sim N(0, 1).$$

Hence

$$P(-1.96 \leq \sqrt{n}(\bar{X} - \mu) \leq 1.96) = 0.95,$$

or

$$P(\bar{X} - 1.96/\sqrt{n} < \mu < \bar{X} + 1.96/\sqrt{n}) = 0.95.$$

So a 95% confidence interval for μ is

$$(\bar{X} - 1.96/\sqrt{n}, \bar{X} + 1.96/\sqrt{n}).$$

Suppose $n = 4$, $\bar{X} = 2.25$. Then a 95% C.I. is $(2.25 - 0.98, 2.25 + 0.98) = (1.27, 3.23)$.

Question: what is $P(1.27 < \mu < 3.23)$? — Note μ is a unknown constant!

Answer: $(1.27, 3.23)$ is one realisation of the random interval $(\bar{X} - 0.98, \bar{X} + 0.98)$ which covers μ with probability 0.95.

Definition. Suppose the joint distribution of $\mathbf{X} = (X_1, \dots, X_n)$ depends some parameter θ .

If $\theta_l(\mathbf{X})$ and $\theta_u(\mathbf{X})$ are two statistics for which

$$P\{\theta_l(\mathbf{X}) < \theta < \theta_u(\mathbf{X})\} = 1 - \alpha,$$

$(\theta_l(\mathbf{X}), \theta_u(\mathbf{X}))$ is called a $100(1 - \alpha)\%$ *confidence interval* for θ .

Remark. (i) $1 - \alpha$ is called the confidence level, which is usually set at 0.95, 0.99 or 0.999. Naturally for given α , we shall search for the interval with the shortest length $\theta_u(\mathbf{X}) - \theta_l(\mathbf{X})$, which gives the **most accurate** estimation.

(ii) We may have $\theta_l(\mathbf{X}) = -\infty$ or $\theta_u(\mathbf{X}) = \infty$, giving a one-sided interval.

(iii) In general we may use non-interval type set $S(\mathbf{X}) \subset \Theta$ as the estimator for θ , i.e.

$$P\{\theta \in S(\mathbf{X})\} = 1 - \alpha.$$

We call $S(\mathbf{X})$ a $100(1 - \alpha)\%$ **confidence set**.

4.2 Method of pivotal functions

Set-up: Observations – \mathbf{X} , Parameter of interest – θ (possible other parameters).

Definition. A function of \mathbf{X} and θ alone is a *pivotal function* if its distribution does not depend on any unknown parameters (i.e. can be calculated numerically).

Examples of pivotal functions: Let X_1, \dots, X_n be a random sample.

1. for exponential distribution with mean θ , $\theta^{-1} \sum_j X_j$ is a pivotal,
2. for $N(\mu, 1)$, $\bar{X} - \mu$ is a pivotal,
3. for $N(\mu, \sigma^2)$, $\sqrt{n}(\bar{X} - \mu)/S$ is a pivotal.

Construction of C.I. based on a pivotal:

Step 1: Find a pivotal function $T = T(\mathbf{X}, \theta)$, and identify its distribution

Step 2: Use the distribution to identify a range of values t_1, t_2 such that

$$P(t_1 < T < t_2) = 1 - \alpha.$$

Step 3: Manipulate the inequalities

$$T > t_1 \quad \text{and} \quad T < t_2$$

to find a set of values for θ . The values included in this set form a $100(1 - \alpha)\%$ C.I. for θ .

Example 2. (Student's t -interval)

(X_1, \dots, X_n) is a random sample from $N(\mu, \sigma^2)$. A pivotal function is

$$T = \sqrt{n}(\bar{X} - \mu)/S \sim t_{n-1},$$

where $S^2 = \frac{1}{n-1} \sum_j (X_j - \bar{X})^2$.

For $n = 10$, we have $P(-2.26 < T < 2.26) = 0.95$. Note that $T > -2.26$ is equivalent to

$$\mu < \bar{X} + 2.26S/\sqrt{n}.$$

(Similar for $T < 2.26$.) A 95% C.I. for μ is $(\bar{X} - 2.26S/\sqrt{n}, \bar{X} + 2.26S/\sqrt{n})$.

Note. (i) In general it is not easy to identify a pivotal function.

(ii) For location parameter μ , $\bar{X} - \mu$ is **likely** to be pivotal.

(iii) The MLE is asymptotically pivotal:

$$\hat{\theta} \sim N(\theta, \mathcal{I}^{-1}(\theta)).$$

Hence, $\sqrt{\mathcal{I}(\hat{\theta})}(\hat{\theta} - \theta) \sim N(0, 1)$ approximately for large n .

4.3 Inverting Tests

One easy way to get a confidence set is to invert a test; see Example 1.

Theorem. Suppose that we have a size α test for the Null Hypothesis $H_0 : \theta = \theta_0$. For each $\theta_0 \in \Theta$, suppose the set $A(\theta_0)$ is the collection of \mathbf{x} for which H_0 is *not rejected*. Then the set

$$S(\mathbf{x}) = \{\theta : \mathbf{x} \in A(\theta), \theta \in \Theta\}$$

is a family of confidence sets for θ at confidence level $1 - \alpha$.

