

# Statistical Modelling I

Kexing Ying

January 11, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Point Estimation</b>	<b>3</b>
2.1	Review . . . . .	3
2.2	Cramér-Rao Lower Bound . . . . .	4
2.3	Asymptotic Properties of Estimators . . . . .	7

# 1 Introduction

In this module, we will consider and analyse the relationship between measurements through the use of statistical models. This is realised in several ways including quantifying distributions, comparing distributions and predicting observations. We shall study these methods through deriving, evaluating and applying estimators, confidence intervals and hypothesis tests based, first on parametric models, and later based on the theory of linear models.

**Definition 1.1** (Statistical Model). A statistical model is a specification of the distribution of  $Y$  up to an unknown parameter  $\theta$ .

**Definition 1.2** (Parameter Space). Given a statistical model  $Y$  up to some parameter  $\theta$ , the set  $\Theta$  of all possible parameter values is called the parameter space.

In this module we will assume  $\Theta \subseteq \mathbb{R}^p$  for some  $p \in \mathbb{N}$  so that we consider *parametric models*. A *semiparametric model* is a statistical model which parameters belong to a more general space, e.g. functions spaces.

As with last years **Probability and Statistics**, we will denote the data  $\mathbf{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$  as a vector and  $\mathbf{Y} = (Y_1, \dots, Y_n)$  a random vector. In this case, the statistical model specifies the joint distribution of  $Y_1, \dots, Y_n$  up to some unknown parameter  $\theta$ . If  $Y_1, \dots, Y_n$  are independent and identically distributed (iid.), then we call it a *random sample*.

Furthermore, in some situations, the random vector  $\mathbf{Y} = (Y_1, \dots, Y_n)$  might be dependent on random or nonrandom values  $x_1, \dots, x_n$ . The  $x_i$ 's are an example of covariates. An example of this could be that, in a clinical trial, some patients are given a treatment while others received a placebo. If we would like to model the outcome of the  $i$ -th patient by their survival time  $Y_i$ , it is clear that as covariate for the  $i$ -th patient, we may use the indicator function for whether or not  $i$  received the treatment as a covariate.

As another example, say we would like to ask whether or not taller people have a higher income. To answer this question we might create a statistical model in which  $Y_i$  is the income,  $x_i$  is the height and

$$Y_i = \beta_0 + x_i\beta_1 + \epsilon_i$$

for  $i = 1, \dots, n$ ,  $\epsilon \sim N(0, \sigma^2)$  iid. and  $\theta = (\beta_0, \beta_1, \sigma^2)$ ,  $\Theta = \mathbb{R}^2 \times [0, \infty)$ .

Having formulated a model, we can draw inferences from the sample. By estimating the unknown parameters we attempt to “fit the model”. Through this, we receive a model that can provide us with point estimates, or better yet, tools that can help us make decisions through some combination of hypothesis tests and confidence intervals.

However, with all statistical models, we have to accept that it will not perfectly reflect reality. But, that is not the point of statistical models anyway. Statistical models are meant to be useful and in general we would like a model to

- agree with the observed data reasonably;
- be relatively simple;
- interpretable, e.g. parameters have a physical interpretation.

With these aims in mind, we might conduct sensitivity analysis in which we discard models that are not adequate for the data through an iterative process.

## 2 Point Estimation

From the introduction, we seen that during the process of fitting the model, we need to estimate  $\theta$  in the statistical model, and furthermore, during the inference process, we need to point estimate, interval estimate or hypothesis test to address our question. We recall from first year that this can be achieved through several methods and we shall quickly review them here.

### 2.1 Review

We recall the following definitions.

**Definition 2.1** (Realisation, Statistic, Estimate, Estimator).

- Data  $y_1, \dots, y_n$  is called a realisation of  $Y_1, \dots, Y_n$ .
- A function  $t$  of observable random variables is called a statistic.
- An estimate of  $\theta$  is  $t(y_1, \dots, y_n)$ .
- An estimator of  $\theta$  is  $T = t(Y_1, \dots, Y_n)$ .

