

# Statistics for Data Science: Week 4

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

`rajita.chandak@epfl.ch`, `myrto.limnios@epfl.ch`



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

- 1 Model phenomenon by **distribution**  $F(y_1, \dots, y_n; \theta)$  on  $\mathcal{Y}^n$ , some  $n \geq 1$ .
- 2 Distributional form is known but  $\theta \in \Theta$  is **unknown**.
- 3 **Observe realisation** of  $(Y_1, \dots, Y_n)^\top \in \mathcal{Y}^n$  from this distribution.
- 4 Use the realisation  $\{Y_1, \dots, Y_n\}$  in order to make **assertions concerning the true value of  $\theta$** , and quantify the uncertainty associated with these assertions.

The first sort of assertion we wish to make is:

- 1 **Point Estimation.** Given realisation  $(Y_1, \dots, Y_n)^\top$  from  $F(y_1, \dots, y_n; \theta)$ , how can we **produce an educated guess for the unknown true parameter  $\theta$** ?

How? With a **point estimator**!

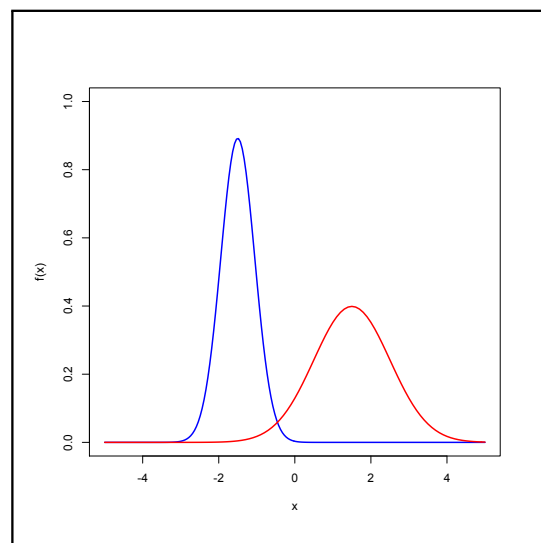
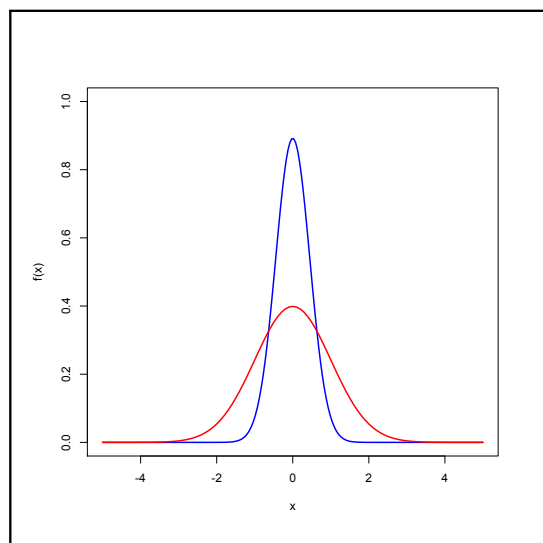
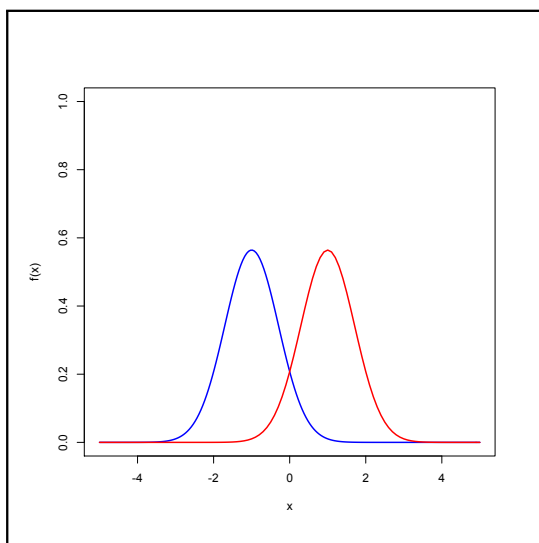
### Definition (Point Estimator)

A statistic with codomain  $\Theta$  is called a *point estimator*, i.e. a point estimator is a statistic  $T : \mathcal{Y}^n \rightarrow \Theta$ .

Since the objective of an estimator is to estimate the  $\theta$  that generated the data, we typically denote it by  $\hat{\theta}(Y_1, \dots, Y_n)$ , or just  $\hat{\theta}$ . Note that  $\theta$  is a deterministic parameter, whereas  $\hat{\theta}$  is a random variable.

But **which** estimator?

- Any statistic taking values in  $\Theta$  could be used!
- Simpler yet, if we are given some  $\hat{\theta}$ , how do we judge its quality?
- Since estimators are *random variables*, every different realisation of the sample  $(Y_1, \dots, Y_n)$  will produce a different realised value for  $\hat{\theta}$ .
- A good estimator should be such that it typically manifests realisations that fall near the true  $\theta$ .
- More precisely, the sampling distribution of an estimator should be concentrated around the true parameter value  $\theta$ .



