

# Chapter 10

## Decision theory

In this chapter, we again denote the observable random variable (the data) by  $X \in \mathcal{X}$ , and its distribution by  $P \in \mathcal{P}$ . The probability model is  $\mathcal{P} := \{P_\theta : \theta \in \Theta\}$ , with  $\theta$  an unknown parameter.

In particular cases, we apply the results with  $X$  being replaced by a vector  $\mathbf{X} = (X_1, \dots, X_n)$ , with  $X_1, \dots, X_n$  i.i.d. with distribution  $P \in \{P_\theta : \theta \in \Theta\}$  (so that  $\mathbf{X}$  has distribution  $\mathbb{P} := \prod_{i=1}^n P \in \{\mathbb{P}_\theta = \prod_{i=1}^n P_\theta : \theta \in \Theta\}$ ).

### 10.1 Decisions and their risk

We give a definition of *risk*, but now somewhat more formal than in Chapter 8.

Let  $\mathcal{A}$  be the *action space*.

- $\mathcal{A} = \mathbb{R}$  corresponds to estimating a real-valued parameter.
- $\mathcal{A} = \{0, 1\}$  corresponds to testing a hypothesis.
- $\mathcal{A} = [0, 1]$  corresponds to randomized tests.
- $\mathcal{A} = \{\text{intervals}\}$  corresponds to confidence intervals.

Given the observation  $X$ , we decide to take a certain action in  $\mathcal{A}$ . Thus, an action is a map  $d : \mathcal{X} \rightarrow \mathcal{A}$ , with  $d(X)$  being the decision taken. If  $\mathcal{A} = \mathbb{R}$  a decision is often called an estimator (denoted by  $T$  for instance). If  $\mathcal{A} = \{0, 1\}$  or  $\mathcal{A} = [0, 1]$ , a decision is often called a test (denoted by  $\phi$  for instance).

A *loss function* (*Verlustfunktion*) is a map

$$L : \Theta \times \mathcal{A} \rightarrow \mathbb{R},$$

with  $L(\theta, a)$  being the loss when the parameter value is  $\theta$  and one takes action  $a$ .

The risk of decision  $d(X)$  is defined as

$$R(\theta, d) := E_\theta L(\theta, d(X)), \quad \theta \in \Theta. \quad \xrightarrow{\text{Expected Loss}}$$

**Example 10.1.1 Estimation** In the case of estimating a parameter of interest  $g(\theta) \in \mathbb{R}$ , the action space is  $\mathcal{A} = \mathbb{R}$  (or a subset thereof). Important loss functions are then

$$L(\theta, a) := w(\theta)|g(\theta) - a|^r,$$

where  $w(\cdot)$  are given non-negative weights and  $r \geq 0$  is a given power. The risk is then

$$R(\theta, d) = w(\theta)E_\theta|g(\theta) - d(X)|^r.$$

A special case is taking  $w \equiv 1$  and  $r = 2$ . Then  $L(\theta, a) = (g(\theta) - a)^2$  is called quadratic loss and

$$R(\theta, d) = E_\theta|g(\theta) - d(X)|^2$$

is called the mean square error.

**Example 10.1.2 Tests** Consider testing the hypothesis

$$H_0 : \theta \in \Theta_0$$

against the alternative

$$H_1 : \theta \in \Theta_1.$$

distinct sets

Here,  $\Theta_0$  and  $\Theta_1$  are given subsets of  $\Theta$  with  $\Theta_0 \cap \Theta_1 = \emptyset$ . As action space, we take  $\mathcal{A} = \{0, 1\}$ , and as loss

$$L(\theta, a) := \begin{cases} 1 & \text{if } \theta \in \Theta_0 \text{ and } a = 1 \\ c & \text{if } \theta \in \Theta_1 \text{ and } a = 0 \\ 0 & \text{otherwise} \end{cases}.$$

Here  $c > 0$  is some given constant. Then

$$R(\theta, d) = \begin{cases} P_\theta(d(X) = 1) & \text{if } \theta \in \Theta_0 \\ cP_\theta(d(X) = 0) & \text{if } \theta \in \Theta_1 \\ 0 & \text{otherwise} \end{cases}.$$

recall it is  
the expectation  
of the  
decision given  
the function

Thus, the risks correspond to the error probabilities (type I and type II errors).

### Note

The best decision  $d$  is the one with the smallest risk  $R(\theta, d)$ . However,  $\theta$  is not known. Thus, if we compare two decision functions  $d_1$  and  $d_2$ , we may run into problems because the risks are not comparable:  $R(\theta, d_1)$  may be smaller than  $R(\theta, d_2)$  for some values of  $\theta$ , and larger than  $R(\theta, d_2)$  for other values of  $\theta$ .

**Example 10.1.3 Estimating the mean**

We revisit Example 5.2.1. Let  $X \in \mathbb{R}$  and let  $g(\theta) = E_\theta X := \mu$ . We take quadratic loss

$$L(\theta, a) := |\mu - a|^2.$$

Assume that  $\text{var}_\theta(X) = 1$  for all  $\theta$ . Consider the collection of decisions

$$d_\lambda(X) := \lambda X, \quad \text{new combination of decisions}$$

where  $0 \leq \lambda \leq 1$ . Then

$$\begin{aligned} R(\theta, d_\lambda) &= \text{var}(\lambda X) + \text{bias}_\theta^2(\lambda X) \\ &= \lambda^2 + (\lambda - 1)^2 \mu^2. \end{aligned}$$

The “optimal” choice for  $\lambda$  would be

$$\lambda_{\text{opt}} := \frac{\mu^2}{1 + \mu^2},$$

because this value minimizes  $R(\theta, d_\lambda)$ . However,  $\lambda_{\text{opt}}$  depends on the unknown  $\mu$ , so  $d_{\lambda_{\text{opt}}}(X)$  is not an estimator.

### Various optimality concepts

We will consider three optimality concepts: *admissibility* (*zulässigkeit*), *minimax* and *Bayes*.

## 10.2 Admissible decisions

**Definition 10.2.1** A decision  $d'$  is called strictly better than  $d$  if

$$R(\theta, d') \leq R(\theta, d), \quad \forall \theta,$$

and

$$\exists \theta : R(\theta, d') < R(\theta, d).$$

When there exists a  $d'$  that is strictly better than  $d$ , then  $d$  is called inadmissible.

### Example 10.2.1 Using only one of the observations

Let, for  $n \geq 2$ ,  $X_1, \dots, X_n$  be i.i.d., with  $g(\theta) := E_\theta(X_i) := \mu$ , and  $\text{var}(X_i) = 1$  (for all  $i$ ). Take quadratic loss  $L(\theta, a) := |\mu - a|^2$ . Consider  $d'(X_1, \dots, X_n) := \bar{X}_n$  and  $d(X_1, \dots, X_n) := X_1$ . Then,  $\forall \theta$ ,

$$R(\theta, d') = \frac{1}{n}, \quad R(\theta, d) = 1,$$

so that  $d$  is inadmissible.

$$\text{Var}(\mu - \bar{X}_n) = \text{Var}\left(\frac{1}{n} \sum_i X_i\right) = \frac{1}{n^2}$$

Note

We note that to show that a decision  $d$  is inadmissible, it suffices to find a strictly better  $d'$ . On the other hand, to show that  $d$  is admissible, one has to verify that there is no strictly better  $d'$ . So in principle, one then has to take all possible  $d'$  into account.

### 10.2.1 Not using the data at all is admissible $\heartsuit$ Interesting!!

Let  $L(\theta, a) := |g(\theta) - a|^r$  and  $d(X) := g(\theta_0)$ , where  $\theta_0$  is some fixed given value.

**Lemma 10.2.1** Assume that  $P_{\theta_0}$  dominates  $P_\theta$ <sup>1</sup> for all  $\theta$ . Then  $d$  is admissible.

<sup>1</sup>Let  $P$  and  $Q$  be probability measures on the same measurable space. Then  $P$  dominates  $Q$  if for all measurable  $B$ ,  $P(B) = 0$  implies  $Q(B) = 0$  ( $Q$  is absolutely continuous with respect to  $P$ ).

**Proof.**

Suppose that  $d'$  is better than  $d$ . Then we have

$$E_{\theta_0} |g(\theta_0) - d'(X)|^r \leq 0.$$

This implies that

$$d'(X) = g(\theta_0), P_{\theta_0}\text{-almost surely.} \quad (10.1)$$

Since by (10.1),

$$P_{\theta_0}(d'(X) \neq g(\theta_0)) = 0$$

the assumption that  $P_{\theta_0}$  dominates  $P_\theta, \forall \theta$ , implies now

$$P_\theta(d'(X) \neq g(\theta_0)) = 0, \forall \theta.$$

That is, for all  $\theta$ ,  $d'(X) = g(\theta_0)$ ,  $P_\theta$ -almost surely, and hence, for all  $\theta$ ,  $R(\theta, d') = R(\theta, d)$ . So  $d'$  is not strictly better than  $d$ . We conclude that  $d$  is admissible.  $\square$

### 10.2.2 A Neyman Pearson test is admissible

Consider testing

$$H_0 : \theta = \theta_0$$

against the alternative

$$H_1 : \theta = \theta_1.$$

Define the risk  $R(\theta, \phi)$  of a test  $\phi$  as the probability of error of first and second kind:

$$R(\theta, \phi) := \begin{cases} E_\theta \phi(X), & \theta = \theta_0 \\ 1 - E_\theta \phi(X), & \theta = \theta_1 \end{cases}.$$

We let  $p_0$  ( $p_1$ ) be the density of  $P_{\theta_0}$  ( $P_{\theta_1}$ ) with respect to some dominating measure  $\nu$  (for example  $\nu = P_{\theta_0} + P_{\theta_1}$ ). A Neyman Pearson test is (see Section 7.1)

$$\phi_{NP} := \begin{cases} 1 & \text{if } p_1/p_0 > c \\ q & \text{if } p_1/p_0 = c \\ 0 & \text{if } p_1/p_0 < c \end{cases}.$$

Here  $0 \leq q \leq 1$ , and  $0 \leq c < \infty$  are given constants.

**Lemma 10.2.2** A Neyman Pearson test is admissible if and only if one of the following two cases hold:

i) its power is strictly less than 1,

or

ii) it has minimal level among all tests with power 1.

**Proof.** Suppose  $R(\theta_0, \phi) < R(\theta_0, \phi_{NP})$ . Then from the Neyman Pearson Lemma, we know that either  $R(\theta_1, \phi) > R(\theta_1, \phi_{NP})$  (i.e., then  $\phi$  is not better than  $\phi_{NP}$ ), or  $c = 0$ . But when  $c = 0$  it holds that  $R(\theta_1, \phi_{NP}) = 0$ , i.e. then  $\phi_{NP}$  has power one.

Similarly, suppose that  $R(\theta_1, \phi) < R(\theta_1, \phi_{NP})$ . Then it follows from the Neyman Pearson Lemma that  $R(\theta_0, \phi) > R(\theta_0, \phi_{NP})$ , because we assume  $c < \infty$ .

So either one better and one worse. You can not ameliorate both and have no admissible

$\Leftrightarrow$  admissibility

higher type I error

$$R_{\theta_1}[\phi_{NP}(X)] - R_{\theta_1}[\phi] \leq c [R_{\theta_1}[\phi] - R_{\theta_1}[\phi_{NP}]]$$

ideally you want the both to be negative

NP Lemma:

lower level

### 10.3 Minimax decisions

**Definition 10.3.1** A decision  $d$  is called **minimax** if

$$\sup_{\theta} R(\theta, d) = \inf_{d'} \sup_{\theta} R(\theta, d').$$

Thus, the minimax criterion concerns the best decision in the worst possible case.

represented by

$$\sup_{\theta}.$$

#### 10.3.1 Minimax Neyman Pearson test

**Lemma 10.3.1** A Neyman Pearson test  $\phi_{NP}$  is minimax, if and only if  $R(\theta_0, \phi_{NP}) = R(\theta_1, \phi_{NP})$ .

**Proof.** Let  $\phi$  be a test, and write for  $j = 0, 1$ ,

$$\text{as the lemma above } r_j := R(\theta_j, \phi_{NP}), \quad r'_j := R(\theta_j, \phi).$$

Suppose that  $r_0 = r_1$  and that  $\phi_{NP}$  is not minimax. Then, for some test  $\phi$ ,

as both equal

$$\max_j r'_j < \max_j r_j.$$

This implies that both

$$r'_0 < r_0, \quad r'_1 < r_1$$

max must be smaller

and by the Neyman Pearson Lemma, this is not possible.

Let  $S = \{(R(\theta_0, \phi), R(\theta_1, \phi)) : \phi : \mathcal{X} \rightarrow [0, 1]\}$ . Note that  $S$  is convex. Thus, if  $r_0 < r_1$ , we can find a test  $\phi$  with  $r_0 < r'_0 < r_1$  and  $r'_1 < r_1$ . So then  $\phi_{NP}$  is not minimax. Similarly if  $r_0 > r_1$ .

more right here less risk here.

so convex  
you trade Type I with Type II

□

### 10.4 Bayes decisions

Suppose the parameter space  $\Theta$  is a measurable space. We can then equip it with a probability measure  $\Pi$ . We call  $\Pi$  the a priori distribution.

**Definition 10.4.1** The **Bayes risk** (with respect to the probability measure  $\Pi$ ) is

$$r(\Pi, d) := \int_{\Theta} R(\vartheta, d) d\Pi(\vartheta).$$

A decision  $d$  is called **Bayes** (with respect to  $\Pi$ ) if

$$r(\Pi, d) = \inf_{d'} r(\Pi, d').$$

So it is the parameter set that minimizes the risk under the current prior probability measure.

Let  $\Pi$  have density  $w := d\Pi/d\mu$  with respect to some dominating measure  $\mu$ .  
We may then write

$$r(\Pi, d) = \int_{\Theta} R(\vartheta, d)w(\vartheta)d\mu(\vartheta) := r_w(d).$$

Thus, the Bayes risk may be thought of as taking a weighted average of the risks. For example, one may want to assign more weight to "important" values of  $\theta$ .

#### 10.4.1 Bayes test

Consider again the testing problem

$$H_0 : \theta = \theta_0$$

against the alternative

$$H_1 : \theta = \theta_1.$$

Let  $L(\theta_0, a) := a$  and  $L(\theta_1, a) := 1 - a$ ,  $w(\theta_0) =: w_0$  and  $w(\theta_1) =: w_1 = 1 - w_0$ .  
Then

$$r_w(\phi) := w_0 R(\theta_0, \phi) + w_1 R(\theta_1, \phi).$$

We take  $0 < w_0 = 1 - w_1 < 1$ .

**Lemma 10.4.1** Bayes test is

$$\phi_{\text{Bayes}} = \begin{cases} 1 & \text{if } p_1/p_0 > w_0/w_1 \\ q & \text{if } p_1/p_0 = w_0/w_1 \\ 0 & \text{if } p_1/p_0 < w_0/w_1 \end{cases}.$$

**Proof.**

$$\begin{aligned} r_w(\phi) &= w_0 \underbrace{\int \phi p_0}_{\text{Type I error}} + w_1 \underbrace{\left(1 - \int \phi p_1\right)}_{\text{Type II error}} \\ &= \int \phi (w_0 p_0 - w_1 p_1) + w_1. \end{aligned}$$

So we choose  $\phi \in [0, 1]$  to minimize  $\phi(w_0 p_0 - w_1 p_1)$ . This is done by taking

$$\phi = \begin{cases} 1 & \text{if } w_0 p_0 - w_1 p_1 < 0 \\ q & \text{if } w_0 p_0 - w_1 p_1 = 0 \\ 0 & \text{if } w_0 p_0 - w_1 p_1 > 0 \end{cases},$$

where for  $q$  we may take any value between 0 and 1. □

Note that

$$2r_w(\phi_{\text{Bayes}}) = 1 - \int |w_1 p_1 - w_0 p_0|. \quad \Leftrightarrow \quad r_w = \frac{1}{2} - \frac{1}{2} \int |w_1 p_1 - w_0 p_0|.$$

In particular, when  $w_0 = w_1 = 1/2$ ,

$$2r_w(\phi_{\text{Bayes}}) = 1 - \int |p_1 - p_0|/2,$$

i.e., the risk is large if the two densities are close to each other.

$$\begin{aligned} 2r_w(\phi_{\text{Bayes}}) &= 2w_1 + 2 \int \frac{w_1 p_1 - w_0 p_0}{w_0 p_0 - w_1 p_0} \\ &\stackrel{+ w_0}{=} 2w_1 + 2 \int \frac{w_1 p_1 - w_0 p_0}{w_1 p_1 - w_0 p_0} \\ &\stackrel{= 1}{=} 1 + \int \frac{w_1 p_1 - w_0 p_0}{w_1 p_1 - w_0 p_0} + 2 \int \frac{w_0 p_0 - w_1 p_0}{w_1 p_1 - w_0 p_0} \\ &\stackrel{\text{decomposed into 2 parts}}{=} 1 + S(w_0 p_0 - w_1 p_0) + 2S(w_1 p_1 - w_0 p_0). \end{aligned}$$

## 10.5 Construction of Bayes estimators

Let  $X$  have distribution  $P \in \mathcal{P} := \{P_\theta : \theta \in \Theta\}$ . Suppose  $\mathcal{P}$  is dominated by a ( $\sigma$ -finite) measure  $\nu$ , and let  $p_\theta = dP_\theta/d\nu$  denote the densities. Let  $\Pi$  be an a priori distribution on  $\Theta$ , with density  $w := d\Pi/d\mu$ . We now think of  $p_\theta$  as the density of  $X$  given the value of  $\theta$ . We write it as

$$p_\theta(x) = p(x|\theta), \quad x \in \mathcal{X}. \quad \text{where unknown}$$

Moreover, we define the marginal density

$$p(\cdot) := \int_{\Theta} p(\cdot|\vartheta)w(\vartheta)d\mu(\vartheta). \quad \begin{array}{l} f(x, \omega) \\ \text{prob of data} \\ \text{given param.} \\ \Rightarrow \text{likelihood.} \end{array}$$

$$p(a \wedge b) = \frac{p(a) \cdot p(b|a)}{p(b)}$$

**Definition 10.5.1** The a posteriori density of  $\theta$  is

$$w(\vartheta|x) = \frac{p(x|\vartheta)w(\vartheta)}{p(x)}, \quad \vartheta \in \Theta, \quad x \in \mathcal{X}. \quad \begin{array}{l} \text{this is} \\ \text{the classical} \\ \text{formula.} \end{array}$$

**Lemma 10.5.1** Given the data  $X = x$ , consider  $\theta$  as a random variable with density  $w(\vartheta|x)$ . Let

$$l(x, a) := E[L(\theta, a)|X = x] = \int_{\Theta} L(\vartheta, a)w(\vartheta|x)d\mu(\vartheta),$$

and

$$d(x) := \arg \min_a l(x, a).$$

Then  $d$  is Bayes decision  $d_{\text{Bayes}}$ .

**Proof.**

$$\text{we are proving } r_w(d') \geq r_w(d) \text{ or } (d_{\text{Bayes}}) - \inf_d r_w(d)$$

$$\begin{aligned} r_w(d') &= \int_{\Theta} R(\vartheta, d')w(\vartheta)d\mu(\vartheta) \\ &= \int_{\Theta} \left[ \int_{\mathcal{X}} L(\vartheta, d'(x))p(x|\vartheta)d\nu(x) \right] w(\vartheta)d\mu(\vartheta) \\ &= \int_{\mathcal{X}} \left[ \int_{\Theta} L(\vartheta, d'(x))w(\vartheta|x)d\mu(\vartheta) \right] p(x)d\nu(x) \\ &= \int_{\mathcal{X}} l(x, d'(x))p(x)d\nu(x) \\ &\geq \int_{\mathcal{X}} l(x, d(x))p(x)d\nu(x) \quad \text{by def. of } d(x) \\ &= r_w(d). \quad \text{which is } \arg \min_a l(x, a) \end{aligned}$$

Notice that changing the order of integration!

General reminder:  
you can always change the order of the integration but remember that the value boundaries will change based on the order of integration.

The ultimate goal is to compute the same area.

**Corollary 10.5.1** The Bayes decision is

$$d_{\text{Bayes}}(X) = \arg \min_{a \in \mathcal{A}} l(X, a),$$

$$f(x|\theta) = g(S(x)) h(x)$$

104

## CHAPTER 10. DECISION THEORY

$$w(\vartheta|x) = \frac{f(\vartheta, x)}{p(x)}$$

where

so if there  
is such decomposition  
as per Le---Schott  
then Risk  
depends on  
sufficient statistics

$$\begin{aligned} l(x, a) &= E(L(\theta, a)|X = x) = \int L(\vartheta, a) w(\vartheta|x) d\mu(\vartheta) \\ &= \int L(\vartheta, a) g_\vartheta(S(x)) w(\vartheta) d\mu(\vartheta) h(x)/p(x). \end{aligned}$$

So

$$d_{\text{Bayes}}(X) = \arg \min_{a \in \mathcal{A}} \int L(\vartheta, a) g_\vartheta(S) w(\vartheta) d\mu(\vartheta),$$

which only depends on the sufficient statistic S.

on stage  
if integrating  
over  
 $\vartheta$ .

### 10.5.1 Bayes test revisited

For the testing problem  
 $H_0 : \theta = \theta_0$   
against the alternative  
 $H_1 : \theta = \theta_1$ ,  
with loss function

This is central!! } As just two cases (two param)  
time risks.

$$l(x, \phi) = E(L|X=x)$$

looks like  $a = \text{decision}$  we have  $L(\theta_0, a) := a, L(\theta_1, a) := 1 - a, a \in \{0, 1\}$ ,  
 $\Rightarrow$  test } loss and test } 1-1 relation

Thus,

$$\arg \min_{\phi} l(\cdot, \phi) = \begin{cases} 1 & \text{if } w_1 p_1 > w_0 p_0 \\ q & \text{if } w_1 p_1 = w_0 p_0 \\ 0 & \text{if } w_1 p_1 < w_0 p_0 \end{cases}.$$

only difference } by def } this new here!!

### 10.5.2 Bayes estimator for quadratic loss

No structure on the weights but using instead the loss function of above.

$$B(x) = p(x|0) - w(0) = p(x) - p(x|x)$$

In the next result, we shall use:

**Lemma 10.5.2** Let  $Z$  be a real-valued random variable. Then

$$\arg \min_{a \in \mathbb{R}} E(Z - a)^2 = EZ. \quad \left. \begin{array}{l} \text{with squared} \\ \text{error Bayesian} \\ \text{estimator} \Rightarrow \text{expected} \\ \text{value} \end{array} \right\}$$

Proof.

$$E(Z - a)^2 = \text{var}(Z) + (a - EZ)^2. \quad \left. \begin{array}{l} \text{var}(Z) \\ + \text{bias}^2 \end{array} \right\} \text{this minimal when } a = EZ$$

Consider the case  $\mathcal{A} = \mathbb{R}$  and  $\Theta \subseteq \mathbb{R}$ . Let  $L(\theta, a) := |\theta - a|^2$ . Then

$$d_{\text{Bayes}}(X) = \underline{E(\theta|X)}.$$

} follows from the lemma above.

For quadratic loss, and for  $T = E(\theta|X)$ , the Bayes risk of an estimator  $T'$  is

$$r_w(T') = E\text{var}(\theta|X) + E(T - T')^2$$

! This result is then used in chapter 11.

(compare with Lemma 5.2.2). This follows from straightforward calculations:

$$r_w(T') = \int R(\vartheta, T') w(\vartheta) d\mu(\vartheta)$$

$$= ER(\theta, T') = E(\theta - T')^2 = E\left[E\left((\theta - T')^2 \mid X\right)\right]$$

and, with  $\theta$  being the random variable,  $\underbrace{\text{var}}_{\text{var}} + \underbrace{\text{bias}^2}_{\text{bias}}$

$$\begin{aligned} E\left((\theta - T')^2 \mid X\right) &= E\left((\theta - T)^2 \mid X\right) + (T - T')^2 \\ &= \text{var}(\theta \mid X) + (T - T')^2. \end{aligned}$$

### 10.5.3 Bayes estimator and the maximum a posteriori estimator

Consider again the case  $\Theta \subseteq \mathbb{R}$ , and  $\mathcal{A} = \Theta$ , and now with loss function  $L(\theta, a) := \mathbf{1}\{| \theta - a | > c\}$  for a given constant  $c > 0$ . Then

$$l(x, a) = \Pi(|\theta - a| > c \mid X = x) = \int_{|\theta - a| > c} \frac{1}{w(\vartheta|x)} w(\vartheta|x) d\vartheta.$$

We note that for  $c \rightarrow 0$

$$\frac{1 - l(x, a)}{2c} = \frac{\Pi(|\theta - a| \leq c \mid X = x)}{2c} \approx w(a|x) = p(x|a) \frac{w(a)}{p(x)}.$$

Thus, for  $c$  small, Bayes rule is approximately  $d_0(x) := \arg \max_{a \in \Theta} p(x|a)w(a)$ . The estimator  $d_0(X)$  is called the *maximum a posteriori estimator* (MAP). If  $w$  is the uniform density on  $\Theta$  (which only exists if  $\Theta$  is bounded), then  $d_0(X)$  is the maximum likelihood estimator.

### 10.5.4 Three worked-out examples

In many examples, it saves a lot of work not to write out complete expressions for the posterior  $w(\vartheta|x)$ . The main interest (at first) is how it depends on  $\vartheta$  and all the other expressions can (at least theoretically, it may be difficult computationally) be found later by using that a density integrates to one. We therefore recall the symbol  $\propto$  (see Section 4.9). For example, we may write

$$\begin{aligned} w(\vartheta|x) &= p(x|\vartheta)w(\vartheta)/p(x) \\ &\propto p(x|\vartheta)w(\vartheta) \end{aligned}$$

because the marginal density  $p(x)$  does not depend on  $\vartheta$ .

As we will see, in the following three examples the posterior is in the same family as the prior. We call them conjugate priors for the distribution concerned.

**Example 10.5.1 Poisson with Gamma prior**

Suppose that given  $\theta$ ,  $X$  has Poisson distribution with parameter  $\theta$ , and that  $\theta$  has the  $\text{Gamma}(k, \lambda)$ -distribution. The density of  $\theta$  is then

$$w(\vartheta) = \lambda^k \vartheta^{k-1} e^{-\lambda\vartheta} / \Gamma(k),$$

where

$$\Gamma(k) = \int_0^\infty e^{-z} z^{k-1} dz.$$

The  $\text{Gamma}(k, \lambda)$  distribution has mean

$$E\theta = \int_0^\infty \vartheta w(\vartheta) d\vartheta = \frac{k}{\lambda}.$$

The a posteriori density is then

$$\begin{aligned} w(\vartheta|x) &= p(x|\vartheta) \frac{w(\vartheta)}{p(x)} \\ &= e^{-\vartheta} \frac{\vartheta^x \lambda^k \vartheta^{k-1} e^{-\lambda\vartheta} / \Gamma(k)}{x!} \\ &= e^{-\vartheta(1+\lambda)} \vartheta^{k+x-1} c(x, k, \lambda), \end{aligned}$$

where  $c(x, k, \lambda)$  is such that

$$\int w(\vartheta|x) d\vartheta = 1.$$

With the  $\propto$  notation

$$\begin{aligned} w(\vartheta|x) &\propto p(x|\vartheta) w(\vartheta) \\ &\propto \underbrace{e^{-\vartheta(1+\lambda)} \vartheta^{k+x-1}}_{\text{constants}}. \end{aligned}$$

Recall (cal)  
drop all  
other terms?  
not depending on  $\vartheta$   
at  $\propto e^{-\vartheta(1+\lambda)} \cdot \vartheta^{k+x-1}$   
 $+ e^{-\vartheta(1+\lambda)} \cdot \vartheta^{k+x-1}$   
 $\cdot \vartheta^{k+x-2}$

You see this saves a lot of writing. We recognize  $w(\vartheta|x)$  as the density of the  $\text{Gamma}(k+x, 1+\lambda)$ -distribution. Bayes estimator with quadratic loss is thus

$$E(\theta|X) = \frac{k+X}{1+\lambda}.$$

The maximum a posteriori estimator is

$$\arg \max_a p(x|a) p(a)$$

$$\frac{k+X-1}{1+\lambda}. \quad \text{Compute the maximum}$$

$$e^{-k-X+1} = 0$$

**Example 10.5.2 Binomial distribution with Beta prior**

Suppose given  $\theta$ ,  $X$  has the  $\text{Binomial}(n, \theta)$ -distribution, and that  $\theta$  is uniformly distributed on  $[0, 1]$ . Then

$$w(\vartheta|x) = \binom{n}{x} \vartheta^x (1-\vartheta)^{n-x} / p(x).$$

$$w(\vartheta|x) = \frac{p(x|\vartheta) p(\vartheta)}{p(x)}$$

$$= \binom{n}{x} \vartheta^x (1-\vartheta)^{n-x} \cdot \frac{1}{1}$$

maximum  
of Gamma

$$\text{Beta} = \frac{1}{\Gamma(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \Gamma(\beta)} \frac{x^{\alpha-1}}{\alpha} (1-x)^{\beta-1}$$

## 10.6. DISCUSSION OF BAYESIAN APPROACH

107

This is the density of the Beta( $x+1, n-x+1$ )-distribution. Thus, with quadratic loss, Bayes estimator is

$$E(\theta|X) = \frac{X+1}{n+2}.$$

More generally, suppose that  $X$  is binomial( $n, \theta$ ) and that  $\theta$  has the Beta( $r, s$ )-prior

$$w(\vartheta) = \frac{\Gamma(r+s)}{\Gamma(r)\Gamma(s)} \vartheta^{r-1} (1-\vartheta)^{s-1}, \quad 0 < \vartheta < 1.$$

Here  $r$  and  $s$  are given positive numbers. The prior expectation is

$$E\theta = \frac{r}{r+s}.$$

Bayes estimator under quadratic loss is the posterior expectation

$$E(\theta|X) = \frac{X+r}{n+r+s}.$$

### Example 10.5.3 Normal distribution with normal prior

Let  $X \sim \mathcal{N}(\theta, 1)$  and  $\theta \sim \mathcal{N}(c, \tau^2)$  for some  $c$  and  $\tau^2$ . We have

$$\begin{aligned} w(\vartheta|x) &= \frac{p(x|\vartheta)w(\vartheta)}{p(x)} \\ &\propto \phi(x-\vartheta)\phi\left(\frac{\vartheta-c}{\tau}\right) \\ &\propto \exp\left[-\frac{1}{2}\left\{(x-\vartheta)^2 + \frac{(\vartheta-c)^2}{\tau^2}\right\}\right] \\ &\propto \exp\left[-\frac{1}{2}\left\{\vartheta - \frac{\tau^2 x + c}{\tau^2 + 1}\right\}^2 \frac{1+\tau^2}{\tau^2}\right]. \end{aligned}$$

We conclude that Bayes estimator for quadratic loss is

$$T_{\text{Bayes}} = E(\theta|X) = \frac{\tau^2 X + c}{\tau^2 + 1}. \quad \left\{ \begin{array}{l} \text{In the case of } c=0 \\ E(\theta|X) = \left(\frac{\tau^2}{\tau^2+1}\right)X \\ = \lambda X \quad 0 \leq \lambda < 1 \end{array} \right.$$

## 10.6 Discussion of Bayesian approach

A main objection against the Bayesian approach is that it is generally subjective. The final estimator depends strongly on the choice of the prior distribution. On the other hand, Bayesian methods are very powerful and often quite natural. The prior may be inspired by or estimated from previous data sets, in which case the above subjectivity problem becomes less pregnant. Furthermore, in complicated models with many unknown parameters, Bayesian methods are a welcome tool for developing sensible algorithms.

so that you have a shrinkage factor, comparable to the one of 10.1.3.

**Credibility sets.** A (frequentist) confidence set for a parameter of interest can be hard to find, and is also less easy to explain to “non-experts”. The Bayesian version of a confidence set is called a *credibility set*, which generally is seen as

an intuitively much clearer concept. For example, in the case of a real-valued parameter  $\theta$ , a  $(1 - \alpha)$ -credibility interval is defined as

$$I := [\hat{\theta}_L(X), \hat{\theta}_R(X)],$$

where the endpoints  $\hat{\theta}_L$  and  $\hat{\theta}_R$  are chosen in such a way that

$$\int_{\hat{\theta}_L(X)}^{\hat{\theta}_R(X)} w(\vartheta|X) d\vartheta = (1 - \alpha).$$

Thus, it is the set which has posterior probability  $(1 - \alpha)$ . A  $(1 - \alpha)$ -credibility set is generally not a  $(1 - \alpha)$ -confidence set, i.e., from a frequentist point of view, its properties are not always clear.

**Pragmatic point of view.** The Bayesian approach is fruitful for the construction of estimators. One can then proceed by studying the frequentist properties of the Bayesian procedure. For example, in the  $\text{Binomial}(n, \theta)$ -model with a uniform prior on  $\theta$ , the Bayes estimator is

$$\hat{\theta}_{\text{Bayes}}(X) = \frac{X + 1}{n + 2}.$$

Given this estimator, one can “forget” we obtained it by Bayesian arguments, and study for example its (frequentist) mean square error.