The proof of this is obvious, because when θ is the true parameter value

$$P_\theta[\theta \in S(\mathbf{X})] = P_\theta[\mathbf{X} \in A(\theta)] = 1 - \alpha.$$

Example 3. Suppose that X_1, X_2, \dots, X_n is a random sample from an exponential distribution with mean $1/\lambda$. The UMPT with size 5% for

$$H_0 : \lambda = \lambda_0 \quad \text{against} \quad H_1 : \lambda > \lambda_0$$

is to reject H_0 when $\sum X_i < \frac{1}{2\lambda_0} \chi_{(2n),0.05}^2$, where $\chi_{(2n),0.05}^2$ is the lower 5% point of the distribution χ_{2n}^2 . The acceptance region here is $\sum X_i > \frac{1}{2\lambda_0} \chi_{(2n),0.05}^2$, which on inversion gives

$$\left(\frac{\chi_{(2n),0.05}^2}{2 \sum X_i}, \infty \right)$$

as the confidence set for λ at the 95% level of confidence.

4.4 Bootstrap confidence intervals

Let X_1, \dots, X_n be sample from unknown distribution of F . We are interested in constructing confidence intervals for some characteristic $\theta = \theta(F)$ of distribution F .

We adopt the so-called *nonparametric bootstrap method* now. It draws a bootstrap sample

$$X_1^*, \dots, X_n^*$$

independently from the uniform distribution over n discrete points $\{X_1, \dots, X_n\}$.

Now in principle the distribution of any statistic based on (X_1^*, \dots, X_n^*) is known, *conditionally on* (X_1, \dots, X_n) .

4.4.1 Bootstrap percentiles

Let $\hat{\theta} = T(X_1, \dots, X_n)$ be an estimator for θ . Define

$$\hat{\theta}^* = T(X_1^*, \dots, X_n^*).$$

Let l_α^* and u_α^* be, respectively, the lower and upper α -point of the distribution $\hat{\theta}^*$, i.e.

$$P\{\hat{\theta}^* \leq l_\alpha^* | X_1, \dots, X_n\} = \alpha,$$

$$P\{\hat{\theta}^* > u_\alpha^* | X_1, \dots, X_n\} = \alpha.$$

Then **the $(1 - \alpha)$ 100%-th bootstrap interval** for θ is defined as

$$(l_{\alpha/2}^*, u_{\alpha/2}^*].$$

In practice we draw B sets of bootstrap samples for some large B , resulting in B bootstrap estimates $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$. Then $l_{\alpha/2}^*$ and $u_{\alpha/2}^*$ are, respectively, the $[B\alpha/2]$ -th smallest and the $[B\alpha/2]$ -th largest values among $\hat{\theta}_1^*, \dots, \hat{\theta}_B^*$.

Note. Both $l_{\alpha/2}^*$ and $u_{\alpha/2}^*$ are statistics depending on $\{X_1, \dots, X_n\}$.

Intuitively, we **must** start with a *good* estimator $\hat{\theta}$. In fact, if $\hat{\theta} = T(X_1, \dots, X_n)$ fulfils some conditions, it holds that as $n \rightarrow \infty$,

$$P(l_{\alpha/2}^* < \theta \leq u_{\alpha/2}^*) \rightarrow 1 - \alpha.$$

4.4.2 Bootstrap- t

Suppose we have a *studentised 'pivot'*

$$V = (\hat{\theta} - \theta) / \hat{\sigma},$$

where $\hat{\theta}$ is an estimator for θ , and $\hat{\sigma}^2$ is an estimator for $\text{Var}(\hat{\theta})$. Define its bootstrap version

$$V^* = (\hat{\theta}^* - \hat{\theta}) / \hat{\sigma}^*.$$

Let

$$\begin{aligned} \alpha/2 &= P\{V^* \leq l_{\alpha/2}^* | X_1, \dots, X_n\} \\ &= P\{V^* > u_{\alpha/2}^* | X_1, \dots, X_n\}. \end{aligned}$$

Then the $100(1 - \alpha)\%$ -th bootstrap- t interval for θ is

$$[\hat{\theta} - \hat{\sigma}u_{\alpha/2}^*, \hat{\theta} - \hat{\sigma}l_{\alpha/2}^*].$$

Remark. The estimator $\hat{\sigma}^2$ may be obtained by bootstrap, then $\hat{\sigma}^{2*}$ will be calculated via a **nested** bootstrap.

Ideally we require that the distribution, or at least the asymptotic distribution, of V does not depend on anything unknown. This will make the approximation

$$P\{\hat{\theta} - \hat{\sigma}u_{\alpha/2}^* \leq \theta < \hat{\theta} - \hat{\sigma}l_{\alpha/2}^*\} \approx \alpha$$

more accurate.

We give a heuristic argument for the above assertion.

If V is an asymptotic pivot in the sense that

$$P(V \leq x) = \Phi(x) + \frac{1}{\sqrt{n}}g(x, \theta) + O\left(\frac{1}{n}\right), \quad (5)$$

similarly in the bootstrap world we have

$$\begin{aligned} P(V^* \leq x | X_1, \dots, X_n) &= \Phi(x) + \frac{1}{\sqrt{n}}g(x, \hat{\theta}) + O_p\left(\frac{1}{n}\right) \\ &= \Phi(x) + \frac{1}{\sqrt{n}}g(x, \theta) + O_p\left(\frac{1}{n}\right). \end{aligned} \quad (6)$$

Hence

$$l_\alpha = l_\alpha^* + O_p\left(\frac{1}{n}\right).$$

In contrast if $\Phi(x)$ depends on θ , then $\Phi(x)$ should be replaced by $\Phi(x, \theta)$ in (5), and by $\Phi(x, \hat{\theta})$ in (6). Now

$$l_\alpha = l_\alpha^* + O_p\left(\frac{1}{\sqrt{n}}\right).$$