## Definition (Mean Squared Error)

Let  $\hat{\theta}$  be an estimator of a parameter  $\theta$  corresponding to a model  $\{F_{\theta} : \theta \in \Theta\}$ ,  $\Theta \subseteq \mathbb{R}^d$ . The mean squared error of  $\hat{\theta}$  is defined as

$$\text{MSE}(\hat{\theta}, \theta) = \mathbb{E} \left[ \left\| \overset{\substack{\text{estimate}}}{\hat{\theta}} - \overset{\substack{\text{true parameter}}}{\theta} \right\|^2 \right].$$

And here's the relation to means and variances:

## Lemma (Bias-Variance Decomposition)

The MSE admits the decomposition

$$\text{MSE}(\hat{\theta}, \theta) = \underbrace{\left\| \mathbb{E}[\hat{\theta}] - \theta \right\|^2}_{\text{bias}^2} + \underbrace{\mathbb{E} \left[ \left\| \hat{\theta} - \mathbb{E}(\hat{\theta}) \right\|^2 \right]}_{\text{variance}}.$$

$$\mathbb{E}[\|X - \mathbb{E}[X]\|^2]$$

low bias, low variance



low bias, high variance



high bias, low variance



$$\|a+b\|^2 = \|a\|^2 + \|b\|^2 + 2a^T b$$

$$a = \hat{\theta} - \mathbb{E}[\hat{\theta}] \quad , \quad b = \mathbb{E}[\hat{\theta}] - \theta$$

Proof.

We expand the MSE after adding and subtracting  $\mathbb{E}[\hat{\theta}]$ :

$$\begin{aligned} \mathbb{E}[\|\hat{\theta} - \theta\|^2] &= \mathbb{E}[\|\underbrace{\hat{\theta} - \mathbb{E}[\hat{\theta}]}_a + \underbrace{\mathbb{E}[\hat{\theta}] - \theta}_b\|^2] \\ &= \mathbb{E}\left[(\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)^T (\hat{\theta} - \mathbb{E}[\hat{\theta}] + \mathbb{E}[\hat{\theta}] - \theta)\right] \\ &= \|\mathbb{E}[\hat{\theta}] - \theta\|^2 + \mathbb{E}[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|^2] + 2\mathbb{E}\left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^T (\mathbb{E}[\hat{\theta}] - \theta)\right] \\ &= \underbrace{\|\mathbb{E}[\hat{\theta}] - \theta\|^2}_{\text{bias}^2} + \underbrace{\mathbb{E}[\|\hat{\theta} - \mathbb{E}[\hat{\theta}]\|^2]}_{\text{var}} + \underbrace{2(\mathbb{E}[\hat{\theta}] - \mathbb{E}[\hat{\theta}])^T}_{=0} (\mathbb{E}[\hat{\theta}] - \theta) \end{aligned}$$

non random  
↓

by linearity of the expectation and since  $(\mathbb{E}[\hat{\theta}] - \theta)$  is deterministic.

□

As foretold, the concentration of an estimator  $\hat{\theta}$  around the true parameter  $\theta$  can always be bounded by the MSE:

## Lemma

Let  $\hat{\theta}$  be an estimator of  $\theta \in \mathbb{R}^p$ . For any  $\epsilon > 0$ ,

$$\mathbb{P}[\underbrace{\|\hat{\theta} - \theta\|}_{\text{marker}} > \epsilon] \leq \frac{\text{MSE}(\hat{\theta}, \theta)}{\epsilon^2} \rightarrow 0$$

*Markov's inequality*

- Note that  $\text{MSE}(\hat{\theta}_n, \theta) \xrightarrow{n \rightarrow \infty} 0 \implies \underbrace{\hat{\theta}_n}_{\text{marker}} \xrightarrow{p} \theta$ .
- When an estimator has this property, we call it **consistent**.