**Complexity regularization.** Here is a “toy” example, where a Bayesian method helps constructing a useful procedure. Let  $X_1, \dots, X_n$  be independent random variables, where  $X_i$  is  $\mathcal{N}(\theta_i, 1)$ -distributed. The  $n$  parameters  $\theta_i$  are all unknown. Thus, there are as many observations as unknowns, a situation where *complexity regularization* is needed. *Complexity regularization* (see Chapter 16) means that in principle, one allows for any parameter value, but that one pays a price for choosing “complex” values. What “complexity” means depends on the situation at hand. We consider in this example the situation where complexity is the opposite of *sparsity*, where the *sparseness* of a vector  $\vartheta$  is defined as its number of non-zero entries. Consider the estimator

$$\hat{\theta} := \arg \min_{\vartheta \in \mathbb{R}^n} \sum_{i=1}^n (X_i - \vartheta_i)^2 + 2\lambda \sum_{i=1}^n |\vartheta_i|,$$

where  $\lambda > 0$  is a regularization parameter. Note that when  $\lambda = 0$ , one has  $\hat{\theta}_i = X_i$  for all  $i$ , whereas on the other extreme, when  $\lambda = \infty$ , one has  $\hat{\theta} \equiv 0$ . The larger  $\lambda$ , the more sparse the estimator will be. In fact, it is easy to verify that for  $i = 1, \dots, n$ ,

$$\hat{\theta}_i = \begin{cases} X_i - \lambda & X_i > \lambda \\ 0 & |X_i| \leq \lambda \\ X_i + \lambda & X_i < -\lambda \end{cases}.$$

This is called the *soft thresholding* estimator. The procedure corresponds to Bayesian maximum a posteriori estimation, with double-exponential (also called Laplace) prior. Indeed, suppose that the prior is  $\theta_1, \dots, \theta_n$  i.i.d. with density

$$w(z) = \underbrace{\frac{1}{\tau\sqrt{2}} \exp\left[-\frac{\sqrt{2}|z|}{\tau}\right]}_{\text{Laplace}}, z \in \mathbb{R},$$

where  $\tau > 0$  is the prior scale parameter ( $\tau^2$  is the variance of this distribution). Given  $X_1, \dots, X_n$ , the posterior distribution of the vector  $\vartheta$  is then

$$(2\pi)^{-n/2} \exp \left[ -\frac{\sum_{i=1}^n (X_i - \vartheta_i)^2}{2} \right] \times (2\pi\tau)^{-n/2} \exp \left[ -\frac{\sqrt{2} \sum_{i=1}^n |\vartheta_i|}{\tau} \right].$$

*seems  
with variance*

Thus,  $\hat{\vartheta}$  with regularization parameter  $\lambda = \sqrt{2}/\tau$  is the maximum a posteriori estimator.

**Bayesian methods as theoretical tool.** In Chapter 11 we will illustrate the fact that Bayesian methods can be exploited as a tool for proving for example frequentist lower bounds. We will see for instance that the Bayesian estimator with constant risk is also the minimax estimator. The idea in such results is to look for “worst possible priors”.

## 10.7 Integrating parameters out \*

Striving at flexible prior distributions one can model them depending on another “hyper-parameter”, say  $\tau$ , i.e., in formula

$$w(\vartheta) := w(\vartheta|\tau).$$

Keeping  $\tau$  fixed and integrating  $\vartheta$  out, the density of  $X$  is then

$$\tilde{p}(x|\tau) := \int p(x|\vartheta)w(\vartheta|\tau)d\mu(\vartheta).$$

One can proceed by estimating  $\tau$ , using for instance maximum likelihood (this is generally computationally quite hard), or the methods of moments. One then obtains a prior  $w(\vartheta|\hat{\tau})$  with estimated parameter  $\hat{\tau}$ . The prior is thus based on the data. The whole procedure is called *empirical Bayes*.

### Example 10.7.1 Poisson with Gamma prior with hyperparameters

Suppose  $X_1, \dots, X_n$  are independent and  $X_i$  has a  $\text{Poisson}(\theta_i)$ -distribution,  $i = 1, \dots, n$ . Assume moreover that  $\theta_1, \dots, \theta_n$  are i.i.d. with  $\text{Gamma}(k, \lambda)$ -distribution, i.e., each has prior density

$$w(z|k, \lambda) = e^{-\lambda z} z^{k-1} \lambda^k / \Gamma(k), \quad z > 0.$$

Both  $k$  and  $\lambda$  are considered as hyper-parameters. Then the density of  $X_1, \dots, X_n$  is

$$\tilde{p}(x_1, \dots, x_n|k, \lambda) \Leftrightarrow \rho(\vartheta) \text{ is the above}$$

$$\begin{aligned} &\propto \int \left( e^{-\sum_{i=1}^n \vartheta_i} \prod_{i=1}^n \vartheta_i^{x_i} e^{-\lambda \sum_{i=1}^n \vartheta_i} \prod_{i=1}^n \vartheta_i^{k-1} \frac{\lambda^k}{\Gamma(k)} \right) d\vartheta_1 \cdots d\vartheta_n \\ &= \prod_{i=1}^n \frac{\Gamma(x_i + k)}{\Gamma(k)} p^k (1-p)^{x_i+k-1}, \end{aligned}$$

where  $p := \lambda/(1 + \lambda)$ . Thus, under  $\tilde{\mathbf{p}}(\cdot|k, \lambda)$ , the observations  $X_1, \dots, X_n$  are independent and  $X_i$  has a negative binomial distribution with parameters  $k$  and  $p$  (check the formula for the negative binomial distribution, see e.g. Example 2.4.1). The mean and variance of the negative binomial distribution can be calculated directly or looked up in a textbook. We then find (for  $i = 1, \dots, n$ ),

$$E(X_i|k, \lambda) = \frac{k(1-p)}{p} = \frac{k}{\lambda}$$

and

$$\text{var}(X_i|k, \lambda) = \frac{k(1-p)}{p^2} = \frac{k(1+\lambda)}{\lambda^2}.$$

We use the method of moments to estimate  $k$  and  $\lambda$ . Let  $\bar{X}_n$  be the sample mean and  $S_n^2 := \sum_{i=1}^n (X_i - \bar{X})^2/(n-1)$  be the sample variance. We solve

$$\frac{\hat{k}}{\hat{\lambda}} = \bar{X}_n, \quad \frac{\hat{k}(1+\hat{\lambda})}{\hat{\lambda}^2} = S_n^2.$$

This yields

$$\hat{k} = \frac{\bar{X}_n^2}{S_n^2 - \bar{X}_n}, \quad \hat{\lambda} = \frac{\bar{X}_n}{S_n^2 - \bar{X}_n}.$$

For given  $k$  and  $\lambda$ , the Bayes estimator of  $\theta_i$  is given in Example 10.5.1. We now insert the estimated values of  $k$  and  $\lambda$  to get an empirical Bayes estimator

$$\hat{\theta}_i = \frac{X_i + \hat{k}}{1 + \hat{\lambda}} = X_i(1 - \bar{X}_n/S_n^2) + \bar{X}_n^2/S_n^2, \quad i = 1, \dots, n.$$

The MLE of  $\theta_i$  is  $X_i$  itself ( $i = 1, \dots, n$ ). We see that the empirical Bayes estimator uses all observations to estimate a particular  $\theta_i$ . The empirical Bayes estimator  $\hat{\theta}_i$  is a convex combination  $(1 - \alpha)X_i + \alpha\bar{X}_n$  of  $X_i$  and  $\bar{X}_n$ , with  $\alpha = \bar{X}_n/S_n^2$  generally close to one if the pooled sample has mean and variance approximately equal, i.e., if the pooled sample is “Poisson-like”.

# Chapter 11

# Proving admissibility and minimaxity

Bayes estimators are quite useful, also for obdurate frequentists. They can be used to construct estimators that are minimax (admissible), or for verification of minimaxity (admissibility).

Let us first recall the definitions. Let  $X \in \mathcal{X}$  have distribution  $P_\theta$ ,  $\theta \in \Theta$ . Let  $T = T(X)$  be a statistic (estimator, decision),  $L(\theta, a)$  be a loss function, and  $R(\theta, T) := E_\theta L(\theta, T(X))$  be the risk of  $T$ .

- $T$  is *minimax* if  $\forall T' \sup_{\theta} R(\theta, T) \leq \sup_{\theta} R(\theta, T')$ . {  
risk is strictly interior.
  - $T$  is *inadmissible* if  $\exists T': \{\forall \theta R(\theta, T') \leq R(\theta, T) \text{ and } \exists \theta R(\theta, T') < R(\theta, T)\}$ .
  - $T$  is *Bayes* (for the prior density  $w$  on  $\Theta$ ) if  $\forall T', r_w(T) \leq r_w(T')$ .

Great example!  
Remember it.

Recall also that Bayes risk for  $w$  is

$$r_w(T) = \int R(\vartheta, T) w(\vartheta) d\mu(\vartheta).$$

Whenever we say that a statistic  $T$  is Bayes, without referring to an explicit prior on  $\Theta$ , we mean that there exists a prior for which  $T$  is Bayes. Of course, if the risk  $R(\theta, T) = R(T)$  does not depend on  $\theta$ , then Bayes risk of  $T$  does not depend on the prior.

Especially in cases where one wants to use the uniform distribution as prior, but cannot do so because  $\Theta$  is not bounded, the notion *extended Bayes* is useful.

**Definition 11.0.1** A statistic  $T$  is called extended Bayes if there exists a sequence of prior densities  $\{w_m\}_{m=1}^{\infty}$  (w.r.t. dominating measures that are allowed to depend on  $m$ ) such that  $r_{w_m}(T) - \inf_{T'} r_{w_m}(T') \rightarrow 0$  as  $m \rightarrow \infty$ .

Idea:  
 You cannot take the uniform (constant) prior belief on the [entire] real line  $\mathbb{R}$ . So you take a big chunk out it and set it to a constant and hope the rest  $\rightarrow 0$ .

## 11.1 Minimaxity



**Lemma 11.1.1** Suppose  $T$  is a statistic with risk  $R(\theta, T) = R(T)$  not depending on  $\theta$ . Then

- (i)  $T$  admissible  $\Rightarrow T$  minimax,
- (ii)  $T$  Bayes  $\Rightarrow T$  minimax,
- and in fact more generally,
- (iii)  $T$  extended Bayes  $\Rightarrow T$  minimax.

Constant risk  
in  $\theta$

**Proof.**

(i)  $T$  is admissible, so for all  $T'$ , either there is a  $\theta$  with  $R(\theta, T') \leq R(T)$ , or  $R(\theta, T') \geq R(T)$  for all  $\theta$ . Hence  $\sup_{\theta} R(\theta, T') \geq R(T)$ .

(ii) Since Bayes implies extended Bayes, this follows from (iii). We nevertheless present a separate proof, as it is somewhat simpler than (iii).

Note first that for any  $T'$ ,

$$r_w(T') = \int R(\vartheta, T') w(\vartheta) d\mu(\vartheta) \quad (11.1)$$

$$\leq \int \underbrace{\sup_{\vartheta} R(\vartheta, T')}_{\text{constant}} w(\vartheta) d\mu(\vartheta) \quad (11.2)$$

$$= \sup_{\vartheta} R(\vartheta, T') \quad (11.3)$$

that is, Bayes risk is always bounded by the supremum risk. Suppose now that  $T'$  is a statistic with  $\sup_{\theta} R(\theta, T') < R(T)$ . Then

$$r_w(T') \leq \sup_{\vartheta} R(\vartheta, T') < R(T) = r_w(T), \quad \text{as again the } \sup_{\vartheta} R(\vartheta, T') \text{ is independent of } \theta \text{ you can } \theta$$

which is in contradiction with the assumption that  $T$  is Bayes.

(iii) Suppose for simplicity that a Bayes decision  $T_m$  for the prior  $w_m$  exists, for all  $m$ , i.e.

$$r_{w_m}(T_m) = \inf_{T'} r_{w_m}(T'), \quad m = 1, 2, \dots$$

By assumption, for all  $\epsilon > 0$ , there exists an  $m$  sufficiently large, such that

$$R(T) = r_{w_m}(T) \leq r_{w_m}(T_m) + \epsilon \leq r_{w_m}(T') + \epsilon \leq \sup_{\theta} R(\theta, T') + \epsilon,$$

because, as we have seen in (11.1), the Bayes risk is bounded by supremum risk. Since  $\epsilon$  can be chosen arbitrary small, this proves (iii).  $\square$

bounded by  
supremum.

**Example 11.1.1** Minimax estimator for Binomial distribution

Consider a Binomial( $n, \theta$ ) random variable  $X$ . Let the prior on  $\theta \in (0, 1)$  be the Beta( $r, s$ ) distribution. Then Bayes estimator for quadratic loss is

$$T = \frac{X + r}{n + r + s}$$

(see Example 10.5.2). Its risk is

$$R(\theta, T) = E_{\theta}(T - \theta)^2$$

here trick  
is to understand  
that  $T_m$   
Bayes estim (convex)  
and  $T$  any  
estim, such  
that by def  
 $r_{w_m}(T) - r_{w_m}(T_m) \leq \epsilon$   
(for  $m$  large  
enough).

cannot  
bc negative  
us bayes decision  
must be minimal  
loss bayes

$$\begin{aligned}
 &= \text{var}_\theta(T) + \text{bias}_\theta^2(T) \\
 &= \frac{n\theta(1-\theta)}{(n+r+s)^2} + \left[ \frac{n\theta+r}{n+r+s} - \frac{(n+r+s)\theta}{n+r+s} \right]^2 \\
 &= \frac{[(r+s)^2 - n]\theta^2 + [n - 2r(r+s)]\theta + r^2}{(n+r+s)^2}.
 \end{aligned}$$

*This can only be constant in  $\theta$  if the coefficients in front of  $\theta^2$  and  $\theta$  are zero.*

$$(r+s)^2 - n = 0, \quad n - 2r(r+s) = 0.$$

Solving for  $r$  and  $s$  gives

$$r = s = \sqrt{n}/2.$$

Plugging these values back in the estimator  $T$  gives

*is risk constant and Bayes*

$$T = \frac{X + \sqrt{n}/2}{n + \sqrt{n}}$$

*such that it always stops the Bayes estim.*

*is minimax* The minimax risk is

$$R(T) = \frac{1}{4(\sqrt{n} + 1)^2}.$$

We can compare this with the supremum risk of the unbiased estimator  $X/n$ :

$$\sup_\theta R(\theta, X/n) = \sup_\theta \frac{\theta(1-\theta)}{n} = \frac{1}{4n}.$$

So for large  $n$ , this does not differ much from the minimax risk.

### 11.1.1 Minimality of the Pitman estimator \*

We consider again the Pitman estimator (see Lemma 9.1.2)

$$T^* = \frac{\int z \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}{\int \mathbf{p}_0(X_1 - z, \dots, X_n - z) dz}.$$

**Lemma 11.1.2**  $T^*$  is extended Bayes (for quadratic loss).

**Proof.** Let  $w_m$  be (the density of) the uniform distribution on the interval  $[-m, m]$ :

$$w_m = \mathbf{l}_{[-m,m]}/2m.$$

The posterior density is then

$$w_m(\vartheta|x) = \frac{p_0(x - \vartheta) \mathbf{l}_{[-m,m]}(\vartheta)}{\int_{-m}^m p_0(x - \vartheta) d\vartheta}.$$

Bayes estimator is thus

$$T_m = \frac{\int_{-m}^m \vartheta p_0(x - \vartheta) d\vartheta}{\int_{-m}^m p_0(x - \vartheta) d\vartheta}.$$

We now compute  $R(\theta, T_m) = E_\theta(T_m - \theta)^2$ . Let

$$T_{a,b}(x) := \frac{\int_a^b z p_0(x-z) dz}{\int_a^b p_0(x-z) dz}.$$

Then for all  $x$ ,  $T_{a,b}(x) \rightarrow T(x)$  as  $a \rightarrow -\infty$  and  $b \rightarrow \infty$ . One can easily verify that also

$$\lim_{a \rightarrow -\infty, b \rightarrow \infty} E_0 T_{a,b}^2(X) \rightarrow E_0 T^2(X).$$

(Note that, for any prior  $w$ ,  $E_0 T^2(X)$  is the Bayes risk  $r_w(T)$  since the risk  $R(\theta, T) = E_0 T^2(X)$  does not depend on  $\theta$ .) Moreover

$$T_{a,b}(X) - \theta = \frac{\int_a^b (z - \theta) p_0(X-z) dz}{\int_a^b p_0(x-z) dz} = \frac{\int_{a-\theta}^{b-\theta} v p_0(X-\theta-v) dv}{\int_{a-\theta}^{b-\theta} p_0(X-\theta-v) dv}.$$

It follows that

$$E_\theta(T_{a,b}(X) - \theta)^2 = E_0 T_{a-\theta, b-\theta}^2(X).$$

Hence,

$$R(\theta, T_m) = E_0 T_{-m-\theta, m-\theta}^2(X).$$

The Bayes risk is

$$r_{w_m}(T_m) = E_{\theta \sim w_m} R(\theta, T_m) = \frac{1}{2m} \int_{-m}^m E_0 T_{-m-\vartheta, m-\vartheta}^2(X) d\vartheta.$$

Hence, for any  $0 < \epsilon < 1$ , we have

$$\begin{aligned} r_{w_m}(T_m) &\geq \inf_{|\vartheta| \leq m(1-\epsilon)} (1-\epsilon) E_0 T_{-m-\vartheta, m-\vartheta}^2(X) \\ &\geq \inf_{a \leq -m\epsilon, b \geq m\epsilon} (1-\epsilon) E_0 T_{a,b}^2(X). \end{aligned}$$

It follows that for any  $0 < \epsilon < 1$ ,

$$\liminf_{m \rightarrow \infty} r_{w_m}(T_m) \geq \liminf_{m \rightarrow \infty} \inf_{a \leq -m\epsilon, b \geq m\epsilon} (1-\epsilon) E_0 T_{a,b}^2(X) = (1-\epsilon) E_0 T^2(X).$$

Hence we have  $r_{w_m}(T_m) \rightarrow E_0 T^2(X)$ , i.e.,  $r_{w_m}(T_m) - r_{w_m}(T) \rightarrow 0$ .  $\square$

**Corollary 11.1.1**  $T^*$  is minimax (for quadratic loss).

## 11.2 Admissibility

In this section, the parameter space is assumed to be an open subset of a topological space, so that we can consider open neighborhoods of members of  $\Theta$ , and continuous functions on  $\Theta$ . We moreover restrict ourselves to statistics  $T$  with  $R(\theta, T) < \infty$ .

$$r_w = \int L(\theta, d) \cdot w(\theta) \cdot d\mu(\theta)$$

**Lemma 11.2.1** Suppose that the statistic  $T$  is Bayes for the prior density  $w$ . Then (i) or (ii) below are sufficient conditions for the admissibility of  $T$ .

- (i) The statistic  $T$  is the unique Bayes decision (i.e.,  $r_w(T) = r_w(T')$  implies that  $\forall \theta, T = T'$   $P_\theta$ -almost surely),
- (ii) For all  $T'$ ,  $R(\theta, T')$  is continuous in  $\theta$ , and moreover, for all open  $U \subset \Theta$ , the prior probability  $\Pi(U) := \int_U w(\vartheta) d\mu(\vartheta)$  of  $U$  is strictly positive.

**Proof.**

- (i) Suppose that for some  $T'$ ,  $R(\theta, T') \leq R(\theta, T)$  for all  $\theta$ . Then also  $r_w(T') \leq r_w(T)$ . Because  $T$  is Bayes, we then must have equality:

$$\underbrace{r_w(T')}_{\text{minimal}} = r_w(T).$$

by def of  
 $r_w(T)$ ; all terms  
equal with  $r_w(T')$   
so that  $r_w(T') < r_w(T)$

So then,  $\forall \theta, T'$  and  $T$  are equal  $P_\theta$ -a.s., and hence,  $\forall \theta, R(\theta, T') = R(\theta, T)$ , so that  $T'$  can not be strictly better than  $T$ .

- (ii) Suppose that  $T$  is inadmissible. Then, for some  $T'$ ,  $R(\theta, T') \leq R(\theta, T)$  for all  $\theta$ , and, for some  $\theta_0$ ,  $R(\theta_0, T') < R(\theta_0, T)$ . This implies that for some  $\epsilon > 0$ , and some open neighborhood  $U \subset \Theta$  of  $\theta_0$ , we have

$$R(\vartheta, T') \leq R(\vartheta, T) - \epsilon, \quad \vartheta \in U. \quad \text{just in der Umgebung}$$

But then

$$\begin{aligned} \underline{r_w(T')} &= \int_U R(\vartheta, T') w(\vartheta) d\nu(\vartheta) + \int_{U^c} R(\vartheta, T') w(\vartheta) d\nu(\vartheta) \\ &\leq \int_U R(\vartheta, T) w(\vartheta) d\nu(\vartheta) - \epsilon \Pi(U) + \int_{U^c} R(\vartheta, T) w(\vartheta) d\nu(\vartheta) \\ &= \underline{r_w(T)} - \epsilon \Pi(U) < r_w(T). \end{aligned}$$

prob of being in  
that subspace.

We thus arrived at a contradiction.  $\Rightarrow T$  is Bayes!

and hence minimal risk.

□

**Lemma 11.2.2** Suppose that  $T$  is extended Bayes, and that for all  $T'$ ,  $R(\theta, T')$  is continuous in  $\theta$ . In fact assume, for all open sets  $U \subset \Theta$ ,

$$\left. \begin{array}{l} \text{extended} \\ \text{bayes for} \\ \text{subsets in Umgebung} \end{array} \right\} \frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \rightarrow 0, \quad \left. \begin{array}{l} \text{if means} \\ \text{the nominator,} \\ \text{converges faster than} \\ \text{the denominator to 0.} \end{array} \right\}$$

as  $m \rightarrow \infty$ . Here  $\Pi_m(U) := \int_U w_m(\vartheta) d\mu_m(\vartheta)$  is the probability of  $U$  under the prior  $\Pi_m$ . Then  $T$  is admissible.

**Proof.** We start out as in the proof of (ii) in the previous lemma. Suppose that  $T$  is inadmissible. Then, for some  $T'$ ,  $R(\theta, T') \leq R(\theta, T)$  for all  $\theta$ , and, for some  $\theta_0$ ,  $R(\theta_0, T') < R(\theta_0, T)$ , so that for some  $\epsilon > 0$ , and some open neighborhood  $U \subset \Theta$  of  $\theta_0$ , we have

$$R(\vartheta, T') \leq R(\vartheta, T) - \epsilon, \quad \vartheta \in U.$$

This would give that for all  $m$ ,

$$\frac{r_w(T) - r_w(T')}{\Pi_m(U)} \geq \left\{ \begin{array}{l} r_{w_m}(T') \leq r_{w_m}(T) - \epsilon \Pi_m(U). \end{array} \right. \quad \text{Integrating over the parameters}$$

Suppose for simplicity that a Bayes decision  $T_m$  for the prior  $w_m$  exists, for all  $m$ , i.e.

$$r_{w_m}(T_m) = \inf_{T'} r_{w_m}(T'), \quad m = 1, 2, \dots$$

Then, for all  $m$ ,

$$r_{w_m}(T_m) \leq r_{w_m}(T') \leq r_{w_m}(T) - \epsilon \Pi_m(U),$$

or

$$\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \geq \epsilon > 0,$$

that is, we arrived at a contradiction.

Because  $\square$   $T_m$  is Bayes!

### 11.2.1 Admissible estimators for the normal mean

Let  $X$  be  $\mathcal{N}(\theta, 1)$ -distributed,  $\theta \in \Theta := \mathbb{R}$  and  $R(\theta, T) := E_\theta(T - \theta)^2$  be the quadratic risk. We consider estimators of the form

$$T = aX + b, \quad a > 0, \quad b \in \mathbb{R}.$$

3 linear or regression  
or smth like  
that.

**Lemma**  $T$  is admissible if and only if one of the following cases hold

- (i)  $a < 1$ ,
- (ii)  $a = 1$  and  $b = 0$ .

prove uniqueness and previous proof.

**Proof.**

$(\Leftarrow)$  (i)

First, we show that  $T$  is Bayes for some prior. It turns out that this works with a normal prior, i.e., we take  $\theta \sim \mathcal{N}(c, \tau^2)$  for some  $c$  and  $\tau^2$  to be specified. In Example 10.5.3 we have seen that Bayes estimator is

$$T_{\text{Bayes}} = E(\theta|X) = \frac{\tau^2 X + c}{\tau^2 + 1}.$$

Taking

$$\frac{\tau^2}{\tau^2 + 1} = a, \quad \frac{c}{\tau^2 + 1} = b,$$

yields  $T = T_{\text{Bayes}}$ .

Next, we check (i) in Lemma 11.2.1, i.e. that  $T$  is unique. Since in view of the calculations in Subsection 10.5.2

$$\text{bayes} = E(\theta|x)$$

$$r_w(T') = E\text{var}(\theta|X) + E(T - T')^2$$

we conclude that if  $r_w(T') = r_w(T)$ , then  $T$  is Bayes and estimator, the same

$$\boxed{E(T - T')^2 = 0} \quad \square$$

Here, the expectation is with  $\theta$  integrated out, i.e., with respect to the measure  $P$  with density

$$p(x) = \int p_\theta(x) w(\theta) d\mu(\theta).$$

Now, we can write  $X = \theta + \epsilon$ , with  $\theta \sim \mathcal{N}(c, \tau^2)$ -distributed, and with  $\epsilon$  a standard normal random variable independent of  $\theta$ . So  $X \sim \mathcal{N}(c, \tau^2 + 1)$ , that is,  $P$  is the  $\mathcal{N}(c, \tau^2 + 1)$ -distribution. Now,  $E(T - T')^2 = 0$  implies  $T = T'$   $P$ -a.s.. Since  $P$  dominates all  $P_\theta$ , we conclude that  $T = T'$   $P_\theta$ -a.s., for all  $\theta$ . So  $T$  is unique, and hence admissible.

( $\Leftarrow$ ) (ii)

In this case,  $T = X$ . We use Lemma 11.2.2. Because  $R(\theta, T) = 1$  for all  $\theta$ , also  $r_w(T) = 1$  for any prior. Let  $w_m$  be the density of the  $\mathcal{N}(0, m)$ -distribution. As we have seen in Example 10.5.3 and also in the previous part of the proof, the Bayes estimator is

$$T_m = \frac{m}{m+1}X. \quad \text{so } E(T) = \theta \rightarrow \text{unbiased}$$

By the bias-variance decomposition, it has risk

$$R(\theta, T_m) = \frac{m^2}{(m+1)^2} + \left( \frac{m}{m+1} - 1 \right)^2 \theta^2 = \frac{m^2}{(m+1)^2} + \frac{\theta^2}{(m+1)^2}.$$

As  $E\theta^2 = m$ , its Bayes risk is

$$r_{w_m}(T_m) = \frac{m^2}{(m+1)^2} + \frac{m}{(m+1)^2} = \frac{m}{m+1}.$$

It follows that

$$r_{w_m}(T) - r_{w_m}(T_m) = 1 - \frac{m}{m+1} = \frac{1}{m+1}.$$

So  $T$  is extended Bayes. But we need to prove the more refined property of Lemma 11.2.2. It is clear that here, we only need to consider open intervals  $U = (u, u+h)$ , with  $u$  and  $h > 0$  fixed. We have

$$\begin{aligned} \Pi_m(U) &= \Phi\left(\frac{u+h}{\sqrt{m}}\right) - \Phi\left(\frac{u}{\sqrt{m}}\right) \\ &= \frac{1}{\sqrt{m}}\phi\left(\frac{u}{\sqrt{m}}\right)h + o(1/\sqrt{m}). \end{aligned}$$

For  $m$  large,

$$\phi\left(\frac{u}{\sqrt{m}}\right) \approx \phi(0) = \frac{1}{\sqrt{2\pi}} > \frac{1}{4} \text{ (say),}$$

so for  $m$  sufficiently large (depending on  $u$ )

$$\phi\left(\frac{u}{\sqrt{m}}\right) \geq \frac{1}{4}.$$

Thus, for  $m$  sufficiently large (depending on  $u$  and  $h$ ), we have

$$\Pi_m(U) \geq \frac{1}{4\sqrt{m}}h.$$

We conclude that for  $m$  sufficiently large

$$\frac{r_{w_m}(T) - r_{w_m}(T_m)}{\Pi_m(U)} \leq \frac{4}{h\sqrt{m}}.$$

Have  
to exercise  
this on  
bayes Risk  
etc.  
Write it down.

As the right hand side converges to zero as  $m \rightarrow \infty$ , this shows that  $X$  is admissible.

( $\Rightarrow$ )

We now have to show that if (i) or (ii) do not hold, then  $T$  is not admissible. This means we have to consider two cases:  $a > 1$  and  $a = 1, b \neq 0$ . In the case  $a > 1$ , we have  $R(\theta, aX + b) \geq \text{var}(aX + b) > 1 = R(\theta, X)$ , so  $aX + b$  is not admissible. When  $a = 1$  and  $b \neq 0$ , it is the bias term that makes  $aX + b$  inadmissible:

$$R(\theta, X + b) = 1 + b^2 \geq 1 = R(\theta, X).$$

Always try this  
Var / Bias decomposition.

$a > 1$   
 $X \sim N(0, 1)$

□

### 11.3 Admissible estimators in exponential families \*

**Lemma 11.3.1** Let  $\theta \in \Theta = \mathbb{R}$  and  $\{P_\theta : \theta \in \Theta\}$  be an exponential family in canonical form:

$$p_\theta(x) = \exp[\theta T(x) - d(\theta)]h(x).$$

Then  $T$  is an admissible estimator of  $g(\theta) := \dot{d}(\theta)$ , under quadratic loss (i.e., under the loss  $L(\theta, a) := |a - g(\theta)|^2$ ).

**Proof.** Recall that

$$\dot{d}(\theta) = E_\theta T, \quad \ddot{d}(\theta) = \text{var}_\theta(T) = I(\theta).$$

(see Section 4.8). Now, let  $T'$  be some estimator, with expectation

$$E_\theta T' := q(\theta).$$

the bias of  $T'$  is

$$b(\theta) = q(\theta) - g(\theta),$$

or

$$q(\theta) = b(\theta) + g(\theta) = b(\theta) + \dot{d}(\theta).$$

This implies

$$\dot{q}(\theta) = \dot{b}(\theta) + I(\theta).$$

By the Cramer Rao lower bound

$$\begin{aligned} R(\theta, T') &= \text{var}_\theta(T') + b^2(\theta) \\ &\geq \frac{[\dot{q}(\theta)]^2}{I(\theta)} + b^2(\theta) = \frac{[\dot{b}(\theta) + I(\theta)]^2}{I(\theta)} + b^2(\theta). \end{aligned}$$

Suppose now that

$$R(\theta, T') \leq R(\theta, T), \forall \theta.$$

Because  $R(\theta, T) = I(\theta)$  this implies

$$\frac{[\dot{b}(\theta) + I(\theta)]^2}{I(\theta)} + b^2(\theta) \leq I(\theta),$$

or

$$I(\theta)\{b^2(\theta) + 2\dot{b}(\theta)\} \leq -[\dot{b}(\theta)]^2 \leq 0.$$

This in turn implies

$$b^2(\theta) + 2\dot{b}(\theta) \leq 0,$$

and hence,  $b(\theta)$  is decreasing and when  $b(\theta) \neq 0$ ,

$$\frac{\dot{b}(\theta)}{b^2(\theta)} \leq -\frac{1}{2},$$

so

$$\frac{d}{d\theta}\left(\frac{1}{b(\theta)}\right) - \frac{1}{2} \geq 0,$$

or

$$\frac{d}{d\theta}\left(\frac{1}{b(\theta)} - \frac{\theta}{2}\right) \geq 0.$$

In other words,  $1/b(\theta) - \theta/2$  is an increasing function.

We will now show that this gives a contradiction, implying that  $b(\theta) = 0$  for all  $\theta$ .

Suppose instead  $b(\theta_0) < 0$  for some  $\theta_0$ . Then also  $b(\vartheta) < 0$  for all  $\vartheta > \theta_0$  since  $b(\cdot)$  is decreasing. It follows that

$$\frac{1}{b(\vartheta)} \geq \frac{1}{b(\theta_0)} + \frac{\vartheta - \theta_0}{2} \rightarrow \infty, \quad \vartheta \rightarrow \infty$$

i.e.,

$$b(\vartheta) \rightarrow 0, \quad \vartheta \rightarrow \infty.$$

This is not possible, as  $b(\theta)$  is a decreasing function.

Similarly, if  $b(\theta_0) > 0$ , take  $\theta_0 \geq \vartheta \rightarrow -\infty$ , to find again

$$b(\vartheta) \rightarrow 0, \quad \vartheta \rightarrow -\infty,$$

which is not possible.