**Example 1.** Let  $Y_1, \dots, Y_n \sim N(\mu, 1)$  iid. for some unknown  $\mu \in \mathbb{R}$ . There are many methods for estimating  $\mu$ .

- the sample mean  $\hat{\mu} = \frac{1}{n} \sum y_i$ ;
- the sample median;
- the  $k$ -trimmed mean where we discard the highest and lowest  $k$  observed  $y_i$  before computing the mean;
- ...

For the sample mean estimate, the corresponding estimator is  $T = \bar{Y} = \frac{1}{n} \sum Y_i$ .

As we can see from the example, there are many possible estimations for the same parameter. To justify the use of a specific estimator, one might use a frequentist's perspective and generate many data and tabulate the results of each estimator. Through this process, one can justify a particular estimator through observed data.

As estimators are random variables, we can formalise this idea by considering properties of its sampling distribution (that is the distribution of the estimator), e.g.

$$\mathbb{P}_\theta(T \in \mathcal{A}), \quad E_\theta(T), \quad \text{Var}_\theta(T), \dots$$

We saw this idea last year in the form of *bias* and *mean square error*. We recall the definitions here.

**Definition 2.2** (Bias). Let  $T$  be an estimator of  $\theta \in \Theta \subseteq \mathbb{R}$ . Then the bias of  $T$  is

$$\text{bias}_\theta(T) = E_\theta(T) - \theta.$$

If  $\text{bias}_\theta(T) = 0$  for all  $\theta \in \Theta$ , then we say  $T$  is unbiased for  $\theta$ .

If the parameter space is higher dimensional, say  $\Theta \subseteq \mathbb{R}^k$ , we may be instead be interested in the value of  $g(\theta)$  for some  $g : \Theta \rightarrow \mathbb{R}$ . Then, we can naturally extend the definition of bias to this by

$$\text{bias}_\theta(T) = E_\theta(T) - g(\theta).$$

**Example 2.** Let  $Y_1, \dots, Y_n \sim N(\mu, \sigma^2)$  iid.  $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times (0, \infty)$ . Then, say if we are in  $\mu$ , we may define  $g : \Theta \rightarrow \mathbb{R} : (\mu, \sigma^2) \mapsto \mu$ .

**Definition 2.3** (Mean Square Error). Let  $T$  be an estimator of  $\theta \in \Theta \subseteq \mathbb{R}$ . Then the mean square error of  $T$  is

$$\text{MSE}_\theta(T) = E_\theta[(T - \theta)^2].$$

In addition to this, we have the standard error of a estimator

**Definition 2.4** (Standard Error). Let  $T$  be an estimator of  $\theta \in \Theta \subseteq \mathbb{R}$ . Then the standard error of  $T$  is

$$\text{SE}_\theta(T) = \sqrt{\text{Var}_\theta(T)}.$$

From last year, we saw the following proposition.

**Proposition 1.** Let  $T$  be an estimator of  $\theta \in \Theta \subseteq \mathbb{R}$ . Then

$$\text{MSE}_\theta(T) = \text{Var}_\theta(T) + (\text{bias}_\theta(T))^2.$$

If we restrict out estimators to be unbiased, often times, we find that the remaining possible estimators well-behaved and we can often find the best estimators by minimising the MSE. However, a biased estimator might have a small MSE than an unbiased estimator (recall sample variance), and it is not necessarily true that such an estimator even exists.

## 2.2 Cramér-Rao Lower Bound

As the mean square error provide us with a method of quantifying how good an estimator is, we are motivated by minimising the mean square error for a family of estimators. That is, if  $\theta \in \Theta$  is a parameter, then is there an estimator  $T$  of  $\theta$  such that for all estimators of  $\theta$ ,  $S$ ,

$$\text{MSE}_\theta(T) \leq \text{MSE}_\theta(S).$$

Unfortunately, the answer to this question is in general, no, however, for unbiased estimators, the answer is often yes. Indeed, if  $T$  is an unbiased estimator, then,

$$\text{MSE}_\theta(T) = \text{Var}_\theta(T) = \text{bias}_\theta(T)^2 = \text{Var}_\theta(T),$$

so it suffices to minimise the variance.

