



HUMBOLDT-UNIVERSITÄT ZU BERLIN

## **Methods of Statistics (M25)**

Lecturer: Prof. Dr. Sven Wang  
WS 24/25

Last Update: Friday 22<sup>nd</sup> November, 2024

# Contents

<b>1</b>	<b>Intro and Disclaimer</b>	<b>4</b>
<b>2</b>	<b>Basic Statistical Concepts</b>	<b>5</b>
<b>3</b>	<b>Parameter Estimation</b>	<b>7</b>
3.1	Maximum Likelihood Principle . . . . .	9
3.2	Bayesian method . . . . .	10
<b>4</b>	<b>Decision Theory</b>	<b>14</b>
4.1	Another optimality concept: Admissibility . . . . .	17
<b>5</b>	<b>Confidence Sets</b>	<b>19</b>
5.1	Hypothesis Testing . . . . .	21
5.1.1	Basic Definitions . . . . .	21
5.1.2	Type I and Type II Errors . . . . .	21
5.1.3	Level and Uniformly Most Powerful Tests . . . . .	21
5.1.4	The Neyman-Pearson Lemma . . . . .	22
5.1.5	Proof of the Neyman-Pearson Lemma . . . . .	22
<b>6</b>	<b>Lecture 6: Neyman-Pearson Lemma and Likelihood Ratio Tests</b>	<b>24</b>
<b>7</b>	<b>Linear Models</b>	<b>29</b>
7.1	Introduction to Linear Regression Models . . . . .	29
7.2	Simple Linear Regression Model . . . . .	29
7.2.1	Statistical Model . . . . .	29
7.2.2	Likelihood Function . . . . .	30
7.2.3	Maximum Likelihood Estimation (MLE) . . . . .	30
7.3	General Linear Models . . . . .	31
7.4	Least Squares Estimation . . . . .	32
7.4.1	Geometric Interpretation . . . . .	32
7.5	Representation for the Least Squares Estimator . . . . .	33
7.6	Optimality of the Least Squares Estimator . . . . .	33
7.7	Conclusion . . . . .	36

<b>8</b>	<b>Lecture 9: Ridge Regression and Its Properties</b>	<b>37</b>
<b>9</b>	<b>Lecture 10: Confidence Sets and Hypothesis Testing in Linear Models</b>	<b>45</b>
9.0.1	The $t$ - and $F$ -Tests . . . . .	48
9.0.2	General Linear Hypothesis Testing Problems . . . . .	49
9.0.3	The $t$ - and $F$ -Tests . . . . .	51
9.0.4	General Linear Hypothesis Testing Problems . . . . .	51

# 1 Intro and Disclaimer

These lecture notes are based on the material presented by Professor Wang during class and are written by students. They may contain errors or omissions. Please refer to the in-person lectures and the literature on Moodle for accurate and authoritative information. If you find an error or want to help, please send an email to:

[said.kassner@student.hu-berlin.de](mailto:said.kassner@student.hu-berlin.de)

[stephensonmonroe@gmail.com](mailto:stephensonmonroe@gmail.com)

[salihiad@hu-berlin.de](mailto:salihiad@hu-berlin.de)

## 2 Basic Statistical Concepts

### Lecture 1

Here is the literature that the class is based on.

**Literature:**

- WS 19/20 R. Altmeyer "*Gliederung Methoden der Statistik*"
- L. Wasserman, *All of Statistics*
- M. Trabs, K. Krenz, M. Jirak and M. Reiss. *Methoden der Statistik*.
- Hastie, Tibshirani, et al., *Elements of Statistical Learning*

Let's start with a (simplest possible) example:

**Example 1 (Polling).** Consider a poll with two answers A and B (representing political parties).

- $N$  = total number of votes
- $M$  = total number of votes supporting party A

**Poll Definitions:**

- $n$  = size of the poll
- $x = (x_1, \dots, x_n)$  = responses, where:

$$x_i = \begin{cases} 0 & \text{if the } i\text{-th person supports B} \\ 1 & \text{if the } i\text{-th person supports A} \end{cases}$$

**Additional Assumptions:**

- $n$ -times, we select a person randomly from the set  $\{1, \dots, N\}$ , and record their (truthful) response.
- Every asked person responds (i.e., no selection bias).
- People can be asked repeatedly.

**Aim of the Poll:** The aim of the poll is to estimate the fraction of party A supporters. This can be written as:

$$\theta = \frac{M}{N} \in [0, 1]$$

An intuitive estimate of  $\theta$  is:

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

**Question:** Is this a good (or best possible) estimator? What properties does it have?

To answer this, we formalize some statistical notions.

**Definition 1 (Sample space).** A *sample space* is a measurable space  $(\mathcal{X}, \mathcal{F})$ , i.e., a set  $\mathcal{X}$  with a  $\sigma$ -algebra  $\mathcal{F}$ , in which our statistical observations take values.

**Definition 2 (Statistical model).** Let  $(\mathcal{X}, \mathcal{F})$  be some sample space and let  $\Theta$  be a set, which we call the *parameter space*. A *statistical model* is a family of probability measures  $\{P_\theta : \theta \in \Theta\}$  on  $(\mathcal{X}, \mathcal{F})$ .

*Remark 1.* Often,  $(\mathcal{X}, \mathcal{F})$  is a "product space." For example, in Example 1,  $\mathcal{X} = \{0, 1\}^n$ , and each  $P_\theta$  is a product distribution, i.e.,  $x_1, \dots, x_n$  are independent, identically distributed. Then we say  $\{P_\theta : \theta \in \Theta\}$  is an *iid statistical model*.