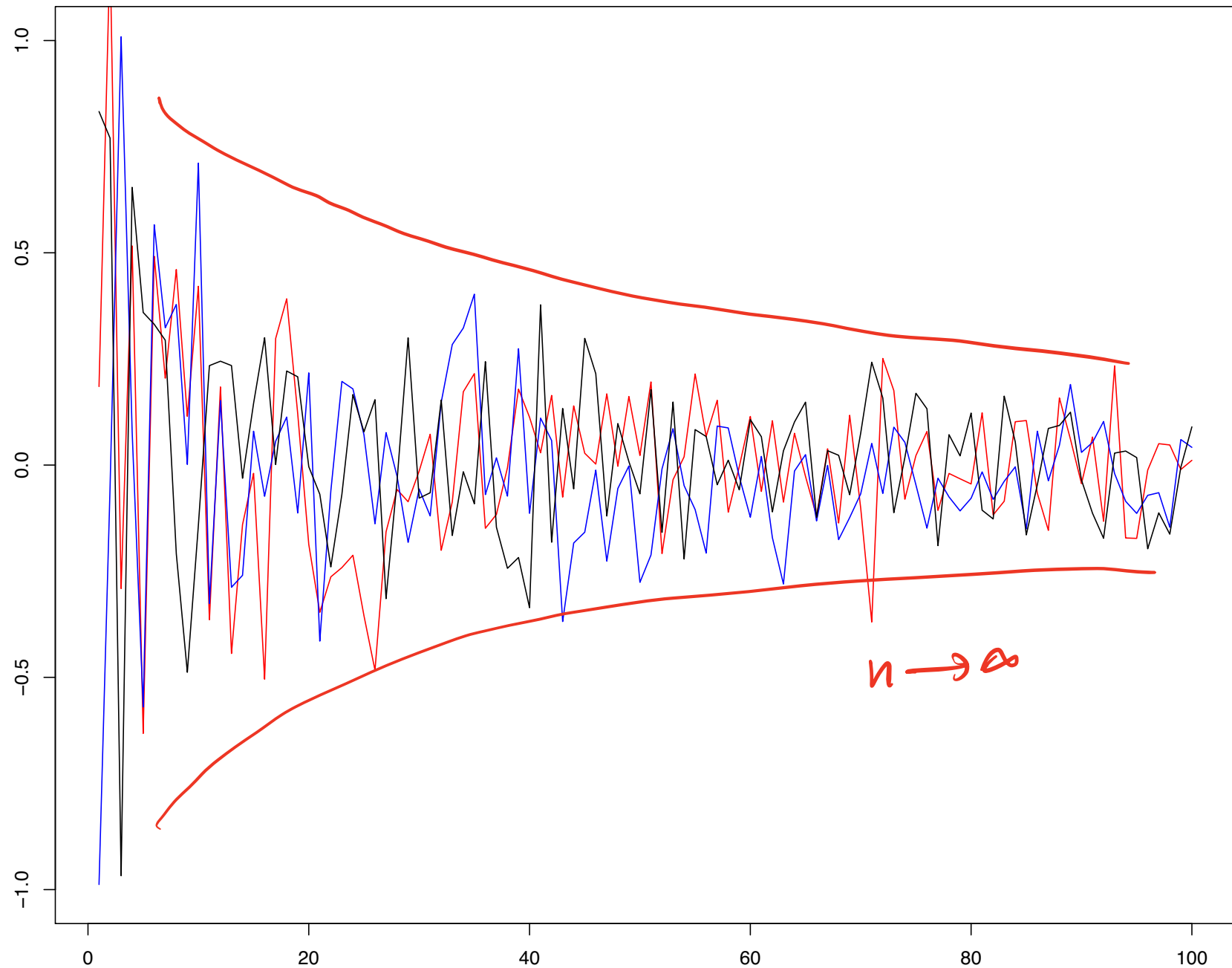
## Definition (Consistency)

An estimator  $\hat{\theta}_n$  of  $\theta$ , constructed on the basis of a sample of size  $n$ , is consistent if  $\hat{\theta}_n \xrightarrow{p} \theta$  as  $n \rightarrow \infty$ .

Note that a vanishing MSE implies consistency, but the converse generally fails.

Consistency of sample mean of sample mean of  $Y_1, \dots, Y_n \sim \mathcal{N}(0, 1)$ , towards the true parameter value 0 (by law of large numbers).

$$\frac{1}{n} \sum Y_i = \hat{\theta} \xrightarrow{?} 0 = \theta$$





Is it always possible to get consistent estimators?

Depends on whether the estimation problem is **well-posed**

## Definition (Identifiability)

A probability model  $\{F_\theta\}_{\theta \in \Theta}$  is called identifiable if for any pair  $\theta_1, \theta_2 \in \Theta$  we have the implication

$$\hat{\theta} \rightarrow \theta \neq \delta, \quad \theta_1 \neq \theta_2 \implies F_{\theta_1} \neq F_{\theta_2}.$$

- Lack of identifiability means that the same model can be produced by more than one parameter.
- In this case we could never distinguish amongst the parameters that give the same model.
- Example: if we have  $N(\mu_1 + \mu_2, \sigma^2)$ , we can never identify each  $\mu_i$ , but only their sum.
- **Henceforth we will tacitly assume identifiability** (and make special mention if it is at stake).

- We can use the MSE to compare estimators or to gauge their performance.
- But is there a *best possible MSE* for a given problem?
- This is a very difficult problem, equivalent to finding a **uniformly optimal estimator**: a statistic  $T_*$  such that

$$\text{MSE}(\hat{\theta}, \theta) \leq \text{MSE}(T, \theta)$$

$$S = \frac{1}{n} \sum y_i$$

for all  $\theta \in \Theta$  and all other estimators  $T$ .