We conclude that  $b(\theta) = 0$  for all  $\theta$ , i.e.,  $T'$  is an unbiased estimator of  $\theta$ . By the Cramer Rao lower bound, we now conclude

$$R(\theta, T') = \text{var}_\theta(T') \geq R(\theta, T) = I(\theta).$$

□

**Example 11.3.1 Admissibility of the sample average for estimating the mean of a normal distribution**

Let  $X$  be  $\mathcal{N}(\theta, 1)$ -distributed, with  $\theta \in \mathbb{R}$  unknown. Then  $X$  is an admissible estimator of  $\theta$ .

**Example 11.3.2 Inadmissibility in the normal distribution with  $\sigma^2$  unknown**

Let  $X$  be  $\mathcal{N}(0, \sigma^2)$ , with  $\sigma^2 \in (0, \infty)$  unknown. Its density is

$$\begin{aligned} p_\theta(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{x^2}{2\sigma^2}\right] \\ &= \exp[\theta T(x) - d(\theta)] h(x), \end{aligned}$$

with

$$T(x) = -x^2/2, \quad \theta = 1/\sigma^2,$$

$$\begin{aligned} d(\theta) &= (\log \sigma^2)/2 = -(\log \theta)/2, \\ \dot{d}(\theta) &= -\frac{1}{2\theta} = -\frac{\sigma^2}{2}, \\ \ddot{d}(\theta) &= \frac{1}{2\theta^2} = \frac{\sigma^4}{2}. \end{aligned}$$

Observe that  $\theta \in \Theta = (0, \infty)$ , which is not the whole real line. So Lemma 11.3.1 cannot be applied. We will now show that  $T$  is not admissible. Define for all  $a > 0$ ,

$$T_a := -aX^2.$$

so that  $T = T_{1/2}$ . We have

$$\begin{aligned} R(\theta, T_a) &= \text{var}_\theta(T_a) + \text{bias}_\theta^2(T_a) \\ &= 2a^2\sigma^4 + [a - 1/2]^2\sigma^4. \end{aligned}$$

Thus,  $R(\theta, T_a)$  is minimized at  $a = 1/6$  giving

$$R(\theta, T_{1/6}) = \sigma^4/6 < \sigma^4/2 = R(\theta, T).$$

## 11.4 Inadmissibility in higher-dimensional settings \*

Let (for  $i = 1, \dots, p$ )  $X_i \sim \mathcal{N}(\theta_i, 1)$  and let  $X_1, \dots, X_p$  be independent. The vector  $\theta := (\theta_1, \dots, \theta_p) \in \mathbb{R}^p$  is unknown. For an estimator  $T = (T_1, \dots, T_p) \in \mathbb{R}^p$ , we define the risk

$$R(\theta, T) := \sum_{i=1}^p E_\theta(T_i - \theta_i)^2.$$

Note that  $R(\theta, X) = p$  where  $X := (X_1, \dots, X_p)$ . One can moreover show (in a similar way as for the case  $p = 1$ ) that  $X$  is minimax, extended Bayes, UMRE and that it reaches the Cramer-Rao lower bound. But for  $p > 2$ ,  $X$  is inadmissible. This follows from the lemma below, which shows that  $X$  can be improved by Stein's estimator. We use the notation  $\|X\|^2 := \sum_{i=1}^p X_i^2$ .

**Definition 11.4.1** Let  $p > 2$  and let  $0 < b < 2(p-2)$  be some constant. Stein's estimator is

$$T^* := \left(1 - \frac{b}{\|X\|^2}\right)X.$$

**Lemma 11.4.1** We have

$$R(\theta, T^*) = p - \left[2b(p-2) - b^2\right]E_\theta \frac{1}{\|X\|^2}.$$

**Proof.** We first calculate

$$\begin{aligned} E_\theta(T_i^* - \theta_i)^2 &= E_\theta \left[ \left(1 - \frac{b}{\|X\|^2}\right) X_i - \theta_i \right]^2 \\ &= E_\theta \left[ (X_i - \theta_i) - \frac{b}{\|X\|^2} X_i \right]^2 \\ &= E_\theta \left[ (X_i - \theta_i)^2 + b^2 \frac{X_i^2}{\|X\|^4} - 2b \frac{X_i(X_i - \theta_i)}{\|X\|^2} \right] \\ &= 1 + b^2 E_\theta \frac{X_i^2}{\|X\|^4} - 2b E_\theta \frac{X_i(X_i - \theta_i)}{\|X\|^2}. \end{aligned}$$

Consider now the expectation in the last term, with  $i = 1$  (say):

$$\begin{aligned} E_\theta \frac{X_1(X_1 - \theta_1)}{\|X\|^2} &= \int \frac{x_1(x_1 - \theta_1)}{\|x\|^2} \prod_{i=1}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= \int \frac{x_1(x_1 - \theta_1)}{\|x\|^2} \phi(x_1 - \theta_1) dx_1 \prod_{i=2}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= - \int \frac{x_1}{\|x\|^2} d\phi(x_1 - \theta_1) \prod_{i=2}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= \int \phi(x_1 - \theta_1) d\left(\frac{x_1}{\|x\|^2}\right) \prod_{i=2}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= \int \phi(x_1 - \theta_1) \left( \frac{1}{\|x\|^2} - 2 \frac{x_1^2}{\|x\|^4} \right) dx_1 \prod_{i=2}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= \int \left( \frac{1}{\|x\|^2} - 2 \frac{x_1^2}{\|x\|^4} \right) \prod_{i=1}^p \left\{ \phi(x_i - \theta_i) dx_i \right\} \\ &= E_\theta \left[ \frac{1}{\|X\|^2} - 2 \frac{X_1^2}{\|X\|^4} \right]. \end{aligned}$$

The same calculation can be done for all other  $i$ . Inserting the result in our formula for  $E_\theta(T_i^* - \theta_i)^2$  gives

$$\begin{aligned} E_\theta(T_i^* - \theta_i)^2 &= 1 + b^2 E_\theta \frac{X_i^2}{\|X\|^4} - 2b E_\theta \left[ \frac{1}{\|X\|^2} - 2 \frac{X_i^2}{\|X\|^4} \right] \\ &= 1 + (b^2 + 4b) E_\theta \frac{X_i^2}{\|X\|^4} - 2b E_\theta \frac{1}{\|X\|^2}. \end{aligned}$$

It follows that

$$\begin{aligned} R(\theta, T^*) &= p + (b^2 + 4b)E_\theta \frac{\sum_{i=1}^p X_i^2}{\|X\|^4} - 2bpE_\theta \frac{1}{\|X\|^2} \\ &= p - \left[ 2b(p-2) - b^2 \right] E_\theta \frac{1}{\|X\|^2}. \end{aligned}$$

□

We thus have the surprising fact that Stein's estimator of  $\theta_i$  uses also the observations  $X_j$  with  $j \neq i$ , even though these observations are independent of  $X_i$  and have a distribution which does not depend on  $\theta_i$ .

Note that  $[2b(p-2) - b^2]$  is maximized for  $b = p-2$ . So the value  $b = p-2$  gives the maximal improvement over  $X$ . Stein's estimator is then

$$T^* = \left[ 1 - \frac{p-2}{\|X\|^2} \right] X.$$

**Remark** It turns out that Stein's estimator is also inadmissible!

**Remark** Let  $g(\theta) := E_\theta 1/\|X\|^2$ . One can show that  $g(0) = 1/(p-2)$ . Moreover,  $g(\theta) \downarrow 0$  as  $\|\theta\| \uparrow \infty$ , so  $R(\theta, T^*) \approx R(\theta, X)$  for  $\|\theta\|$  large.

**Remark** Let us take an empirical Bayesian point of view. Suppose  $\theta_1, \dots, \theta_p$  are i.i.d. with the  $\mathcal{N}(0, \tau^2)$ -distribution. If  $\tau^2$  is known, Bayes estimator is

$$T_{i, \text{Bayes}} = \frac{\tau^2}{1 + \tau^2} X_i, \quad i = 1, \dots, p$$

(see Example 5.2.1). Given  $\theta_i$ ,  $X_i \sim \mathcal{N}(\theta_i, 1)$  ( $i = 1, \dots, p$ ). So unconditionally,  $X_i \sim \mathcal{N}(0, 1 + \tau^2)$  ( $i = 1, \dots, p$ ). Thus, unconditionally,  $X_1, \dots, X_p$  are identically distributed, each having the  $\mathcal{N}(0, \sigma^2)$ -distribution with  $\sigma^2 = 1 + \tau^2$ . As estimator of the variance  $\sigma^2$  we may use the sample version  $\hat{\sigma}^2 := \sum_{i=1}^p X_i^2/p = \|X\|^2/p$  (we need not center with the sample average as the unconditional mean of the  $X_i$  is known to be zero). That is, we estimate  $\tau^2$  by

$$\hat{\tau}^2 := \hat{\sigma}^2 - 1 = \|X\|^2/p - 1.$$

This leads to the empirical Bayes estimator

$$T_{i, \text{emp. Bayes}} := \frac{\hat{\tau}^2}{1 + \hat{\tau}^2} X = \left[ 1 - \frac{p}{\|X\|^2} \right] X.$$

This shows that when  $p > 4$ , then Stein's estimator with  $b = p$  is an empirical Bayes estimator.

# Chapter 12

## The linear model

Consider  $n$  independent observations  $Y_1, \dots, Y_n$ . This time we do not assume that they are identically distributed. Let  $X \in \mathbb{R}^{n \times p}$  be a given matrix with (non-random) entries  $\{x_{i,j} : i = 1, \dots, n, j = 1, \dots, p\}$ . The matrix  $X$  is considered as (fixed) input and the vector  $Y = (Y_1, \dots, Y_n)^T$  as (random) output. One also calls the columns of  $X$  the co-variables. The matrix  $X$  is the *design matrix*. We assume it to be non-random, that is, we consider the case of fixed design.

### 12.1 Definition of the least squares estimator

We aim at predicting  $Y$  given  $X$  and decide to do this by linear approximation: we look for the best linear approximation of  $Y_i$  given  $x_{i,1}, \dots, x_{i,p}$ . We measure the fit using the residual sum of squares. This means that we minimize

$$\sum_{i=1}^n \left( Y_i - a - \sum_{j=1}^p x_{i,j} b_j \right)^2.$$

over  $a \in \mathbb{R}$  and  $b = (b_1, \dots, b_p)^T \in \mathbb{R}^p$ .

To simplify the expressions, we rename the quantities involved as follows. Define for all  $i$ ,  $x_{i,p+1} := 1$  and define  $b_{p+1} := a$ . Then for all  $i$ ;  $a + \sum_{j=1}^p x_{i,j} b_j = \sum_{j=1}^{p+1} x_{i,j} b_j$ . In other words, if we put in the matrix  $X$  a column containing only 1's then we may omit the constant  $a$ . Thus, putting the column of only 1's in front and replacing  $p+1$  by  $p$ , we let

$$X := \begin{pmatrix} 1 & x_{1,2} & \cdots & x_{1,p} \\ 1 & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,2} & \cdots & x_{n,p} \end{pmatrix}$$

Then we minimize

$$\sum_{i=1}^n \left( Y_i - \sum_{j=1}^{p+1} x_{i,j} b_j \right)^2.$$

over  $b = (b_1, \dots, b_p)^T \in \mathbb{R}^p$ .

Let us denote the Euclidean norm of a vector  $v \in \mathbb{R}^n$  by<sup>1</sup>

$$\|v\|_2 := \sqrt{\sum_{i=1}^n v_i^2}.$$

Write

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}.$$

Then

$$\sum_{i=1}^n \left( Y_i - \sum_{j=1}^{p+1} x_{i,j} b_j \right)^2 = \|Y - Xb\|_2^2.$$

**Definition 12.1.1** Suppose  $X$  has rank  $p$ . One calls

$$\hat{\beta} := \arg \min_{b \in \mathbb{R}^p} \|Y - Xb\|_2^2$$

the least squares estimator.

We say that the least squares  $\hat{\beta}$  is obtained by (linear) regression of  $Y$  on  $X$ .

The distance between  $Y$  and the space  $\{Xb : b \in \mathbb{R}^p\}$  spanned by the columns of  $X$  is minimized by projecting  $Y$  on this space. In fact, one has

$$\frac{1}{2} \frac{\partial}{\partial b} \|Y - Xb\|_2^2 = -X^T(Y - Xb).$$

It follows that  $\hat{\beta}$  is a solution of the so-called normal equations

$$X^T(Y - X\hat{\beta}) = 0$$

or

$$X^T Y = X^T X \hat{\beta}.$$

If  $X$  has rank  $p$ , the matrix  $X^T X$  has an inverse  $(X^T X)^{-1}$  and we get

full rank and  
thus invertible

$$\hat{\beta} = (X^T X)^{-1} X^T Y.$$

The projection of  $Y$  on  $\{Xb : b \in \mathbb{R}^p\}$  is

$$\underbrace{X(X^T X)^{-1} X^T Y}_{\text{projection}} \quad \left\{ \begin{array}{l} \text{best estimator} \\ \text{is the projection} \\ \text{according to the} \\ \text{LSSE} \end{array} \right.$$

Recall that a projection is a linear map of the form  $PP^T$  such that  $P^T P = I$ . We can write  $X(X^T X)^{-1} X^T := PP^T$ .<sup>2</sup>

\$\left\{ \begin{array}{l} \text{so "in"} \\ \text{projection} \\ p^{-1} = p^T \end{array} \right.\$

<sup>1</sup>We sometimes omit the subscript “2”

<sup>2</sup>Write the singular value decomposition of  $X$  as  $X = P\phi Q^T$ , where  $\phi = \text{diag}(\phi_1, \dots, \phi_p)$  contains the singular values and where  $P^T P = I$  and  $Q^T Q = I$ .

**Example with  $p = 1$**

For  $p = 1$

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}.$$

Then

$$X^T X = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix},$$

for  $\tilde{\beta}^1 = \frac{1}{n} \cdot \frac{2x^2}{\sum_{i=1}^n x_i^2} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$

$$(X^T X)^{-1} = \left( n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 \right)^{-1} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix}$$

$$= \left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^{-1} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}.$$

Moreover

$$X^T Y = \begin{pmatrix} n \bar{Y} \\ \sum_{i=1}^n x_i Y_i \end{pmatrix}.$$

We now let (changing notation:  $\hat{\alpha} := \hat{\beta}_1$ ,  $\hat{\beta} := \hat{\beta}_2$ )

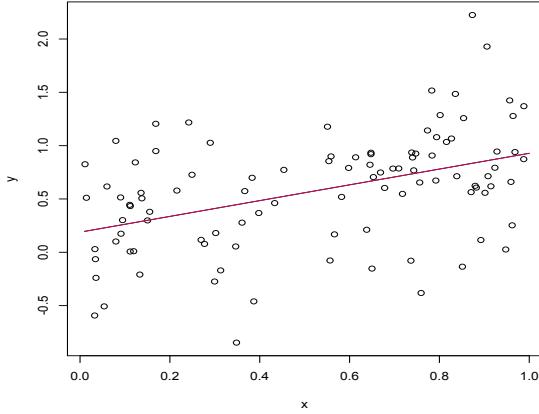
$$\begin{aligned} \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} &= (X^T X)^{-1} X^T Y \\ &= \left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^{-1} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \begin{pmatrix} n \bar{Y} \\ \sum_{i=1}^n x_i Y_i \end{pmatrix} \\ &= \left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^{-1} \begin{pmatrix} \sum_{i=1}^n x_i^2 \bar{Y} - \bar{x} \sum_{i=1}^n x_i Y_i \\ -n \bar{x} \bar{Y} + \sum_{i=1}^n x_i Y_i \end{pmatrix} \\ &= \left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^{-1} \begin{pmatrix} \sum_{i=1}^n (x_i - \bar{x})^2 \bar{Y} - \bar{x} (\sum_{i=1}^n x_i Y_i - n \bar{x} \bar{Y}) \\ \sum_{i=1}^n x_i Y_i - n \bar{x} \bar{Y} \end{pmatrix}. \end{aligned}$$

Here we used that  $\sum_{i=1}^n x_i^2 = \sum_{i=1}^n (x_i - \bar{x})^2 + n \bar{x}^2$ . We can moreover write

$$\sum_{i=1}^n x_i Y_i - n \bar{x} \bar{Y} = \sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y}).$$

Thus

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} \bar{Y} - \hat{\beta} \bar{x} \\ \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{pmatrix}.$$



Simulated data with  $Y = .3 + .6 \times x + \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, \frac{1}{4})$ ,  $\hat{\alpha} = .19$ ,  $\hat{\beta} = .740$

## 12.2 Intermezzo: the $\chi^2$ distribution

Let  $Z_1, \dots, Z_p$  be i.i.d.  $\mathcal{N}(0, 1)$ -distributed. Define the  $p$ -vector

$$Z := \begin{pmatrix} Z_1 \\ \vdots \\ Z_p \end{pmatrix}.$$

Then  $Z$  is  $\mathcal{N}(0, I)$ -distributed, with  $I$  the  $p \times p$  identity matrix. The  $\chi^2$ -distribution with  $p$  degrees of freedom is defined as the distribution of

$$\|Z\|_2^2 := \sum_{j=1}^p Z_j^2.$$

Notation:  $\|Z\|_2^2 \sim \chi_p^2$ .

For a symmetric positive definite matrix  $\Sigma$ , one can define the square root  $\Sigma^{1/2}$  as a symmetric positive definite matrix satisfying

$$\Sigma^{1/2} \Sigma^{1/2} = \Sigma. \quad \text{Similar to Cholesky but not entirely equal}$$

Its inverse is denoted by  $\Sigma^{-1/2}$  (which is the square root of  $\Sigma^{-1}$ ). If  $Z \in \mathbb{R}^p$  is  $\mathcal{N}(0, \Sigma)$ -distributed, the transformed vector

$$\tilde{Z} := \Sigma^{-1/2} Z$$

is  $\mathcal{N}(0, I)$ -distributed. It follows that

$$\underline{Z^T \Sigma^{-1} Z} = \tilde{Z}^T \tilde{Z} = \|\tilde{Z}\|_2^2 \sim \chi_p^2.$$

## 12.3 Distribution of the least squares estimator

**Definition 12.3.1** For  $f = EY$  we let  $\hat{\beta}^* := (X^T X)^{-1} X^T f$  and we call  $X\hat{\beta}^*$  the best linear approximation of  $f$ . (No covariance between the error terms)

**Lemma 12.3.1** Suppose  $E\epsilon\epsilon^T = \sigma^2 I$  where  $\epsilon := Y - f$ . Then

i)  $E\hat{\beta} = \beta^*$ ,  $\text{Cov}(\hat{\beta}) = \sigma^2(X^T X)^{-1}$ ,

ii)  $E\|X(\hat{\beta} - \beta^*)\|_2^2 = \sigma^2 p$ ,

iii)  $E\|X\hat{\beta} - f\|_2^2 = \underbrace{\sigma^2 p}_{\text{estimation error}} + \underbrace{\|X\beta^* - f\|_2^2}_{\text{misspecification error}}$ .

it's always a trade off.  
the longer the model the  
better you can approximate  
but the worst you can  
estimate.



**Proof.**

i) By straightforward computation  $\hat{\beta} = (X^T X)^{-1} X^T f$

$$\hat{\beta} - \beta^* = \underbrace{(X^T X)^{-1} X^T}_{:=A} \epsilon.$$

We therefore have

$$E(\hat{\beta} - \beta^*) = AE\epsilon = 0,$$

and the covariance matrix of  $\hat{\beta}$  is

$$\begin{aligned} \text{Cov}(\hat{\beta}) &= \text{Cov}(A\epsilon) \\ &= A \underbrace{\text{Cov}(\epsilon)}_{=\sigma^2 I} A^T \\ &= \sigma^2 A A^T = \sigma^2 (X^T X)^{-1}. \end{aligned}$$

ii) Define the projection  $PP^T := X(X^T X)^{-1} X^T$ . Then

$$\|X(\hat{\beta} - \beta^*)\|_2^2 = \|PP^T \epsilon\|_2^2 =: \sum_{j=1}^p V_j^2,$$

where  $V := P^T \epsilon$ ,

$$\begin{aligned} EV &= P^T E\epsilon = 0, \\ \text{and} \quad \text{Cov}(V) &= (P^T) \text{Cov}(\epsilon) P = \sigma^2 I. \end{aligned}$$

It follows that

$$E \sum_{j=1}^p V_j^2 = \sum_{j=1}^p EV_j^2 = \sigma^2 p.$$

iii) It holds by Pythagoras' rule for all  $b$

$$\|Xb - f\|_2^2 = \|X(b - \beta^*)\|_2^2 + \|X\beta^* - f\|_2^2$$

since  $X\beta^* - f$  is orthogonal to  $X$  ( $X^T(X\beta^* - f) = 0$ )

And hence to  $X(\hat{\beta} - \beta^*)$  ( $X^T X(\hat{\beta} - \beta^*) = 0$ )

□

**Lemma 12.3.2** Suppose  $\epsilon := Y - f \sim \mathcal{N}(0, \sigma^2 I)$ . Then we have

i)  $\hat{\beta} - \beta^* \sim \mathcal{N}(0, \sigma^2(X^T X)^{-1})$ ,

ii)  $\frac{\|X(\hat{\beta} - \beta^*)\|_2^2}{\sigma^2} \sim \chi_p^2$  where  $\chi_p^2$  is  $\chi^2$ -distributed with  $p$  degrees of freedom (see Section 12.2 for a definition).

Recall

$$Xb = Y$$

so that  $Y - E(Y)$

on same axis

$$\xrightarrow{\sim} X\beta^* - f$$

$$\xrightarrow{\sim} X(b - \beta^*)$$

and  $f = E(Y)$  on same axis as  $Y$ .

**Proof.**

i) Since  $\hat{\beta}$  is a linear function of the multivariate normal  $\epsilon$ , the least squares estimator  $\hat{\beta}$  is also multivariate normal.

ii) Define the projection  $PP^T := X(X^T X)^{-1} X^T$ . Then

$$\|X(\hat{\beta} - \beta^*)\|_2^2 = \|PP^T \epsilon\|_2^2 := \sum_{j=1}^p V_j^2.$$

Now  $V := P^T \epsilon$  has i.i.d.  $\mathcal{N}(0, \sigma^2)$  entries.  $\square$

**Remark** The misspecification error  $\|X\beta^* - f\|_2^2$  comes from the possible misspecification of the linear model. That is,  $f$  need not be a linear combination of the columns of  $X$ . One sometimes also calls  $\|X\beta^* - f\|_2^2$  the approximation error. The estimation error is here the variance term  $\sigma^2 p$ .

**Remark** More generally, many estimators are approximately normally distributed (for example the sample median) and many test statistics have approximately a  $\chi^2$  null-distribution (for example the  $\chi^2$  goodness-of-fit statistic). This phenomenon occurs because many models can in a certain sense be approximated by the linear model and many minus log-likelihoods resemble the least squares loss function (see Chapter 14). Understanding the linear model is a first step towards understanding a wide range of more complicated models.

**Corollary 12.3.1** Suppose the linear model is well-specified: for some  $\beta \in \mathbb{R}^p$

$$EY = X\beta.$$

Assume  $\epsilon := Y - EY \sim \mathcal{N}(0, \sigma^2)$ , where  $\sigma^2 := \sigma_0^2$  is known. Then a test for  $H_0: \beta = \beta_0$ ,  $x_{\beta_0} = \epsilon(y)$  as again  $y = x_{\beta_0} + \epsilon$  is:

reject  $H_0$  when  $\|X(\hat{\beta} - \beta_0)\|_2^2 / \sigma_0^2 > G_p^{-1}(1 - \alpha)$ , where  $G_p$  is the distribution function of a  $\chi_p^2$ -distributed random variable.

**Remark** When  $\sigma^2$  is unknown one may estimate it using the estimator

$$\hat{\sigma}^2 = \frac{\|\hat{\epsilon}\|_2^2}{n-p},$$

and taking the expectation yields the result.

where  $\hat{\epsilon} := Y - X\hat{\beta}$  is the vector of residuals. Under the assumptions of the previous corollary (but now with possibly unknown  $\sigma^2$ ) the test statistic  $\|X(\hat{\beta} - \beta_0)\|_2^2 / (\hat{\sigma}^2 p)$  has a so-called  $F$ -distribution with  $p$  and  $n - p$  degrees of freedom.

Recall  $F$ -distribution:  $\frac{X_m}{Y_n}$  distribution where  $X$  and  $Y$  distributed  $\chi^2$  variables

## 12.4 Intermezzo: some matrix algebra

Let  $z \in \mathbb{R}^p$  be a vector and  $B \in \mathbb{R}^{q \times p}$  be a  $q \times p$ -matrix, ( $p \geq q$ ) with rank  $q$ . Moreover, let  $V \in \mathbb{R}^{p \times p}$  be a positive definite  $p \times p$ -matrix.

**Lemma 12.4.1** We have

$$\max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T a\} = z^T z - z^T B^T (BB^T)^{-1} Bz.$$

Max w.r.t.  
nonzero  
vectors

**Proof.** We use Lagrange multipliers  $\lambda \in \mathbb{R}^p$ . We have

$$\frac{\partial}{\partial a} \{2a^T z - a^T a + \underbrace{2a^T B^T \lambda}_{\text{Notice that } Ba=0 \text{ makes anything with } a^T a \text{ zero!}}\} = z - a + B^T \lambda.$$

Hence for

$$a_* := \arg \max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T a\},$$

we have

$$z - a_* + B^T \lambda = 0,$$

or

$$a_* = z + B^T \lambda.$$

The restriction  $Ba_* = 0$  gives

$$Bz + BB^T \lambda = 0.$$

So

$$\lambda = -(BB^T)^{-1}Bz.$$

Inserting this in the solution  $a^*$  gives

$$a_* = z - B^T(BB^T)^{-1}Bz.$$

Now

$$\begin{aligned} a_*^T a_* &= \left( z^T - z^T B^T (BB^T)^{-1} B \right) \left( z - B^T (BB^T)^{-1} Bz \right) \\ &= z^T z - z^T B^T (BB^T)^{-1} Bz. \end{aligned}$$

So

$$2a_*^T z - a_*^T a_* = z^T z - z^T B^T (BB^T)^{-1} Bz.$$

□

**Lemma 12.4.2** *We have*

$$\max_{a \in \mathbb{R}^p: Ba=0} \{2a^T z - a^T V a\} = z^T V^{-1} z - z^T V^{-1} B^T \left( BV^{-1} B^T \right)^{-1} BV^{-1} z.$$

**Proof.** Make the transformation  $b := V^{1/2}a$ , and  $y := V^{-1/2}z$ , and  $C = BV^{-1/2}$ . Then

$$\begin{aligned} \max_{a: Ba=0} \{2a^T z - a^T V a\} &= \max_{b: Cb=0} \{2b^T y - b^T b\} \\ &= y^T y - y^T C^T (CC^T)^{-1} C y \\ &= z^T V^{-1} z - z^T V^{-1} B^T \left( BV^{-1} B^T \right)^{-1} BV^{-1} z. \end{aligned}$$

□

**Corollary 12.4.1** *Let  $L(a) := 2a^T z - a^T V a$ . The difference between the unrestricted maximum and the restricted maximum of  $L(a)$  is*

$$\max_a L(a) - \max_{a: Ba=0} L(a) = z^T V^{-1} B^T \left( BV^{-1} B^T \right)^{-1} BV^{-1} z.$$

## 12.5 Testing a linear hypothesis

In this section we assume the model

$$Y = X\beta + \epsilon,$$

where  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ . We want to test the hypothesis

$$H_0 : B\beta = 0 ,$$

where  $B \in \mathbb{R}^{q \times p}$  is a given  $q \times p$  matrix.

Let

$$\hat{\beta}_0 := \arg \min_{b \in \mathbb{R}^p: Bb=0} \|Y - Xb\|_2^2$$

be the least squares estimator under the restriction  $B\hat{\beta} = 0$ .

**Lemma 12.5.1** Under  $H_0$

$$\frac{\|Y - X\hat{\beta}_0\|_2^2 - \|Y - X\hat{\beta}\|_2^2}{\sigma^2}$$

has a  $\chi_q^2$ -distribution.

**Proof.** Since  $\|Y - Xb\|^2 = \epsilon^T Y - 2\epsilon^T X(b - \beta) + (b - \beta)^T X^T X(b - \beta)$  we have under  $H_0$  (write ' $\tilde{b}$ ' :=  $b - \beta$ )

$$\hat{\beta}_0 - \beta = \arg \max_{\tilde{b} \in \mathbb{R}^p: B\tilde{b}=0} \left\{ 2\epsilon^T X\tilde{b} - \tilde{b}^T X^T X\tilde{b} \right\}.$$

Therefore, invoking Corollary 12.4.1

$$\begin{aligned} & \underbrace{\|Y - X\hat{\beta}_0\|_2^2}_{\text{unrestricted}} - \underbrace{\|Y - X\hat{\beta}\|_2^2}_{\text{restricted}} \\ &= \underbrace{\epsilon^T X(X^T X)^{-1} B^T}_{:= Z^T} \left( B(X^T X)^{-1} B^T \right)^{-1} \underbrace{B(X^T X)^{-1} X^T \epsilon}_{:= Z} \\ &= Z^T \left( B(X^T X)^{-1} B^T \right)^{-1} Z. \end{aligned}$$

The  $q$ -vector

$$Z := B(X^T X)^{-1} X^T \epsilon$$

has a multivariate normal distribution with mean zero and covariance matrix

$$\sigma^2 \left( B(X^T X)^{-1} X^T \right) \left( B(X^T X)^{-1} X^T \right)^T = \sigma^2 B(X^T X)^{-1} B^T.$$

It follows that under  $H_0$

$$\frac{\|Y - X\hat{\beta}_0\|_2^2 - \|Y - X\hat{\beta}\|_2^2}{\sigma^2}$$

has a  $\chi_q^2$ -distribution. □

# Chapter 13

## Asymptotic theory



In this and subsequent chapters, the observations  $X_1, \dots, X_n$  are considered as the first  $n$  of an infinite sequence of i.i.d. random variables  $X_1, \dots, X_n, \dots$  with values in  $\mathcal{X}$  and with distribution  $P$ . We say that the  $X_i$  are i.i.d. copies, of some random variable  $X \in \mathcal{X}$  with distribution  $P$ . We let  $\underline{P} = P \times P \times \dots$  be the distribution of the whole sequence  $\{X_i\}_{i=1}^\infty$ .

The model class for  $P$  is

$$\mathcal{P} := \{P_\theta : \theta \in \Theta\}.$$

When  $P = P_\theta$ , we write  $\underline{P} = \underline{P}_\theta = P_\theta \times P_\theta \times \dots$ . The parameter of interest is

$$\gamma := g(\theta) \in \mathbb{R}^p,$$

where  $g : \Theta \rightarrow \mathbb{R}^p$  is a given function. We let

$$\Gamma := \{g(\theta) : \theta \in \Theta\}$$

be the parameter space for  $\gamma$ .

An estimator

$$T_n(X_1, \dots, X_n)$$

based on the data  $X_1, \dots, X_n$ , is some function  $T_n(\cdot)$  evaluated at the data  $X_1, \dots, X_n$ . We often write shorthand

$$\underline{T}_n = T_n(X_1, \dots, X_n).$$

We assume the estimator  $T_n$  is defined for all  $n$ , i.e., we actually consider a sequence of estimators  $\{T_n\}_{n=1}^\infty$ . We are interested in estimators  $\underline{T}_n \in \Gamma$  of  $\gamma$ .

**Remark** Under the i.i.d. assumption, it is natural to assume that each  $T_n$  is a symmetric function of the data, that is

$$\underline{T}_n(X_1, \dots, X_n) = T_n(X_{\pi_1}, \dots, X_{\pi_n})$$

for all permutations  $\pi$  of  $\{1, \dots, n\}$ . In that case, one can write  $T_n$  in the form  $T_n = Q(\hat{P}_n)$ , where  $\hat{P}_n$  is the empirical distribution (see also Subsection 2.4.1).

Not depend  
on the order →  
just depends on the probability of each observation

### 13.1 Types of convergence

**Definition 13.1.1** Let  $\{Z_n\}_{n=1}^{\infty}$  and  $Z$  be  $\mathbb{R}^p$ -valued random variables defined on the same probability space<sup>1</sup>. We say that  $Z_n$  converges in probability to  $Z$  if for all  $\epsilon > 0$ ,

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

Notation:  $Z_n \xrightarrow{\mathbb{P}} Z$ .

**Remark** Chebyshev's inequality can be a tool to prove convergence in probability. It says that for all increasing functions  $\psi : [0, \infty) \rightarrow [0, \infty)$ , one has

$$\mathbb{P}(\|Z_n - Z\| \geq \epsilon) \leq \frac{\mathbb{E}\psi(\|Z_n - Z\|)}{\psi(\epsilon)}. \quad \left. \begin{array}{l} \text{Normalized difference} \\ \text{of expected value} \\ \text{and variable} \end{array} \right\}$$

It follows that proving that the right hand side  $\rightarrow 0$  as  $n \rightarrow \infty$  proves convergence in probability

**Definition 13.1.2** Let  $\{Z_n\}_{n=1}^{\infty}$  and  $Z$  be  $\mathbb{R}^p$ -valued random variables. We say that  $Z_n$  converges in distribution to  $Z$ , if for all continuous and bounded functions  $f$ ,

$$\lim_{n \rightarrow \infty} \mathbb{E}f(Z_n) = \mathbb{E}f(Z).$$

Notation:  $Z_n \xrightarrow{\mathcal{D}} Z$ .