*Remark 2 (Back to Example 1).* Here:

- $\mathcal{X} = \{0, 1\}^n$
- $\Theta = [0, 1]$
- $\mathcal{F} = \mathcal{P}(\{0, 1\}^n)$
- $P_\theta = (\text{Bernoulli}(\theta))^{\otimes n}$

*Remark 3.* If every person could only be asked once, we would have  $P_\theta = \text{Hypergeometric}(N, M, n)$ , which "converges" to the Bernoulli model as  $N, M \rightarrow \infty$ . We might have to discretize  $\Theta$  and take  $\theta = \frac{M}{N}$ . (Exercise: Think about it!)

# 3 Parameter Estimation

Assume that  $\Theta \subseteq \mathbb{R}^p$ , for  $p \geq 1$ . This is the setting of parametric statistics. [Assume  $\Theta$  is measurable.]

**Definition 3 (Estimator).** An estimator for  $\theta \in \Theta$  is any measurable function:

$$\hat{\theta} : (\mathcal{X}, \mathcal{F}) \rightarrow \Theta.$$

Any function that, based on some data  $x \in \mathcal{X}$ , outputs a guess / estimate  $\hat{\theta}(x) \in \Theta$ .

## Lecture 2

**Last time:** Statistical model = family of probability measures on  $(\mathcal{X}, \mathcal{F})$  indexed by  $\theta \in \Theta$ .

**Sample space:**  $(\mathcal{X}, \mathcal{F})$

**Estimator:** = measurable function  $(\mathcal{X}, \mathcal{F}) \rightarrow \Theta$

Now, what are some desirable properties we would like to have?

**Definition 4 (Unbiased estimator).** Let  $\Theta = \mathbb{R}^p$  (measurable),  $p \geq 1$ . An estimator  $\hat{\theta}$  is unbiased if

$$\mathbb{E}_{\theta}[\hat{\theta}] = \mathbb{E}_{\mathbb{P}_{\theta}}[\hat{\theta}] = \theta, \text{ for all } \theta \in \Theta.$$

Where  $\mathbb{E}_{\theta}[\cdot] = \mathbb{E}_{\mathbb{P}_{\theta}}[\cdot]$  denotes expectation under the law  $\mathbb{P}_{\theta}$ .

In more explicit terms:

$$\mathbb{E}_{x \sim \mathbb{P}_{\theta}}[\hat{\theta}(x)] = \theta \quad \forall \theta$$

*Remark 4 (Unbiasedness).* Unbiasedness means "no systematic errors." However, we'd also like a "good"  $\hat{\theta}$  to be concentrated around the data-generating parameter.

**Definition 5 (Consistent estimator).** Let  $(\mathbb{P}_{\theta}^n)_{\theta \in \Theta}$  be a sequence of statistical models ( $n \geq 1$ ), on the same parameter space  $\Theta$  not depending on  $n \geq 1$ .

Let  $\hat{\theta}_n$  be a sequence of estimators. Then  $\hat{\theta}_n$  is called consistent if for every  $\theta \in \Theta$ ,

$$\hat{\theta}_n \xrightarrow{\mathbb{P}_{\theta}^n} \theta$$

or explicitly, for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta^n(|\hat{\theta}_n - \theta| > \epsilon) = 0.$$

**Back to Example 1:**

- $X_i = \{0, 1\}^n$
- $\Theta = [0, 1]$
- $\mathbb{P}_\theta^n = \text{Bernoulli}(\theta)^{\otimes n}$
- $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$

**Unbiasedness:**

Let  $\theta \in \Theta$ , then

$$\mathbb{E}_\theta \left[ \frac{1}{n} \sum_{i=1}^n X_i \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\theta[X_i] = \frac{1}{n} \sum_{i=1}^n \theta = \theta.$$

Thus,  $\mathbb{E}_\theta[\hat{\theta}_n] = \theta$ .

**Consistency:**

- We could use the Weak Law of Large Numbers (WLLN).
- Alternatively,

$$\text{Var}_\theta(\hat{\theta}_n) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}_\theta(X_i) = \frac{1}{n^2} \sum_{i=1}^n \theta(1 - \theta) = \frac{\theta(1 - \theta)}{n}$$

which tends to zero as  $n \rightarrow \infty$ .

It follows: For every  $\epsilon > 0$ ,

$$\mathbb{P}_\theta^n(|\hat{\theta}_n - \theta| > \epsilon) = \mathbb{P}_\theta^n(|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]| > \epsilon)$$

By Markov's inequality:

$$\mathbb{P}_\theta^n(|\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n]| > \epsilon) \leq \frac{\mathbb{E}_\theta \left[ (\hat{\theta}_n - \mathbb{E}[\hat{\theta}_n])^2 \right]}{\epsilon^2} = \frac{\text{Var}_\theta(\hat{\theta}_n)}{\epsilon^2} = \frac{\theta(1 - \theta)}{n\epsilon^2}$$

which tends to zero as  $n \rightarrow \infty$ . Thus,

$$(\hat{\theta}_n : n \geq 1) \text{ is consistent.} \quad \square$$



### 3.1 Maximum Likelihood Principle

Is there another way to motivate  $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n X_i$ ? Yes, it turns out it is the maximum likelihood estimator, i.e.,

MLE = "parameter which assigns the highest probability to the observed data."