**Theorem 1** (Cramér-Rao Lower Bound). Suppose  $T = T(X)$  is an unbiased estimator for  $\theta \in \Theta \subseteq \mathbb{R}$  based on  $X = (X_1, \dots, X_n)$  with joint pdf  $f_\theta(x)$ . Then under mild regularity conditions (which is elaborated on in the proof below),

$$\text{Var}_\theta(T) \geq \frac{1}{I(\theta)},$$

where

$$I(\theta) = E_{\theta} \left[ \left\{ \frac{\partial}{\partial \theta} \log f_{\theta}(X) \right\}^2 \right],$$

and we call  $I(\theta)$  the *Fisher information* of the sample.

By computing, we find the Fisher information to equal the following.

$$I(\theta) = -E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log f_{\theta}(X) \right].$$

Indeed,

$$\begin{aligned} E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log f_{\theta}(X) \right] &= E_{\theta} \left[ \frac{\partial}{\partial \theta} \frac{f'_{\theta}(X)}{f_{\theta}(X)} \right] \\ &= E_{\theta} \left[ -\frac{f'_{\theta}(X)}{f_{\theta}^2(X)} f'_{\theta}(X) + \frac{f''_{\theta}(X)}{f_{\theta}(X)} \right] \\ &= E_{\theta} \left[ -\left( \frac{\partial}{\partial \theta} \log f_{\theta}(X) \right)^2 \right] + E_{\theta} \left[ \frac{f''_{\theta}(X)}{f_{\theta}(X)} \right]. \end{aligned}$$

So, the result follows as,

$$\begin{aligned} E_{\theta} \left[ \frac{f''_{\theta}(X)}{f_{\theta}(X)} \right] &= \int_{x \in A} \frac{f''_{\theta}(x)}{f_{\theta}(x)} f_{\theta}(x) dx \\ &= \int_{x \in A} f''_{\theta}(x) dx = \frac{\partial^2}{\partial \theta^2} \int_{x \in A} f_{\theta}(x) dx = 0, \end{aligned}$$

where we denoted  $A$  as the support of  $f_{\theta}$ . This is a useful identity whenever the second derivative is easy to compute.

**Corollary 1.1.** Suppose  $X_1, \dots, X_n$  is a random sample. Then,  $f_{\theta}^{(1)}$  is the pdf of single observation, then

$$I_f(\theta) = nI_{f_{\theta}^{(1)}}(\theta).$$

*Proof.* Since a random sample is iid.  $f_{\theta}(x) = \prod f_{\theta}^{(1)}$  and so

$$I_f(\theta) = -E_{\theta} \left[ \frac{\partial^2}{\partial \theta^2} \log f_{\theta}(X) \right] = \sum_{i=1}^n -E_{\theta} \left( \frac{\partial^2}{\partial \theta^2} \log f_{\theta}^{(1)}(X_i) \right) = nI_{f_{\theta}^{(1)}}(\theta).$$

□

From this, we can conclude that the Fisher information is proportional to the sample size.

**Example 3.** Let us find the Fisher information for the random sample  $X_1, \dots, X_n \sim \text{Bern}(\theta)$ .

By the above corollary, we have  $I_f(\theta) = nI_{f_{\theta}^{(1)}}(\theta)$ . So, since the pmf of a Bernoulli random variable is  $f_{\theta}^{(1)}(x) = \theta^x(1-\theta)^{1-x}$ , we have

$$\frac{\partial}{\partial \theta} \log f_{\theta}^{(1)}(x) = \frac{x}{\theta} - \frac{1-x}{1-\theta} = \frac{x-\theta}{\theta(1-\theta)},$$

hence,

$$I_{f_\theta^{(1)}}(\theta) = E \left[ \left( \frac{x - \theta}{\theta(1 - \theta)} \right)^2 \right] = \frac{1}{\theta^2(1 - \theta)^2} \text{Var}(X) = \frac{1}{\theta(1 - \theta)}.$$

Thus, the Fisher information of the random sample is just  $I_f(\theta) = n/\theta(1 - \theta)$ .