- To see this, let  $T = c$  be a trivial (constant) estimator and observe that for any non-trivial estimator  $S$  we have  $\text{MSE}(S, \theta) > \text{MSE}(T, \theta)$  at  $\theta = c$ .
- So if we want to do well for all  $\theta$  we can't do perfectly for any specific  $\theta$ .
- Here's a simpler question to ask instead (ruling out trivial estimators):  
 $E[\hat{\theta}] = \theta$

Among unbiased estimators (bias zero), can we make the MSE arbitrarily small?

- If so, **how?** (what is the crucial ingredient at play?)

- Rephrasing, we are asking whether there is fundamental lower bound for the variance of an unbiased estimator of  $\theta$ .
- We will concentrate on **1-dimensional parameters** for simplicity.

## The Question

For  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  with joint density/frequency  $f(\mathbf{y}; \theta)$  depending on an unknown  $\theta \in \mathbb{R}$ , does there exist some function  $\Lambda(\theta) > 0$  such that

$$\text{var}[\hat{\theta}] \stackrel{=}{\geq} \Lambda(\theta), \quad \forall \theta$$

for any estimator  $\hat{\theta}$  such that  $\mathbb{E}[\hat{\theta}] = \theta$ ?

- At a next step we can ask if this bound is achievable.

Let's assume we can interchange differentiation and integration in the form Leibniz Rule

$\frac{d}{d\theta} \int S(\mathbf{y}) f(\mathbf{y}; \theta) d\mathbf{y} \stackrel{!}{=} \int S(\mathbf{y}) \frac{f(\mathbf{y}; \theta)}{f(\mathbf{y}; \theta)} \frac{d}{d\theta} f(\mathbf{y}; \theta) d\mathbf{y} = \int S(\mathbf{y}) f(\mathbf{y}; \theta) \frac{d}{d\theta} \log f(\mathbf{y}; \theta) d\mathbf{y}$   
whenever an integral such as the one on the left hand side presents itself.

① Setting  $U = \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta)$  and  $S(\mathbf{y}) = 1$ , this gives that

$$\mathbb{E}[U] = \int f(\mathbf{y}; \theta) \frac{\partial}{\partial \theta} \log f(\mathbf{y}; \theta) d\mathbf{y} = \frac{d}{d\theta} \int f(\mathbf{y}; \theta) d\mathbf{y} \stackrel{\frac{\partial}{\partial \theta} 1}{=} 0$$

② Therefore  $\text{var}[U] = \mathbb{E}[U^2] = \mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right)^2 \right]$

③ For  $\hat{\theta}$  unbiased, our “interchange ansatz” with  $S(\mathbf{y}) = \hat{\theta}(\mathbf{y})$  gives  $\hat{\theta} \Rightarrow \mathbb{E}[\hat{\theta}] = \theta$   
 $\text{cov}(\hat{\theta}, U) = \mathbb{E}[\hat{\theta} U] - \underbrace{\mathbb{E}[\hat{\theta}] \mathbb{E}[U]}_{=0} = \int \underbrace{\hat{\theta}(\mathbf{y})}_{S(\mathbf{y})} f(\mathbf{y}; \theta) \frac{d}{d\theta} \log f(\mathbf{y}; \theta) d\mathbf{y} = \frac{d}{d\theta} \underbrace{\mathbb{E}[\hat{\theta}]}_{=\theta} \stackrel{\frac{\partial}{\partial \theta} \theta}{=} 1$

Now the Cauchy-Schwartz inequality gives

$$\text{var}(\hat{\theta}) \geq \frac{\text{cov}^2(\hat{\theta}, U)}{\text{var}(U)} \Rightarrow \text{var}(\hat{\theta}) \geq \frac{1}{\mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right)^2 \right]} = \frac{1}{\mathcal{I}(\theta)}$$

$(\mathbb{E}[XY])^2 \leq \mathbb{E}[X^2] \mathbb{E}[Y^2]$   
 $X = \hat{\theta}, Y = U$

$$\frac{\partial}{\partial \theta} \int f(y; \theta) = \int \frac{\partial}{\partial \theta} f(y; \theta) = \int \frac{f(y; \theta)}{f(y; \theta)} \frac{\partial}{\partial \theta} f(y; \theta)$$

$$\frac{\partial}{\partial \theta} \log f(y; \theta) = \frac{1}{f(y; \theta)} \cdot \frac{\partial}{\partial \theta} f(y; \theta)$$

$$= \int f(y) \frac{\partial}{\partial \theta} \log(f(y; \theta))$$

In summary, we have established:

## Cramér-Rao Lower Bound

Given sufficient regularity, any unbiased estimator  $\hat{\theta}(\mathbf{Y})$  of finite variance satisfies:

$$\text{var}[\hat{\theta}(\mathbf{Y})] \geq \frac{1}{\mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right)^2 \right]} = \frac{1}{\mathcal{I}_n(\theta)}$$

The quantity  $\mathcal{I}_n(\theta)$  is fundamental, and called **Fisher information**.