**In our example**, each  $\mathbb{P}_\theta^n$  has a probability density (likelihood)

$$\begin{aligned}\mathbb{P}_\theta^n(x) &= \prod_{i=1}^n \mathbb{P}_\theta(x_i) = \prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \\ &= \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}.\end{aligned}$$

Fixing  $x \in \{0, 1\}^n$  and maximizing in  $\theta \in [0, 1]$  gives the following:

- If  $\sum_{i=1}^n x_i = 0$ , then  $\hat{\theta}_n = 0$  is the maximizer.
- If  $\sum_{i=1}^n x_i = n$ , then  $\hat{\theta}_n = 1$  is the maximizer.
- If  $\hat{\theta}_n \in \{1, \dots, n-1\}$ , then writing  $S_n = \sum_{i=1}^n x_i$  gives:

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathbb{P}_\theta^n(x) &= S_n \theta^{S_n-1} (1 - \theta)^{n-S_n-1} - (n - S_n) \theta^{S_n} (1 - \theta)^{n-S_n-1} = 0 \\ &\Leftrightarrow S_n(1 - \theta) - \theta(n - S_n) \\ &\Leftrightarrow \theta = \frac{S_n}{n}. \quad \square\end{aligned}$$

**Definition 6 (Dominated statistical model & MLE).** A model  $(\mathbb{P}_\theta^n)_{\theta \in \Theta}$  is called dominated if there exists a measure  $\mu$  on  $(\mathcal{X}, \mathcal{F})$  such that for every  $\theta \in \Theta$ ,  $\mathbb{P}_\theta \ll \mu$  or equivalently (by Radon-Nikodym), for all  $\theta \in \Theta$ , there is a probability density  $\frac{d\mathbb{P}_\theta}{d\mu}$  of  $\mathbb{P}_\theta$  with respect to  $\mu$ .

The MLE is defined as any  $\hat{\theta} \in \Theta$  that maximizes the function

$$\theta \mapsto \frac{d\mathbb{P}_\theta}{d\mu}(x) = \mathbb{P}_\theta(x).$$

*Remark 5 (Caveats).* • MLE might not be unique.

- MLE might not exist.
- It's not always clear that some selection  $\hat{\theta}(x) \in \arg \max_{\theta} \mathbb{P}_\theta(x)$  is a measurable function of  $x \in \mathcal{X}$ . However, there are measurable selection theorems that permit a measurable choice of  $\hat{\theta}$  under very general conditions.

*Remark 6.* In all the models we study, we will work with the Lebesgue measure (for continuous data) or the counting measure (for discrete data).

**Example 2 (Normal model).** Consider random samples  $X_1, \dots, X_n \sim N(\mu, \sigma^2)$  for some unknown  $\mu \in \mathbb{R}$ ,  $\sigma^2 > 0$ , and let  $\lambda$  denote the Lebesgue measure on  $\mathbb{R}^n$ .

$$\theta = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty).$$

Then the likelihood is:

$$\begin{aligned} L(\mu, \sigma^2 \mid x) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \left(\frac{X_i - \mu}{\sigma}\right)^2\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \|X - \mu \cdot 1_n\|^2\right), \end{aligned}$$

where by  $1_n$  we denote the vector of ones of dimension  $n$ .

Here, the MLE is given as:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i \quad [\text{Sample mean}], \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2 \quad [\text{Sample variance}].$$

$$\mathbb{E}_\theta[\hat{\mu}] = \mu, \quad \mathbb{E}_\theta[\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2, \quad \text{so } \hat{\sigma}^2 \text{ is biased.}$$

$\Rightarrow$  MLE is not always unbiased, and not always a "good" method.

## 3.2 Bayesian method

**Motivation** In Bayesian statistics, a key element is the prior distribution, which we denote by  $\pi$ , reflecting our "beliefs" about the parameter  $\theta \in \Theta$  before observing data ( $\pi$  is a probability measure on  $\Theta$ ).

A prior  $\pi$ , together with a model ( $P_\theta : \theta \in \Theta$ ), gives rise to a joint probability distribution for the pair  $(\theta, x) \in \Theta \times \mathcal{X}$ .

The Bayesian approach bases statistical inference on the posterior distribution of  $\theta$  conditioned on  $x$ .

**Joint probability:**

$$(\theta, x) \mapsto \pi(\theta)P_\theta(x)$$

conditional distribution of  $x \mid \theta$

**Posterior:**

$$\pi(\theta | x) = \frac{\pi(\theta)P_\theta(x)}{\int_{\Theta} \pi(\theta)P_\theta(x) d\theta}$$

*Remark 7.* Bayesian methods automatically generate "error bars" because the posterior is not an estimator but a whole probability distribution.

## Lecture 3

**Definition 7 (Prior, Posterior, Bayes' Rule).** Let  $\mathcal{F}_\Theta$  be a  $\sigma$ -algebra on  $\Theta$ , and suppose

$$\{P_\theta : \theta \in \Theta\}$$

is a dominated statistical model with densities  $p_\theta(x)$ , and assume that

$$(\theta, x) \mapsto p_\theta(x)$$

is "jointly measurable" (i.e., w.r.t.  $\sigma(\mathcal{F}_\Theta \times \mathcal{F})$ ).

Let  $\Pi$  be a prior distribution on  $\Theta$ , with density  $\pi(\theta)$  w.r.t. measure  $\nu(\cdot)$ . Then, define the posterior density

$$\pi(\theta | x) := \frac{p_\theta(x)\pi(\theta)}{\int_{\Theta} p_{\tilde{\theta}}(x) d\Pi(\tilde{\theta})}.$$

The corresponding probability measure  $\Pi(\cdot | x)$  is called the posterior distribution:

$$\begin{aligned} \Pi(B | x) &= \int_B \pi(\theta | x) d\nu(\theta), \quad B \in \mathcal{F}_\Theta. \\ &= \frac{\int_B p_\theta(x) \pi(\theta) d\nu(\theta)}{\int_{\Theta} p_{\tilde{\theta}}(x) d\Pi(\tilde{\theta})}, \\ &= \frac{\int_B p_\theta(x) d\Pi(\theta)}{\int_{\Theta} p_{\tilde{\theta}}(x) d\Pi(\tilde{\theta})}, \end{aligned}$$

*Remark 8.* Think of  $\Theta \subseteq \mathbb{R}^p$ ,  $\nu(\cdot)$  as a Lebesgue measure,  $\pi(\cdot)$  as a Lebesgue density.