**Remark** Convergence in probability implies convergence in distribution, but not the other way around.

#### Example 13.1.1 The central limit theorem (CLT)

Let  $X_1, X_2, \dots$  be i.i.d. real-valued random variables with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_n := \sum_{i=1}^n X_i/n$  be the average of the first  $n$ . Then by the central limit theorem (CLT),

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma^2),$$

No matter  
the underlying  
list of  
 $x_1, x_2, \dots$

that is

$$\mathbb{P}\left(\sqrt{n}(\bar{X}_n - \mu) \leq z\right) \rightarrow \Phi(z), \quad \forall z.$$

just normalizing  
by the variance

by portmanteau Theorem below

The following theorem says that for convergence in distribution, one actually can do with one-dimensional random variables. We omit the proof.

**Theorem 13.1.1 (Cramér-Wold device)** Let  $(\{Z_n\}, Z)$  be a collection of  $\mathbb{R}^p$ -valued random variables. Then

$$Z_n \xrightarrow{\mathcal{D}} Z \Leftrightarrow a^T Z_n \xrightarrow{\mathcal{D}} a^T Z \quad \forall a \in \mathbb{R}^p.$$

<sup>1</sup>Let  $(\Omega, \mathcal{A}, \mathbb{P})$  be a probability space, and  $X : \Omega \rightarrow \mathbb{R}^p$  and  $Y : \Omega \rightarrow \mathbb{R}^q$  be two measurable maps. Then  $X$  and  $Y$  are called random variables, and they are defined on the same probability space  $\Omega$ .

**Example 13.1.2 Multivariate CLT**

Let  $X_1, X_2, \dots$  be i.i.d. copies of a random variable  $X = (X^{(1)}, \dots, X^{(p)})^T$  in  $\mathbb{R}^p$ . Assume  $EX := \mu = (\mu_1, \dots, \mu_p)^T$  and  $\Sigma := \text{Cov}(X) := EXX^T - \mu\mu^T$  exist. Then for all  $a \in \mathbb{R}^p$ ,

$$Ea^T X = a^T \mu, \quad \text{var}(a^T X) = a^T \Sigma a.$$

Define

$$\bar{X}_n = (\bar{X}_n^{(1)}, \dots, \bar{X}_n^{(p)})^T.$$

By the 1-dimensional CLT, for all  $a \in \mathbb{R}^p$ ,

$$\sqrt{n}(a^T \bar{X}_n - a^T \mu) \xrightarrow{\mathcal{D}} \mathcal{N}(0, a^T \Sigma a). \quad \text{multivariate CLT}$$

The Cramér-Wold device therefore gives the  $p$ -dimensional CLT

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma). \quad \text{univariate CLT}$$

We state the Portmanteau Theorem:

**Theorem 13.1.2** Let  $(\{Z_n\}, Z)$  be a collection of  $\mathbb{R}^p$ -valued random variables. Denote the distribution of  $Z$  by  $Q$  and let  $G = Q(Z \leq \cdot)$  be its distribution function. The following statements are equivalent:

- (i)  $Z_n \xrightarrow{\mathcal{D}} Z$  (i.e.,  $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z)$   $\forall f$  bounded and continuous).
- (ii)  $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z)$   $\forall f$  bounded and Lipschitz.<sup>2</sup> means that the slopes of the function is bounded.
- (iii)  $\mathbb{E}f(Z_n) \rightarrow \mathbb{E}f(Z)$   $\forall f$  bounded and  $Q$ -a.s. continuous.
- (iv)  $\mathbb{P}(Z_n \leq z) \rightarrow G(z)$  for all  $G$ -continuity points  $z$ . almost surely

### 13.1.1 Stochastic order symbols

Let  $\{Z_n\}$  be a collection of  $\mathbb{R}^p$ -valued random variables, and let  $\{r_n\}$  be strictly positive random variables. We write

$$Z_n = \mathcal{O}_{\mathbb{P}}(1)$$

( $Z_n$  is bounded in probability) if

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(\|Z_n\| > M) = 0. \quad \text{Not exploding in limit.}$$

This corresponds in fact of rewriting the expression like  
is  $O(g(x))$   $\Leftrightarrow \lim_{n \rightarrow \infty} \frac{\|Z_n\|}{g(r_n)} \leq M \text{ a.s.}$

Given that:  $\mathbb{P}(\|Z_n\| \geq M g(x)) = 0$  as  $n \rightarrow \infty$ ,  $M \rightarrow \infty$

This is also called uniform tightness of the sequence  $\{Z_n\}$ . We write  $Z_n = \mathcal{O}_{\mathbb{P}}(r_n)$  if  $Z_n/r_n = \mathcal{O}_{\mathbb{P}}(1)$ .

If  $Z_n$  converges in probability to zero we write this as

$$Z_n = o_{\mathbb{P}}(1).$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|Z_n - 0| > r_n) = 0$$

Moreover,  $Z_n = o_{\mathbb{P}}(r_n)$  ( $Z_n$  is of small order  $r_n$  in probability) if  $Z_n/r_n = o_{\mathbb{P}}(1)$ .

<sup>2</sup>A real-valued function  $f$  on (a subset of)  $\mathbb{R}^p$  is Lipschitz if for a constant  $C_L$  and all  $(z, \tilde{z})$  in the domain of  $f$ ,  $|f(z) - f(\tilde{z})| \leq C_L \|z - \tilde{z}\|$ .

### 13.1.2 Some implications of convergence

**Lemma 13.1.1** Suppose that  $Z_n$  converges in distribution. Then  $Z_n = O_P(1)$ .

**Proof.** To simplify, take  $p = 1$  (Cramér-Wold device). Let  $Z_n \xrightarrow{D} Z$ , where  $Z$  has distribution function  $G$ . Then for every  $G$ -continuity point  $M$ ,

$$\mathbb{P}(Z_n > M) \rightarrow 1 - G(M), \quad \text{by Portmanteau above.}$$

and for every  $G$ -continuity point  $-M$ ,

$$\mathbb{P}(Z_n \leq -M) \rightarrow G(-M).$$

Since  $1 - G(M)$  as well as  $G(-M)$  converge to zero as  $M \rightarrow \infty$ , the result follows.  $\square$

**Example 13.1.3** Averages differ from their means by an order  $1/\sqrt{n}$  in probability

Let  $X_1, X_2, \dots$  be i.i.d. copies of a random variable  $X \in \mathbb{R}$  with  $EX = \mu$  and  $\text{var}(X) < \infty$ . Then by the CLT,

$$\bar{X}_n - \mu = O_P(1/\sqrt{n}) \Leftrightarrow f(\bar{X}_n - \mu) = O_P(1)$$

**Theorem 13.1.3 (Slutsky's Theorem)** Let  $(\{Z_n, A_n\}, Z)$  be a collection of  $\mathbb{R}^p$ -valued random variables, and  $a \in \mathbb{R}^p$  be a vector of constants. Assume that

$$Z_n \xrightarrow{D} Z, \quad A_n \xrightarrow{P} a. \quad \text{Then}$$

$$A_n^T Z_n \xrightarrow{D} a^T Z.$$

**Proof.** Take a bounded Lipschitz function  $f$ , say

$$|f| \leq C_B, \quad |f(z) - f(\tilde{z})| \leq C_L \|z - \tilde{z}\|. \quad \text{Lip. cond.}$$

Then  $\|f(A_n^T Z_n) - f(a^T Z)\| \rightarrow 0$  since  $A_n \xrightarrow{P} a$  and  $f$  is bounded and Lipschitz.

$$\begin{aligned} & \left| \mathbb{E}f(A_n^T Z_n) - \mathbb{E}f(a^T Z) \right| \\ & \leq \left| \mathbb{E}f(A_n^T Z_n) - \mathbb{E}f(a^T Z_n) \right| + \left| \mathbb{E}f(a^T Z_n) - \mathbb{E}f(a^T Z) \right|. \end{aligned}$$

$$\begin{aligned} & \left| \mathbb{E}f(A_n^T Z_n) - \mathbb{E}f(a^T Z_n) \right| \\ & \leq \|a - b\| \cdot \|f(A_n^T Z_n) - f(a^T Z_n)\| \\ & \leq \|a - b\| \cdot \|f\| \cdot \|A_n^T Z_n - a^T Z_n\| \\ & \leq \|a - b\| \cdot \|f\| \cdot \|A_n - a\| \cdot \|Z_n - Z\| \end{aligned}$$

Because the function  $z \mapsto f(a^T z)$  is bounded and Lipschitz (with Lipschitz constant  $\|a\|C_L$ ), we know that the second term goes to zero. As for the first term, we argue as follows. Let  $\epsilon > 0$  and  $M > 0$  be arbitrary. Define  $S_n :=$

$\{\|Z_n\| \leq M, \|A_n - a\| \leq \epsilon\}$ . Then

$$\begin{aligned} & \left| \mathbb{E}f(A_n^T Z_n) - \mathbb{E}f(a^T Z_n) \right| \\ & \leq \mathbb{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \mathbf{1}_{\{S_n\}} \\ & = \mathbb{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \mathbf{1}_{\{S_n^c\}} \\ & + \mathbb{E} \left| f(A_n^T Z_n) - f(a^T Z_n) \right| \mathbf{1}_{\{S_n^c\}} \end{aligned}$$



$$\begin{aligned} & \text{thus by Jensen's} \\ & \text{inequality} \\ & \mathbb{E}f(g(x)) \geq g(\mathbb{E}x) \end{aligned}$$

Absolute value convex so  $g(x) \geq g(\mathbb{E}x)$

Conv. in distribution

### 13.2. CONSISTENCY AND ASYMPTOTIC NORMALITY

135

Now

$$\leq C_L \epsilon M + 2\bar{C}_B \mathbb{P}(S_n^c). \quad (13.1)$$

Thus, both terms in (13.1) can be made arbitrary small by appropriately choosing  $\epsilon$  small and  $n$  and  $M$  large.  $\square$

## 13.2 Consistency and asymptotic normality

**Definition 13.2.1** A sequence of estimators  $\{T_n\}$  of  $\gamma = g(\theta)$  is called consistent if

$$T_n \xrightarrow{\text{P}_\theta} \gamma. \quad \left\{ \begin{array}{l} \text{Under the probability} \\ \text{measure of the true parameter.} \end{array} \right.$$

**Definition 13.2.2** A sequence of estimators  $\{T_n\}$  of  $\gamma = g(\theta)$  is called asymptotically normal with asymptotic covariance matrix  $V_\theta$ , if

$$\sqrt{n}(T_n - \gamma) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, V_\theta).$$

**Example 13.2.1** Consistency and asymptotic normality of the average  
Suppose  $\mathcal{P}$  is the location model

$$\mathcal{P} = \left\{ P_{\mu, F_0}(X \leq \cdot) := F_0(\cdot - \mu), \mu \in \mathbb{R}, F_0 \in \mathcal{F}_0 \right\}.$$

The parameter is then  $\theta = (\mu, F_0)$  and  $\Theta = \mathbb{R} \times \mathcal{F}_0$ . We assume for all  $F_0 \in \mathcal{F}_0$

$$\text{would be } \int x dF_0(x) = 0, \sigma_{F_0}^2 := \int x^2 dF_0(x) < \infty.$$

Let  $g(\theta) := \mu$  and  $T_n := (X_1 + \dots + X_n)/n =: \bar{X}_n$ . Then, by the law of large numbers,  $T_n$  is a consistent estimator of  $\mu$  and, by the central limit theorem,

$$\sqrt{n}(T_n - \mu) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, \sigma_{F_0}^2).$$

## 13.3 Asymptotic linearity

As we will show, for many estimators asymptotic normality is a consequence of asymptotic linearity, that is, the estimator is approximately an average, to which we can apply the CLT.

**Definition 13.3.1** The sequence of estimators  $\{T_n\}$  of  $\gamma = g(\theta) \in \mathbb{R}^p$  is called asymptotically linear if for a function  $l_\theta : \mathcal{X} \rightarrow \mathbb{R}^p$ , with  $E_\theta l_\theta(X) = 0$  and

$$E_\theta l_\theta(X) l_\theta^T(X) =: V_\theta < \infty, \rightarrow \text{Mean 0 and finite variance}$$

it holds that

$$T_n - \gamma = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbb{P}_\theta}(1/\sqrt{n}).$$

i.e. the statistic must be on average correct.

<sup>3</sup>In the one-dimensional case ( $p = 1$ ) we thus have  $V_\theta = E_\theta l_\theta^2(X)$ .

Cause taking the limit the only linear influence term remains and this fully describes  $T_n - \gamma$

affine function as  $T_n(X_1, \dots, X_n)$  so that  $\frac{1}{n} \sum_{i=1}^n l_\theta(X_i)$  corresponds adding to

**Remark.** We then call  $l_\theta$  the *influence function* of (the sequence)  $T_n$ . Roughly speaking,  $l_\theta(x)$  approximately measures the influence of an additional observation  $x$  (compare with the influence function as defined in Section 8.4).

→ **Example 13.3.1 Influence function of the sample average**

Assuming the entries of  $X$  have finite variance, the estimator  $T_n := \bar{X}_n$  is a linear and hence asymptotically linear estimator of the mean  $\mu$ , with influence function

$$l_\theta(x) = x - \mu.$$

↳ we proved before that  
in fact  $T_n - \mu$  is of order  $\frac{1}{\sqrt{n}}$

so  
that

**Example 13.3.2 Influence function of the sample variance**

Let  $X$  be real-valued, with  $E_\theta X =: \mu$ ,  $\text{var}_\theta(X) =: \sigma^2$  and  $\kappa =: E_\theta(X - \mu)^4$  (assumed to exist). The sample variance is

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

Let  $\hat{\sigma}_n^2$  be the estimator

$\left(\frac{1}{n-1} - \frac{1}{n}\right) \sum$   
One sees that

$$\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

convergence in Prob. immediate see

$$S^2 - \hat{\sigma}_n^2 = \mathcal{O}_P(1/n) = o_P(1/\sqrt{n}).$$

We rewrite  $\hat{\sigma}^2$  as

$$\begin{aligned} \hat{\sigma}_n^2 &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + (\bar{X}_n - \mu)^2 - \frac{2}{n} \sum_{i=1}^n (X_i - \mu)(\bar{X}_n - \mu) \\ &= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X}_n - \mu)^2. \end{aligned}$$

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

why  $\frac{1}{n}$ ?

solving the squares of above and inserting it

Because by the CLT,  $(\bar{X}_n - \mu) = \mathcal{O}_{P_\theta}(n^{-1/2})$ , we get

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + \mathcal{O}_{P_\theta}(1/n).$$

This is the influence

hence given the definition asymptotic linear.

So asymptotically one does not notice that  $\mu$  is estimated, and  $\hat{\sigma}_n^2$  (and also  $S^2$ ) is asymptotically linear with influence function

$$l_\theta(x) = (x - \mu)^2 - \sigma^2.$$

As  $T_n - \gamma$  and the above

The asymptotic variance is

$$V_\theta = E_\theta \left( (X - \mu)^2 - \sigma^2 \right)^2 = \kappa - \sigma^4.$$

derived

$$\begin{aligned} E_\theta (X - \mu)^4 + \sigma^4 - &\underbrace{E_\theta 2(X - \mu)^2 \sigma^2}_{- \sigma^4} \\ &- \sigma^4 \end{aligned}$$

## 13.4 The $\delta$ -technique

**Theorem 13.4.1** Let  $(\{T_n\}, Z)$  be a collection of random variables in  $\mathbb{R}^p$ ,  $c \in \mathbb{R}^p$  be a nonrandom vector, and  $\{r_n\}$  be a nonrandom sequence of positive numbers, with  $r_n \downarrow 0$ . Moreover, let  $\dot{h}: \mathbb{R}^p \rightarrow \mathbb{R}$  be differentiable at  $c$ , with derivative  $\dot{h}(c) \in \mathbb{R}^p$ . Suppose that

$$(T_n - c)/r_n \xrightarrow{\mathcal{D}} Z.$$

Then

$$\left( h(T_n) - h(c) \right)/r_n \xrightarrow{\mathcal{D}} \dot{h}(c)^T Z \quad \left\{ \begin{array}{l} \text{Note that when} \\ \text{there is convergence} \\ \text{to dist; } Z = O(\mathbf{1}). \end{array} \right.$$

and in fact

$$h(T_n) - h(c) = \dot{h}(c)^T (T_n - c) + o_P(r_n). \quad \left\{ \begin{array}{l} \text{asymptotically} \\ \text{linear} \end{array} \right.$$

Better to view  
on the other  
downloaded paper

**Proof.** By Slutsky's Theorem,

$$\dot{h}(c)^T (T_n - c)/r_n \xrightarrow{\mathcal{D}} \dot{h}(c)^T Z.$$

convergence in distribution

Since  $(T_n - c)/r_n$  converges in distribution, we know that  $\|T_n - c\|/r_n = O_P(1)$ . Hence,  $\|T_n - c\| = O_P(r_n)$ . The result follows now from

$$h(T_n) - h(c) = \dot{h}(c)^T (T_n - c) + o(\|T_n - c\|) = \dot{h}(c)^T (T_n - c) + o_P(r_n).$$

first order approximation

□

**Corollary 13.4.1** Let  $T_n$  be an asymptotically normal estimator of  $\gamma = g(\theta) \in \mathbb{R}^p$  with asymptotic covariance matrix

$$V_\theta.$$

Suppose  $h$  is differentiable at  $\gamma$ . Then  $h(T_n)$  is an asymptotically normal estimator of  $h(\gamma)$  with asymptotic variance<sup>4</sup>

$$\dot{h}(\gamma)^T V_\theta \dot{h}(\gamma).$$

$\Rightarrow$  Asymptotically normal.

$$\dot{h}(T_n - \gamma) \xrightarrow{\mathcal{D}} N(0, V_\theta)$$

it follows:

$$\frac{h(T_n) - h(\gamma)}{\sqrt{r_n}} \xrightarrow{\mathcal{D}} \dot{h}(\gamma) \cdot N(0, V_\theta)$$

$$\Rightarrow \xrightarrow{\mathcal{D}} N(0, \dot{h}(\gamma)^T V_\theta \dot{h}(\gamma))$$

If moreover  $T_n$  is an asymptotically linear estimator of  $\gamma$ , with influence function

$$l_\theta$$

then  $h(T_n)$  is an asymptotically linear estimator of  $h(\gamma)$  with influence function

$$\dot{h}(\gamma)^T l_\theta.$$

Check  
at it  
again

**Example 13.4.1** Asymptotic linear estimator of the parameter of the exponential distribution

Let  $X_1, \dots, X_n$  be a sample from the Exponential( $\theta$ ) distribution, with  $\theta > 0$ .

<sup>4</sup>For  $p = 1$  the asymptotic variance of  $h(T_n)$  is thus  $\dot{h}^2(\gamma) V_\theta$ .

always by def  
and not asymptotically!

## CHAPTER 13. ASYMPTOTIC THEORY

by def at CLT.

Then  $\bar{X}_n$  is a linear estimator of  $E_\theta X = 1/\theta := \gamma$ , with influence function  $l_\theta(x) = x - 1/\theta$ . The variance of  $\sqrt{n}(T_n - 1/\theta)$  is  $1/\theta^2 = \gamma^2$ . By Theorem 13.4.1,  $1/\bar{X}_n$  is an asymptotically linear estimator of  $\theta$ . In this case,  $h(\gamma) = 1/\gamma$ , so that  $h(\gamma) = -1/\gamma^2$ . The influence function of  $1/\bar{X}_n$  is thus

$$h(\gamma)l_\theta(x) = -\frac{1}{\gamma^2}(x - \gamma) = -\theta^2(x - 1/\theta).$$

The asymptotic variance of  $1/\bar{X}_n$  is

$$[h(\gamma)]^2 \gamma^2 = \frac{1}{\gamma^2} = \theta^2. \quad \text{from the corollary, 13.1.}$$

So

$$\sqrt{n}\left(\frac{1}{\bar{X}_n} - \theta\right) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, \theta^2).$$

solve like  
this in the  
exercise!

$$h(\bar{X}_n) = \frac{1}{\bar{X}_n}$$

### Example 13.4.2 Two-dimensional asymptotic linearity of the sample average and sample variance

Consider again Example 13.3.2. Let  $X$  be real-valued, with  $E_\theta X := \mu$ ,  $\text{var}_\theta(X) := \sigma^2$  and  $\kappa := E_\theta(X - \mu)^4$  (assumed to exist). Define moreover, for  $r = 1, 2, 3, 4$ , the  $r$ -th moment  $\mu_r := E_\theta X^r$ . We again consider the estimator

$$\hat{\sigma}_n^2 := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

We have

$$\hat{\sigma}_n^2 = h(T_n),$$

where  $T_n = (T_{n,1}, T_{n,2})^T$ , with

$$T_{n,1} = \bar{X}_n, \quad T_{n,2} = \frac{1}{n} \sum_{i=1}^n X_i^2,$$

and

$$h(t) = t_2 - t_1^2, \quad t = (t_1, t_2)^T. \quad \begin{cases} \text{def of} \\ \text{variance} \end{cases}$$

The estimator  $T_n$  has influence function

$$l_\theta(x) = \begin{pmatrix} x - \mu_1 \\ x^2 - \mu_2 \end{pmatrix}. \quad \begin{cases} \text{linear estim.} \\ \text{not even asymp. linear} \end{cases}$$

By the 2-dimensional CLT,

$$\sqrt{n}\left(T_n - \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}\right) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, V_\theta),$$

with

$$V_\theta = \begin{pmatrix} \mu_2 - \mu_1^2 & \mu_3 - \mu_1\mu_2 \\ \mu_3 - \mu_1\mu_2 & \mu_4 - \mu_2^2 \end{pmatrix}. \quad \begin{cases} \mu_i \text{ are the} \\ \text{central moments} \end{cases}$$

It holds that

$$h\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}\right) = \begin{pmatrix} 2\mu_1 \\ 1 \end{pmatrix}, \quad \begin{cases} \text{de-w.r.t. this} \\ \text{or w.r.t. this} \end{cases}$$

so that  $\hat{\sigma}_n^2$  has influence function

$$\begin{aligned} h^T \left( \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \right) l_\theta(x) &= \begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix}^T \begin{pmatrix} x - \mu_1 \\ x^2 - \mu_2 \end{pmatrix} \\ &= (x - \mu)^2 - \sigma^2 \end{aligned}$$

(invoking  $\mu_1 = \mu$ ). After some calculations, one finds moreover that

$$\begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix}^T V_\theta \begin{pmatrix} -2\mu_1 \\ 1 \end{pmatrix} = \kappa - \sigma^4,$$

i.e., the  $\delta$ -method gives the same result as the ad hoc method in Example 13.3.2, as it of course should.



# Chapter 14

## M-estimators

Recall the maximum likelihood estimator as defined in Section 2.4.3. In this chapter we introduce a general class of estimators of which the MLE is a special case. They are defined as minimizers of some empirical risk function.

Let, for each  $\gamma \in \Gamma$ , be defined some loss function  $\rho_\gamma(X)$ . These are for instance constructed as in Chapter 10: we let  $L(\theta, a)$  be the loss when taking action  $a$ . Then, we fix some decision  $d(x)$ , and rewrite

$$L(\theta, d(x)) := \rho_\gamma(x),$$

assuming the loss  $L$  depends only on  $\theta$  via the parameter of interest  $\gamma = q(\theta)$ .

We now require that the *theoretical risk*

$$\mathcal{R}(c) := E_\theta \rho_c(X)$$

is minimized at the value  $c = \gamma$  i.e.,

$$\gamma = \arg \min_{c \in \Gamma} E_\theta \rho_c(X) = \arg \min_{c \in \Gamma} \mathcal{R}(c). \quad (14.1)$$

Alternatively, given  $\rho_c$ , one may view (14.1) as the *definition* of  $\gamma$ .

If  $c \mapsto \rho_c(x)$  is differentiable for all  $x$ , we write

$$\psi_c(x) := \dot{\rho}_c(x) := \frac{\partial}{\partial c} \rho_c(x).$$

Then, assuming we may interchange differentiation and taking expectations<sup>1</sup>, we have

$$\dot{\mathcal{R}}(\gamma) = 0,$$

where  $\dot{\mathcal{R}}(c) = E_\theta \psi_c(X)$ .

Define now the *empirical risk*

$$\hat{\mathcal{R}}_n(c) := \frac{1}{n} \sum_{i=1}^n \rho_c(X_i), \quad c \in \Gamma.$$

<sup>1</sup>If  $|\partial \rho_c / \partial c| \leq H(\cdot)$  where  $E_\theta H(X) < \infty$ , then it follows from the dominated convergence theorem that  $\partial [E_\theta \rho_c(X)] / \partial c = E_\theta [\partial \rho_c(X) / \partial c]$  or otherwise put,  $\dot{\mathcal{R}}(c) = E_\theta \psi(X)$ .

**Definition 14.0.1** The M-estimator  $\hat{\gamma}_n$  of  $\gamma$  is defined as

$$\hat{\gamma}_n := \arg \min_{c \in \Gamma} \frac{1}{n} \sum_{i=1}^n \rho_c(X_i) = \arg \min_{c \in \Gamma} \hat{\mathcal{R}}_n(c).$$

estimator minimizing  
emp risk

The “M” in “M-estimator” stands for Minimizer (or - take minus signs - Maximizer).

If  $\rho_c(x)$  is differentiable in  $c$  for all  $x$ , we generally can define  $\hat{\gamma}_n$  as the solution of putting the derivatives

$$\dot{\hat{\mathcal{R}}}_n(c) = \frac{\partial}{\partial c} \frac{1}{n} \sum_{i=1}^n \rho_c(X_i) = \frac{1}{n} \sum_{i=1}^n \psi_c(X_i)$$

→ assuming

to zero. This is called the Z-estimator. The “Z” in “Z-estimator” stands for Zero.

like this  $\psi_c$  will find your optimum

interchange  
diff.  
and  
integration.

**Definition 14.0.2** The Z-estimator  $\hat{\gamma}_n$  of  $\gamma$  is defined as a solution of the equations

$$\check{\mathcal{R}}_n(\hat{\gamma}_n) = 0$$

where  $\dot{\check{\mathcal{R}}}_n(c) = \frac{1}{n} \sum_{i=1}^n \psi_c(X_i)$ .

**Remark** A solution  $\hat{\gamma}_n \in \Gamma$  is then assumed to exist.

**Example 14.0.1** The least squares estimator

Let  $X \in \mathbb{R}$ , and let the parameter of interest be the mean  $\mu = E_\theta X$ . Assume  $X$  has finite variance  $\sigma^2$ . Then

$$\mu = \arg \min_c E_\theta(X - c)^2,$$

as (recall), by the bias-variance decomposition

$$E_\theta(X - c)^2 = \sigma^2 + (\mu - c)^2.$$

So in this case, we can take

$$\rho_c(x) = (x - c)^2.$$

Clearly

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

is minimized at  $c = \bar{X}_n := \sum_{i=1}^n X_i / n$ . See also Section 2.3.

## 14.1 MLE as special case of M-estimation

Suppose  $\Theta \subset \mathbb{R}^p$  and that the densities  $p_\theta = dP_\theta / d\nu$  exist w.r.t. some  $\sigma$ -finite measure  $\nu$ .

**Definition 14.1.1** The quantity

$$K(\tilde{\theta}|\theta) = E_{\theta} \log \left( \frac{p_{\theta}(X)}{p_{\tilde{\theta}}(X)} \right)$$

is called the Kullback Leibler information, or the relative entropy.

**Remark** Some care has to be taken, not to divide by zero! This can be handled e.g., by assuming that the support  $\{x : p_{\theta}(x) > 0\}$  does not depend on  $\theta$  (see also Condition I in the CRLB of Chapter 5).

Take for all  $\tilde{\theta} \in \Theta$

*risk measure*  $L(\theta, x)$

$$\rho_{\tilde{\theta}}(x) = -\log p_{\tilde{\theta}}(x).$$

As  $\log \left( \frac{p_{\theta}(x)}{p_{\tilde{\theta}}(x)} \right)$

i.e. here  
- likelihood  
= risk.

Then

$$\mathcal{R}(\tilde{\theta}) = -E_{\theta} \log p_{\tilde{\theta}}(X).$$

One easily sees that

$$K(\tilde{\theta}|\theta) = \mathcal{R}(\tilde{\theta}) - \mathcal{R}(\theta).$$

Let us restate Lemma 2.4.1 and reprove it in a slightly different manner.

**Lemma 14.1.1** The function  $\mathcal{R}(\tilde{\theta}) = -E_{\theta} \log p_{\tilde{\theta}}(X)$  is minimized at  $\tilde{\theta} = \theta$ :

$$\theta = \arg \min_{\tilde{\theta}} \mathcal{R}(\tilde{\theta}).$$

Proof  $\theta$  is the minimizer for the M-estimator

**Proof.** We will show that

$$K(\tilde{\theta}|\theta) \geq 0. \quad \text{from the fact that } K(\tilde{\theta}|\theta) = R(\tilde{\theta}) - R(\theta)$$

This follows from Jensen's inequality. Since the log-function is concave,

$$\begin{aligned} K(\tilde{\theta}|\theta) &= -E_{\theta} \log \left( \frac{p_{\tilde{\theta}}(X)}{p_{\theta}(X)} \right) \\ &\geq -\log \left( E_{\theta} \left( \frac{p_{\tilde{\theta}}(X)}{p_{\theta}(X)} \right) \right) \\ &= -\log 1 = 0. \end{aligned}$$

*notice  $\log(\cdot)$  is convex but  $-\log(\cdot)$  is concave*

□

With  $\rho_{\tilde{\theta}}(x) = -\log p_{\tilde{\theta}}(x)$  we find  $\psi_{\tilde{\theta}}(x) := \dot{\rho}_{\tilde{\theta}}(x) = -s_{\tilde{\theta}}(x)$ . Recall that  $s_{\theta}$  is the score function

$$s_{\theta} = \dot{p}_{\theta}/p_{\theta},$$

see Definition 4.7.1. We have seen moreover in Lemma 4.7.1 that  $E_{\theta} s_{\theta}(X) = 0$ . This is just another way to see that  $\theta$  is a solution of the equation

$$\dot{\mathcal{R}}(\theta) = 0,$$

where  $\dot{\mathcal{R}}(\tilde{\theta}) = E_{\theta} \psi_{\tilde{\theta}}(X)$  with  $\psi_{\tilde{\theta}} = -\dot{\rho}_{\tilde{\theta}}/p_{\tilde{\theta}}$ .

Proof of  $\theta$  being the minimizer for the  $\hat{\theta}$ -estimator.

With  $\rho_{\tilde{\theta}}(x) = -\log p_{\tilde{\theta}}(x)$  the M-estimator is the maximum likelihood estimator

$$\begin{aligned}\hat{\theta} &= \arg \min_{\tilde{\theta} \in \Theta} \mathcal{L}_X(\tilde{\theta}) \\ &= \arg \min_{\tilde{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n \left( -\log p_{\tilde{\theta}}(X_i) \right).\end{aligned}$$

## 14.2 Consistency of M-estimators



Note that  $\gamma$  minimizes a theoretical expectation, whereas the M-estimator  $\hat{\gamma}_n$  minimizes the empirical average. Likewise,  $\gamma$  is a solution of putting a theoretical expectation to zero, whereas the Z-estimator  $\hat{\gamma}_n$  is the solution of putting an empirical average to zero.