In summary, we have established:

## Cramér-Rao Lower Bound

Given sufficient regularity, any unbiased estimator  $\hat{\theta}(\mathbf{Y})$  of finite variance satisfies:

$$\text{var}[\hat{\theta}(\mathbf{Y})] \geq \frac{1}{\mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right)^2 \right]} = \frac{1}{\mathcal{I}_n(\theta)}$$

The quantity  $\mathcal{I}_n(\theta)$  is fundamental, and called **Fisher information**.

- If  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$  has iid entries, we have  $f(\mathbf{y}; \theta) = \prod_{i=1}^n f(y_i; \theta)$  and so

$$\mathcal{I}_n(\theta) = n\mathcal{I}_1(\theta).$$

- By further interchanges of integration/differentiation it typically holds that

$$\mathcal{I}_n(\theta) = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f(\mathbf{Y}; \theta) \right]$$

- The deeper meaning of all this will become clearer when we study **likelihood**.

Is the Cramér-Rao bound **tight (achievable)** though?

if  $\text{var}[\hat{\theta}] \stackrel{?}{=} \frac{1}{\mathcal{I}_n(\theta)}$

then  $\text{var}[\hat{\theta}] = \frac{\text{cov}^2 \left[ \hat{\theta}, \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right]}{\text{var} \left[ \frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) \right]}$

which occurs if and only if  $\frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta)$  is a linear function of  $\hat{\theta}$  (correlation 1):

$$\frac{\partial}{\partial \theta} \log f(\mathbf{Y}; \theta) = A(\theta) \hat{\theta}(\mathbf{Y}) + B(\theta)$$

Solving this differential equation yields, for all  $\mathbf{y}$ ,

$$\log f(\mathbf{y}; \theta) = \underbrace{A^*(\hat{\theta}) + B^*(\theta)} + S(\mathbf{y})$$

so that  $\text{var}_{\theta}(\hat{\theta})$  attains the lower bound if and only if the density (frequency) of  $\mathbf{Y}$  is a one-parameter exponential family with sufficient statistic  $\hat{\theta}$ .



So **what ingredients** go into pushing towards this lower bound?

The **Rao-Blackwell Theorem** tells us that sufficiency is key:

### Theorem (Rao-Blackwell Theorem)

Let  $\hat{\theta}$  be an unbiased estimator of  $\theta$  with finite variance, and let  $T$  be sufficient for  $\theta$ . Then  $\hat{\theta}^* := \mathbb{E}[\hat{\theta} | T]$  is also an unbiased estimator of  $\theta$  and

$$\text{var}(\hat{\theta}^*) \leq \text{var}(\hat{\theta}).$$

Equality is attained if and only if  $\mathbb{P}_{\theta}[\hat{\theta}^* = \hat{\theta}] = 1$ .

So **what ingredients** go into pushing towards this lower bound?

The **Rao-Blackwell Theorem** tells us that sufficiency is key:

### Theorem (Rao-Blackwell Theorem)

*Let  $\hat{\theta}$  be an unbiased estimator of  $\theta$  with finite variance, and let  $T$  be sufficient for  $\theta$ . Then  $\hat{\theta}^* := \mathbb{E}[\hat{\theta} | T]$  is also an unbiased estimator of  $\theta$  and*

$$\frac{1}{I(\theta)} \stackrel{?}{=} \text{var}(\hat{\theta}^*) \leq \text{var}(\hat{\theta}).$$

*Equality is attained if and only if  $\mathbb{P}_{\theta}[\hat{\theta}^* = \hat{\theta}] = 1$ .*

### Comments:

- Throwing away irrelevant aspects of the data improves estimation quality.
- These irrelevant aspects contribute to the variation of the estimator (as they have sampling variation of their own), but without furnishing any useful information on the parameter
- $\hat{\theta}^* = \mathbb{E}[\hat{\theta} | T]$  is called a “Rao-Blackwellised” version of  $\hat{\theta}$ .