**Exception:**  $\Theta = \{0, 1\}$  in hypothesis testing. Then, we'd take  $\nu(\cdot)$  to be the counting measure.

From the posterior, we can derive several estimators:

- **Maximum-a-posterior (MAP) estimator:**

$$\hat{\theta}_{\text{MAP}}(x) = \operatorname{argmax}_{\theta \in \Theta} \pi(\theta | x).$$

- **Posterior mean:** Say  $\Theta \subseteq \mathbb{R}^p$  convex

$$\hat{\theta}(x) = \int_{\Theta} \theta \pi(\theta | x) d\nu(\theta) \in \mathbb{R}^p.$$

**Back to Example 1.1:** *Binomial model:*  $\mathcal{X} = \{0, 1, \dots, n\}$ ,  $p_{\theta} = \text{Bin}(n, \theta)$ ,  $\theta \in \Theta = [0, 1]$ .

*Prior (uniform):*  $\Pi = \text{Unif}(0, 1)$ .

We know:

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{MLE}} \quad (\text{for the uniform prior}),$$

$$\hat{\theta}_{\text{MAP}} = \frac{X}{n}.$$

- **Posterior mean:**

$$\pi(\theta|x) = \frac{p_{\theta}(x)}{\int p_{\bar{\theta}}(x)d\bar{\theta}} \propto \binom{n}{k} \theta^x (1 - \theta)^{n-x}.$$

- **Binomial distribution:**

$$\text{Bin}(n, p)(k) = \binom{n}{k} p^k (1 - p)^{n-k},$$

where  $k \in \{0, \dots, n\}$  is the number of successes, and  $p$  is the probability of success, and  $n$  would be interpreted as the "number of coin flips".

$$\pi(\theta|x) \propto \theta^x (1 - \theta)^{n-x}.$$

and

$$\int_0^1 \pi(\theta | x) d\nu = 1$$

We conclude that  $\pi(\theta|x)$  is a **Beta-distribution** on  $[0, 1]$ ,

$$\text{Beta}(x + 1, n - x + 1).$$

The mean is given by:

$$\hat{\theta} = \frac{x + 1}{n + 2}.$$

*Remark 9 (Beta distribution).* The Beta distribution is defined as:

$$\text{Beta}(a, b), \quad a, b \geq 0.$$

The probability density function of the Beta distribution is given by:

$$P_{\text{Beta}(a,b)}(x) = x^a(1-x)^b.$$

**Definition 8 (Conjugate Bayesian models).** Let  $(P_\theta : \theta \in \Theta)$  be a statistical model. Then, some family  $\mathcal{D}$  of p.m.s on  $\Theta$  is called *conjugate* if

$$\Pi \in \mathcal{D} \implies \Pi(\cdot|x) \in \mathcal{D} \quad \text{for all } x \in \mathcal{X}.$$

**Examples:**

- $(\text{Bin}(n, \theta)) : \theta \in [0, 1], \quad \mathcal{D} = \text{Beta}(a, b), \quad a, b \geq 0.$
- $(\mathcal{N}(\mu, \sigma^2)) : \mu \in \mathbb{R}, \quad \mathcal{D} = \{\mathcal{N}(\mu, n^2), \mu \in \mathbb{R}, n^2 > 0\}, \quad \sigma^2 \text{ known.}$

## 4 Decision Theory

Here suppose that  $\Theta \subseteq \mathbb{R}^p$ .

**Definition 9 (Loss function).** A function  $\ell : \Theta \times \mathbb{R}^p \rightarrow [0, \infty)$  is a *loss function* if for every  $\theta \in \Theta$ ,  $\ell(\theta, \cdot)$  is measurable. Given some estimator  $\hat{\theta}$ , the expected loss is

$$\mathcal{R}(\theta, \hat{\theta}) = \mathbb{E}_{\theta} [\ell(\theta, \hat{\theta})].$$

**Example:** Take  $\ell(\theta, \hat{\theta}) = \|\theta - \hat{\theta}\|_{\mathbb{R}^p}^2$ . Then,

$$\mathcal{R}(\theta, \hat{\theta}) = \mathbb{E}_{\theta} [\|\theta - \hat{\theta}\|_{\mathbb{R}^p}^2]$$

is the mean squared error (MSE).

**Proposition 1 (Bias-Variance Decomposition).** Let  $\hat{\theta} \in L^2(\mathbb{P}_{\theta})$ . Then it holds that:

$$R(\hat{\theta}, \theta) = (\mathbb{E}_{\theta}[\hat{\theta}] - \theta)^2 + \text{Var}_{\theta}(\hat{\theta}).$$

*Proof.* We have

$$R(\hat{\theta}, \theta) = \mathbb{E}_{\theta}[(\hat{\theta} - \mathbb{E}_{\theta}[\hat{\theta}] + \mathbb{E}_{\theta}[\hat{\theta}] - \theta)^2]$$

Expanding the squared term:

$$= \mathbb{E}_{\theta}[(\hat{\theta} - \mathbb{E}_{\theta}[\hat{\theta}])^2] + (\mathbb{E}_{\theta}[\hat{\theta}] - \theta)^2 + 2\mathbb{E}_{\theta}[(\hat{\theta} - \mathbb{E}_{\theta}[\hat{\theta}])(\mathbb{E}_{\theta}[\hat{\theta}] - \theta)].$$