With the Fisher information, we can apply the Cramér-Rao lower bound theorem allowing us to conclude that an unbiased estimator  $T$  for  $\theta$  has variance  $\text{Var}(T) \geq \theta(1 - \theta)/n = \text{Var}(\bar{X})$ . This allows us to conclude that the sample mean  $\bar{X}$  minimises the mean square error among unbiased estimators for  $\theta$ .

Let us now prove the Cramér-Rao lower bound theorem.

*Proof.* (Cramér-Rao lower bound theorem). Let us first specify the regularity conditions for the Cramér-Rao lower bound theorem.

- Assume that the set  $A := \text{supp} f_\theta = \{x \in \mathbb{R}^n \mid f_\theta(x) > 0\}$  is independent of  $\theta$ .
- $\Theta$  is an open interval in  $\mathbb{R}$ .
- For all  $\theta \in \Theta$  there exists  $\frac{\partial f_\theta}{\partial \theta}$ .
- Differentiation and integration commutes (for the specific cases where it is used).

As we saw last year, the space of random variables form an inner product space with the inner product

$$\langle X, Y \rangle = E[XY],$$

and so, the Cauchy-Schwarz inequality applies. That is for all random variables  $X, Y$

$$[E(XY)]^2 \leq E(X^2)E(Y^2).$$

So, we have

$$\begin{aligned} \text{Var}_\theta(T) I_f(\theta) &= E_\theta[(T - E_\theta T)^2] E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f_\theta(X) \right)^2 \right] \\ &\geq \left( E_\theta \left[ (T - E_\theta(T)) \frac{\partial}{\partial \theta} \log f_\theta(X) \right] \right)^2. \end{aligned}$$

Thus, it suffices to show that the expectation on the right hand side evaluates to 1.

$$\begin{aligned} E_\theta \left[ (T - E_\theta(T)) \frac{\partial}{\partial \theta} \log f_\theta(X) \right] &= E_\theta \left[ (T - E_\theta(T)) \frac{\frac{\partial}{\partial \theta} f_\theta(X)}{f_\theta(X)} \right] \\ &= \int_{x \in A} (T(x) - E_\theta(T)) \frac{\frac{\partial}{\partial \theta} f_\theta(x)}{f_\theta(x)} f_\theta(x) dx \\ &= \int_{x \in A} T(x) \frac{\partial}{\partial \theta} f_\theta(x) dx - \int_{x \in A} E_\theta(T) \frac{\partial}{\partial \theta} f_\theta(x) dx \\ &= \frac{\partial}{\partial \theta} \int_{x \in A} T(x) f_\theta(x) dx - E_\theta(T) \frac{\partial}{\partial \theta} \int_{x \in A} f_\theta(x) dx \\ &= \frac{\partial}{\partial \theta} E_\theta(T) - 0 = \frac{\partial}{\partial \theta} \theta = 1. \end{aligned}$$

□

## 2.3 Asymptotic Properties of Estimators

While the Cramér-Rao lower bound theorem provides us with a lower bound for the variance for non-biased estimators, as we have previously seen, it is not always true that there exists an unbiased estimator. So, rather than giving up, we instead study the estimators as the sample size becomes large.

As we have seen, evaluating an estimator  $T = T(X_1, \dots, X_n)$  of  $\theta$  depends on its *sampling distribution*. From the sampling distribution, one can possibly find properties about the estimator such as is bias, mean square error and so on. However, it is not necessarily true that an estimator has a closed form. Indeed, often times, the estimator is defined as a solution to some equation.

To simplify this, one often consider  $T_n = T_n(X_1, \dots, X_n)$  as a sequence of random variables indexed by  $n \in \mathbb{N}$  and consider the stochastic convergence of the variables in question. We recall from last term's probability course, there are three different notions of convergence for random variables,

- convergence in probability;
- convergence almost surely (almost everywhere);
- convergence in distribution.

Let us quickly define them here again.