## Proof.

Since  $T$  is sufficient for  $\theta$ ,  $\mathbb{E}[\hat{\theta} | T = t] = h(t)$  is independent of  $\theta$ , so that  $\hat{\theta}^*$  is well-defined as a statistic (depends only on  $\mathbf{Y}$  and not  $\theta$ ). Then,

$$\mathbb{E}[\hat{\theta}^*] = \mathbb{E}[\mathbb{E}[\hat{\theta} | T]] = \mathbb{E}[\hat{\theta}] = \theta.$$

*defn, property of conditional exp.*

Furthermore, from the law of total variance, we have

$$\text{var}(\hat{\theta}) = \text{var}[\mathbb{E}(\hat{\theta} | T)] + \mathbb{E}[\text{var}(\hat{\theta} | T)] \geq \text{var}[\mathbb{E}(\hat{\theta} | T)] = \text{var}(\hat{\theta}^*)$$

*by assumption.*  
 *$\hat{\theta}^*$  by defn.*  
 *$\geq 0$*   
 *$= 0?$*

In addition, note that

$$\text{var}(\hat{\theta} | T) := \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta} | T])^2 | T] = \mathbb{E}[(\hat{\theta} - \hat{\theta}^*)^2 | T]$$

*defn.*  
*defn.*

so that  $\mathbb{E}[\text{var}(\hat{\theta} | T)] = \mathbb{E}(\hat{\theta} - \hat{\theta}^*)^2 > 0$  unless if  $\mathbb{P}(\hat{\theta}^* = \hat{\theta}) = 1$ . □

Suppose that  $\hat{\theta}$  is an unbiased estimator of  $g(\theta)$  and  $T, S$  are  $\theta$ -sufficient.

- What is the relationship between  $\text{var}(\underbrace{\mathbb{E}[\hat{\theta}|T]}_{\hat{\theta}_T^*}) \stackrel{?}{\underset{\leq}{\geq}} \text{var}(\underbrace{\mathbb{E}[\hat{\theta}|S]}_{\hat{\theta}_S^*})$
- Intuition suggests that whichever of  $T, S$  carries the least irrelevant information (in addition to the relevant information) should “win”
  - More formally, if  $T = h(S)$  then we should expect that  $\hat{\theta}_T^*$  dominate  $\hat{\theta}_S^*$ .

Suppose that  $\hat{\theta}$  is an unbiased estimator of  $g(\theta)$  and  $T, S$  are  $\theta$ -sufficient.

- What is the relationship between  $\text{var}(\underbrace{\mathbb{E}[\hat{\theta}|T]}_{\hat{\theta}_T^*}) \stackrel{?}{\underset{\leq}{\geq}} \text{var}(\underbrace{\mathbb{E}[\hat{\theta}|S]}_{\hat{\theta}_S^*})$
- Intuition suggests that whichever of  $T, S$  carries the least irrelevant information (in addition to the relevant information) should “win”
  - More formally, if  $T = h(S)$  then we should expect that  $\hat{\theta}_T^*$  dominate  $\hat{\theta}_S^*$ .

## Proposition

For  $\hat{\theta}$  an unbiased estimator of  $\theta$  and  $T, S$  two  $\theta$ -sufficient statistics, define

$$\hat{\theta}_T^* := \mathbb{E}[\hat{\theta}|T] \quad \& \quad \hat{\theta}_S^* := \mathbb{E}[\hat{\theta}|S].$$

Then, the following implication holds

$$T = h(S) \implies \text{var}(\hat{\theta}_T^*) \leq \text{var}(\hat{\theta}_S^*)$$

- Essentially this means that the best possible “Rao-Blackwellisation” is achieved by conditioning on a minimally sufficient statistic.

## Proof.

Recall the *tower property* of conditional expectation: if  $Y = f(X)$ , then

$$\mathbb{E}[Z|Y] = \mathbb{E}\{\mathbb{E}(Z|X)|Y\}.$$

Since  $T = f(S)$  we have

$$\begin{aligned} \hat{\theta}_T^* &\stackrel{\text{defn.}}{=} \mathbb{E}[\hat{\theta}|T] \\ &\stackrel{\text{tower property}}{=} \mathbb{E}[\mathbb{E}(\hat{\theta}|S)|T] \\ &\stackrel{\text{minimal Suff.}}{\rightarrow} \mathbb{E}[\hat{\theta}_S^*|T] \end{aligned}$$

The conclusion now follows from the Rao-Blackwell theorem. □

$$\mathbb{E} [\|\hat{\theta} - \theta\|^2]$$
$$\max (\hat{\theta} - \theta)$$

- So now we have a means to judge the quality of an estimator
  - In certain cases, we even know what's the best performance we can hope for.
  - And (minimal) **sufficiency can help** us approach it.
  - But how can we actually come up with an estimator in the first place?
  - We need **general methods that can be applied in any model context** to yield an estimator.
  - Preferably methods that yield **good** estimators relative to our performance measures/bounds.
- ↪ The main focus will be on a key method called **maximum likelihood**.

**Motivation:** recall our understanding of statistics as “inverse probability”.

→ For the moment, consider the discrete case for simplicity.

## Probability Perspective

Given a parameter  $\theta \in \Theta$ , then for any  $(y_1, \dots, y_n)^\top \in \mathcal{Y}^n$ , we can evaluate

$$\underline{(y_1, \dots, y_n)} \mapsto \mathbb{P}_\theta[Y_1 = y_1, \dots, Y_n = y_n]$$

that is, how the probability varies as a function of the sample (=the result).