Since  $\mathbb{E}_{\theta}[\hat{\theta} - \mathbb{E}_{\theta}[\hat{\theta}]] = 0$ , the last term vanishes, leaving us with:

$$R(\hat{\theta}, \theta) = \text{Var}_{\theta}(\hat{\theta}) + (\mathbb{E}_{\theta}[\hat{\theta}] - \theta)^2. \quad \square$$

**Definition 10 (Minimax Risk).** Given an estimator  $\hat{\theta}$  in a model  $(\mathbb{P}_{\theta} : \theta \in \Theta)$ , the maximal risk of  $\hat{\theta}$  is

$$\sup_{\theta \in \Theta} R(\theta, \hat{\theta}).$$

The minimax risk of a model  $(\mathbb{P}_{\theta} : \theta \in \Theta)$  is given as

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}),$$

where the infimum is taken over all estimators. An estimator is called minimax if

$$\sup_{\theta \in \Theta} R(\theta, \hat{\theta}) = \inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}).$$

**Definition 11 (Bayes Risk).** Given a prior  $\pi$  on  $\Theta$ , the  $\pi$ -Bayes risk of a decision rule  $\delta$  for the loss function  $L$  is defined as

$$R_{\Pi}(\delta) = \mathbb{E}_{\Pi}[R(\delta, \theta)] = \int_{\Theta} R(\delta, \theta) \pi(\theta) d\theta = \int_{\Theta} \int_{\mathcal{X}} L(\delta(x), \theta) \pi(\theta) p_{\theta}(x) dx d\theta.$$

A  $\Pi$ -Bayes decision rule  $\hat{\theta}_{\Pi}$  is any decision rule that minimizes  $R_{\Pi}(\hat{\theta})$ .

SW:  $\ell$  instead of  $L$  below,  $p_{\theta}(x)$  instead of  $f(x, \theta)$  Note: Has been corrected.

**Definition 12 (Posterior Risk).** For a Bayesian model, the posterior risk  $R_{\Pi}$  is defined as the average loss under the posterior distribution for some observation  $x \in \mathcal{X}$ :

$$R_{\Pi(\cdot|x)}(\delta) = \mathbb{E}_{\Pi}[\ell(\delta(x), \theta)|x].$$

Here, the notation  $\mathbb{E}_{\Pi}[\cdot|x]$  stands for the expectation under the posterior distribution.

**Proposition 2 (Bayes Risk and Posterior Risk).** An estimator  $\delta$  that minimizes the  $\Pi$ -posterior risk  $R_{\Pi}$  also minimizes the  $\pi$ -Bayes risk  $R_{\pi}$ .

*Proof.* The  $\pi$ -Bayes risk can be rewritten as

$$\begin{aligned} R_{\pi}(\delta) &= \int_{\Theta} \mathbb{E}_{\theta}[\ell(\delta(X), \theta)] \pi(\theta) d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} \ell(\delta(x), \theta) \pi(\theta) p_{\theta}(x) dx d\theta \\ &= \int_{\mathcal{X}} \int_{\Theta} \ell(\delta(x), \theta) \frac{p_{\theta}(x) \pi(\theta)}{\int_{\Theta} p_{\theta'}(x) \pi(\theta') d\theta'} \times \underbrace{\int_{\Theta} p_{\theta'}(x) \pi(\theta') d\theta'}_{=: n(x) \geq 0} dx d\theta \\ &= \int_{\mathcal{X}} \mathbb{E}_{\Pi}[\ell(\delta(x), \theta)|x] n(x) dx. \end{aligned}$$

[Notation  $n(x)$  motivated by the word ‘normalising constant’].

Let  $\delta_{\Pi}$  be a decision rule that minimizes the posterior risk, i.e., such that for all  $x \in \mathcal{X}$ ,

$$\mathbb{E}_{\Pi}[\ell(\delta_{\Pi}(x), \theta)|x] \leq \mathbb{E}_{\Pi}[\ell(\delta(x), \theta)|x].$$

Multiplying by  $n(x) \geq 0$  and integrating on both sides over  $\mathcal{X}$  yields the desired result.  $\square$

**Example 3.** For the quadratic risk with the squared-loss, the posterior risk is minimized by taking  $\delta(X) = \mathbb{E}_\Pi[\theta|X]$ , by minimizing the quadratic function in  $\delta$ . Other losses will give other ways to minimize the posterior risk, and other Bayes decision rules.

**Proposition 3.** Let  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\theta)_{\theta \in \Theta})$  be a statistical model and let  $\hat{\theta}$  be an estimator. Then we have

$$\sup_{\theta \in \Theta} R(\hat{\theta}, \theta) = \sup_{\Pi} \int_{\Theta} R(\hat{\theta}, \theta) \Pi(d\theta),$$

where the supremum is taken over all prior distributions  $\Pi$ .

*Proof.* Obviously, we have

$$\int_{\Theta} R(\hat{\theta}, \theta) \Pi(d\theta) \leq \sup_{\theta \in \Theta} R(\hat{\theta}, \theta).$$

On the other hand, by using the prior distributions  $\delta_\theta$  (Dirac measure on  $\theta \in \Theta$ ), we obtain

$$\sup_{\Pi} \int_{\Theta} R(\hat{\theta}, \theta) \Pi(d\theta) \geq \int_{\Theta} R(\hat{\theta}, \theta) \delta_\theta(d\theta) = R(\hat{\theta}, \theta). \quad \square$$

[Note: In the following we use the notation  $\delta$  for decision rules while on the blackboard we used  $\hat{\theta}$  or  $\tilde{\theta}$ . If you want to adjust this please contact me so that I can give you access.]

