

Example 2 (cont'd)

Indeed, let $Y_i \sim N(\mu, \sigma^2)$. Then \bar{Y} and S^2 are independent and

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

It follows that $\text{var}(S^2) = 2(\sigma^2)^2/(n-1)$. Furthermore, combining this with the results of Example 1, we have that

$$\frac{\frac{\sqrt{n}(\bar{Y}-\mu)}{\sigma}}{\sqrt{\frac{(n-1)S^2}{\sigma^2}/(n-1)}} = \frac{\sqrt{n}(\bar{Y}-\mu)}{S} \sim t_{n-1}.$$

Quantities s.t. $(n-1)S^2/\sigma^2$ or $\sqrt{n}(\bar{Y}-\mu)/S^2$, which are not statistics but have a known distribution, are called pivotal quantities.

They play a key role in the development of confidence intervals and hypothesis testing as we will see in L5 and L6.

Example 3

Let Y_1, \dots, Y_n be an iid random sample with $Y_i \sim \text{Unif}(0, \theta)$, with $\theta > 0$. For the statistics \bar{Y} and $Y_{(n)}$, let's compute some of their features.

For the sample average we have

$$\begin{aligned} E(\bar{Y}) &= E(Y_1) = \theta/2, \quad \text{var}(\bar{Y}) = E((\bar{Y})^2) - E(\bar{Y})^2 = \frac{\text{var}(\text{Unif})}{n} = \theta^2/(12n). \end{aligned}$$

For the maximum we have

$$\begin{aligned} E(X_{(n)}) &= \int_{-\infty}^{+\infty} t f_{X_{(n)}}(t) dt = \int_0^\theta t \cdot n F(t)^{n-1} f(t) dt \\ &= \int_0^\theta t n (t/\theta)^{n-1} (1/\theta) dt = \theta n / (n+1), \end{aligned}$$

The larger the sample size, the closer we are to θ

and

$$\text{var}(Y_{(n)}) = \theta^2 n / (n+1)^2.$$

Note that in the case of $Y_{(n)}$ we had to use it's pdf in order to compute the two moments.

Some statistics target θ the parameter of a distribution (or a component of it).

In that case, we call them estimators and denote them by a Greek letter with a hat, e.g. $\hat{\theta}$.

An estimator thus is a function of the random sample, and can tell us something useful about the parameter θ .

For example, the sample average is useful when we want to learn about the population average μ , the sample median Q_2 is useful for learning about the population median and so on.

→ hist is an estimate of the pdf

We now look at methods for computing estimators and then we will see methods for comparing estimators.

Method of Moments

Useful when we want to estimate θ that can be expressed as a function of moments of Y .

Is one of the oldest statistical estimation methods, dating back to 1936.

The method consists in equating sample moments, e.g. \bar{Y} , $\overline{Y^2}$, $\overline{Y^3}$, ... with the corresponding population moments, e.g. $E(Y)$, $E(Y^2)$, $E(Y^3)$ and solving these equations in terms of the parameter θ .

Example 4

Suppose Y_1, \dots, Y_n is an iid sample from some distribution F_θ and let $E(Y_1) = \mu$ be the unknown parameter.

Equating the sample moment with the corresponding population moment leads to

$$\underline{\bar{Y} = E(Y) = \mu.}$$

So \bar{Y} is the method of moment estimator for μ , i.e. $\hat{\mu}_{MM} = \bar{Y}$ (reads: the method of moments estimator for μ is \bar{Y}).

Example 5

Let Y_1, \dots, Y_n be an iid random sample with $E(Y_1) = \mu$ and $\text{var}(Y_1) = \sigma^2$, with μ, σ^2 unknown.

First note that $\sigma^2 = E(Y^2) - E(Y)^2 = E(Y^2) - \mu^2$. Equating the first two moments leads to

estimator for μ

$$\bar{Y} = E(Y) = \mu,$$

$$E(Y^2) = \sigma^2 + \mu^2 = \overline{Y^2}.$$

"method of moments"

estimator for σ^2

We conclude that

$$\hat{\sigma}_{MM}^2 = \overline{Y^2} - \bar{Y}^2 = \hat{\sigma}^2,$$

(we have encountered this estimator previously.) and $\hat{\mu}_{MM} = \bar{Y}$.

So $(\bar{Y}, \hat{\sigma}^2)$ is the method of moments estimator for (μ, σ^2) .

A signal+noise problem

MM. not useful for not i.i.d. samples \Rightarrow used as a starting point.

In many situations, we have measurements Y_i , which can be thought of as the sum of a physical signal $g_i(\beta)$ and a noise ϵ_i , i.e.

$$Y_i = g_i(\beta) + \epsilon_i,$$

where ϵ_i is some unknown component whereas $g_i(\beta)$ is the signal, which may depend on some unknown parameter β .

For instance, when measuring the resistance of an electronic circuit with a multimeter, we observe a realisation of Y_i , say $y_i = 12 \Omega$.

This value depends, in part, on the real resistance of the equipment $g_i(\beta)$ (e.g. the sum of the resistance of all its component) and for the other part, on the accuracy of the instrument; see also L2 for another example.