By the law of large numbers, averages converge to expectations. So the M-estimator (Z-estimator) does make sense. However, consistency and further properties are not immediate, because we actually need convergence the averages to expectations over a range of values  $c \in \Gamma$  simultaneously. This is the topic of empirical process theory.   
 many possible parametrizations.

consistency  
 $T_n \rightarrow \theta$  as  $n \rightarrow \infty$

We will borrow the notation from empirical process theory. For a function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , we let

$$P_\theta f := E_\theta f(X), \quad \hat{P}_n f := \frac{1}{n} \sum_{i=1}^n f(X_i).$$

Then, by the law of large numbers, if  $P_\theta |f| < \infty$ ,  $f$  bounded func

$$|(\hat{P}_n - P_\theta)f| \rightarrow 0, \quad \mathbb{P}_\theta\text{-a.s.}$$

With this new notation we have

$$\hat{R}_n(c) = \hat{P}_n \rho_c, \quad R(c) = P_\theta \rho_c.$$

**Theorem 14.2.1** Suppose the uniform convergence

$$\begin{aligned}\text{Then } &\xrightarrow{\text{this is uniform convergence}} \sup_{c \in \Gamma} |\hat{R}_n(c) - R(c)| \rightarrow 0, \quad \mathbb{P}_\theta\text{-a.s.} \\ &\xrightarrow{\quad} R(\hat{\gamma}_n) \rightarrow R(\gamma), \quad \mathbb{P}_\theta\text{-a.s.}\end{aligned}$$

**Proof.** The uniform convergence implies

$$\begin{aligned}0 &\leq P_\theta(\rho_{\hat{\gamma}_n} - \rho_\gamma) \xrightarrow{\text{minimizer}} R(\hat{\gamma}) - R(\gamma) \xrightarrow{\text{minimizes}} \hat{P}_n \rho_{\hat{\gamma}_n} \neq P_\theta \rho_\gamma. \\ &= -(\hat{P}_n - P_\theta)(\rho_{\hat{\gamma}_n} - \rho_\gamma) + (\hat{P}_n \rho_{\hat{\gamma}_n} - P_\theta \rho_\gamma) \xrightarrow{\text{minimizes}} 0 \\ &\leq -(\hat{P}_n - P_\theta)(\rho_{\hat{\gamma}_n} - \rho_\gamma) \xrightarrow{\text{Cauchy Schwartz}} \\ &\leq |(\hat{P}_n - P_\theta)\rho_{\hat{\gamma}_n}| + |(\hat{P}_n - P_\theta)\rho_\gamma| \\ &\leq \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta)\rho_c| + |(\hat{P}_n - P_\theta)\rho_\gamma| \\ &\leq 2 \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta)\rho_c|.\end{aligned}$$

Uniform convergence is stronger than pointwise convergence. It means that a sequence of functions  $f_n$  converges to a limiting function  $f$  on a set  $E$  if given any positive number  $\epsilon$ , a number  $N$  can be found such that each function  $f_n, f_{N+1}, \dots$  differ from  $f$  by no more than  $\epsilon$  at every point in  $E$ .

which finally  $\hat{\gamma}_n = \gamma$   
implies consistency of the estim.

so that  $\sup_{c \in \Gamma} |\hat{R}_n(c) - R(c)| = 0$   $\xrightarrow{\text{as this smaller equal first condition}}$   
implies  $\hat{R}_n(c) = R(c)$

through this

## 14.2. CONSISTENCY OF M-ESTIMATORS

this is why well separation is necessary

We will need that convergence of to the minimum value also implies convergence of the  $\arg \min$ , i.e., convergence of the location of the minimum. To this end, we present the following definition.

**Definition** The minimizer  $\gamma$  of  $\mathcal{R}(c)$  is called well-separated if for all  $\epsilon > 0$ ,

$$\inf \left\{ \mathcal{R}(c) : c \in \Gamma, \|c - \gamma\| \geq \epsilon \right\} \geq \mathcal{R}(\gamma).$$

so no other  $c \in \Gamma$  reaches the minimum.

If  $\gamma$  is well-separated,  $\mathcal{R}(\hat{\gamma}_n) \rightarrow \mathcal{R}(\gamma)$   $\mathbb{P}_\theta$ -a.s.. implies  $\hat{\gamma}_n \rightarrow \gamma$   $\mathbb{P}_\theta$ -a.s..

So the general strategy is the following:  
 1. Prove uniform convergence of  $\hat{\rho}_n(x)$ .  
 2. If  $\mathcal{R}(c)$  is well separated then you obtain asymptotic consistency as  $\hat{\gamma}_n \rightarrow \gamma$ .

Proof of uniform convergence.



In the next lemma, we give sufficient conditions for the uniform in  $c$  convergence of the empirical risk  $\hat{\mathcal{R}}_n(c)$  to the theoretical risk  $\mathcal{R}(c)$ . Consistency of the M-estimator is then a consequence, as was shown in Theorem 14.2.1. (For consistency the assumption of a compact parameter space  $\Gamma$  can often be omitted if  $c \mapsto \rho_c$  is convex. We skip the details.)

**Lemma 14.2.1** Suppose that  $\Gamma$  is compact, that  $c \mapsto \rho_c(x)$  is continuous for all  $x$ , and that

$$P_\theta \left( \sup_{c \in \Gamma} |\rho_c| \right) < \infty.$$

Then we have the uniform convergence

$$\sup_{c \in \Gamma} |(\hat{\mathcal{R}}_n - P_\theta)\rho_c| \rightarrow 0, \quad \mathbb{P}_\theta\text{-a.s.} \quad (14.2)$$

**Proof.** Define for each  $\delta > 0$  and  $c \in \Gamma$ ,

$$w(\cdot, \delta, c) := \sup_{\tilde{c} \in \Gamma: \|\tilde{c} - c\| < \delta} |\rho_{\tilde{c}} - \rho_c|. \quad \text{Maximum risk.}$$

Then for all  $x$ , as  $\delta \downarrow 0$ ,

$$w(x, \delta, c) \rightarrow 0.$$

So also, by dominated convergence (look at wikipedia page in Italian).

$$P_\theta w(\cdot, \delta, c) \rightarrow 0.$$

Hence, for all  $\epsilon > 0$ , there exists a  $\delta_c$  such that

$\delta_c \rightarrow 0$ . { the one of convergence for which was shown. }  $P_\theta w(\cdot, \delta_c, c) \leq \epsilon.$  Let  $B_c := \{\tilde{c} \in \Gamma : \|\tilde{c} - c\| < \delta_c\}$ . exists this ball!

Then  $\{B_c : c \in \Gamma\}$  is a covering of  $\Gamma$  by open sets. Since  $\Gamma$  is compact, there exists finite sub-covering

$$B_{c_1} \dots B_{c_N}.$$

For  $c \in B_{c_j}$ ,

$$|\rho_c - \rho_{c_j}| \leq w(\cdot, \delta_{c_j}, c_j). \quad \text{by def of } w(\cdot, \delta, c)$$

cause this supremum by definition.

*signs  
versus*

It follows that

$$\begin{aligned} \sup_{c \in \Gamma} |(\hat{P}_n - P_\theta) \rho_c| &\leq \max_{1 \leq j \leq N} |(\hat{P}_n - P_\theta) \rho_{c_j}| \\ &+ \max_{1 \leq j \leq N} (\hat{P}_n w(\cdot, \delta_{c_j}, c_j) + P_\theta w(\cdot, \delta_{c_j}, c_j)) \\ &\rightarrow 2 \max_{1 \leq j \leq N} P_\theta w(\cdot, \delta_{c_j}, c_j) \leq \underline{2\epsilon}, \text{ } \mathbb{P}_\theta\text{-a.s..} \end{aligned}$$

$\hat{P}_n \nearrow P$   
 $\rho_c \nearrow \rho_c$

$\hat{P}_n \nearrow P$   
 $\rho_c \nearrow \rho_c$

*there together two in expectation*

*the expectation right on hand side*

*integrate with pts  
on the phone.*

□

### Example 14.2.1 Consistency of the MLE in the logistic location family

The above theorem directly uses the definition of the M-estimator, and does not rely on having an explicit expression available. Here is an example where an explicit expression is indeed not possible. Consider the logistic location family, where the densities are

$$p_\theta(x) = \frac{e^{x-\theta}}{(1+e^{x-\theta})^2}, \quad x \in \mathbb{R},$$

where  $\theta \in \Theta \subset \mathbb{R}$  is the location parameter. Take

$$\rho_\theta(x) := -\log p_\theta(x) = \theta - x + 2 \log(1 + e^{x-\theta}).$$

Then  $\hat{\theta}_n$  is the MLE. It is a solution of

$$\frac{2}{n} \sum_{i=1}^n \frac{e^{X_i - \hat{\theta}_n}}{1 + e^{X_i - \hat{\theta}_n}} = 1.$$

*This one is exponent.*

*Nevertheless*  
we know  
 $\hat{\theta}_n \rightarrow \theta$ .  
given the previous  
theorem.

*This expression cannot be made into an explicit expression for  $\hat{\theta}_n$ . However, we do note the caveat that in order to be able to apply the above consistency theorem, we need to assume that  $\Theta$  is compact. This problem can be circumvented by using the result below for Z-estimators.*

To prove consistency of a Z-estimator of a one-dimensional parameter is relatively easy.

Recall that  $\psi_c = \dot{\rho}_c$  and  $\dot{\mathcal{R}}(c) = P_\theta \psi_c := E_\theta \psi_c(X)$ ,  $c \in \Gamma$ . Recall further that  $\dot{\mathcal{R}}(\gamma) = 0$  since  $\gamma$  is defined as the minimizer of  $\mathcal{R}(\cdot)$ .

**Theorem 14.2.2** Suppose that  $\Gamma \subset \mathbb{R}$  and that  $\psi_c(x)$  is continuous in  $c$  for all  $x$ . Assume moreover that

$$P_\theta |\psi_c| < \infty, \quad \forall c,$$

and that  $\exists \delta > 0$  such that

$$\dot{\mathcal{R}}(c) > 0, \quad \gamma < c < \gamma + \delta,$$

$$\dot{\mathcal{R}}(c) < 0, \quad \gamma - \delta < c < \gamma.$$

✓

*One dimensional case*

*Consistency  
for Z-estimator*

Then for  $n$  large enough,  $\mathbb{P}_\theta$ -a.s., there is a solution  $\hat{\gamma}_n$  of  $\dot{\mathcal{R}}_n(\hat{\gamma}_n) = 0$ , and this solution  $\hat{\gamma}_n$  is consistent.

**Proof.** Let  $0 < \epsilon < \delta$  be arbitrary. By the law of large numbers,  $\mathbb{P}_\theta$ -a.s. for  $n$  sufficiently large,

$$\dot{\mathcal{R}}_n(\gamma + \epsilon) > 0, \quad \dot{\mathcal{R}}_n(\gamma - \epsilon) < 0.$$

The continuity of  $c \mapsto \psi_c$  implies that then  $\dot{\mathcal{R}}_n(\hat{\gamma}_n) = 0$  for some  $|\hat{\gamma}_n - \gamma| < \epsilon$ .

□

### 14.3 Asymptotic normality of M-estimators

For a function  $f : \mathcal{X} \rightarrow \mathbb{R}^p$  we let  $P_\theta f := E_\theta f(X) \in \mathbb{R}^p$  (whenever it exists). Moreover, we let

$$P_\theta f f^T = E_\theta f(X) f^T(X) \in \mathbb{R}^{p \times p}$$

(whenever it exists). The covariance matrix of the vector  $f(X)$  is thus

$$\Sigma := P_\theta f f^T - (P_\theta f)(P_\theta f)^T$$

The CLT says that

$$\sqrt{n}(\hat{P}_n - P_\theta)f \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, \Sigma).$$

With this notation we can translate Definition 13.3.1 as:  $T_n$  is an asymptotically linear estimator of  $\gamma$  if

$$T_n - \gamma = \hat{P}_n l_\theta + o_{\mathbb{P}_\theta}(1/\sqrt{n}),$$

with  $P_\theta l_\theta = 0$  and  $V_\theta := P_\theta l_\theta l_\theta^T < \infty$ .

Definition Denote now

$$\nu_n(c) := \sqrt{n}(\hat{P}_n - P_\theta)\psi_c = \sqrt{n}\left(\dot{\mathcal{R}}_n(c) - \dot{\mathcal{R}}(c)\right), \quad c \in \Gamma.$$

**Definition 14.3.1** The stochastic process

$$\{\nu_n(c) : c \in \Gamma\}$$

is called the empirical process indexed by  $c$ . The empirical process is called asymptotically continuous at  $\gamma$  if for all (possibly random) sequences  $\{\gamma_n\}$  in  $\Gamma$ , with  $\|\gamma_n - \gamma\| = o_{\mathbb{P}_\theta}(1)$ , we have

$$|\nu_n(\gamma_n) - \nu_n(\gamma)| = o_{\mathbb{P}_\theta}(1). \quad \begin{matrix} \text{convergence} \\ \text{probability} \\ \xrightarrow{\mathbb{P}_\theta} 0 \end{matrix}$$

i.e. if  
 $y_n \xrightarrow{\mathbb{P}_\theta} y$   
so does  
 $\nu_n(\hat{\mathcal{R}}(y_n))$

For verifying asymptotic continuity, there are various tools, which involve complexity assumptions on the map  $c \mapsto \psi_c$ . This goes beyond the scope of these notes. But let us see what asymptotic continuity can bring up.

Recall that  $\dot{\mathcal{R}}(c) = P_\theta \psi_c$ . We assume that

$$M_\theta := \frac{\partial}{\partial c^T} \dot{\mathcal{R}}(c) \Big|_{c=\gamma}$$

exists. It is a  $p \times p$  matrix. We require it to be of full rank, which amounts to assuming that  $\gamma$ , as a solution to  $\dot{\mathcal{R}}(\gamma) = 0$ , is well-identified.

$R(c) = P_\theta \psi_c$   
 $L(c, d(x))$   
This is of  $\dim R$ ,  
so that

$$\text{Recall } \dot{\mathcal{R}}(c) = P_\theta \psi_c = \dot{\phi}_c(x)$$

*Note to the student!*

**Theorem 14.3.1** Let  $\hat{\gamma}_n$  be the Z-estimator of  $\gamma$ . Suppose that  $\hat{\gamma}_n$  is a consistent estimator of  $\gamma$ , and that  $\nu_n$  is asymptotically continuous at  $\gamma$ . Suppose moreover  $M_\theta^{-1}$  exists, and also

$$\underline{P_\theta \psi_\gamma \psi_\gamma^T} \quad \left. \begin{array}{l} \text{P. 1. 1. p} \\ \text{Jacobi} \end{array} \right\}$$

Then  $\hat{\gamma}_n$  is asymptotically linear with influence function

$$l_\theta = -M_\theta^{-1} \psi_\gamma.$$

**Proof.** By definition,

$$\dot{\mathcal{R}}_n(\hat{\gamma}_n) = 0, \quad \dot{\mathcal{R}}(\gamma) = 0.$$

So we have

$$\begin{aligned} 0 &= \dot{\mathcal{R}}_n(\hat{\gamma}_n) \\ &= \nu_n(\hat{\gamma}_n)/\sqrt{n} + \dot{\mathcal{R}}(\hat{\gamma}_n) \quad \xrightarrow{\text{error: this}} \quad \left[ \begin{array}{l} \text{as all of the terms involved} \\ \text{contain } \mathcal{R}(\gamma) \text{ and } \dot{\mathcal{R}}(\gamma), \\ \text{which by def} = 0. \end{array} \right] \\ &= \nu_n(\hat{\gamma}_n)/\sqrt{n} + \dot{\mathcal{R}}(\hat{\gamma}_n) - \dot{\mathcal{R}}(\gamma) \\ &=: \underline{(i)} + \underline{(ii)}. \quad \xrightarrow{\text{by def}} \end{aligned}$$

For the first term, we use the asymptotic continuity of  $\nu_n$  at  $\gamma$ :

$$\begin{aligned} (i) &= \nu_n(\hat{\gamma}_n)/\sqrt{n} \\ &= \nu_n(\gamma)/\sqrt{n} + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \\ &= \dot{\mathcal{R}}_n(\gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}). \end{aligned} \quad \left. \begin{array}{l} \text{idea: } \nu(\gamma) - \nu(\hat{\gamma}_n) = \sigma_{\mathbf{P}_\theta}(1) \\ \text{asymptotic linearity.} \end{array} \right\}$$

For the second term, we use the differentiability of  $\dot{\mathcal{R}}(c) = P_\theta \psi_c$  at  $c = \gamma$ :

$$\begin{aligned} (ii) &= \dot{\mathcal{R}}(\hat{\gamma}_n) - \dot{\mathcal{R}}(\gamma) \\ &= M_\theta(\hat{\gamma}_n - \gamma) + o(\|\hat{\gamma}_n - \gamma\|). \quad \xrightarrow{\text{taylor expansion. first term.}} \end{aligned}$$

So we arrive at

$$0 = \dot{\mathcal{R}}_n(\gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}) + M_\theta(\hat{\gamma}_n - \gamma) + o(\|\hat{\gamma}_n - \gamma\|).$$

Because, by the CLT,  $\dot{\mathcal{R}}_n(\gamma) = \mathcal{O}_{\mathbf{P}_\theta}(1/\sqrt{n})$ , this implies  $\|\hat{\gamma}_n - \gamma\| = \mathcal{O}_{\mathbf{P}_\theta}(1/\sqrt{n})$ . Hence

$$0 = \dot{\mathcal{R}}_n(\gamma) + M_\theta(\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}),$$

or

$$\begin{aligned} M_\theta(\hat{\gamma}_n - \gamma) &= -\dot{\mathcal{R}}_n(\gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \\ &= -\hat{P}_n \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \end{aligned}$$

or

$$(\hat{\gamma}_n - \gamma) = -\hat{P}_n M^{-1} \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}).$$

Important  $\square$

Not true actually same issue as with notation on the left  
Very tricky business actually everything unchanged.  
i.e. our  $R$  and the other  $R_n$

**Corollary 14.3.1** Under the conditions of Theorem 14.3.1

$$\sqrt{n}(\hat{\gamma}_n - \gamma) \xrightarrow{D_\theta} \mathcal{N}(0, V_\theta),$$

with

$$V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1}. \quad \begin{cases} \text{for the variance of the} \\ \text{linear term } [-P_n M^{-1} \psi_\gamma] \end{cases}$$

Asymptotic linearity can also be established directly, under rather restrictive assumptions, see Theorem 14.3.2 coming up next. We assume quite a lot of smoothness for the functions  $\psi_c$  (namely, derivatives that are Lipschitz), so that asymptotic linearity can be proved by straightforward arguments. We stress however that such smoothness assumptions are by no means necessary.

**Theorem 14.3.2** Let  $\hat{\gamma}_n$  be the Z-estimator of  $\gamma$ , and suppose that  $\hat{\gamma}_n$  is a consistent estimator of  $\gamma$ . Suppose that, for all  $c$  in a neighborhood  $\{c \in \Gamma : \|c - \gamma\| < \epsilon\}$ , the map  $c \mapsto \psi_c(x)$  is differentiable for all  $x$ , with derivative

$$\dot{\psi}_c(x) = \frac{\partial}{\partial c^T} \psi_c(x)$$

(a  $p \times p$  matrix). Assume moreover that, for all  $c$  and  $\tilde{c}$  in a neighborhood of  $\gamma$ , and for all  $x$ , we have, in matrix-norm<sup>2</sup>,

$$\|\dot{\psi}_c(x) - \dot{\psi}_{\tilde{c}}(x)\| \leq H(x)\|c - \tilde{c}\|, \quad \begin{cases} \text{Lipschitz} \\ \text{derivative.} \end{cases}$$

where  $H : \mathcal{X} \rightarrow \mathbb{R}$  satisfies

$$P_\theta H < \infty.$$

Then

$$M_\theta := \left. \frac{\partial}{\partial c^T} \dot{\mathcal{R}}(c) \right|_{c=\gamma} = P_\theta \dot{\psi}_\gamma. \quad (14.3)$$

Assuming  $M_\theta^{-1}$  and  $J_\theta := E_\theta \psi_\gamma \psi_\gamma^T$  exist, the influence function of  $\hat{\gamma}_n$  is

$$l_\theta = -M_\theta^{-1} \psi_\gamma.$$

X  
 No  
 need  
 to  
 know  
 the  
 proof.

**Proof.** Result (14.3) follows from the dominated convergence theorem.

By the mean value theorem,

$$\begin{aligned} 0 &= \dot{\mathcal{R}}_n(\hat{\gamma}) \\ &= \hat{P}_n \dot{\psi}_{\tilde{\gamma}_n} \\ &= \hat{P}_n \psi_\gamma + \hat{P}_n \dot{\psi}_{\tilde{\gamma}_n(\cdot)}(\tilde{\gamma}_n - \gamma) \quad \begin{cases} \text{mean} \\ \text{value} \end{cases} \quad \begin{cases} \text{theorem} \\ \text{theorem} \end{cases} \\ &= \dot{\mathcal{R}}_n(\gamma) + \hat{P}_n \dot{\psi}_{\tilde{\gamma}_n(\cdot)}(\tilde{\gamma}_n - \gamma) \end{aligned}$$

where for all  $x$ ,  $\|\tilde{\gamma}_n(x) - \gamma\| \leq \|\tilde{\gamma}_n - \gamma\|$ . Thus

$$0 = \dot{\mathcal{R}}_n(\gamma) + \hat{P}_n \dot{\psi}_\gamma(\tilde{\gamma}_n - \gamma) + \hat{P}_n (\dot{\psi}_{\tilde{\gamma}_n(\cdot)} - \dot{\psi}_\gamma)(\tilde{\gamma}_n - \gamma),$$

<sup>2</sup>For a matrix  $A$ ,  $\|A\| := \sup_{v \neq 0} \|Av\|/\|v\|$ .

$$\dot{\hat{R}}_n(\gamma) + \hat{P}_n \psi_\gamma (\hat{\gamma}_n - \gamma) = \hat{P}_n (\dot{\psi}_{\hat{\gamma}_n} - \dot{\psi}_\gamma) (\hat{\gamma}_n - \gamma)$$

so that

*if taking the norm  
and considering Lipschitz derivatives*

$$\left\| \dot{\hat{R}}_n(\gamma) + \hat{P}_n \psi_\gamma (\hat{\gamma}_n - \gamma) \right\| \leq \left( \hat{P}_n H \right) \|\hat{\gamma}_n - \gamma\|^2 = \mathcal{O}_{\mathbf{P}_\theta}(1) \|\hat{\gamma}_n - \gamma\|^2,$$

where in the last inequality, we used  $\hat{P}_n H < \infty$ . Now, by the law of large numbers,

$$\hat{P}_n \dot{\psi}_\gamma = P_\theta \dot{\psi}_\gamma + o_{\mathbf{P}_\theta}(1) = M_\theta + o_{\mathbf{P}_\theta}(1).$$

Thus

$$\left| \dot{\hat{R}}_n(\gamma) + M_\theta (\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(\|\hat{\gamma}_n - \gamma\|) \right| = \mathcal{O}_{\mathbf{P}_\theta}(\|\hat{\gamma}_n - \gamma\|^2).$$

Because  $\dot{\hat{R}}_n(\gamma) = \mathcal{O}_{\mathbf{P}_\theta}(1/\sqrt{n})$  by the CLT, this ensures that  $\|\hat{\gamma}_n - \gamma\| = \mathcal{O}_{\mathbf{P}_\theta}(1/\sqrt{n})$ . It follows that

$$\left| \dot{\hat{R}}_n(\gamma) + M_\theta (\hat{\gamma}_n - \gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \right| = \mathcal{O}_{\mathbf{P}_\theta}(1/n).$$

Hence

$$M_\theta (\hat{\gamma}_n - \gamma) = -\dot{\hat{R}}_n(\gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n})$$

and so

$$\begin{aligned} \hat{\gamma}_n - \gamma &= -M_\theta^{-1} \dot{\hat{R}}_n(\gamma) + o_{\mathbf{P}_\theta}(1/\sqrt{n}) \\ &= -\hat{P}_n M_\theta^{-1} \psi_\gamma + o_{\mathbf{P}_\theta}(1/\sqrt{n}). \end{aligned}$$

□

**Note** The asymptotic normality follows again from the asymptotic linearity established in Theorem 14.3.2.

## 14.4 Asymptotic normality of the MLE

In this section, we show that, under regularity conditions, the MLE is asymptotically normal with asymptotic covariance matrix the inverse of the Fisher-information matrix  $I(\theta)$ . We use that maximum likelihood estimation is a special case of M-estimation and apply the results of the previous section. In order to do so we need to assume regularity conditions.

Let  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  be dominated by a  $\sigma$ -finite dominating measure  $\nu$ , and write the densities as  $p_\theta = dP_\theta/d\nu$ . Suppose that  $\Theta \subset \mathbb{R}^p$ . Assume that the support of  $p_\theta$  does not depend on  $\theta$  (Condition I in Section 5.5). As loss we take minus the log-density:

$$\rho_\theta := -\log p_\theta.$$

*the index  
must enter the  
relation*

The MLE is

$$\hat{\theta}_n := \arg \max_{\theta \in \Theta} \hat{P}_n \log p_\theta.$$

We suppose that the score function

$$s_\theta = \frac{\partial}{\partial \theta} \log p_\theta = \frac{\dot{p}_\theta}{p_\theta}$$

exists, and that we may interchange differentiation and integration, so that the score has mean zero:

$$\overbrace{P_\theta s_\theta} = \int \dot{p}_\theta d\nu = \frac{\partial}{\partial \theta} \int p_\theta d\nu = \frac{\partial}{\partial \theta} 1 = 0.$$

Recall that the Fisher-information matrix is

$$\boxed{I(\theta) := P_\theta s_\theta s_\theta^T.}$$

Now, it is clear that  $\psi_\theta = -s_\theta$ , and, assuming derivatives exist and that again we may change the order of differentiation and integration,

$$\boxed{M_\theta = P_\theta \dot{\psi}_\theta = -P_\theta \dot{s}_\theta,}$$

and (see also Lemma 4.7.1)

$$\begin{aligned} \underline{P_\theta \dot{s}_\theta} &= P_\theta \left( \frac{\ddot{p}_\theta}{p_\theta} - \frac{\dot{p}_\theta \dot{p}_\theta^T}{p_\theta p_\theta} \right) \text{ This is} \\ &= \left( \frac{\partial^2}{\partial \theta \partial \theta^T} 1 \right) - P_\theta s_\theta s_\theta^T \\ &= 0 - I(\theta) = \boxed{-I(\theta)} \end{aligned}$$

Hence, in this case,  $M_\theta = -I(\theta)$ , and the influence function is

$$l_\theta = I(\theta)^{-1} s_\theta.$$

So the asymptotic covariance matrix of the MLE  $\hat{\theta}_n$  is

$$\boxed{I(\theta)^{-1} \left( P_\theta s_\theta s_\theta^T \right) I(\theta)^{-1} = I(\theta)^{-1}.}$$

It follows that

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

## 14.5 Two further examples of M-estimation

In this section we examine the  $\alpha$ -quantile and the Huber estimator.

### Example 14.5.1 Asymptotic normality of the $\alpha$ -quantile

*In this example, the parameter of interest is the  $\alpha$ -quantile. We will consider a loss function which does not satisfy regularity conditions, but nevertheless leads to an asymptotically linear, and hence asymptotically normal, estimator.*



Let  $\mathcal{X} := \mathbb{R}$ . The distribution function of  $X$  is denoted by  $F$ . Let  $0 < \alpha < 1$  be given. The  $\alpha$ -quantile of  $F$  is  $\gamma = F^{-1}(\alpha)$  (assumed to exist). We moreover assume that  $F$  has density  $f$  with respect to Lebesgue measure, and that  $f(x) > 0$  in a neighborhood of  $\gamma$ . As loss function we take

$$\rho_c(x) := \rho(x - \gamma), \quad \text{here is where it enters the relation.}$$

where

$$\rho(x) := (1 - \alpha)|x|l\{x < 0\} + \alpha|x|l\{x > 0\}.$$

We now first check that for  $\mathcal{R}(c) := P_\theta \rho_c$

$$\arg \min_c \mathcal{R}(c) = F^{-1}(\alpha) := \gamma.$$

We have

$$\dot{\rho}(x) = \alpha l\{x > 0\} - (1 - \alpha)l\{x < 0\}.$$

Note that  $\dot{\rho}$  does not exist at  $x = 0$ . This is one of the irregularities in this example.

It follows that

$$\psi_c(x) = -\alpha l\{x > c\} + (1 - \alpha)l\{x < c\}.$$

Hence

$$\dot{\mathcal{R}}(c) = P_\theta \psi_c = -\alpha + F(c).$$

} as times  
- 1 more  
ableitung!

(the fact that  $\psi_c$  is not defined at  $x = c$  can be shown not to be a problem, roughly because a single point has probability zero, as  $F$  is assumed to be continuous).

So

$$\boxed{\dot{\mathcal{R}}(\gamma) = 0, \text{ for } \gamma = F^{-1}(\alpha)}.$$

We now derive  $M_\theta$ , which is a scalar in this case:

$$\begin{aligned} M_\theta &= \frac{d}{dc} \dot{\mathcal{R}}(c) \Big|_{c=\gamma} \\ &= \frac{d}{dc} (-\alpha + F(c)) \Big|_{c=\gamma} \\ &= f(\underline{\gamma}) = f(F^{-1}(\alpha)). \end{aligned}$$

The influence function is thus <sup>3</sup>

$$l_\theta(x) = -M_\theta^{-1} \psi_\gamma(x) = \frac{1}{f(\gamma)} \left\{ -l\{x < \gamma\} + \alpha \right\}.$$

<sup>3</sup>Note that in the special case  $\alpha = 1/2$  (where  $\gamma$  is the median), this becomes

$$l_\theta(x) = \begin{cases} -\frac{1}{2f(\gamma)} & x < \gamma \\ +\frac{1}{2f(\gamma)} & x > \gamma \end{cases}.$$

We conclude that, for  $\hat{\mathcal{R}}_n(c) := \hat{P}_n \rho_c$  and

$$\hat{\gamma}_n := \arg \min_c \hat{\mathcal{R}}_n(c),$$

which we write as the sample quantile  $\hat{\gamma}_n = \hat{F}_n^{-1}(\alpha)$  (or an approximation thereof up to order  $o_{\mathbf{P}_\theta}(1/\sqrt{n})$ ), one has

$$\sqrt{n}(\hat{F}_n^{-1}(\alpha) - F^{-1}(\alpha)) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}\left(0, \frac{\alpha(1-\alpha)}{f^2(F^{-1}(\alpha))}\right).$$

### Example 14.5.2 Asymptotic normality of the Huber estimator

In this example, we illustrate that the Huber-estimator is asymptotically linear and hence asymptotically normal. Let again  $\mathcal{X} = \mathbb{R}$  and  $F$  be the distribution function of  $X$ . We let the parameter of interest be the a location parameter. The Huber loss function is

$$\rho_c(x) = \rho(x - c),$$

with

$$\rho(x) = \begin{cases} x^2 & |x| \leq k \\ k(2|x| - k) & |x| > k \end{cases}.$$

We define  $\gamma$  as

$$\gamma := \arg \min_c P_\theta \rho_c.$$

It holds that

$$\dot{\rho}(x) = \begin{cases} 2x & |x| \leq k \\ +2k & x > k \\ -2k & x < -k \end{cases}.$$

Therefore,

$$\psi_c(x) = \begin{cases} -2(x - c) & |x - c| \leq k \\ -2k & x - c > k \\ +2k & x - c < -k \end{cases}.$$

One easily derives that

$$\begin{aligned} \dot{\mathcal{R}}(c) := P_\theta \psi_c &= -2 \int_{-k+c}^{k+c} x dF(x) + 2c[F(k+c) - F(-k+c)] \\ &= 2k[1 - F(k+c)] + 2kF(-k+c). \end{aligned}$$

So

$$\underline{M_\theta} = \frac{d}{dc} \dot{\mathcal{R}}(c) \Big|_{c=\gamma} = 2[F(k+\gamma) - F(-k+\gamma)].$$

The influence function of the Huber estimator is

$$l_\theta(x) = \frac{1}{[F(k+\gamma) - F(-k+\gamma)]} \begin{cases} x - \gamma & |x - \gamma| \leq k \\ +k & x - \gamma > k \\ -k & x - \gamma < -k \end{cases}.$$

For  $k \rightarrow 0$ , this corresponds to the influence function of the median.

## 14.6 Asymptotic relative efficiency

In this section, we assume that the parameter of interest is real-valued:

$$\gamma \in \Gamma \subseteq \mathbb{R}.$$

**Definition 14.6.1** Let  $T_{n,1}$  and  $T_{n,2}$  be two estimators of  $\gamma$ , that satisfy

$$\sqrt{n}(T_{n,j} - \gamma) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, V_{\theta,j}), \quad j = 1, 2.$$

Then

$$e_{2:1} := \frac{V_{\theta,1}}{V_{\theta,2}}$$

is called the asymptotic relative efficiency of  $T_{n,2}$  with respect to  $T_{n,1}$ .

If  $e_{2:1} > 1$ , the estimator  $T_{n,2}$  is asymptotically more efficient than  $T_{n,1}$ . An asymptotic  $(1-\alpha)$ -confidence interval for  $\gamma$  based on  $T_{n,2}$  is then narrower than the one based on  $T_{n,1}$ .

**Example 14.6.1** Asymptotic relative efficiency of sample mean and sample median

Let  $\mathcal{X} = \mathbb{R}$ , and  $F$  be the distribution function of  $X$ . Suppose that  $F$  is symmetric around the parameter of interest  $\mu$ . In other words,

$$F(\cdot) = F_0(\cdot - \mu),$$

where  $F_0$  is symmetric around zero. We assume that  $F_0$  has finite variance  $\sigma^2$ , and that it has density  $f_0$  w.r.t. Lebesgue measure, with  $f_0(0) > 0$ . Take  $T_{n,1} := \bar{X}_n$ , the sample mean, and  $T_{n,2} := \hat{F}_n^{-1}(1/2)$ , the sample median. Then  $V_{\theta,1} = \sigma^2$  and  $V_{\theta,2} = 1/(4f_0^2(0))$  (the latter being derived in Example 14.5.1). So

$$e_{2:1} = 4\sigma^2 f_0^2(0).$$

Whether the sample mean is the winner, or rather the sample median, depends thus on the distribution  $F_0$ . Let us consider three cases.