**Proposition 4.** Let  $\pi$  be a prior on  $\Theta$  such that

$$R_\pi(\delta_\pi) = \sup_{\theta \in \Theta} R(\delta_\pi, \theta),$$

where  $\delta_\pi$  is a  $\pi$ -Bayes rule. Then it holds that

1. The rule  $\delta_\pi$  is minimax.
2. If  $\delta_\pi$  is unique Bayes, then it is unique minimax.

*Proof.* Let  $\delta$  be any decision rule. Then

$$\sup_{\theta \in \Theta} R(\delta, \theta) \geq \mathbb{E}_\pi[R(\delta, \theta)],$$

$$\int_{\Theta} R(\delta, \theta) \pi(\theta) d\theta \geq \mathbb{E}_\pi[R(\delta, \theta)],$$

$$\int_{\Theta} R(\delta, \theta) \pi(\theta) d\theta = R_\pi(\delta_\pi) = \sup_{\theta \in \Theta} R(\delta_\pi, \theta).$$

Taking the infimum over  $\delta$  gives the result.

2. If  $\delta_\pi$  is unique Bayes, the second inequality is strict for any  $\delta' \neq \delta_\pi$ . □



**Corollary 1.** *If a (unique) Bayes rule  $\delta_\pi$  has constant risk in  $\theta$ , then it is (unique) minimax.*

*Proof.* If a Bayes rule  $\delta_\pi$  has constant risk, then

$$R_\pi(\delta_\pi) = \mathbb{E}_\pi[R(\delta_\pi, \theta)] = \sup_{\theta \in \Theta} R(\delta_\pi, \theta),$$

where  $R(\delta_\pi, \theta)$  is constant in  $\theta$ . Uniqueness of the Bayes rule implies uniqueness of the minimax rule, as in part 2 of the former proposition.  $\square$

**Example 4.** Hence, if the maximal risk of a Bayes rule  $\delta_\pi$  equals the Bayes risk, then  $\pi$  is least favorable, and the corresponding Bayes rule is minimax.

- In a  $\text{Bin}(n, \theta)$  model, let  $\pi_{a,b}$  be a  $\text{Beta}(a, b)$  prior on  $\theta \in [0, 1]$ . Then the unique Bayes rule for  $\pi_{a,b}$  over the quadratic risk is the posterior mean  $\delta_{a,b} = \bar{\theta}_{a,b}$ . Trying to solve the equation

$$R(\delta_{a,b}, \theta) = \text{const.} \quad \forall \theta \in [0, 1]$$

we can find a prior  $\pi_{a^*, b^*}$  and a corresponding Bayes rule  $\delta_{\pi_{a^*, b^*}}$  of constant risk. It is therefore unique minimax, and different from the MLE (see Examples sheet).

- In a  $\mathcal{N}(\theta, 1)$  model,  $\bar{X}_n$  is minimax, as proved later.

## 4.1 Another optimality concept: Admissibility

**Definition 13.** A decision rule  $\delta$  is *inadmissible* if there exists  $\delta'$  such that

$$R(\delta', \theta) \leq R(\delta, \theta) \quad \forall \theta \in \Theta \quad \text{and} \quad R(\delta', \theta) < R(\delta, \theta) \quad \text{for some } \theta \in \Theta.$$

*Remark 10.* • The intuition is that there is no reason to choose an inadmissible estimator or decision rule: it would always be better to choose another estimator that dominates it.

- Admissibility is not the only criterion to evaluate an estimator: In most cases, a constant estimator will be admissible for the quadratic risk, but it is often not reasonable.

**Proposition 5.** 1. *A unique Bayes rule is admissible.*

2. *If  $\delta$  is admissible and has constant risk, then it is minimax.*

Proof may be done in the Examples sheet.

**Definition 14.** For a vector  $X \in \mathbb{R}^p$ , the *James–Stein estimator* is defined as

$$\delta^{JS}(X) = \left(1 - \frac{p-2}{\|X\|^2}\right) X.$$

In a Gaussian model  $X \sim \mathcal{N}(\theta, I_p)$  for  $\theta \in \mathbb{R}^p$  (with a single observation, to simplify notation), the risk of the MLE is given by

$$R(\hat{\theta}_{\text{MLE}}, \theta) = \mathbb{E}_\theta[\|X - \theta\|^2] = \sum_{j=1}^p \mathbb{E}_\theta[(X_j - \theta_j)^2] = p.$$

For  $X \sim \mathcal{N}(\theta, I_p)$  with  $p \geq 3$ , the risk of the James–Stein estimator satisfies for all  $\theta \in \mathbb{R}^p$

$$R(\delta^{JS}, \theta) < p.$$

## 5 Confidence Sets

**Definition 15 (Confidence Set).** Let  $(\mathcal{X}, \mathcal{F}, (\mathbb{P}_\theta)_{\theta \in \Theta})$  be a statistical model. For a given  $\alpha \in [0, 1]$ , consider sets  $C_{1-\alpha}(x) \subseteq \Theta$  for each  $x \in \mathcal{X}$ . Then  $C_{1-\alpha}(x)$  is called a random confidence set at level  $1 - \alpha$  (or with coverage probability  $1 - \alpha$ ) if

$$\forall \theta \in \Theta : \mathbb{P}_\theta(\theta \in C_{1-\alpha}) = \mathbb{P}_\theta(\{x \in \mathcal{X} : \theta \in C_{1-\alpha}(x)\}) \geq 1 - \alpha.$$