## Statistics Perspective

Given a sample  $(y_1, \dots, y_n)^\top \in \mathcal{Y}^n$ , then for any  $\theta \in \Theta$  we can calculate

$$\underline{\theta} \mapsto \mathbb{P}_\theta[Y_1 = y_1, \dots, Y_n = y_n]$$

that is, how the probability varies as a function of  $\theta$  (=the model).

**Intuition:** we imagine that, having our sample, the values of  $\theta$  that are most plausible are those that render the observed sample most probable...



This motivates the following definition...

### Definition (Likelihood)

Let  $(Y_1, \dots, Y_n)$  be a sample of random variables with joint density/frequency  $f(\underbrace{y_1, \dots, y_n}_{\uparrow \quad \uparrow}; \theta)$ , where  $\theta \in \mathbb{R}^p$ . The likelihood of  $\theta$  is defined as

$$L(\theta) = \underbrace{f(Y_1, \dots, Y_n; \theta)}_{\text{likelihood}}.$$

If  $(Y_1, \dots, Y_n)^\top$  has i.i.d. entries, each with density/frequency  $f(y_i; \theta)$  then,

$$L(\theta) = \prod_{i=1}^n f(Y_i; \theta)$$

This motivates the following definition...

### Definition (Likelihood)

Let  $(Y_1, \dots, Y_n)$  be a sample of random variables with joint density/frequency  $f(y_1, \dots, y_n; \theta)$ , where  $\theta \in \mathbb{R}^p$ . The likelihood of  $\theta$  is defined as

$$L(\theta) = f(\underbrace{Y_1}, \dots, \underbrace{Y_n}, \theta).$$

If  $(Y_1, \dots, Y_n)^\top$  has i.i.d. entries, each with density/frequency  $f(y_i; \theta)$  then,

$$L(\theta) = \prod_{i=1}^n \underbrace{f(Y_i; \theta)}$$

... and the following estimation method

### Definition (Maximum Likelihood Estimator)

In the same context, a maximum likelihood estimator (MLE) of  $\hat{\theta}$  is an estimator such that

$$L(\theta) \leq L(\hat{\theta}), \quad \forall \theta \in \Theta.$$

$$\hat{\theta}_1, \hat{\theta}_2 : L(\hat{\theta}_1) = L(\hat{\theta}_2)$$

Many comments are in order:

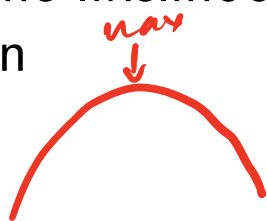
- When there exists a unique maximum, we speak of **the** MLE  $\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta)$
- The likelihood is a random function.
- It is the joint density/frequency of the sample, but viewed as a function of  $\theta$ .
- It is NOT “the probability of  $\theta$ ”
- $L(\theta)$  is the answer to the question *how does the joint density/probability of the sample vary as we vary  $\theta$ ?*
  - In the discrete case it is exactly “the probability of observing our sample” as a function of  $\theta$ .
  - In the continuous case, since  $F(\mathbf{y} + \varepsilon/2; \theta) - F(\mathbf{y} - \varepsilon/2; \theta) \approx \|\varepsilon\| f(\mathbf{y}; \theta)$  as  $\|\varepsilon\| \downarrow 0$ , we can view  $\|\varepsilon\| \times L(\theta)$  as being the “probability of observing something in the neighbourhood of our sample”, as a function of  $\theta$ .

- If a sufficient statistic  $T$  exists for  $\theta$  then Fisher-Neyman factorisation implies

$$\prod f = \prod g(T; \theta) h(Y) \quad L(\theta) = g(T(\mathbf{Y}); \theta) \underbrace{h(\mathbf{Y})}_{\propto g(T(\mathbf{Y}); \theta)}$$

i.e. **any** MLE depends on data **only through a sufficient statistic**.

- Since the sufficient statistic was arbitrary, if a minimally sufficient statistic exists, the MLE will have used an estimator that has achieved the maximal sufficient reduction of the data.
- MLE's are also *equivariant*. If  $g : \Theta \rightarrow \Theta'$  is a bijection, and if  $\hat{\theta}$  is the MLE of  $\theta$ , then  $g(\hat{\theta})$  is the MLE of  $g(\theta)$  (you can take the hat out:  $g(\hat{\theta}) = \widehat{g(\theta)}$ )
- When the likelihood is differentiable in  $\theta$ , its maximum  $L(\theta)$  must solve the equation



$$\nabla_{\theta} L(\theta) = 0,$$

$$\frac{\partial}{\partial \theta} L(\theta) = 0$$

- But before declaring a solution as an MLE, we must verify it to be a maximum (and not a minimum!).