Linear regression

The signal $g_i(\beta)$ may depend on some known features x_{i1}, \dots, x_{ip} , through the linear function \rightarrow *not suitable for all applications*

$$g_i(\beta) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, \quad i = 1, \dots, n,$$

with unknown parameters $\beta = (\beta_0, \dots, \beta_p)$.

The aim is to understand the impact of x_{ij} 's on Y_i , i.e. β_i .

We have the pair $y_i, g_i(\beta)$, where y_i is what we measure and $g_i(\beta)$ is how the system should behave according to modeller's view.

This is commonly known as a linear regression problem, and one of the points is how to estimate β using $y_i, g_i(\beta)$ for all i .

Method of Least Squares

For any fixed β , the deviances $y_i - g_i(\beta)$, tell's us by how much our model $g_i(\beta)$ misses the observed value. It seems intuitive then to look for a β that leads to smallest deviances.

The method of least squares consists in estimating β by the value that leads to smallest sum of squared deviances.

That is, the least squares estimator is defined by

$$\hat{\beta}_{LS} = \arg \inf_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - g_i(\beta))^2.$$

Example 6

Let y_1, \dots, Y_n be counts of bacteria in a culture of cells, measured at time points t_1, \dots, t_n . Aim: study bacteria growth rate.

A possible model for this problem is

$$\begin{aligned} Y_i &= g(t_i; \beta) + \epsilon_i, \\ g(t_i; \theta) &= \beta_0 + \theta_1 t_i, \quad i = 1, \dots, n, \end{aligned}$$

where $\beta = (\beta_0, \beta_1) \in \mathbb{R}^2$ is unknown.

In words: (we assume) the counts follow a linear equation with time, and we wish to learn about the parameters of this line.

To find the LS estimator we solve in β_0, β_1 the system

$$\frac{d}{d\beta} \sum_{i=1}^n (y_i - g_i(\theta))^2 = 0, \quad \rightarrow \text{to find the minimum}$$

Example 6 (cont'd)

This system is

$$\begin{aligned}\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1}) &= 0 \\ \sum_{i=1}^n x_{i1} (y_i - \beta_0 - \beta_1 x_{i1}) &= 0\end{aligned}$$

and the solution is

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \frac{s_{y,x}}{s_x^2} \bar{x}, \\ \hat{\beta}_1 &= \frac{s_{y,x}}{s_x^2}.\end{aligned}$$

Handwritten annotations: A blue arrow points from $s_{y,x}$ in the first equation to the word "covariance". A blue arrow points from s_x^2 in the second equation to the word "variance".

and

The vector $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ as above is the least squares estimator for β .

Method of Maximum Likelihood

Let Y_1, \dots, Y_n be an iid random sample from with $Y_i \sim F_\theta$ and pdf f and unknown parameter θ .

The Maximum Likelihood Estimator (MLE) is defined by

$$\hat{\theta} = \arg \inf_{\theta \in \Theta} L(\theta).$$

Under standard regularity conditions, the MLE is also defined as the solution to the likelihood equation

$$\frac{d\ell(\theta)}{d\theta} = 0.$$

Note that θ may be d -dimensional vector, in which case the likelihood equation consists in d simultaneous equations.

Example 7

Let Y_1, \dots, Y_n be an iid random sample with $Y_i \sim \text{Ber}(\theta)$, with θ unknown. Let's compute the maximum likelihood estimator for θ .

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1-y_i} \\ &= \theta^{\sum_i y_i} (1 - \theta)^{n - \sum_i y_i}, \end{aligned}$$

Since $L(\theta) > 0$, for all θ , we can apply the log to get the log-likelihood function. So we solve the likelihood equation $d\ell(\theta)/d\theta = 0$, i.e.

$$\frac{\sum_i y_i}{\theta} - \frac{n - \sum_i y_i}{1 - \theta} = 0,$$

to get the solution $\hat{\theta} = \bar{y}$

Furthermore, $d^2\ell(\theta)/d\theta^2$ at $\theta = \hat{\theta}$ is negative, so $\hat{\theta}$ is a local maximum, thus it's the MLE of θ . In this case, the MLE coincides with the MME.

Example 8

Let Y_1, \dots, Y_m be a random vector with distribution $\text{Mn}(n, \theta_1, \dots, \theta_m)$ with $0 < \theta_i < 1$ for all i , $\sum_i \theta_i = 1$ and $n = \sum_i Y_i$.

For example, in a sample of Chinese population of Hong Kong in 1937, blood types occur with the following frequencies, where M and N are red cell antigens

	Blood Type			
	M	MN	N	Total
Frequency	342	500	187	1029

Intuitively, the estimated probability of each blood type is the ratio of the observed frequency divided by n , i.e. $\hat{\theta}_i = y_i/n$ for all i . This is the MLE of θ .

Indeed, the log-likelihood is

$$\ell(\theta) = \log n! - \sum_{i=1}^3 \log y_i! + \sum_{i=1}^3 y_i \log \theta_i.$$