**Definition 2.5** (Convergence in Probability). Let  $(X_n)_{n=1}^{\infty}$  be a sequence of random variables. Then,  $(X_n)$  converges to the random variable  $X$  in probability if for all  $\epsilon > 0$ ,

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

**Definition 2.6** (Convergence Almost Surely). Let  $(X_n)_{n=1}^{\infty}$  be a sequence of random variables. Then,  $(X_n)$  converges to the random variable  $X$  almost surely if

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1,$$

where  $\lim_{n \rightarrow \infty} X_n = X$  is denoting the event

$$\{\omega \in \Omega \mid X_n(\omega) \rightarrow X(\omega)\}.$$

With other words,  $X_n \rightarrow X$  almost surely, if the set of points  $\omega$  such that  $X_n(\omega)$  does not converge to  $X(\omega)$  has measure 0.

**Definition 2.7.** Let  $(X_n)_{n=1}^{\infty}$  be a sequence of random variables. Then,  $(X_n)$  converges to the random variable  $X$  with cdf  $F_X$  in distribution if

$$\lim_{n \rightarrow \infty} \mathbb{P}(X_n \leq x) = F_X(x),$$

for all  $x$  at which  $F_X$  is continuous.

We also recall the following chain of implications,

$$X_n \rightarrow_{\text{as}} X \implies X_n \rightarrow_{\text{p}} X \implies X_n \rightarrow_{\text{d}} X.$$

If  $X = c$  is a constant, then

$$X_n \rightarrow_{\text{p}} X \iff X_n \rightarrow_{\text{d}} X.$$

We apply this notion onto estimators.

**Definition 2.8** (Consistency). A sequence of estimators  $(T_n)_{n=1}^{\infty}$  for  $g(\theta)$  is called (weakly) consistent if for all  $\theta \in \Theta$ ,  $\epsilon > 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|T_n - g(\theta)| > \epsilon) = 0.$$

While it is possible to prove consistency for certain estimators, it is often a non-trivial task. Instead, we often look at whether a sequence of estimators are asymptotically unbiased and prove the a special class of these estimators are consistent.

**Definition 2.9** (Asymptotically Unbiased Estimators). A sequence of estimators  $(T_n)_{n=1}^{\infty}$  for  $g(\theta)$  is called asymptotically unbiased if for all  $\theta \in \Theta$ ,

$$E_{\theta}(T_n) \rightarrow g(\theta).$$

We see that  $E_{\theta}(T_n)$  is simply a value and this is simply the convergence for real sequences. Before moving on to prove results about estimators, let us quickly recall the Markov inequality.

**Proposition 2.** Let  $X$  be a random variable with  $X(\Omega) \subseteq [0, \infty)$ , then for all  $a \in \mathbb{R}$ ,

$$\mathbb{P}(|X| \geq a) \leq \frac{E(|X|)}{a}.$$

*Proof.* See first year notes. □

**Lemma 2.1** (MSE Consistency). Let  $(T_n)_{n=1}^{\infty}$  be asymptotically unbiased for  $g(\theta)$  for all  $\theta \in \Theta$ . Then, if  $\text{Var}_{\theta}(T_n) \rightarrow 0$  as  $n \rightarrow \infty$ ,  $T_n$  is consistent for  $g(\theta)$ .

*Proof.* Let  $\epsilon > 0$ , then, by Markov's inequality,

$$\begin{aligned} \mathbb{P}_{\theta}(|T_n - g(\theta)| \geq \epsilon) &= \mathbb{P}_{\theta}((T_n - g(\theta))^2 \geq \epsilon^2) \leq \frac{1}{\epsilon^2} E_{\theta}(T_n - g(\theta))^2 \\ &= \frac{\text{MSE}_{\theta}(T_n)}{\epsilon^2} = \frac{1}{\epsilon^2} (\text{Var}_{\theta}(T_n) + (E_{\theta}(T_n) - g(\theta))^2). \end{aligned}$$

Since the right hand side tends to 0 as  $n \rightarrow \infty$ , so is the left hand side. □

Thus, to show that a sequence of estimators is consistent, it suffices to show that it is asymptotically unbiased and its variance tends to 0.