- If the likelihood is twice differentiable in  $\theta$ , we can verify this by checking

$$-\nabla_{\theta}^2 L(\theta)|_{\theta=\hat{\theta}} \succ 0, \quad \frac{\partial^2}{\partial \theta^2} L(\theta)$$

i.e that minus the Hessian is positive definite. In one dimension, this reduces to the standard second derivative criterion.

- To solve  $\nabla_{\theta} L(\theta) = 0$  when the  $Y_i$  are independent, we must painfully calculate the derivative of an  $n$ -fold product.  $\frac{\partial}{\partial \theta} (f(\theta) + g(\theta)) = \frac{\partial}{\partial \theta} f(\theta) + \frac{\partial}{\partial \theta} g(\theta)$
- To avoid this, we focus instead on the **loglikelihood**  $\ell(\theta) := \log L(\theta)$  instead. Maximisation of  $\ell$  is equivalent to maximisation of  $L$  by monotonicity.
- When the  $Y_i$  are independent,  $\ell$  has the advantage of being a sum rather than a product

$$\ell(\theta) = \log \left( \prod_{i=1}^n f_{Y_i}(Y_i; \theta) \right) = \sum_{i=1}^n \log f_{Y_i}(Y_i; \theta).$$

- Of course, under twice differentiability, verification of a maximum can be checked again by whether or not

$$\nabla_{\theta} \ell(\theta)|_{\theta=\hat{\theta}} = 0 \quad \& \quad -\nabla_{\theta}^2 \ell(\theta)|_{\theta=\hat{\theta}} \succ 0.$$

## Example (MLE for Bernoulli trials)

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Bernoulli}(p)$ . The likelihood is

$\hat{p}_{MLE}$

$$L(p) = \prod_{i=1}^n f(Y_i; p) = \prod_{i=1}^n \underbrace{p^{Y_i} (1-p)^{1-Y_i}}_{\text{density of Ber}(p)} \stackrel{\text{algebra}}{=} p^{\sum_{i=1}^n Y_i} (1-p)^{n - \sum_{i=1}^n Y_i}$$

giving loglikelihood

$\downarrow$   
property of log

$$\ell(p) = \log p \sum_{i=1}^n Y_i + \log(1-p) \left( n - \sum_{i=1}^n Y_i \right).$$

This is twice differentiable in  $p$  and we calculate

$$\frac{d}{dp} \ell(p) = p^{-1} \sum_{i=1}^n Y_i - (1-p)^{-1} \left( n - \sum_{i=1}^n Y_i \right) \stackrel{!}{=} 0$$

## Example (MLE for Bernoulli trials, continued)

Solving

$$p^{-1} \sum_{i=1}^n Y_i - (1-p)^{-1} \left( n - \sum_{i=1}^n Y_i \right) \stackrel{\downarrow}{=} 0,$$

$$\hat{p}_{MLE} = \frac{1}{n} \sum Y_i$$

we get the unique root  $\frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}$ . Calling this  $\hat{p}$ , we now verify that

$$\frac{d^2}{dp^2} \ell(p) = -\stackrel{\downarrow}{p^2} \sum_{i=1}^n Y_i - (1-\stackrel{\downarrow}{p})^{-2} \left( n - \sum_{i=1}^n Y_i \right),$$

which is a negative expression, since  $0 \leq \sum_{i=1}^n Y_i \leq n$  and  $p \in (0, 1)$ . Thus

$$\hat{p} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$$

is the unique MLE of  $p$ .



## Example (MLE for exponential distribution)

Let  $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Exp}(\lambda)$ . The likelihood is

$$L(\lambda) = \prod_{i=1}^n f(Y_i; \lambda) = \prod_{i=1}^n \lambda e^{-\lambda Y_i} = \lambda^n \exp \left\{ -\lambda \sum_{i=1}^n Y_i \right\}.$$

and the log likelihood is

$$\ell(\lambda) = n \log \lambda - \lambda \sum_{i=1}^n Y_i.$$

This is twice differentiable in  $\lambda$  and we calculate

$$\frac{d}{d\lambda} \ell(\lambda) = n\lambda^{-1} - \sum_{i=1}^n Y_i = 0$$



## Example (MLE for exponential distribution, continued)

Setting  $\ell'(\lambda) = 0$  we get a unique root

$$\left(\frac{1}{n} \sum_{i=1}^n Y_i\right)^{-1} = 1/\bar{Y}.$$

Call this  $\hat{\lambda}$ , and note that

$$\frac{d^2}{d\lambda^2} \ell(\lambda) = -\frac{n}{\lambda^2}$$

is always negative, since  $\lambda > 0$ . Thus

$$\hat{\lambda} = \left(\frac{1}{n} \sum_{i=1}^n Y_i\right)^{-1} = 1/\bar{Y}$$

is the unique MLE of  $\lambda$ .