**Note:** The following example was only started in Lecture 4 and may be fully covered in Lecture 5.

**Example 5.** Consider the statistical model  $(\{0, 1\}^n, \mathcal{P}(\{0, 1\}^n), (\mathbb{P}_p)_{p \in [0, 1]})$  with  $\mathbb{P}_p = \text{Ber}(p)^{\otimes n}$  and independent observations  $X_k \sim \text{Ber}(p)$ ,  $k \in \{0, \dots, n\}$ . We are looking for a confidence interval  $C_{1-\alpha}$  around  $\hat{p} = \bar{X}_n = \frac{1}{n} \sum_{k=1}^n X_k$ , i.e.,

$$C_{1-\alpha} = [\bar{X}_n - a, \bar{X}_n + b] \quad (C_{1-\alpha}(x) = [\bar{X}_n(x) - a(x), \bar{X}_n(x) + b(x)])$$

should satisfy (where  $a$  and  $b$  might be random) that

$$1 - \alpha \leq \mathbb{P}_p(p \in C_{1-\alpha}) = \mathbb{P}_p(\bar{X}_n - a \leq p \leq \bar{X}_n + b) = \mathbb{P}_p(-b \leq \bar{X}_n - p \leq a).$$

Let  $t \mapsto F_p^n(t) := \mathbb{P}_p(\bar{X}_n - p \leq t)$  be the distribution function. Then,

$$\begin{aligned} \mathbb{P}_p(-b \leq \bar{X}_n - p \leq a) &= \mathbb{P}_p(\bar{X}_n - p \leq a) - \mathbb{P}_p(\bar{X}_n - p < -b) \\ &= F_p^n(a) - F_p^n(-b) + R_n, \end{aligned}$$

where  $R_n = \mathbb{P}_p(\bar{X}_n - p = -b)$ . Choose  $a, b$  as quantiles of  $\mathbb{P}_p^n$ , i.e.,  $a = (F_p^n)^{-1}(1 - \alpha/2)$  and  $-b = (F_p^n)^{-1}(\alpha/2)$  (with quantile function  $t \mapsto (F_p^n)^{-1}(t) := \inf\{t \in \mathbb{R} : F_p^n(t) \geq t\}$ ).

However,  $F_p^n$  and thus  $a, b$  are unknown. Consider two possibilities:

**Normal Approximation.** It holds that  $\mathbb{E}_p^n[X_k] = p$ ,  $\sigma := \text{Var}_p^n(X_k) = p(1-p)$ . By the central limit theorem, we have

$$\frac{\sqrt{n}}{\sigma} (\bar{X}_n - p) = \frac{1}{\sqrt{n}} \sum_{k=1}^n \frac{X_k - p}{\sigma} \xrightarrow{d} N(0, 1), \quad n \rightarrow \infty.$$

For  $Z \sim N(0, 1)$ , it holds that

$$\begin{aligned} F_p^n(a) &= \mathbb{P}_p^n(\bar{X}_n - p \leq a) = \mathbb{P}_p^n\left(\frac{\sqrt{n}}{\sigma}(\bar{X}_n - p) \leq \frac{\sqrt{n}}{\sigma}a\right) \approx \mathbb{P}(|Z| \leq \frac{\sqrt{n}}{\sigma}a) \\ &= \Phi\left(\frac{\sqrt{n}}{\sigma}a\right) = \Phi(z_\beta) \end{aligned}$$

for  $a := \frac{\sigma}{\sqrt{n}}z_\beta$  (where  $z_\beta$  is the  $\beta$ -quantile of the  $N(0, 1)$ -distribution, i.e.,  $\Phi(z_\beta) = \beta$ ). In particular,  $R_n = o(1)$  (i.e.,  $R_n \rightarrow 0$  as  $n \rightarrow \infty$ ). For  $a = b$  (since the  $N(0, 1)$ -distribution is symmetric) and because  $\Phi(-x) = 1 - \Phi(x)$ , it follows that

$$\mathbb{P}_p^n(p \in C_{1-\alpha}) = F_p^n(a) - F_p^n(-a) + R_n \approx \Phi(z_\beta) - (1 - \Phi(z_\beta)) + o(1) = 2\Phi(z_\beta) - 1 + o(1).$$

For  $\beta = 1 - \alpha/2$ ,  $C_{1-\alpha}$  is an asymptotically correct confidence interval. However,  $p$  and therefore  $\sigma$  and  $a$  are unknown. Solutions:

- Estimate  $\sigma = p(1-p) \leq 1/4$  to widen the confidence interval.
- For the empirical variance  $\hat{\sigma}^2 = \frac{1}{n} \sum_{k=1}^n (X_k - \bar{X}_n)^2$ , we have  $\hat{\sigma}^2 \rightarrow \sigma^2$  almost surely (by the law of large numbers). Using Slutsky's lemma (Lemma: For random variables  $(X_n, Y_n)_{n \geq 1}$  with  $X_n \xrightarrow{d} X$ ,  $Y_n \xrightarrow{p} c \in \mathbb{R}$  (where  $c$  is deterministic), it holds that  $X_n + Y_n \xrightarrow{d} X + c$  and  $X_n \cdot Y_n \xrightarrow{d} c \cdot X$ ), it follows that

$$\frac{\sqrt{n}}{\hat{\sigma}} (\bar{X}_n - p) = \frac{\sigma}{\hat{\sigma}} \cdot \frac{\sqrt{n}}{\sigma} (\bar{X}_n - p) \xrightarrow{d} N(0, 1), \quad n \rightarrow \infty.$$

From this, we derive  $a = \frac{\hat{\sigma}}{\sqrt{n}} z_{1-\alpha/2}$  (randomly chosen).

$$\begin{aligned} \mathbb{P}_p^n(p \in C_{1-\alpha}) &= \mathbb{P}_p^n(|\bar{X}_n - p| \leq a) = \mathbb{P}_p^n\left(\left|\frac{\sqrt{n}}{\hat{\sigma}}(\bar{X}_n - p)\right| \leq z_{1-\alpha/2}\right) \\ &\approx \mathbb{P}(|Z| \leq z_{1-\alpha/2}) = 2\Phi(z_{1-\alpha/2}) - 1 = 1 - \alpha. \end{aligned}$$