However, while consistency is a nice property for a sequence of estimators to have, it is a very minimal requirement. So, in order to derive hypothesis tests and confidence intervals,



we also need the sampling distribution of  $T_n$ . As we have seen previously, the sample mean estimators  $T_n$  for a normal distribution  $N(\theta, 1)$  has distribution  $T_n \sim N(\theta, 1/n)$ , and so, by centring and scaling, we have

$$\sqrt{n}(T_n - \theta) \sim N(0, 1),$$

for all  $n \geq 1$ . This means we can work with the CDF of  $T_n$  allowing us to easily analyse the behaviours of these estimators. However, this is in general not the case and in most cases, we cannot derive easily the distributions of the estimators. Nonetheless, often, we may approximate their distribution with a normal distribution.

**Definition 2.10** (Asymptotically Normal). A sequence of estimators  $T_n$  for  $\theta \in \mathbb{R}$  is asymptotically normal if, for some  $\sigma^2(\theta)$ ,

$$\sqrt{n}(T_n - \theta) \rightarrow N(0, \sigma^2(\theta)),$$

in distribution.

From last term, we recall the central limit theorem (CLT).

**Theorem 2** (Central Limit Theorem). Let  $Y_1, \dots, Y_n$  be iid. random variables with  $E(Y_i) = \mu$  and  $\text{Var}(Y_i) = \sigma^2 < \infty$ . Then the sequence  $\sqrt{n}(\bar{Y} - \mu)$  converges in distribution to a  $N(0, \sigma^2)$  distribution.

*Proof.* See the *probability for statistics* course. □

The central limit theorem allows us to conclude that a large class of estimators are asymptotically normal. Indeed, sample means and estimators which can be written as a combination of sample means under weak conditions are certainly asymptotically normal. However, we would also consider other estimators.

**Lemma 2.2** (Slutsky's lemma). Let  $X_n, X$  and  $Y_n$  be random variables (or random vectors). If  $X_n \rightarrow X$  in distribution and  $Y_n \rightarrow c$  in probability for some constant  $c$ , then

- $X_n + Y_n \rightarrow X + c$  in distribution;
- $Y_n X_n \rightarrow cX$  in distribution;
- $Y_n^{-1} X_n \rightarrow c^{-1} X$  in distribution if  $c \neq 0$ .

*Proof.* See the *probability for statistics* course. □

Another useful result for determining whether or not a sequence of estimators are asymptotically normal is the  $\delta$ -method. The  $\delta$ -method allows us to consider whether or not the transformation of an asymptotically normal estimator remains asymptotically normal.

**Theorem 3** ( $\delta$ -Method). Suppose that  $T_n$  is an asymptotically normal estimator of  $\theta$  with

$$\sqrt{n}(T_n - \theta) \rightarrow_d N(0, \sigma^2(\theta)),$$

and  $g : \Theta \subseteq \mathbb{R} \rightarrow \mathbb{R}$  is a differentiable function with  $g'(\theta) \neq 0$ . Then

$$\sqrt{n}(g(T_n) - g(\theta)) \rightarrow_d N(0, g'(\theta)^2 \sigma^2(\theta)).$$

*Proof.* Since  $g$  is differentiable,

$$g(T_n) = g(\theta) + g'(\theta)(T_n - \theta) + o((T_n - \theta)^2),$$

where  $R$  is the remainder. So

$$\sqrt{n}(g(T_n) - g(\theta)) = g'(\theta)\sqrt{n}(T_n - \theta) + o((T_n - \theta)^2).$$

Thus, assuming the remainder is negligible, we have

$$\sqrt{n}(g(T_n) - g(\theta)) \rightarrow_d N(0, g'(\theta)^2 \sigma^2(\theta)).$$

□

Lastly, a useful result we will often use (perhaps implicitly) is the continuous mapping theorem. Alike sequential continuity for metric spaces, the continuous mapping theorem will allow us to preserve stochastic convergence under continuous mappings.

**Theorem 4** (Continuous Mapping Theorem). Let  $g : \mathbb{R}^k \rightarrow \mathbb{R}^m$  be continuous at every point of a set  $C$  such that  $\mathbb{P}(X \in C) = 1$ . Then if  $X_n \rightarrow X$ , then  $g(X_n) \rightarrow g(X)$  for all three notions of convergence, i.e. convergence in distribution, in probability and almost surely.