**Case i** Let  $F_0$  be the standard normal distribution, i.e.,  $F_0 = \Phi$ . Then  $\sigma^2 = 1$  and  $f_0(0) = 1/\sqrt{2\pi}$ . Hence

$$e_{2:1} = \frac{2}{\pi} \approx 0.64.$$

So  $\bar{X}_n$  is the winner. Note that  $\bar{X}_n$  is the MLE in this case.

**Case ii** Let  $F_0$  be the Laplace distribution, with variance  $\sigma^2$  equal to one. This distribution has density

$$f_0(x) = \frac{1}{\sqrt{2}} \exp[-\sqrt{2}|x|], \quad x \in \mathbb{R}.$$

So we have  $f_0(0) = 1/\sqrt{2}$ , and hence

$$e_{2:1} = 2.$$

at least.

You can then calculate the number of obs. you need more using  $V_{\theta,2}$  for having the same CI length:

i.e.:

$$\frac{\text{length}_1}{\text{length}_2} = \frac{2\sqrt{\frac{1}{n_1}} \cdot \Phi(1-\alpha)}{2\sqrt{\frac{1}{n_2}} \cdot \Phi(1-\alpha)} \stackrel{?}{=} 1$$

$$\Leftrightarrow \frac{V_{\theta,1}}{V_{\theta,2}} = 1 \Leftrightarrow \frac{V_{\theta,1}}{V_{\theta,2}} = \frac{n_2}{n_1}$$

Thus, the sample median, which is the MLE for this case, is the winner.

*Some idea  
wilcoxon test,  
It dist. not  
normal you  
are better  
off using  
the median!*

Case iii Suppose

$$F_0 = (1 - \eta)\Phi + \eta\Phi(\cdot/3). \quad \begin{array}{l} \text{mixed dist.} \\ \text{= How many obs do I need then?} \end{array}$$

This means that the distribution of  $X$  is a mixture, with mixing probabilities  $1 - \eta$  and  $\eta$ , of two normal distributions, one with unit variance, and one with variance  $3^2 = 9$ . Otherwise put, associated with  $X$  is an unobservable label  $Y \in \{0, 1\}$ . Given  $Y = 1$ , the random variable  $X$  is  $\mathcal{N}(\mu, 1)$ -distributed. Given  $Y = 0$ , the random variable  $X$  has a  $\mathcal{N}(\mu, 3^2)$  distribution. Moreover,  $P(Y = 1) = 1 - P(Y = 0) = 1 - \eta$ . Hence

$$\sigma^2 := \text{var}(X) = (1 - \eta)\text{var}(X|Y = 1) + \eta\text{var}(X|Y = 0) = (1 - \eta) + 9\eta = 1 - 8\eta.$$

It furthermore holds that

$$f_0(0) = (1 - \eta)\phi(0) + \frac{\eta}{3}\phi(0) = \frac{1}{\sqrt{2\pi}} \left(1 - \frac{2\eta}{3}\right).$$

It follows that

$$e_{2:1} = \frac{2}{\pi} \left(1 - \frac{2\eta}{3}\right)^2 (1 + 8\eta).$$

Let us now further compare the results with the  $\alpha$ -trimmed mean. Because  $F$  is symmetric, it turns out that the  $\alpha$ -trimmed mean has the same influence function as the Huber-estimator with  $k = F^{-1}(1 - \alpha)$ :

$$l_\theta(x) = \frac{1}{F_0(k) - F(-k)} \begin{cases} x - \mu, & |x - \mu| \leq k \\ +k, & x - \mu > k \\ -k, & x - \mu < -k \end{cases}.$$

(This can be seen from Example 15.3.2 ahead which is not part of the exam). The influence function is used to compute the asymptotic variance  $V_{\theta,\alpha}$  of the  $\alpha$ -trimmed mean:

$$V_{\theta,\alpha} = \frac{\int_{F_0^{-1}(\alpha)}^{F_0^{-1}(1-\alpha)} x^2 dF_0(x) + 2\alpha(F_0^{-1}(1 - \alpha))^2}{(1 - 2\alpha)^2}.$$

From this, we then calculate the asymptotic relative efficiency of the  $\alpha$ -trimmed mean w.r.t. the mean. Note that the median is the limiting case with  $\alpha \rightarrow 1/2$ .

Table: Asymptotic relative efficiency of  $\alpha$ -trimmed mean over mean

	$\alpha = 0.05$	$0.125$	$0.5$
$\eta = 0.00$	0.99	0.94	0.64
0.05	1.20	1.19	0.83
0.25	1.40	1.66	1.33

## 14.7 Asymptotic pivots

Again throughout this section, enough regularity is assumed, such as existence of derivatives and interchanging integration and differentiation.



Recall the definition of an asymptotic pivot (see Section 6.2). It is a function  $Z_n(\gamma) := Z_n(X_1, \dots, X_n, \gamma)$  of the data  $X_1, \dots, X_n$  and the parameter of interest  $\gamma = g(\theta) \in \mathbb{R}^p$ , such that its asymptotic distribution does not depend on the unknown parameter  $\theta$ , i.e., for a random variable  $Z$ , with distribution  $Q$  not depending on  $\theta$ ,

$$\boxed{Z_n(\gamma) \xrightarrow{\mathcal{D}_\theta} Z, \forall \theta.}$$

An asymptotic pivot can be used to construct approximate  $(1 - \alpha)$ -confidence intervals for  $\gamma$ , and tests for  $H_0 : \gamma = \gamma_0$  with approximate level  $\alpha$ .

Consider now an asymptotically normal estimator  $T_n$  of  $\gamma$ , which is asymptotically unbiased and has asymptotic covariance matrix  $V_\theta$ , that is

$$\sqrt{n}(T_n - \gamma) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, V_\theta), \forall \theta.$$

This still depends on the parameter that affects the variance.

(assuming such an estimator exists). Then, depending on the situation, there are various ways to construct an asymptotic pivot.

full rank

### 1<sup>st</sup> asymptotic pivot

If the asymptotic covariance matrix  $V_\theta$  is non-singular, and depends only on the parameter of interest, say  $V_\theta = V(\gamma)$  (for example, if  $\gamma = \theta$ ), then an asymptotic pivot is

$$Z_{n,1}(\gamma) := n(T_n - \gamma)^T V(\gamma)^{-1} (T_n - \gamma). \quad \begin{cases} \text{use wild device} \\ \text{for the proof} \end{cases}$$

The asymptotic distribution is the  $\chi^2$ -distribution with  $p$  degrees of freedom.

by the lemma

very smart

### 2<sup>nd</sup> asymptotic pivot

If, for all  $\theta$ , one has a consistent estimator  $\hat{V}_n$  of  $V_\theta$ , then an asymptotic pivot is

$$Z_{n,2}(\gamma) := n(T_n - \gamma)^T \hat{V}_n^{-1} (T_n - \gamma).$$

The asymptotic distribution is again the  $\chi^2$ -distribution with  $p$  degrees of freedom. This follows from Slutsky's lemma.

Proof:

$$\sqrt{n}(T_n - \gamma) \xrightarrow{\mathcal{D}} N(0, I)$$

$$\hat{V}_n \xrightarrow{\mathcal{P}} V_\theta$$

By Slutsky:

$$\hat{V}_n^{-1} \xrightarrow{\mathcal{D}} V_\theta^{-1}$$

Then by the lemma above:

$$Z_{n,2}(\gamma) \xrightarrow{\mathcal{D}} \chi_p^2$$

### Estimators of the asymptotic variance

o If  $\hat{\theta}_n$  is a consistent estimator of  $\theta$  and if  $\theta \mapsto V_\theta$  is continuous, one may insert  $\hat{V}_n := V_{\hat{\theta}_n}$ .

o If  $T_n = \gamma_n$  is the M-estimator of  $\gamma$ ,  $\gamma$  being the solution of  $P_\theta \psi_\gamma = 0$ , then (under regularity) the asymptotic covariance matrix is

$$V_\theta = M_\theta^{-1} J_\theta M_\theta^{-1},$$

where

$$\boxed{J_\theta = P_\theta \psi_\gamma \psi_\gamma^T,}$$

Let  $z_n \xrightarrow{\mathcal{D}} N(0, I)$   
 Lemma:  $\|z_n\|^2 \xrightarrow{\mathcal{D}} \chi_p^2$   
 Proof: Assume  $f = I$  (by wild device)  
 Let  $f$  be bounded & continuous: Then  $\mathbb{E}[f(z_n)] \rightarrow \mathbb{E}[f(I)]$   
 bounded & continuous:  
 $\rightarrow \mathbb{E}[f(z_n)^2] = \mathbb{E}[f(I)^2]$   
 The function then here has mult with the covariance matrix

and

$$\begin{aligned}\underline{M}_\theta &:= \widehat{\mathcal{R}}(\gamma) \\ &:= \frac{\partial}{\partial c^T} \dot{\mathcal{R}}(c) \Big|_{c=\gamma} \\ &= \frac{\partial}{\partial c^T} P_\theta \psi_c \Big|_{c=\gamma} \\ &= \underline{P_\theta \dot{\psi}_\gamma}.\end{aligned}$$

Then one may estimate  $J_\theta$  and  $M_\theta$  by

$$\boxed{\hat{J}_n := \hat{P}_n \psi_{\hat{\gamma}_n} \psi_{\hat{\gamma}_n}^T = \frac{1}{n} \sum_{i=1}^n \psi_{\hat{\gamma}_n}(X_i) \psi_{\hat{\gamma}_n}^T(X_i),}$$

and

$$\begin{aligned}\hat{M}_n &:= \ddot{\mathcal{R}}_n(\hat{\gamma}_n) \\ &= \hat{P}_n \dot{\psi}_{\hat{\gamma}_n} \\ &= \frac{1}{n} \sum_{i=1}^n \dot{\psi}_{\hat{\gamma}_n}(X_i),\end{aligned}$$

respectively. Under some regularity conditions,

$$\boxed{\hat{V}_n := \hat{M}_n^{-1} \hat{J}_n \hat{M}_n^{-1}.}$$

is a consistent estimator of  $V_\theta$ <sup>4</sup>.

## 14.8 Asymptotic pivot based on the MLE

Suppose that  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  has  $\Theta \subset \mathbb{R}^p$ , and that  $\mathcal{P}$  is dominated by some  $\sigma$ -finite measure  $\nu$ . Let  $p_\theta := dP_\theta/d\nu$  denote the densities, and let

$$\hat{\theta}_n := \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_\vartheta(X_i)$$

be the MLE.

Assume enough regularity such as existence of derivatives and interchanging integration and differentiation. Recall that  $\hat{\theta}_n$  is an M-estimator with loss function  $\rho_\vartheta = -\log p_\vartheta$ , and hence under regularity conditions,  $\psi_\vartheta = \dot{\rho}_\vartheta$  is minus

---

<sup>4</sup>From most algorithms used to compute the M-estimator  $\hat{\gamma}_n$ , one easily can obtain  $\hat{M}_n$  and  $\hat{J}_n$  as output. Recall e.g. that the Newton-Raphson algorithm is based on the iterations

$$\hat{\gamma}_{\text{new}} = \hat{\gamma}_{\text{old}} - \left( \ddot{\mathcal{R}}_n(\hat{\gamma}_{\text{old}}) \right)^{-1} \dot{\mathcal{R}}_n(\hat{\gamma}_{\text{old}}).$$

the score function  $s_\theta := \dot{p}_\theta/p_\theta$ . The asymptotic variance of the MLE is  $I^{-1}(\theta)$ , where  $I(\theta) := P_\theta s_\theta s_\theta^T$  is the Fisher information:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, I^{-1}(\theta)), \forall \theta$$

(see Section 14.4). Thus, in this case

$$\underline{Z_{n,1}(\theta)} = n(\hat{\theta}_n - \theta)^T I(\theta)(\hat{\theta}_n - \theta),$$

and, with  $\hat{I}_n$  being a consistent estimator of  $I(\theta)$

$$\boxed{Z_{n,2}(\theta) = n(\hat{\theta}_n - \theta)^T \hat{I}_n(\hat{\theta}_n - \theta).}$$

Note that one may take

$$\hat{I}_n := -\frac{1}{n} \sum_{i=1}^n \dot{s}_{\hat{\theta}_n}(X_i) = -\frac{\partial^2}{\partial \theta \partial \theta^T} \left. \frac{1}{n} \sum_{i=1}^n \log p_\theta(X_i) \right|_{\theta=\hat{\theta}_n}$$

as estimator of the Fisher information<sup>5</sup>.

### 3rd asymptotic pivot

Define now the twice log-likelihood ratio

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) := 2 \sum_{i=1}^n \left[ \log p_{\hat{\theta}_n}(X_i) - \log p_\theta(X_i) \right].$$

It turns out that the log-likelihood ratio is indeed an asymptotic pivot. A practical advantage is that it is self-normalizing: one does not need to explicitly estimate asymptotic (co-)variances.

**Lemma 14.8.1** Under regularity conditions,  $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta)$  is an asymptotic pivot for  $\theta$ . Its asymptotic distribution is again the  $\chi^2$ -distribution with  $p$  degrees of freedom:

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \xrightarrow{\mathcal{D}_\theta} \chi_p^2 \quad \forall \theta.$$

**Sketch of the proof.** We have by a two-term Taylor expansion

$$\left. \begin{aligned} 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) &= 2n\hat{P}_n \left[ \log p_{\hat{\theta}_n} - \log p_\theta \right] \\ &\approx 2n(\hat{\theta}_n - \theta)^T \hat{P}_n s_\theta + n(\hat{\theta}_n - \theta)^T \hat{P}_n \dot{s}_\theta(\hat{\theta}_n - \theta) \\ &\approx 2n(\hat{\theta}_n - \theta)^T \hat{P}_n s_\theta - n(\hat{\theta}_n - \theta)^T I(\theta)(\hat{\theta}_n - \theta), \end{aligned} \right\} \begin{aligned} \text{Given } \hat{\theta}_n - \theta &\approx \frac{1}{n} \sum_{i=1}^n \log p_{\hat{\theta}_n}(X_i) \\ &= I(\theta)^{-1} \hat{P}_n s_\theta \end{aligned}$$

where in the second step, we used  $\hat{P}_n \dot{s}_\theta \approx P_\theta \dot{s}_\theta = -I(\theta)$ . (If you like you may compare this two-term Taylor expansion with the one in the sketch of proof of Le Cam's 3rd Lemma (which is not part of the exam)). With the remainder terms of the two-term Taylor expansion being asymptotically negligible, we are

<sup>5</sup>In other words (as for general M-estimators), the algorithm (e.g. Newton Raphson) for calculating the maximum likelihood estimator  $\hat{\theta}_n$  generally also provides an estimator of the Fisher information as by-product.

mathematically speaking dealing with a situation as in Lemma 12.3.2 where the least squares estimator is studied<sup>6</sup>. The MLE  $\hat{\theta}_n$  is asymptotically linear with influence function  $l_\theta = I(\theta)^{-1}s_\theta$ :

$$\hat{\theta}_n - \theta = I(\theta)^{-1}\hat{P}_ns_\theta + o_{\mathbb{P}_\theta}(n^{-1/2}).$$

Hence,

$$\underline{2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta)} \approx n(\hat{P}_ns_\theta)^T I(\theta)^{-1}(\hat{P}_ns_\theta).$$

The result now follows from

$$\sqrt{n}\hat{P}_ns_\theta \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, I(\theta)).$$

insert  
in the  
above

□

## 14.9 MLE for the multinomial distribution

Let  $X_1, \dots, X_n$  be i.i.d. copies of  $X$ , where  $X \in \{1, \dots, k\}$  is a label, with

$$P_\theta(X = j) := \pi_j, \quad j = 1, \dots, k.$$

where the probabilities  $\pi_j$  are positive and add up to one:  $\sum_{j=1}^k \pi_j = 1$ , but are assumed to be otherwise unknown. Then there are  $p := k - 1$  unknown parameters, say  $\theta = (\pi_1, \dots, \pi_{k-1})$ . Define  $N_j := \#\{i : X_i = j\}$ . Note that  $(N_1, \dots, N_k)$  has a multinomial distribution with parameters  $n$  and  $(\pi_1, \dots, \pi_k)$ .

Given the constraint:  
 $\partial f = p-1$

**Lemma 14.9.1** For each  $j = 1, \dots, k$ , the MLE of  $\pi_j$  is

$$\hat{\pi}_j = \frac{N_j}{n}.$$

Not  
to know

**Proof.** The log-densities can be written as

$$\log p_\theta(x) = \sum_{j=1}^k \mathbf{1}\{x = j\} \log \pi_j,$$

so that

$$\sum_{i=1}^n \log p_\theta(X_i) = \sum_{j=1}^k N_j \log \pi_j.$$

Putting the derivatives with respect to  $\theta = (\pi_1, \dots, \pi_{k-1})$ , (with  $\pi_k = 1 - \sum_{j=1}^{k-1} \pi_j$ ) to zero gives,

$$\frac{N_j}{\hat{\pi}_j} - \frac{N_k}{\hat{\pi}_k} = 0.$$

each variable appears  
twice 1 time and 1 time  
here at k

Hence

$$\hat{\pi}_j = N_j \frac{\hat{\pi}_k}{N_k}, \quad j = 1, \dots, k,$$

<sup>6</sup> $\hat{P}_ns_\theta$  takes the role of  $X^T\epsilon/n$  and  $I(\theta)$  takes the role of  $X^TX/n$

and thus

$$1 = \sum_{j=1}^k \hat{\pi}_j = n \frac{\hat{\pi}_k}{N_k},$$

$\downarrow N_j$

yielding

$$\hat{\pi}_k = \frac{N_k}{n},$$

and hence

$$\hat{\pi}_j = \frac{N_j}{n}, \quad j = 1, \dots, k.$$

□

We now first calculate  $Z_{n,1}(\theta)$ . For that, we need to find the Fisher information  $I(\theta)$ .

**Lemma 14.9.2** *The Fisher information is*

$$I(\theta) = \begin{pmatrix} \frac{1}{\pi_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\pi_{k-1}} \end{pmatrix} + \frac{1}{\pi_k} \boldsymbol{\iota} \boldsymbol{\iota}^T,$$

<sup>7</sup>

where  $\boldsymbol{\iota}$  is the  $(k-1)$ -vector  $\boldsymbol{\iota} := (1, \dots, 1)^T$ .

**Proof.** We have

$$s_{\theta,j}(x) = \frac{1}{\pi_j} \mathbf{1}\{x = j\} - \frac{1}{\pi_k} \mathbf{1}\{x = k\}.$$

So

$$\begin{aligned} (I(\theta))_{j_1, j_2} &= E_\theta \left( \frac{1}{\pi_{j_1}} \mathbf{1}\{X = j_1\} - \frac{1}{\pi_k} \mathbf{1}\{X = k\} \right) \left( \frac{1}{\pi_{j_2}} \mathbf{1}\{X = j_2\} - \frac{1}{\pi_k} \mathbf{1}\{X = k\} \right) \\ &= \begin{cases} \frac{1}{\pi_k} & j_1 \neq j_2 \\ \frac{1}{\pi_j} + \frac{1}{\pi_k} & j_1 = j_2 = j \end{cases}. \end{aligned}$$

□

We thus find

$$\begin{aligned} Z_{n,1}(\theta) &= n(\hat{\theta}_n - \theta)^T I(\theta)(\hat{\theta}_n - \theta) \\ &= n \begin{pmatrix} \hat{\pi}_1 - \pi_1 \\ \vdots \\ \hat{\pi}_{k-1} - \pi_{k-1} \end{pmatrix}^T \left[ \begin{pmatrix} \frac{1}{\pi_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{\pi_{k-1}} \end{pmatrix} + \frac{1}{\pi_k} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} \right] \begin{pmatrix} \hat{\pi}_1 - \pi_1 \\ \vdots \\ \hat{\pi}_{k-1} - \pi_{k-1} \end{pmatrix} \\ &= n \sum_{j=1}^{k-1} \frac{(\hat{\pi}_j - \pi_j)^2}{\pi_j} + n \frac{1}{\pi_k} \left( \sum_{j=1}^{k-1} (\hat{\pi}_j - \pi_j) \right)^2 \end{aligned}$$

---

<sup>7</sup>To invert such a matrix, one may apply the formula  $(A + bb^T)^{-1} = A^{-1} - \frac{A^{-1}bb^TA^{-1}}{1+b^TA^{-1}b}$ .

$$\begin{aligned}
&= n \sum_{j=1}^k \frac{(\hat{\pi}_j - \pi_j)^2}{\pi_j} \\
&= \sum_{j=1}^k \frac{(N_j - n\pi_j)^2}{n\pi_j}.
\end{aligned}$$

This is called the Pearson's chi-square

$$\sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}.$$

A version of  $Z_{n,2}(\theta)$  is to replace, for  $j = 1, \dots, k$ ,  $\pi_j$  by  $\hat{\pi}_j$  in the expression for the Fisher information. This gives

$$Z_{n,2}(\theta) = \sum_{j=1}^k \frac{(N_j - n\pi_j)^2}{N_j}.$$

This is called the Pearson's chi-square

$$\sum \frac{(\text{observed} - \text{expected})^2}{\text{observed}}.$$

Finally, the log-likelihood ratio pivot is

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) = 2 \sum_{j=1}^k N_j \log \left( \frac{\hat{\pi}_j}{\pi_j} \right).$$

The approximation  $\log(1+x) \approx x - x^2/2$  shows that  $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) \approx Z_{n,2}(\theta)$ :

$$\begin{aligned}
2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta) &= -2 \sum_{j=1}^k N_j \log \left( 1 + \frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right) \\
&\approx -2 \sum_{j=1}^k N_j \left( \frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right) + \sum_{j=1}^k N_j \left( \frac{\pi_j - \hat{\pi}_j}{\hat{\pi}_j} \right)^2 \\
&= Z_{n,2}(\theta).
\end{aligned}$$

The three asymptotic pivots  $Z_{n,1}(\theta)$ ,  $Z_{n,2}(\theta)$  and  $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta)$  are each asymptotically  $\chi_{k-1}^2$ -distributed under  $\mathbb{P}_\theta$ .

## 14.10 Likelihood ratio tests

For the simple hypothesis

$$H_0 : \theta = \theta_0,$$

we can use  $2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta_0)$  as test statistic: reject  $H_0$  if

$$\boxed{2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\theta_0) > G_p^{-1}(1 - \alpha),}$$

where  $G_p$  is the distribution function of the  $\chi^2_p$ -distribution. This test has approximately level  $\alpha$  as was shown in Lemma 14.8.1.

Consider now the hypothesis

$$H_0 : R(\theta) = 0,$$

where

$$R(\theta) = \begin{pmatrix} R_1(\theta) \\ \vdots \\ R_q(\theta) \end{pmatrix}$$

are  $q$  restrictions on  $\theta$  (with  $R : \mathbb{R}^p \rightarrow \mathbb{R}^q$  a given function<sup>8</sup>).

Let  $\hat{\theta}_n$  be the unrestricted MLE, that is

$$\hat{\theta}_n = \arg \max_{\vartheta \in \Theta} \sum_{i=1}^n \log p_\vartheta(X_i).$$

Moreover, let  $\hat{\theta}_n^0$  be the restricted MLE, defined as

$$\hat{\theta}_n^0 = \arg \max_{\vartheta \in \Theta : R(\vartheta)=0} \sum_{i=1}^n \log p_\vartheta(X_i).$$

*restrict.  
↓ & all var / all directions*

Define the  $(q \times p)$ -matrix

$$\dot{R}(\theta) = \frac{\partial}{\partial \vartheta^T} R(\vartheta) \Big|_{\vartheta=\theta}.$$

*relax  
the restrictions  
a bit*

We assume  $\dot{R}(\theta)$  has rank  $q$ .

Let

$$\mathcal{L}_n(\hat{\theta}_n) - \mathcal{L}_n(\hat{\theta}_n^0) = \sum_{i=1}^n \left[ \log p_{\hat{\theta}_n}(X_i) - \log p_{\hat{\theta}_n^0}(X_i) \right]$$

be the log-likelihood ratio for testing  $H_0 : R(\theta) = 0$ .

**Lemma 14.10.1** Under regularity conditions, and if  $H_0 : R(\theta) = 0$  holds, we have

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \xrightarrow{\mathcal{D}_\theta} \chi_q^2.$$

**Sketch of the proof.** Let

$$\mathbf{Z}_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n s_\theta(X_i).$$

As in the sketch of the proof of Lemma 14.8.1, we can use a two-term Taylor expansion to show for any sequence  $\vartheta_n$  satisfying  $\vartheta_n = \theta + \mathcal{O}_{\mathbf{P}_\theta}(n^{-1/2})$ , that

$$2 \sum_{i=1}^n \left[ \log p_{\vartheta_n}(X_i) - \log p_\theta(X_i) \right]$$

<sup>8</sup>The notation  $R$  is used here for the restrictions. It has nothing to do with the risk  $\mathcal{R}$

$$= \underline{2\sqrt{n}(\vartheta_n - \theta)^T \mathbf{Z}_n - n(\vartheta_n - \theta)^2 I(\theta)(\vartheta_n - \theta)} + o_{\mathbf{P}_\theta}(1).$$

Here, we also again use that  $\sum_{i=1}^n \dot{s}_{\vartheta_n}(X_i)/n = -I(\theta) + o_{\mathbf{P}_\theta}(1)$ . Moreover, by a one-term Taylor expansion, and invoking that  $R(\theta) = 0$ ,

$$R(\vartheta_n) = \dot{R}(\theta)(\vartheta_n - \theta) + o_{\mathbf{P}_\theta}(n^{-1/2}).$$

Insert Corollary 12.4.1 with  $z := \mathbf{Z}_n$ ,  $B := \dot{R}(\theta)$ , and  $V = I(\theta)$ . This gives

$$\begin{aligned} & 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \\ &= 2 \sum_{i=1}^n \left[ \log p_{\hat{\theta}_n}(X_i) - \log p_\theta(X_i) \right] - 2 \sum_{i=1}^n \left[ \log p_{\hat{\theta}_n^0}(X_i) - \log p_\theta(X_i) \right] \\ &= \mathbf{Z}_n^T I(\theta)^{-1} \dot{R}^T(\theta) \left( \dot{R}(\theta) I(\theta)^{-1} \dot{R}(\theta)^T \right)^{-1} \dot{R}(\theta) I(\theta)^{-1} \mathbf{Z}_n + o_{\mathbf{P}_\theta}(1) \\ &:= \mathbf{Y}_n^T W^{-1} \mathbf{Y}_n + o_{\mathbf{P}_\theta}(1), \end{aligned}$$

where  $\mathbf{Y}_n$  is the  $q$ -vector

$$\mathbf{Y}_n := \dot{R}(\theta) I(\theta)^{-1} \mathbf{Z}_n,$$

and where  $W$  is the  $(q \times q)$ -matrix

$$W := \dot{R}(\theta) I(\theta)^{-1} \dot{R}(\theta)^T.$$

We know that

$$\mathbf{Z}_n \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, I(\theta)).$$

Hence

$$\mathbf{Y}_n \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(0, W),$$

so that

$$\mathbf{Y}_n^T W^{-1} \mathbf{Y}_n \xrightarrow{\mathcal{D}_\theta} \chi_q^2.$$

□

**Corollary 14.10.1** *From the sketch of the proof of Lemma 14.10.1, one sees that moreover (under regularity),*

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n(\hat{\theta}_n - \hat{\theta}_n^0)^T I(\theta)(\hat{\theta}_n - \hat{\theta}_n^0),$$

and also

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n(\hat{\theta}_n - \hat{\theta}_n^0)^T I(\hat{\theta}_n^0)(\hat{\theta}_n - \hat{\theta}_n^0).$$

## 14.11 Contingency tables

Let  $X$  be a bivariate label, say  $X \in \{(j, k) : j = 1, \dots, r, k = 1, \dots, s\}$ . For example, the first index may correspond to sex ( $r = 2$ ) and the second index to the color of the eyes ( $s = 3$ ). The probability of the combination  $(j, k)$  is

$$\pi_{j,k} := P_\theta(X = (j, k)).$$

Let  $X_1, \dots, X_n$  be i.i.d. copies of  $X$ , and

$$N_{j,k} := \#\{X_i = (j, k)\}.$$

From Section 14.9, we know that the (unrestricted) MLE of  $\pi_{j,k}$  is equal to

$$\hat{\pi}_{j,k} := \frac{N_{j,k}}{n}.$$

We now want to test whether the two labels are independent. The null-hypothesis is

$$H_0 : \pi_{j,k} = (\pi_{j,+}) \times (\pi_{+,k}) \quad \forall (j, k).$$

Here

$$\pi_{j,+} := \sum_{k=1}^s \pi_{j,k}, \quad \pi_{+,k} := \sum_{j=1}^r \pi_{j,k}.$$

One may check that the restricted MLE is

$$\hat{\pi}_{j,k}^0 = (\hat{\pi}_{j,+}) \times (\hat{\pi}_{+,k}),$$

where

$$\hat{\pi}_{j,+} := \sum_{k=1}^s \hat{\pi}_{j,k}, \quad \hat{\pi}_{+,k} := \sum_{j=1}^r \hat{\pi}_{j,k}.$$

The log-likelihood ratio test statistic is thus

$$\begin{aligned} 2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) &= 2 \sum_{j=1}^r \sum_{k=1}^s N_{j,k} \left[ \log \left( \frac{N_{j,k}}{n} \right) - \log \left( \frac{N_{j,+} N_{+,k}}{n^2} \right) \right] \\ &= 2 \sum_{j=1}^r \sum_{k=1}^s N_{j,k} \log \left( \frac{n N_{j,k}}{N_{j,+} N_{+,k}} \right). \end{aligned}$$

Its approximation as given in Corollary 14.10.1 is

$$2\mathcal{L}_n(\hat{\theta}_n) - 2\mathcal{L}_n(\hat{\theta}_n^0) \approx n \sum_{j=1}^r \sum_{k=1}^s \frac{(N_{j,k} - N_{j,+} N_{+,k}/n)^2}{N_{j,+} N_{+,k}}.$$

This is Pearson's chi-squared test statistic for testing independence.

To find out what the value of  $q$  is in this example, we first observe that the unrestricted case has  $p = rs - 1$  free parameters. Under the null-hypothesis,

there remain  $(r-1)+(s-1)$  free parameters. Hence, the number of restrictions is

$$q = \left( rs - 1 \right) - \left( (r-1) + (s-1) \right) = (r-1)(s-1).$$

Thus, under  $H_0 : \pi_{j,k} = (\pi_{j,+}) \times (\pi_{+,k}) \forall (j, k)$ , we have

$$n \sum_{j=1}^r \sum_{k=1}^s \frac{(N_{j,k} - N_{j,+}N_{+,k}/n)^2}{N_{j,+}N_{+,k}} \xrightarrow{\mathcal{D}_\theta} \chi^2_{(r-1)(s-1)}.$$



# Chapter 15

## Abstract asymptotics \*

In Subsection 2.4.1 we discussed so-called plug-in estimators. The idea is that an estimator  $\hat{\gamma}_n$  can (typically) be written as some functional  $Q$  of the empirical distribution  $\hat{P}_n$ . When  $P$  is the “true” distribution and  $\gamma = Q(P)$ , then the point is studying how close  $\hat{\gamma}_n(\hat{P}_n)$  is to  $\gamma = Q(P)$  when  $\hat{P}_n$  is close to  $P$ . This is the topic of the first part of this chapter.

In Section 5.5 we obtained the Cramér Rao lower bound. A somewhat disappointing result was that it can only be reached within exponential families (see Lemma 5.6.1). We have also seen in Section 14.4 that the MLE is *asymptotically unbiased* and reaches *asymptotically* the CRLB (its asymptotic covariance matrix is  $I(\theta)^{-1}$ , where  $I(\theta)$  is the Fisher information for estimating  $\theta$ ). The topic of the second part of this chapter is to show (for the one-dimensional case for simplicity) that  $I(\theta)$  is indeed asymptotically the efficient variance.

One may now wonder why the inverse of  $I(\theta)$  is there. Think of it in this way. The Fisher information was obtained by looking at derivatives of the map

$$\theta \mapsto \log p_\theta.$$

But what actually plays a role in the inverse map

$$P \mapsto \theta$$

or, in case  $\gamma = g(\theta)$  is the parameter of interest, the map

$$P \mapsto \gamma.$$

Then remember that the derivative of the inverse of a function (say  $f : \mathbb{R} \mapsto \mathbb{R}$ ) is the inverse of its derivative. In our case the mapping  $P \mapsto \gamma$  is rather abstract, so studying its derivatives, as done in the first part of this chapter, requires some new notions.

## 15.1 Plug-in estimators \*

When  $\mathcal{X}$  is Euclidean space, one can define the distribution function  $F(x) := P_\theta(X \leq x)$  and the empirical distribution function

$$\hat{F}_n(x) = \frac{1}{n} \#\{X_i \leq x, 1 \leq i \leq n\}.$$

This is the distribution function of a probability measure that puts mass  $1/n$  at each observation. For general  $\mathcal{X}$ , we define likewise the empirical distribution  $\hat{P}_n$  as the distribution that puts mass  $1/n$  at each observation, i.e., more formally

$$\hat{P}_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i},$$

where  $\delta_x$  is a point mass at  $x$ . Thus, for (measurable) sets  $A \subset \mathcal{X}$ ,

$$\hat{P}_n(A) = \frac{1}{n} \#\{X_i \in A, 1 \leq i \leq n\}.$$

For (measurable) functions  $f : \mathcal{X} \rightarrow \mathbb{R}^r$  (for some  $r \in \mathbb{N}$ ), we write, as in previous sections,

$$\hat{P}_n f := \frac{1}{n} \sum_{i=1}^n f(X_i) = \int f d\hat{P}_n.$$

Thus, for sets,

$$\hat{P}_n(A) = \hat{P}_n \mathbf{1}_A.$$

Again, as in previous sections, we use the same notations for expectations under  $P_\theta$ :

$$P_\theta f := E_\theta f(X) = \int f dP_\theta,$$

so that

$$P_\theta(A) = P_\theta \mathbf{1}_A.$$

The parameter of interest is denoted as

$$\gamma = g(\theta) \in \mathbb{R}^p.$$

It can often be written in the form

$$\gamma = Q(P_\theta),$$

where  $Q$  is some functional on (a superset of) the model class  $\mathcal{P}$ . Assuming  $Q$  is also defined at the empirical measure  $\hat{P}_n$ , the plug-in estimator of  $\gamma$  is now

$$T_n := Q(\hat{P}_n).$$

Conversely,

**Definition 15.1.1** If a statistic  $T_n$  can be written as  $T_n = Q(\hat{P}_n)$ , then it is called a Fisher-consistent estimator of  $\gamma = g(\theta)$ , if  $Q(P_\theta) = g(\theta)$  for all  $\theta \in \Theta$ .

We will also encounter modifications, where

$$T_n = Q_n(\hat{P}_n),$$

and for  $n$  large,

$$Q_n(P_\theta) \approx Q(P_\theta) = g(\theta).$$

**Example 15.1.1 Plug-in estimator of (functions of) the mean**

Consider a given  $f : \mathcal{X} \rightarrow \mathbb{R}^r$  and a given  $h : \mathbb{R}^r \rightarrow \mathbb{R}^p$ . Let  $\gamma := h(P_\theta f)$ . The plug-in estimator is then  $T_n = h(\hat{P}_n f)$ .

**Example 15.1.2 M- and Z-estimators are plug-in estimators**

The M-estimator

$$\hat{\gamma}_n := \arg \min_{c \in \Gamma} \hat{P}_n \rho_c$$

is a plug-in estimator of

$$\gamma = \arg \min_{c \in \Gamma} P_\theta \rho_c.$$

Similarly, the Z-estimator  $\hat{\gamma}_n$  as solution of

$$\hat{P}_n \psi_c \Big|_{c=\hat{\gamma}_n} = 0$$

is a plug-in estimator of

$$P_\theta \psi_c \Big|_{c=\gamma} = 0.$$

**Example 15.1.3 The  $\alpha$ -trimmed mean as plug-in estimator**

Let  $\mathcal{X} = \mathbb{R}$  and consider the  $\alpha$ -trimmed mean

$$T_n := \frac{1}{n - 2[n\alpha]} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} X_{(i)}.$$

What is its theoretical counterpart? Because the  $i$ -th order statistic  $X_{(i)}$  can be written as

$$X_{(i)} = \hat{F}_n^{-1}(i/n),$$

and in fact

$$X_{(i)} = \hat{F}_n^{-1}(u), \quad i/n \leq u < (i+1)/n,$$

we may write, for  $\alpha_n := [n\alpha]/n$ ,

$$\begin{aligned} T_n &= \frac{n}{n - 2[n\alpha]} \frac{1}{n} \sum_{i=[n\alpha]+1}^{n-[n\alpha]} \hat{F}_n^{-1}(i/n) \\ &= \frac{1}{1 - 2\alpha_n} \int_{\alpha_n + 1/n}^{1 - \alpha_n} \hat{F}_n^{-1}(u) du := Q_n(\hat{P}_n). \end{aligned}$$

Replacing  $\hat{F}_n$  by  $F$  gives

$$\begin{aligned} Q_n(F) &= \frac{1}{1-2\alpha_n} \int_{\alpha_n+1/n}^{1-\alpha_n} F^{-1}(u) du \\ &\approx \frac{1}{1-2\alpha} \int_{\alpha}^{1-\alpha} F^{-1}(u) du \\ &= \frac{1}{1-2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x) \\ &:= Q(P_\theta). \end{aligned}$$

#### Example 15.1.4 Histogram as plug-in estimator of a density

Let  $\mathcal{X} = \mathbb{R}$ , and suppose  $X$  has density  $f$  w.r.t., Lebesgue measure. Suppose  $f$  is the parameter of interest. We may write

$$f(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h}.$$

Replacing  $F$  by  $\hat{F}_n$  here does not make sense. Thus, this is an example where  $Q(P) = f$  is only well defined for distributions  $P$  that have a density  $f$ . We may however slightly extend the plug-in idea, by using the estimator

$$\hat{f}_n(x) := \frac{\hat{F}_n(x+h_n) - \hat{F}_n(x-h_n)}{2h_n} := Q_n(\hat{P}_n),$$

with  $h_n$  “small” ( $h_n \rightarrow 0$  as  $n \rightarrow \infty$ ).

## 15.2 Consistency of plug-in estimators \*

We first present the uniform convergence of the empirical distribution function to the theoretical one.

Such uniform convergence results hold also in much more general settings (see also (14.2) in the proof of consistency for M-estimators).

**Theorem 15.2.1 (Glivenko-Cantelli)** Let  $\mathcal{X} = \mathbb{R}$ . We have

$$\sup_x |\hat{F}_n(x) - F(x)| \rightarrow 0, \quad \mathbb{P}_\theta - \text{a.s.}$$

**Proof.** We know that by the law of large numbers, for all  $x$

$$|\hat{F}_n(x) - F(x)| \rightarrow 0, \quad \mathbb{P}_\theta - \text{a.s.},$$

so also for all finite collection  $a_1, \dots, a_N$ ,

$$\max_{1 \leq j \leq N} |\hat{F}_n(a_j) - F(a_j)| \rightarrow 0, \quad \mathbb{P}_\theta - \text{a.s.}$$

Let  $\epsilon > 0$  be arbitrary, and take  $a_0 < a_1 < \dots < a_{N-1} < a_N$  in such a way that

$$F(a_j) - F(a_{j-1}) \leq \epsilon, \quad j = 1, \dots, N$$

where  $F(a_0) := 0$  and  $F(a_N) := 1$ . Then, when  $x \in (a_{j-1}, a_j]$ ,

$$\hat{F}_n(x) - F(x) \leq \hat{F}_n(a_j) - F(a_{j-1}) \leq F_n(a_j) - F(a_j) + \epsilon,$$

and

$$\hat{F}_n(x) - F(x) \geq \hat{F}_n(a_{j-1}) - F(a_j) \geq \hat{F}_n(a_{j-1}) - F(a_{j-1}) - \epsilon,$$

so

$$\sup_x |\hat{F}_n(x) - F(x)| \leq \max_{1 \leq j \leq N} |\hat{F}_n(a_j) - F(a_j)| + \epsilon \rightarrow \epsilon, \text{ I}\mathbb{P}_{\theta}\text{-a.s..}$$

□

#### Example 15.2.1 Consistency of the sample median

Let  $\mathcal{X} = \mathbb{R}$  and let  $F$  be the distribution function of  $X$ . We consider estimating the median  $\gamma := F^{-1}(1/2)$ . We assume  $F$  to continuous and strictly increasing. The sample median is

$$T_n := \hat{F}_n^{-1}(1/2) := \begin{cases} X_{((n+1)/2)} & n \text{ odd} \\ [X_{(n/2)} + X_{(n/2+1)}]/2 & n \text{ even} \end{cases}.$$

So

$$\hat{F}_n(T_n) = \frac{1}{2} + \begin{cases} 1/(2n) & n \text{ odd} \\ 0 & n \text{ even} \end{cases}.$$

It follows that

$$\begin{aligned} |F(T_n) - F(\gamma)| &\leq |\hat{F}_n(T_n) - F(T_n)| + |\hat{F}_n(T_n) - F(\gamma)| \\ &= |\hat{F}_n(T_n) - F(T_n)| + |\hat{F}_n(T_n) - \frac{1}{2}| \\ &\leq |\hat{F}_n(T_n) - F(T_n)| + \frac{1}{2n} \rightarrow 0, \text{ I}\mathbb{P}_{\theta}\text{-a.s..} \end{aligned}$$

So  $\hat{F}_n^{-1}(1/2) = T_n \rightarrow \gamma = F^{-1}(1/2)$ ,  $\text{I}\mathbb{P}_{\theta}\text{-a.s.}$ , i.e., the sample median is a consistent estimator of the population median.

### 15.3 Asymptotic normality of plug-in estimators $\star$

Let  $\gamma := Q(P) \in \mathbb{R}^p$  be the parameter of interest. The idea in this subsection is to apply a  $\delta$ -method, but now in a nonparametric framework. The parametric  $\delta$ -method says that if  $\hat{\theta}_n$  is an asymptotically linear estimator of  $\theta \in \mathbb{R}^p$ , and if  $\gamma = g(\theta)$  is some function of the parameter  $\theta$ , with  $g$  being differentiable at  $\theta$ , then  $\hat{\gamma}$  is an asymptotically linear estimator of  $\gamma$ . Now, we write  $\gamma = Q(P)$  as a function of the probability measure  $P$  (with  $P = P_\theta$ , so that  $g(\theta) = Q(P_\theta)$ ). We let  $P$  play the role of  $\theta$ , i.e., we use the probability measures themselves as parameterization of  $\mathcal{P}$ . We then have to redefine differentiability in an abstract setting, namely we differentiate w.r.t.  $P$ .

#### Definition 15.3.1

- The influence function of  $Q$  at  $P$  is

$$l_P(x) := \lim_{\epsilon \downarrow 0} \frac{Q((1-\epsilon)P + \epsilon\delta_x) - Q(P)}{\epsilon}, \quad x \in \mathcal{X},$$

whenever the limit exists.

- o The map  $Q$  is called Gâteaux differentiable at  $P$  if for all probability measures  $\tilde{P}$ , we have

$$\lim_{\epsilon \downarrow 0} \frac{Q((1-\epsilon)P + \epsilon\tilde{P}) - Q(P)}{\epsilon} = E_{\tilde{P}}l_P(X).$$

- o Let  $d$  be some (pseudo-)metric on the space of probability measures. The map  $Q$  is called Fréchet differentiable at  $P$ , with respect to the metric  $d$ , if

$$Q(\tilde{P}) - Q(P) = E_{\tilde{P}}l_P(X) + o(d(\tilde{P}, P)).$$

**Remark 1** In line with the notation introduced previously, we write for a function  $f : \mathcal{X} \rightarrow \mathbb{R}^r$  and a probability measure  $\tilde{P}$  on  $\mathcal{X}$

$$\tilde{P}f := E_{\tilde{P}}f(X).$$

**Remark 2** If  $Q$  is Fréchet or Gâteaux differentiable at  $P$ , then

$$Pl_P(:= E_Pl_P(X)) = 0.$$

**Remark 3** If  $Q$  is Fréchet differentiable at  $P$ , and if moreover

$$d((1-\epsilon)P + \epsilon\tilde{P}, P) = o(\epsilon), \quad \epsilon \downarrow 0,$$

then  $Q$  is Gâteaux differentiable at  $P$ :

$$\begin{aligned} Q((1-\epsilon)P + \epsilon\tilde{P}) - Q(P) &= ((1-\epsilon)P + \epsilon\tilde{P})l_P + o(\epsilon) \\ &= \epsilon\tilde{P}l_P + o(\epsilon). \end{aligned}$$

The following result says that Fréchet differentiable functionals are generally asymptotically linear.

**Lemma 15.3.1** Suppose that  $Q$  is Fréchet differentiable at  $P$  with influence function  $l_P$ , and that

$$d(\hat{P}_n, P) = \mathcal{O}_{\mathbf{P}}(n^{-1/2}). \tag{15.1}$$

Then

$$Q(\hat{P}_n) - Q(P) = \hat{P}_n l_P + o_{\mathbf{P}}(n^{-1/2}).$$

**Proof.** This follows immediately from the definition of Fréchet differentiability.

□

**Corollary 15.3.1** Assume the conditions of Lemma 15.3.1, with influence function  $l_P$  satisfying  $V_P := Pl_P l_P^T < \infty$ . Then

$$\sqrt{n}(Q(\hat{P}_n) - Q(P)) \xrightarrow{\mathcal{D}_P} \mathcal{N}(0, V_P).$$

**An example where (15.1) holds**

Suppose  $\mathcal{X} = \mathbb{R}$  and that we take

$$d(\tilde{P}, P) := \sup_x |\tilde{F}(x) - F(x)|.$$

Then indeed  $d(\hat{P}_n, P) = O_P(n^{-1/2})$ . This follows from Donsker's theorem, which we state here without proof:

**Theorem 15.3.1** (*Donsker's theorem*) Suppose  $F$  is continuous. Then

$$\sup_x \sqrt{n} |\hat{F}_n(x) - F(x)| \xrightarrow{\mathcal{D}} Z,$$

where the random variable  $Z$  has distribution function

$$G(z) = 1 - 2 \sum_{j=1}^{\infty} (-1)^{j+1} \exp[-2j^2 z^2], \quad z \geq 0.$$



Fréchet differentiability is generally quite hard to prove, and often not even true. We will sketch how to establish Gâteaux differentiability in two examples.

**Example 15.3.1 Asymptotic linearity of the Z-estimator**

We consider again the asymptotic linearity of the Z-estimator. Throughout in this example, we assume enough regularity. Let  $\gamma$  be defined by the equation

$$P\psi_\gamma = 0.$$

Let  $P_\epsilon := (1 - \epsilon)P + \epsilon\tilde{P}$ , and let  $\gamma_\epsilon$  be a solution of the equation

$$P_\epsilon\psi_{\gamma_\epsilon} = 0.$$

We assume that as  $\epsilon \downarrow 0$ , also  $\gamma_\epsilon \rightarrow \gamma$ . It holds that

$$(1 - \epsilon)P\psi_{\gamma_\epsilon} + \epsilon\tilde{P}\psi_{\gamma_\epsilon} = 0,$$

so

$$P\psi_{\gamma_\epsilon} + \epsilon(\tilde{P} - P)\psi_{\gamma_\epsilon} = 0,$$

and hence

$$P(\psi_{\gamma_\epsilon} - \psi_\gamma) + \epsilon(\tilde{P} - P)\psi_{\gamma_\epsilon} = 0.$$

Assuming differentiability of  $c \mapsto P\psi_c$ , we obtain

$$\begin{aligned} P(\psi_{\gamma_\epsilon} - \psi_\gamma) &= \left( \frac{\partial}{\partial c^T} P\psi_c \Big|_{c=\gamma} \right) (\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|) \\ &:= M_P(\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|). \end{aligned}$$

Moreover, again under regularity

$$\begin{aligned} (\tilde{P} - P)\psi_{\gamma_\epsilon} &= (\tilde{P} - P)\psi_\gamma + (\tilde{P} - P)(\psi_{\gamma_\epsilon} - \psi_\gamma) \\ &= (\tilde{P} - P)\psi_\gamma + o(1) = \tilde{P}\psi_\gamma + o(1). \end{aligned}$$

It follows that

$$M_P(\gamma_\epsilon - \gamma) + o(|\gamma_\epsilon - \gamma|) + \epsilon(\tilde{P} - P)\psi_\gamma + o(\epsilon) = 0,$$

or, assuming  $M_P$  to be invertible,

$$(\gamma_\epsilon - \gamma)(1 + o(1)) = -\epsilon M_P^{-1}\tilde{P}\psi_\gamma + o(\epsilon),$$

which gives

$$\frac{\gamma_\epsilon - \gamma}{\epsilon} \rightarrow -M_P^{-1}\tilde{P}\psi_\gamma.$$

The influence function is thus (as already seen in Subsection 14.3)

$$l_P = -M_P^{-1}\psi_\gamma.$$

**Example 15.3.2 Asymptotic linearity of the  $\alpha$ -trimmed mean**  
The  $\alpha$ -trimmed mean is a plug-in estimator of

$$\gamma := Q(P) = \frac{1}{1-2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} x dF(x).$$

Using partial integration, may write this as

$$(1-2\alpha)\gamma = (1-\alpha)F^{-1}(1-\alpha) - \alpha F^{-1}(\alpha) - \int_\alpha^{1-\alpha} v dF^{-1}(v).$$

The influence function of the quantile  $F^{-1}(v)$  is

$$q_v(x) = -\frac{1}{f(F^{-1}(v))} \left( \mathbb{I}\{x \leq F^{-1}(v)\} - v \right)$$

(see Example 14.5.1), i.e., for the distribution  $P_\epsilon = (1-\epsilon)P + \epsilon\tilde{P}$ , with distribution function  $F_\epsilon = (1-\epsilon)F + \epsilon\tilde{F}$ , we have

$$\begin{aligned} \lim_{\epsilon \downarrow 0} \frac{F_\epsilon^{-1}(v) - F^{-1}(v)}{\epsilon} &= \tilde{P}q_v \\ &= -\frac{1}{f(F^{-1}(v))} \left( \tilde{F}(F^{-1}(v)) - v \right). \end{aligned}$$

Hence, for  $P_\epsilon = (1-\epsilon)P + \epsilon\tilde{P}$ ,

$$(1-2\alpha) \lim_{\epsilon \downarrow 0} \frac{Q((1-\epsilon)P + \epsilon\tilde{P}) - Q(P)}{\epsilon}$$

$$\begin{aligned} &= (1-\alpha)\tilde{P}q_{1-\alpha} - \alpha\tilde{P}q_\alpha - \int_\alpha^{1-\alpha} v d\tilde{P}q_v \\ &= \int_\alpha^{1-\alpha} \frac{1}{f(F^{-1}(v))} \left( \tilde{F}(F^{-1}(v)) - v \right) dv \\ &= \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \frac{1}{f(u)} \left( \tilde{F}(u) - F(u) \right) dF(u) \\ &= \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \left( \tilde{F}(u) - F(u) \right) du \\ &= (1-2\alpha)\tilde{P}l_P, \end{aligned}$$

where

$$l_P(x) = -\frac{1}{1-2\alpha} \int_{F^{-1}(\alpha)}^{F^{-1}(1-\alpha)} \left( \mathbb{1}\{x \leq u\} - F(u) \right) du.$$

We conclude that, under regularity conditions, the  $\alpha$ -trimmed mean is asymptotically linear with the above influence function  $l_P$ , and hence asymptotically normal with asymptotic variance  $P l_P^2$ .

## 15.4 Asymptotic Cramer Rao lower bound \*

Let  $X$  have distribution  $P \in \{P_\theta : \theta \in \Theta\}$ . We assume for simplicity that  $\Theta \subset \mathbb{R}$  and that  $\theta$  is the parameter of interest. Let  $T_n$  be an estimator of  $\theta$ .

Throughout this section, we take certain, sometimes unspecified, regularity conditions for granted.



In particular, we assume that  $\mathcal{P}$  is dominated by some  $\sigma$ -finite measure  $\nu$ , and that the Fisher-information

$$I(\theta) := E_\theta s_\theta^2(X)$$

exists for all  $\theta$ . Here,  $s_\theta$  is the score function

$$s_\theta := \frac{d}{d\theta} \log p_\theta = \dot{p}_\theta / p_\theta,$$

with  $p_\theta := dP_\theta / d\nu$ .

Recall now that if  $T_n$  is an unbiased estimator of  $\theta$ , then by the Cramer Rao lower bound,  $1/I(\theta)$  is a lower bound for its variance (under regularity Conditions I and II, see Section 5.5).

**Definition 15.4.1** Suppose that

$$\sqrt{n}(T_n - \theta) \xrightarrow{\mathcal{D}_\theta} \mathcal{N}(b_\theta, V_\theta), \quad \forall \theta.$$

Then  $b_\theta$  is called the asymptotic bias, and  $V_\theta$  the asymptotic variance. The estimator  $T_n$  is called asymptotically unbiased if  $b_\theta = 0$  for all  $\theta$ . If  $T_n$  is asymptotically unbiased and moreover  $V_\theta = 1/I(\theta)$  for all  $\theta$ , and some regularity conditions holds, then  $T_n$  is called asymptotically efficient.

**Remark 1** The assumptions in the above definition, are for all  $\theta$ . Clearly, if one only looks at one fixed given  $\theta_0$ , it is easy to construct a super-efficient estimator, namely  $T_n = \theta_0$ . More generally, to avoid this kind of super-efficiency, one does not only require conditions to hold for all  $\theta$ , but in fact uniformly in  $\theta$ , or for all sequences  $\{\theta_n\}$ . The regularity one needs here involves the idea that one actually needs to allow for sequences  $\theta_n$  the form  $\theta_n = \theta + h/\sqrt{n}$ . In fact, the regularity requirement is that also, for all  $h$ ,

$$\sqrt{n}(T_n - \theta_n) \xrightarrow{\mathcal{D}_{\theta_n}} \mathcal{N}(0, V_\theta).$$

To make all this mathematically precise is quite involved. We refer to van der Vaart (1998). A glimpse is given in Le Cam's 3<sup>rd</sup> Lemma, see the next subsection.

**Remark 2** Note that when  $\theta = \theta_n$  is allowed to change with  $n$ , this means that distribution of  $X_i$  can change with  $n$ , and hence  $X_i$  can change with  $n$ . Instead of regarding the sample  $X_1, \dots, X_n$  are the first  $n$  of an infinite sequence, we now consider for each  $n$  a new sample, say  $X_{1,1}, \dots, X_{n,n}$ .

**Remark 3** We have seen that the MLE  $\hat{\theta}_n$  generally is indeed asymptotically unbiased with asymptotic variance  $V_\theta$  equal to  $1/I(\theta)$ , i.e., under regularity assumptions, the MLE is asymptotically efficient.

For asymptotically linear estimators, with influence function  $l_\theta$ , one has asymptotic variance  $V_\theta = E_\theta l_\theta^2(X)$ . The next lemma indicates that generally  $1/I(\theta)$  is indeed a lower bound for the asymptotic variance.

**Lemma 15.4.1** *Suppose asymptotic linearity:*

$$T_n - \theta = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbf{P}_\theta}(n^{-1/2}),$$

where  $E_\theta l_\theta(X) = 0$ ,  $E_\theta l_\theta^2(X) := V_\theta < \infty$ . Assume moreover that

$$E_\theta l_\theta(X) s_\theta(X) = 1. \quad (15.2)$$

Then

$$V_\theta \geq \frac{1}{I(\theta)}.$$

**Proof.** This follows from the Cauchy-Schwarz inequality:

$$\begin{aligned} 1 &= |\text{cov}_\theta(l_\theta(X), s_\theta(X))|^2 \\ &\leq \text{var}_\theta(l_\theta(X)) \text{var}_\theta(s_\theta(X)) = V_\theta I(\theta). \end{aligned}$$

□

It may look like a coincidence when in a special case, equality (15.2) indeed holds. But actually, it is true in quite a few cases. This may at first seem like magic.

We consider two examples. To simplify the expressions, we again write short-hand

$$P_\theta f := E_\theta f(X).$$

**Example 15.4.1 Equation (15.2) for Z-estimators**

This example examines the Z-estimator of  $\theta$ . Then we have, for  $P = P_\theta$ ,

$$P\psi_\theta = 0.$$

The influence function is

$$l_\theta = -\psi_\theta/M_\theta,$$

where

$$M_\theta := \frac{d}{d\theta} P\psi_\theta.$$

Under regularity, we have

$$M_\theta = P\dot{\psi}_\theta = \int \dot{\psi}_\theta p_\theta d\nu, \quad \dot{\psi}_\theta = \frac{d}{d\theta} \psi_\theta.$$

We may also write

$$M_\theta = - \int \psi_\theta \dot{p}_\theta d\nu, \quad \dot{p}_\theta = \frac{d}{d\theta} p_\theta.$$

This follows from the chain rule

$$\frac{d}{d\theta} \psi_\theta p_\theta = \dot{\psi}_\theta p_\theta + \psi_\theta \dot{p}_\theta,$$

and (under regularity)

$$\int \frac{d}{d\theta} \psi_\theta p_\theta d\nu = \frac{d}{d\theta} \int \psi_\theta p_\theta d\nu = \frac{d}{d\theta} P\psi_\theta = \frac{d}{d\theta} 0 = 0.$$

Thus

$$Pl_\theta s_\theta = -M_\theta^{-1} P\psi_\theta s_\theta = -M_\theta^{-1} \int \psi_\theta \dot{p}_\theta d\nu = 1,$$

that is, (15.2) holds.

#### Example 15.4.2 Equation (15.2) for plug-in estimators

We consider now the plug-in estimator  $Q(\hat{P}_n)$ . Suppose that  $Q$  is Fisher consistent (i.e.,  $Q(P_\theta) = \theta$  for all  $\theta$ ). Assume moreover that  $Q$  is Fréchet differentiable with respect to the metric  $d$ , at all  $P_\theta$ , and that

$$d(P_{\tilde{\theta}}, P_\theta) = \mathcal{O}(|\tilde{\theta} - \theta|).$$

Then, by the definition of Fréchet differentiability

$$h = Q(P_{\theta+h}) - Q(P_\theta) = P_{\theta+h} l_\theta + o(|h|) = (P_{\theta+h} - P_\theta) l_\theta + o(|h|),$$

or, as  $h \rightarrow 0$ ,

$$\begin{aligned} 1 &= \frac{(P_{\theta+h} - P_\theta) l_\theta}{h} + o(1) = \frac{\int l_\theta (p_{\theta+h} - p_\theta) d\nu}{h} + o(1) \\ &\rightarrow \int l_\theta \dot{p}_\theta d\nu = P_\theta (l_\theta s_\theta). \end{aligned}$$

So (15.2) holds.

## 15.5 Le Cam's 3<sup>rd</sup> Lemma \*

The following example serves as a motivation to consider sequences  $\theta_n$  depending on  $n$ . It shows that pointwise asymptotics can be very misleading.

**Example 15.5.1 Hedges-Lehmann example of super-efficiency**

Let  $X_1, \dots, X_n$  be i.i.d. copies of  $X$ , where  $X = \theta + \epsilon$ , and  $\epsilon$  is  $\mathcal{N}(0, 1)$ -distributed. Consider the estimator

$$T_n := \begin{cases} \bar{X}_n, & \text{if } |\bar{X}_n| > n^{-1/4} \\ \bar{X}_n/2, & \text{if } |\bar{X}_n| \leq n^{-1/4} \end{cases}.$$

Then

$$\sqrt{n}(T_n - \theta) \xrightarrow{\mathcal{D}_\theta} \begin{cases} \mathcal{N}(0, 1), & \theta \neq 0 \\ \mathcal{N}(0, \frac{1}{4}), & \theta = 0 \end{cases}.$$

So the pointwise asymptotics show that  $T_n$  can be more efficient than the sample average  $\bar{X}_n$ . But what happens if we consider sequences  $\theta_n$ ? For example, let  $\theta_n = h/\sqrt{n}$ . Then, under  $\mathbb{P}_{\theta_n}$ ,  $\bar{X}_n = \bar{\epsilon}_n + h/(\sqrt{n}) = \mathcal{O}_{\mathbb{P}_{\theta_n}}(n^{-1/2})$ . Hence,  $\mathbb{P}_{\theta_n}(|\bar{X}_n| > n^{-1/4}) \rightarrow 0$ , so that  $\mathbb{P}_{\theta_n}(T_n = \bar{X}_n) \rightarrow 0$ . Thus,

$$\begin{aligned} \sqrt{n}(T_n - \theta_n) &= \sqrt{n}(T_n - \theta_n)l\{T_n = \bar{X}_n\} + \sqrt{n}(T_n - \theta_n)l\{T_n = \bar{X}_n/2\} \\ &\xrightarrow{\mathcal{D}_{\theta_n}} \mathcal{N}\left(-\frac{h}{2}, \frac{1}{4}\right). \end{aligned}$$

The asymptotic mean square error  $\text{AMSE}_\theta(T_n)$  is defined as the asymptotic variance + asymptotic squared bias:

$$\text{AMSE}_{\theta_n}(T_n) = \frac{1+h^2}{4}.$$

The  $\text{AMSE}_\theta(\bar{X}_n)$  of  $\bar{X}_n$  is its normalized non-asymptotic mean square error, which is

$$\text{AMSE}_{\theta_n}(\bar{X}_n) = \text{MSE}_{\theta_n}(\bar{X}_n) = 1.$$

So when  $h$  is large enough, the asymptotic mean square error of  $T_n$  is larger than that of  $\bar{X}_n$ .

Le Cam's 3<sup>rd</sup> lemma shows that asymptotic linearity for all  $\theta$  implies asymptotic normality, now also for sequences  $\theta_n = \theta + h/\sqrt{n}$ . The asymptotic variance for such sequences  $\theta_n$  does not change. Moreover, if (15.2) holds for all  $\theta$ , the estimator is also asymptotically unbiased under  $\mathbb{P}_{\theta_n}$ .

**Lemma 15.5.1** (Le Cam's 3<sup>rd</sup> Lemma) Suppose that for all  $\theta$ ,

$$T_n - \theta = \frac{1}{n} \sum_{i=1}^n l_\theta(X_i) + o_{\mathbb{P}_\theta}(n^{-1/2}),$$

where  $P_\theta l_\theta = 0$ , and  $V_\theta := P_\theta l_\theta^2 < \infty$ . Then, under regularity conditions,

$$\sqrt{n}(T_n - \theta_n) \xrightarrow{\mathcal{D}_{\theta_n}} \mathcal{N}\left(\{P_\theta(l_\theta s_\theta) - 1\}h, V_\theta\right).$$

We will present a sketch of the proof of this lemma. For this purpose, we need the following auxiliary lemma.

**Lemma 15.5.2** (*Auxiliary lemma*) Let  $Z \in \mathbb{R}^2$  be  $\mathcal{N}(\mu, \Sigma)$ -distributed, where

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \quad \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{1,2} \\ \sigma_{1,2} & \sigma_2^2 \end{pmatrix}.$$

Suppose that

$$\mu_2 = -\sigma_2^2/2.$$

Let  $Y \in \mathbb{R}^2$  be  $\mathcal{N}(\mu + a, \Sigma)$ -distributed, with

$$a = \begin{pmatrix} \sigma_{1,2} \\ \sigma_2^2 \end{pmatrix}.$$

Let  $\phi_Z$  be the density of  $Z$  and  $\phi_Y$  be the density of  $Y$ . Then we have the following equality for all  $z = (z_1, z_2) \in \mathbb{R}^2$ :

$$\phi_Z(z)e^{z_2} = \phi_Y(z).$$

**Proof.** The density of  $Z$  is

$$\phi_Z(z) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left[-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right].$$

Now, one easily sees that

$$\Sigma^{-1}a = \begin{pmatrix} 0 \\ 1 \end{pmatrix}.$$

So

$$\begin{aligned} \frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu) &= \frac{1}{2}(z - \mu - a)^T \Sigma^{-1} (z - \mu - a) \\ &\quad + a^T \Sigma^{-1} (z - \mu) - \frac{1}{2}a^T \Sigma^{-1} a \end{aligned}$$

and

$$\begin{aligned} a^T \Sigma^{-1} (z - \mu) - \frac{1}{2}a^T \Sigma^{-1} a &= \begin{pmatrix} 0 \\ 1 \end{pmatrix}^T (z - \mu) - \frac{1}{2} \begin{pmatrix} 0 \\ 1 \end{pmatrix}^T a \\ &= z_2 - \mu_2 - \frac{1}{2}\sigma_2^2 = z_2. \end{aligned}$$

□

**Sketch of proof of Le Cam's 3<sup>rd</sup> Lemma.** Set

$$\Lambda_n := \sum_{i=1}^n \left[ \log p_{\theta_n}(X_i) - \log p_{\theta}(X_i) \right].$$

Then under  $\mathbb{P}_{\theta}$ , by a two-term Taylor expansion,

$$\begin{aligned} \Lambda_n &\approx \frac{h}{\sqrt{n}} \sum_{i=1}^n s_{\theta}(X_i) + \frac{h^2}{2} \frac{1}{n} \sum_{i=1}^n \dot{s}_{\theta}(X_i) \\ &\approx \frac{h}{\sqrt{n}} \sum_{i=1}^n s_{\theta}(X_i) - \frac{h^2}{2} I(\theta), \end{aligned}$$

as

$$\frac{1}{n} \sum_{i=1}^n \dot{s}_\theta(X_i) \approx E_\theta \dot{s}_\theta(X) = -I(\theta).$$

We moreover have, by the assumed asymptotic linearity, under  $\mathbb{P}_\theta$ ,

$$\sqrt{n}(T_n - \theta) \approx \frac{1}{\sqrt{n}} \sum_{i=1}^n l_\theta(X_i).$$

Thus,

$$\begin{pmatrix} \sqrt{n}(T_n - \theta) \\ \Lambda_n \end{pmatrix} \xrightarrow{\mathcal{D}_\theta} Z,$$

where  $Z \in \mathbb{R}^2$ , has the two-dimensional normal distribution:

$$Z = \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ -\frac{h^2}{2} I(\theta) \end{pmatrix}, \begin{pmatrix} V_\theta & hP_\theta(l_\theta s_\theta) \\ hP_\theta(l_\theta s_\theta) & h^2 I(\theta) \end{pmatrix} \right).$$

Thus, we know that for all bounded and continuous  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , one has

$$\mathbb{E}_\theta f(\sqrt{n}(T_n - \theta), \Lambda_n) \rightarrow \mathbb{E}f(Z_1, Z_2).$$

Now, let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be bounded and continuous. Then, since

$$\prod_{i=1}^n p_{\theta_n}(X_i) = \prod_{i=1}^n p_\theta(X_i) e^{\Lambda_n},$$

we may write

$$\mathbb{E}_{\theta_n} f(\sqrt{n}(T_n - \theta)) = \mathbb{E}_\theta f(\sqrt{n}(T_n - \theta)) e^{\Lambda_n}.$$

The function  $(z_1, z_2) \mapsto f(z_1) e^{z_2}$  is continuous, but not bounded. However, one can show that one may extend the Portmanteau Theorem to this situation. This then yields

$$\mathbb{E}_\theta f(\sqrt{n}(T_n - \theta)) e^{\Lambda_n} \rightarrow \mathbb{E}f(Z_1) e^{Z_2}.$$

Now, apply the auxiliary Lemma, with

$$\mu = \begin{pmatrix} 0 \\ -\frac{h^2}{2} I(\theta) \end{pmatrix}, \Sigma = \begin{pmatrix} V_\theta & hP_\theta(l_\theta s_\theta) \\ hP_\theta(l_\theta s_\theta) & h^2 I(\theta) \end{pmatrix}.$$

Then we get

$$\mathbb{E}f(Z_1) e^{Z_2} = \int f(z_1) e^{z_2} \phi_Z(z) dz = \int f(z_1) \phi_Y(z) dz = \mathbb{E}f(Y_1),$$

where

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} hP_\theta(l_\theta s_\theta) \\ \frac{h^2}{2} I(\theta) \end{pmatrix}, \begin{pmatrix} V_\theta & hP_\theta(l_\theta s_\theta) \\ hP_\theta(l_\theta s_\theta) & h^2 I(\theta) \end{pmatrix} \right),$$

so that

$$Y_1 \sim \mathcal{N}(hP_\theta(l_\theta s_\theta), V_\theta).$$

So we conclude that

$$\sqrt{n}(T_n - \theta) \xrightarrow{\mathcal{D}_{\theta_n}} Y_1 \sim \mathcal{N}(hP_\theta(l_\theta s_\theta), V_\theta).$$

Hence

$$\sqrt{n}(T_n - \theta_n) = \sqrt{n}(T_n - \theta) - h \xrightarrow{\mathcal{D}_{\theta_n}} \mathcal{N}(h\{P_\theta(l_\theta s_\theta) - 1\}, V_\theta).$$

□



# Chapter 16

## Complexity regularization

Suppose again the framework where we have i.i.d. copies  $X_1, \dots, X_n$  of a population random variable  $X$ , and that we model the distribution  $P$  of  $X$  as  $P \in \mathcal{P} = \{P_\theta : \theta \in \Theta\}$ . If  $\Theta$  is finite-dimensional, one can regard its dimension,  $p$  say, as a measure of the complexity of the parameter space  $\Theta$ . If  $p$  is larger than  $n$ , there are more parameters than observations and one may intuitively already see that this is not a statistically favourable situation. It is in a way comparable to a system with more unknowns than equations: an ill-posed system. A model with very many parameters may fit the data very well, but will have little predictive power. With too many parameters, there is a danger of overfitting. Complexity regularization can be roughly seen as a way to deal with a complex parameter space by letting the data decide which sub-model yields a good trade-off between the approximation error and estimation error.

*Set of all of the possible probabilities deriving from  $\theta$ . In other words, set of all non-zero probabilities that fulfill normalization.*

We note that even when the parameter space is  $\infty$ -dimensional one need not always apply complexity regularization. The dimension of  $\Theta$  is in itself not always the best description of complexity. If  $\Theta$  is a metric space, one can apply its so-called *entropy* as a measure of complexity. The details go beyond the scope of these lecture notes. We instead will give the non-parametric and the high-dimensional regression problem as prototype examples.

### 16.1 Non-parametric regression

Consider  $n$  real-valued response variables  $Y_1, \dots, Y_n$  which depend on some fixed co-variables  $x_1, \dots, x_n$  (with  $x_i$  in some space  $\mathcal{X}$  for all  $i$ ) in the following manner:

$$Y_i = f(x_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\epsilon_1, \dots, \epsilon_n$  is unobservable noise, and where  $f$  is an unknown function. Suppose we think of using the least squares estimator for estimating  $f$ . Let  $Y := (Y_1, \dots, Y_n)^T$ . With some abuse of notation, we identify a function  $f$ :

$\mathcal{X} \rightarrow \mathbb{R}$  with the vector

$$\mathbf{f} := (\mathbf{f}(x_1), \dots, \mathbf{f}(x_n))^T \in \mathbb{R}^n$$

The least squares estimator is

$$\hat{\mathbf{f}}_{\text{LS}} := \arg \min_{\mathbf{f} \in \mathbb{R}^n} \|\mathbf{Y} - \mathbf{f}\|_2^2.$$

Clearly, if all  $x_i$  are distinct, then  $\hat{\mathbf{f}}_{\text{LS}} = \mathbf{Y}$  and one has a perfect fit

$$\|\mathbf{Y} - \hat{\mathbf{f}}_{\text{LS}}\|_2^2 = 0.$$

This is a typical instance of overfitting. The estimator  $\hat{\mathbf{f}}_{\text{LS}}$  just reproduces the data and has no predictive power.

We need a model for  $f$ , say  $f \in \mathcal{F}$  with  $\mathcal{F}$  some class of functions.

## 16.2 Smoothness classes

Suppose  $\mathcal{X}$  is some interval in  $\mathbb{R}$ , say  $\mathcal{X} = [0, 1]$ . It depends on the situation but a reasonable assumption on  $f$  may be that it is not too wiggly. One may formulate that mathematically for example by saying that  $f$  is differentiable with derivative  $f'$  with  $|f'|$  within bounds. One may for example measure the roughness of  $f$  by the (Sobolev) semi-norm of its derivative

$$\sqrt{\int_0^1 |f'(x)|^2 dx}.$$

*smoothness assumption.*

↑

Notice that here taking the interval between  $[0, 1]$  result from the definition of the problem  $X \in [0, 1]$ .

One way to go is now to do least squares over all  $f$  under the restriction that  $\int_0^1 |f'(x)|^2 dx \leq M^2$  where  $M$  is a given constant. It turns out that a more flexible approach is to apply the Lagrangian version of this. Then, choose a tuning parameter  $\lambda \geq 0$  and define the estimator  $\hat{f}$  as

$$\hat{f} := \arg \min_{\mathbf{f}} \left\{ \|\mathbf{Y} - \mathbf{f}\|_2^2 + \lambda^2 \int_0^1 |f'(x)|^2 dx \right\}.$$

↓  
how do you choose that?

This is called (a version of) Tikhonov regularization. One may view

$$\lambda^2 \int_0^1 |f'(x)|^2 dx$$

as a penalty for choosing a too wiggly function. The penalty regularizes the function. The tuning parameter  $\lambda$  controls the amount of regularization: the larger  $\lambda$  the more regular the estimator will be. The choice of  $\lambda$  is not an easy point. From theoretical point of view there are some guidelines. (Choosing  $\lambda^2$  of order  $n^{1/3}$  trades off approximation error and estimation error. One may alternatively invoke Bayesian arguments to choose  $\lambda$ .) In practice one may apply cross-validation.

So far we described smoothness in terms of first derivatives being bounded. One may also use higher order derivatives if one believes the unknown function  $f$  has these. Let  $f^{(m)}$  denote the derivative of order  $m$  of the function  $f : [0, 1] \rightarrow \mathbb{R}$ . A possible penalty is then

$$\lambda^2 \int_0^1 |f^{(m)}(x)|^2 dx.$$

The tuning parameter  $\lambda$  can be chosen smaller for higher values of  $m$ . (A value of order  $n^{\frac{1}{2m+1}}$  gives a trade-off between approximation error and estimation error.) The resulting estimator is called a *smoothing spline*.

With a quadratic penalty, the penalized least squares estimator  $\hat{f}$  is not difficult to compute as it is a minimizer of a quadratic function. Nevertheless, explicit expressions are typically not available. In the next section, we provide explicit expressions for the continuous version, as a “curiosity”.

### 16.3 A continuous version with explicit solution ★

We examine the continuous version of the previous section. The problem can be explicitly solved. Suppose we observe a function  $y : [0, 1] \rightarrow \mathbb{R}$  and we aim at smoothing it using the penalty of the previous section. We let

$$\hat{f} = \arg \min_f \left\{ \int_0^1 |y(x) - f(x)|^2 dx + \lambda^2 \int_0^1 |f'(x)|^2 dx \right\}. \quad (16.1)$$

**Lemma 16.3.1** *Let  $\hat{f}$  be given in (16.1). Then*

$$\hat{f}(x) = \frac{C}{\lambda} \cosh\left(\frac{x}{\lambda}\right) + \frac{1}{\lambda} \int_0^x y(u) \sinh\left(\frac{u-x}{\lambda}\right) du,$$

where

$$C = Y(1) - \left\{ \frac{1}{\lambda} \int_0^1 Y(u) \sinh\left(\frac{1-u}{\lambda}\right) du \right\} / \sinh\left(\frac{1}{\lambda}\right),$$

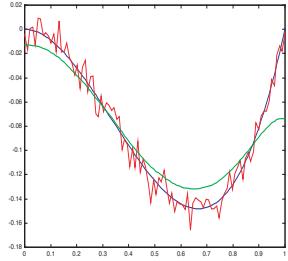
with

$$Y(x) = \int_0^x y(u) du.$$

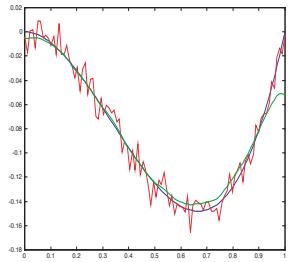
The proof is calculus of variations and will be skipped.

Yet, with this explicit expression it is very easy to carry out the estimation procedure. Here is an example.

## Numerical example



Denoised, lambda=0.1  
Error=2.8119e-04



Denoised, lambda=0.05  
Error=7.8683e-05

In black is the unknown function  $f$  (as this is a simulation we know the unknown  $f$ ). The red wiggly function is the observed function  $y$ . The green function is the estimator  $\hat{f}$ . We see in the second figure that by decreasing the tuning parameter  $\lambda$  the estimator  $\hat{f}$  is closer to the unknown truth  $f$ .



## 16.4 Estimating a function of bounded variation

Suppose again  $\mathcal{X} = [0, 1]$ . Instead of taking the penalty

$$\lambda^2 \int_0^1 |f'(x)|^2 dx$$

one may alternatively choose

$$\lambda^2 \int_0^1 |f'(x)| dx$$

i.e. without the square. This seems like a minor modification but it makes a huge difference. One may further relax the assumption of differentiability and define the total variation of the function  $f$  as

$$TV(f) := \sum_{i=2}^n |f(x_i) - f(x_{i-1})|,$$

where we assume that the  $x_i$  are ordered:  $x_1 < \dots < x_n$ . This leads to the estimator

$$\hat{f} := \arg \min_f \left\{ \|Y - f\|_2^2 + \lambda^2 \text{TV}(f) \right\}.$$

!

In analogy, if  $Y$  denotes the heights of mountains at locations where a road is to be built, one tries to even out the mountains and valleys using the total variation penalty (so that the road will not have steep slopes) and at the same time move little earth measured in terms of the least squares loss. This can be done in an iterative manner by filling the valley where the slope is steepest, or digging away the mountain with steepest slope. The estimator is locally adaptive: by increasing  $\lambda$  only local changes are made. This is in contrast with the estimator of Section 16.2 or its continuous version in Section 16.3. One can see in the numerical example of Section 16.3 that changing  $\lambda$  has a non-local effect.

One may write the estimator as solution of a linear least squares problem with an  $\ell_1$ -penalty on the coefficients. This will relate the problem with that of Section 16.5. The results there will help to understand why removing the square in the penalty drastically changes the estimator.

Let us define  $f(x_0) =: 0$ . Then for  $i = 1, \dots, n$

$$\begin{aligned} f(x_i) &= \sum_{j=1}^i f(x_j) - f(x_{j-1}) \\ &= \sum_{i=1}^n \underbrace{\left( f(x_j) - f(x_{j-1}) \right)}_{:= b_j} \underbrace{l\{j \leq i\}}_{:= \xi_{i,j}} \\ &= \sum_{i=1}^n b_j \xi_{i,j}. \end{aligned}$$

Putting the coefficients  $b_j$ ,  $j = 1, \dots, n$  in a vector  $b = (b_1, \dots, b_n)^T$  and the  $\xi_{i,j}$  in a matrix

$$X := \begin{pmatrix} \xi_{1,1} & \cdots & \xi_{1,n} \\ \vdots & \ddots & \vdots \\ \xi_{n,1} & \cdots & \xi_{n,n} \end{pmatrix}$$

we see that the total variation penalized estimator is

$$\hat{f} = \arg \min_{\substack{f \\ f=Xb}} \left\{ \|Y - f\|_2^2 + \lambda^2 \|b\|_1 \right\}$$

where  $\|b\|_1 = \sum_{j=1}^n |b_j|$  denotes the  $\ell_1$ -norm of the vector  $b \in \mathbb{R}^n$ . This rewriting makes clear that there are as many parameters as there are observations, namely  $n$ . The penalty takes care that despite this fact one will not overfit the data (when  $\lambda$  is not too small).

We end this section with a small trip to the case where  $\mathcal{X}$  is two-dimensional, say  $\mathcal{X} = [0, 1]^2$ . Then the unknown  $f$  is an image, say, and the observations are

$$Y_{i_1, i_2} = f(x_{i_1, i_2}) + \epsilon_{i_1, i_2}.$$

Consider  $n = m^2$  observations and say they are on a regular grid

$$\underline{x_{i_1, i_2}} = \left( \frac{i_1}{m}, \frac{i_2}{m} \right), \underline{i_1 = 1, \dots, m}, \underline{i_2 = 1, \dots, m}.$$

Now, how can we model smoothness of an image? Again, one may opt to work with the squares of derivatives, a counterpart of  $\int |f'(x)|^2 dx$  for the one-dimensional case. However, as in the one-dimensional case, taking squares has non-local effects. The penalized least squares reconstruction of the image will then look blurred. For example, if the image is a landscape with lakes and rivers, there are sharp boundaries which will be blurred by taking a penalty based on squared derivatives. An alternative is again a total variation penalty. There are several definitions around for total variation in dimension larger than 1. One possibility in our 2-dimensional case is to define

$$\text{TV}(f) := \sum_{i_1=2}^m \sum_{i_2=2}^m |\Delta f(x_{i_1, i_2})|,$$

where

$$\Delta f(x_{i_1, i_2}) := f\left(\frac{i_1}{m}, \frac{i_2}{m}\right) - f\left(\frac{i_1-1}{m}, \frac{i_2}{m}\right) - f\left(\frac{i_1}{m}, \frac{i_2-1}{m}\right) + f\left(\frac{i_1-1}{m}, \frac{i_2-1}{m}\right).$$

Think about  
the geometrical  
construction of  
this

The image reconstruction algorithm is

$$\hat{f} := \arg \min_f \left\{ \sum_{i_1=1}^m \sum_{i_2=1}^m \left( Y_{i_1, i_2} - f(x_{i_1, i_2}) \right)^2 + \lambda^2 \text{TV}(f) \right\}.$$

This estimator can again be written as a least squares estimator of a linear function with an  $\ell_1$ -penalty on the coefficients.

## 16.5 The ridge and Lasso penalty

In the linear model one has data  $(x_1, Y_1), \dots, (x_n, Y_n)$  with  $x_i \in \mathbb{R}^p$  a  $p$ -dimensional row vector and  $Y_i \in \mathbb{R}$  ( $i = 1, \dots, n$ ) and one wants to find a good linear approximation using the least squares loss function

$$b \mapsto \sum_{i=1}^n \left( Y_i - \sum_{j=1}^p x_{i,j} b_j \right)^2,$$

Define the design matrix

$$X := \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} x_{1,1} & \cdots & x_{1,p} \\ \vdots & \ddots & \vdots \\ x_{n,1} & \cdots & x_{n,p} \end{pmatrix}$$

and the vector of responses

$$Y := \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}.$$

Set the vector of coefficients to  $b := (b_1, \dots, b_p)^T$ . Then

$$\sum_{i=1}^n \left( Y_i - \sum_{j=1}^p x_{i,j} b_j \right)^2 = \|Y - Xb\|_2^2.$$

If  $p \geq n$  and  $X$  has rank  $n$ , minimizing this over all  $b \in \mathbb{R}^p$  gives a “perfect” solution  $\hat{\beta}_{\text{LS}}$  with  $X\hat{\beta}_{\text{LS}} = Y$ . This solution just reproduces the data and is therefore of no use. We say that it overfits.

**Definition 16.5.1** The ridge regression estimator is

$$\hat{\beta}_{\text{ridge}} := \arg \min_{b \in \mathbb{R}^p} \left\{ \|Y - Xb\|_2^2 + \lambda^2 \|b\|_2^2 \right\},$$

where  $\lambda > 0$  is a regularization parameter.

**Definition 16.5.2** The Lasso<sup>1</sup> estimator is

$$\hat{\beta}_{\text{Lasso}} := \arg \min_{b \in \mathbb{R}^p} \left\{ \|Y - Xb\|_2^2 + 2\lambda \|b\|_1 \right\},$$

where  $\lambda > 0$  is a regularization parameter and  $\|b\|_1 := \sum_{j=1}^p |b_j|$  is the  $\ell_1$ -norm of  $b$ .

**Note** Recall the definition of the Bayesian MAP estimator as given in Subsection 10.5.3. Consider the model  $Y = X\beta + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ . The ridge regression estimator is the MAP estimator using as prior  $\beta_1, \dots, \beta_p$  i.i.d.  $\sim [\mathcal{N}(0, \tau^2)]$  with  $\tau = \sigma/\lambda$ . The Lasso estimator is the MAP using as prior  $\beta_1, \dots, \beta_p$  i.i.d.  $\sim [\text{Laplace}(0, \tau^2)]$  where the standard deviation  $\tau$  is  $\tau = \sqrt{2}\sigma^2/\lambda$ . See also Section 10.6.

**Remark** As  $\lambda$  grows the ridge estimator shrinks the coefficients. They will however not be set exactly to zero. The coefficients of the Lasso estimator shrink as well, and some - or even many - are set exactly to zero. The ridge estimator can be useful if  $p$  is moderately large. For very large  $p$  the Lasso is often preferred. The idea is that one should not try to estimate something when the signal is below the noise level. Instead, then one should simply put it to zero.



**Remark** Both ridge estimator and Lasso are biased. As  $\lambda$  increases the bias increases, but the variance decreases.

*↳ This is possibly the entire point of Lasso and Ridge.*

**Remark** The regularization parameter  $\lambda$  is for example chosen by using “cross validation” or (information) theoretic or Bayesian arguments. Below, we will see that for the Lasso a choice of order  $\sqrt{n \log p}$  is theoretically justified.

<sup>1</sup>This is relatively recent methodology, introduced as Lasso by Tibshirani, R., 1996: Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B* 58, 267-288

We now investigate further expressions for both ridge estimator and Lasso.

**Lemma 16.5.1** *The ridge estimator  $\hat{\beta}_{\text{ridge}}$  is given by*

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda^2 I)^{-1} X^T Y.$$

**Proof.** We have

$$\begin{aligned} \frac{1}{2} \frac{\partial}{\partial b} \left\{ \|Y - Xb\|_2^2 + \lambda^2 \|b\|_2^2 \right\} &= -X^T(Y - Xb) + \lambda^2 b \\ &= -X^T Y + \left( X^T X + \lambda^2 I \right) b. \end{aligned}$$

The estimator  $\hat{\beta}_{\text{ridge}}$  puts this to zero.  $\square$

As Hat  
means  
the inverse  
thereof.

**Corollary 16.5.1** *Suppose orthonormal design:  $X^T X = nI$  (thus  $p \leq n$  necessarily). Then*

$$\hat{\beta}_{\text{ridge}} = X^T Y / (n + \lambda^2).$$

After some calculations, as in Example 5.2.1, one sees that when  $\epsilon_1, \dots, \epsilon_n$  are i.i.d. with mean zero and variance  $\sigma^2$ , then

$$\mathbb{E} \|X\hat{\beta}_{\text{ridge}} - f\|_2^2 = \underbrace{\left[ \frac{\lambda^2/n}{1 + \lambda^2/n} \right]^2 \|X\beta^*\|_2^2}_{\text{(bias)}^2} + \underbrace{\left[ \frac{1}{1 + \lambda^2/n} \right]^2 p\sigma^2}_{\text{variance}} + \underbrace{\|X\beta^* - f\|_2^2}_{\text{misspecification error}},$$

where  $X\beta^*$  is the projection of  $f$  on the space spanned by the columns of  $X$ . We see that in order to trade off bias and variance, we have to know the variance  $\sigma^2$ , but what is worse, we also have to know the  $(\text{bias})^2 \|X\beta^*\|_2^2$ . But the bias is unknown as  $f$  is unknown. Thus, we are facing the same problems as in Example 5.2.1.

For the Lasso estimator there is no simple expression in general. We therefore only consider the special case of orthonormal design.

**Lemma 16.5.2** *Suppose  $X^T X = nI$  (thus  $p \leq n$  necessarily). Define  $Z := X^T Y$ . Then for  $j = 1, \dots, p$*

$$\hat{\beta}_{\text{Lasso},j} = \begin{cases} Z_j/n - \lambda/n & Z_j \geq \lambda \\ 0 & |Z_j| \leq \lambda \\ Z_j/n + \lambda/n & Z_j \leq -\lambda \end{cases}.$$

**Proof.** Write  $\hat{\beta}_{\text{Lasso}} =: \hat{\beta}$  for short. We can write

$$\|Y - Xb\|_2^2 = \|Y\|_2^2 - 2b^T X^T Y + nb^T b = -2b^T Z + nb^T b.$$

Thus for each  $j$  we minimize

$$-2b_j Z_j + nb_j^2 + 2\lambda|b_j|.$$

If  $\hat{\beta}_j > 0$  it must be a solution of putting the derivative of the above expression to zero:

$$-Z_j + n\hat{\beta}_j + \lambda = 0,$$

or

$$\hat{\beta}_j = Z_j/n - \lambda/n.$$

Similarly, if  $\hat{\beta}_j < 0$  we must have

$$-Z_j + n\hat{\beta}_j - \lambda = 0.$$

Otherwise  $\hat{\beta}_j = 0$ .  $\square$

### Some notation

- o For a vector  $z \in \mathbb{R}^p$  we let  $\|z\|_\infty := \max_{1 \leq j \leq p} |z_j|$  be its  $\ell_\infty$ -norm.
- o We let  $X_1, \dots, X_p$  denote the columns of  $X$ .
- o For a subset  $S \subset \{1, \dots, p\}$  we let  $X\beta_S^*$  be the best linear approximation of  $f := EY$  using the variables in  $S$ , i.e.,  $X\beta_S^*$  is the projection in  $\mathbb{R}^n$  of  $f$  on the linear space  $\{\sum_{j \in S} X_j b_{S,j} : b_S \in \mathbb{R}^{|S|}\}$ .

In the next theorem we again assume orthogonal design. For general design, one needs so-called restricted eigenvalues or compatibility conditions between  $\ell_2$ -norms and  $\ell_1$ -norms.

**Theorem 16.5.1** Consider again fixed design with  $X^T X = nI$ . Let  $f = EY$  and  $\epsilon = Y - f$ . Fix some level  $\alpha \in (0, 1)$  and suppose that for some  $\lambda_\alpha$  it holds that

$$\mathbb{P}(\|X^T \epsilon\|_\infty > \lambda_\alpha) \leq \alpha.$$

Then for  $\lambda > \lambda_\alpha$  we have with probability at least  $1 - \alpha$

$$\|X\hat{\beta}_{\text{Lasso}} - f\|_2^2 \leq \min_S \left\{ \underbrace{\frac{(\lambda + \lambda_\alpha)^2}{n}|S|}_{\text{estimation error}} + \underbrace{\|X\beta_S^* - f\|_2^2}_{\text{approximation error}} \right\}.$$

**Proof.** Write  $\hat{\beta} := \hat{\beta}_{\text{Lasso}}$  and  $f = X\beta$ . On the set where  $\|X^T \epsilon\|_\infty \leq \lambda_\alpha$  we have

$$-n|\beta_j| > \lambda + \lambda_\alpha \Rightarrow n|\hat{\beta}_j - \beta_j| \leq \lambda + \lambda_\alpha,$$

$$-n|\beta_j| \leq \lambda + \lambda_\alpha \Rightarrow |\hat{\beta}_j - \beta_j| \leq |\beta_j|.$$

So with probability at least  $(1 - \alpha)$ ,

$$\begin{aligned} \|X\hat{\beta}_{\text{Lasso}} - f\|_2^2 &\leq \frac{(\lambda + \lambda_\alpha)^2}{n} \left( \#\{j : n|\beta_j| > \lambda + \lambda_\alpha\} \right) + \sum_{n|\beta_j| \leq \lambda + \lambda_\alpha} n\beta_j^2 \\ &= \min_S \left\{ \frac{(\lambda + \lambda_\alpha)^2}{n}|S| + \|X\beta_S^* - f\|_2^2 \right\}. \end{aligned}$$



$\square$

If we compare the result of Theorem 16.5.1 with result iii) of Lemma 12.3.1 concerning the ordinary least squares estimator, we see that the Lasso exhibits

an automatic trade-off between approximation error and estimation error. This is called *adaptation*. As we will see below, the parameter  $\lambda$  can typically be chosen of order  $\sqrt{n \log p}$ . Therefore the price paid for not knowing a priori which subset of the coefficients is relevant is of order  $\log p$ . This is generally considered as being a relatively low price.

**Note** One may write

$$\|X\beta_S^* - f\|_2^2 = \|X\beta_S^* - X\beta^*\|_2^2 + \|X\beta^* - f\|_2^2$$

where  $X\beta^*$  is the projection of  $f$  on the space spanned by the columns of  $X$ . Thus, the “approximation error” actually consists of two terms. The second term is the misspecification error: it vanishes when the linear model is well-specified.

**Corollary 16.5.2** Suppose that  $f = X\beta$  where  $\beta$  has  $s := \#\{j : \beta_j \neq 0\}$  non-zero components. Then under the conditions of the above theorem, with probability at least  $1 - \alpha$

$$\|X(\hat{\beta}_{\text{Lasso}} - \beta)\|_2^2 \leq \frac{(\lambda + \lambda_\alpha)^2}{n} s.$$

The above corollary tells us that the Lasso estimator adapts to favourable situations where  $\beta$  has many zeroes (i.e. where  $\beta$  is sparse).

To complete the story, we need to study a bound for  $\lambda_\alpha$ . It turns out that for many types of error distributions, one can take  $\lambda_\alpha$  of order  $\sqrt{n \log p}$ . We show this for the case of i.i.d.  $\mathcal{N}(0, \sigma^2)$  noise. For that purpose, we start with a bound for the tails of a standard normal random variable.

**Lemma 16.5.3** Suppose  $Z \sim \mathcal{N}(0, 1)$ . Then for all  $t > 0$

$$\mathbb{P}(Z \geq \sqrt{2t}) \leq \exp[-t].$$

**Proof.** First check that for all  $u > 0$

$$E \exp[uZ] = \exp[u^2/2].$$

Then by Chebyshevs inequality

$$\mathbb{P}(Z > \sqrt{2t}) \leq \exp[u^2/2 - u\sqrt{2t}].$$

Now choose  $u = \sqrt{2t}$ .  $\square$

**Corollary 16.5.3** Let  $Z_1, \dots, Z_p$  be  $p$  (possibly dependent) standard normal random variables. Then for all  $t$

$$\begin{aligned} \mathbb{P}\left(\max_{1 \leq j \leq p} |Z_j| \geq \sqrt{2(\log(2p) + t)}\right) &\leq \sum_{j=1}^p \mathbb{P}(|Z_j| \geq \sqrt{2(\log(2p) + t)}) \\ &\leq 2p \exp[-(\log(2p) + t)] = \exp[-t]. \end{aligned}$$

**Remark** In the above corollary we used  $\mathbb{P}(A \cup B) \leq \mathbb{P}(A) + \mathbb{P}(B)$ . This is called the *union bound*.

**Corollary 16.5.4** Let  $\epsilon_1, \dots, \epsilon_n$  be i.i.d.  $\mathcal{N}(0, \sigma^2)$  and let  $X = (X_1, \dots, X_p)$  with  $X_1, \dots, X_p$  be fixed vectors in  $\mathbb{R}^n$  with  $\|X_j\|_2 = n$  for all  $j$ . Then for  $0 < \alpha < 1$

$$\mathbb{P}\left(\|X^T \epsilon\|_\infty \geq \sigma \sqrt{2n(\log(2p/\alpha))}\right) \leq \alpha.$$

**Remark.** The value  $\alpha = \frac{1}{2}$  thus gives a bound for the median of  $\|X \hat{\beta}_{\text{Lasso}} - f\|_2^2$ . In the case of Gaussian errors one may use “concentration of measure” to deduce that  $\|X \hat{\beta}_{\text{Lasso}} - f\|_2^2$  is “concentrated” around its median.

## 16.6 Conclusion

We have seen in this chapter that several concepts from classical statistics also play their role in the modern version of statistics where the parameter is possibly high dimensional. The classical least squares methodology keeps its prominent place, but now it is equipped with a regularization penalty. More generally, M-estimators (e.g. maximum likelihood estimators) can also be used in high dimensions, applying again some regularization technique? The bias-variance decomposition continues to play its role too, for instance leading to guidelines for choosing tuning parameters.

Shrinkage estimators play an important role in high-dimensional statistics. This is also related to the result of Section 11.4 where we have seen that the sample average in dimension higher than 2 is inadmissible as it can be improved by a shrinkage estimator.

Complexity regularization can typically be seen as a Bayesian MAP approach. One may also use the a posteriori mean as estimator etc. Today, Bayesian approaches to high-dimensional and non-parametric problems are very important and successful.

Complexity regularization is typically invoked for the construction of adaptive estimators. An adaptive estimator mimics the situation where we knew beforehand the complexity of the underlying target to be estimated. To evaluate the performance of an adaptive estimator, one uses as benchmark the case where the (hopefully low) complexity of the target is indeed known. Thus the benchmark comes from classical statistical theory.



# Chapter 17

## Literature

- J.O. Berger (1985) *Statistical Decision Theory and Bayesian Analysis* Springer  
A fundamental book on Bayesian theory.
- P.J. Bickel, K.A. Doksum (2001) *Mathematical Statistics, Basic Ideas and Selected Topics* Volume I, 2<sup>nd</sup> edition, Prentice Hall  
Quite general, and mathematically sound.
- D.R. Cox and D.V. Hinkley (1974) *Theoretical Statistics* Chapman and Hall  
Contains good discussions of various concepts and their practical meaning.  
Mathematical development is sketchy.
- A. DasGupta (2011) *Probability for Statistics and Machine Learning*, Springer  
Contains all the probability theory background needed. (Look out for the upcoming book *Statistical Theory, a Comprehensive Course* by the same author.)
- J.G. Kalbfleisch (1985) *Probability and Statistical Inference* Volume 2, Springer  
Treats likelihood methods.
- L.M. Le Cam (1986) *Asymptotic Methods in Statistical Decision Theory* Springer  
Treats decision theory on a very abstract level.
- E.L. Lehmann (1983) *Theory of Point Estimation* Wiley  
A “klassiker”. The lecture notes partly follow this book
- E.L. Lehmann (1986) *Testing Statistical Hypothesis* 2<sup>nd</sup> edition, Wiley  
Goes with the previous book.
- J.A. Rice (1994) *Mathematical Statistics and Data Analysis* 2<sup>nd</sup> edition, Duxbury Press  
A more elementary book.

- M.J. Schervish (1995) *Theory of Statistics* Springer  
Mathematically exact and quite general. Also good as reference book.
- R.J. Serfling (1980) *Approximation Theorems of Mathematical Statistics* Wiley  
Treats asymptotics.
- A.W. van der Vaart (1998) *Asymptotic Statistics* Cambridge University Press  
Treats modern asymptotics and e.g. semiparametric theory
- L. Wasserman (2004) *All of Statistics. A Concise Course in Statistical Inference* Springer.  
Contains a wide range of topics in mathematical statistics and machine learning.