## 5.1 Hypothesis Testing

### 5.1.1 Basic Definitions

Let  $(P_\theta : \theta \in \Theta)$  be a statistical model, and let  $\Theta = \Theta_0 \cup \Theta_1$  be a partition:

- A **statistical test** is a measurable function of the data  $\varphi : (\mathcal{X}, \mathcal{F}) \rightarrow [0, 1]$ .
- If  $\varphi(x) \in \{0, 1\}$  for all  $x \in \mathcal{X}$ , then  $\varphi$  is a **non-randomized test**; otherwise, it is **randomized**.
- $H_0 : \theta \in \Theta_0$  is the **null hypothesis**.
- $H_1 : \theta \in \Theta_1$  is the **alternative hypothesis**.
- The map  $\theta \rightarrow \beta_\varphi(\theta) = P_\theta(\varphi = 1)$  is called the **power function** of a test  $\varphi$ .

### 5.1.2 Type I and Type II Errors

- For  $\theta \in \Theta_0$ ,  $\beta_\varphi(\theta)$  represents the **Type I error** (wrongly rejecting the null).
- For  $\theta \in \Theta_1$ ,  $1 - \beta_\varphi(\theta)$  represents the **Type II error** (failing to reject the alternative when it is true).

$$1 - \beta_\varphi(\theta) = 0 \quad \Theta_0 \quad \Theta_1 \quad \Theta$$

**Note:**

$$1 - P_\theta(\varphi = 1) = P_\theta(\varphi = 0) = P_\theta \text{ (wrongly accepting the null)}$$

### 5.1.3 Level and Uniformly Most Powerful Tests

**Definition 16 (Level).** A test  $\varphi : \mathcal{X} \rightarrow [0, 1]$  has **level**  $\alpha \in [0, 1]$  if

$$\sup_{\theta \in \Theta_0} \beta_\varphi(\theta) \leq \alpha.$$

**Definition 17 (Uniformly Most Powerful Test).** Given a level  $\alpha \in (0, 1)$ ,  $\varphi : \mathcal{X} \rightarrow [0, 1]$  is called **uniformly most powerful (UMP)** if, for every other test  $\varphi'$  of level  $\alpha$  and all  $\theta \in \Theta_1$ ,

$$\beta_\varphi(\theta) \geq \beta_{\varphi'}(\theta).$$

### 5.1.4 The Neyman-Pearson Lemma

The Neyman-Pearson Lemma provides a basis for constructing the most powerful tests for simple hypotheses:

**Theorem 1 (Neyman-Pearson Lemma).** *Let  $\Theta_0 = \{\theta_0\}$  and  $\Theta_1 = \{\theta_1\}$  be simple hypotheses:*

1. **Existence:** *There exists a test  $\varphi$  and a constant  $k \in [0, \infty)$  such that  $P_{\theta_0}(\varphi = 1) = \alpha$ , with*

$$\varphi(x) = \begin{cases} 1, & \text{if } \frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} > k \\ 0, & \text{if } \frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} < k \end{cases}$$

*Here,  $p_{\theta_1}$  and  $p_{\theta_0}$  are densities with respect to some dominated measure  $\mu$ .*

2. **Sufficiency:** *If  $\varphi$  satisfies  $P_{\theta_0}(\varphi = 1) = \alpha$  and the above form, then  $\varphi$  is a UMP level  $\alpha$  test.*
3. **Necessity:** *If  $\varphi_k$  is UMP for level  $\alpha$ , then it must be of the form shown above.*

### 5.1.5 Proof of the Neyman-Pearson Lemma

1. Define the likelihood ratio  $r(x) = \frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} \in [0, \infty)$ . Let  $F_0$  be the CDF of  $r(x)$  under  $P_{\theta_0}$ .

$$F_0(t) = P_{\theta_0}(r(x) \leq t).$$

Define  $\alpha(t) = 1 - F_0(t) = P_{\theta_0}(r(x) > t)$  and note:

- $\alpha$  is right-continuous:

$$\lim_{\epsilon \rightarrow 0} \alpha(t + \epsilon) = P_{\theta_0}(r(x) > t).$$

- $\alpha$  is non-increasing.
- $\alpha$  has left limits.

$\alpha$  is **cadlag**: It is continuous from the right and has a left limit.

There exists  $k \in [0, \infty)$  such that  $\alpha \leq \alpha(k^-)$  and  $\alpha \geq \alpha(k)$ .

We define the test

$$\varphi(x) = \begin{cases} 1 & \text{if } r(x) > k, \\ \gamma & \text{if } r(x) = k, \\ 0 & \text{if } r(x) < k. \end{cases}$$

Set  $\gamma = \frac{\alpha - \alpha(k)}{\alpha(k^-) - \alpha(k)}$ .

The level of  $\varphi$  is

$$\begin{aligned} E_{\theta_0}[\varphi(x)] &= P_{\theta_0}(\varphi(x) = 1). \\ &= P_{\theta_0}(r(x) > k) + P_{\theta_0}(r(x) = k) \cdot \gamma = \alpha. \end{aligned}$$

# 6 Lecture 6: Neyman-Pearson Lemma and Likelihood Ratio Tests

## Neyman-Pearson Lemma

### Power of a Test

The **power** of a test is defined as:

$$E_{\theta_1}[\varphi] = P_{\theta_1}(\varphi = 1)$$

### Likelihood Ratio Test

The **likelihood ratio** is given by:

$$\Lambda(x) = \frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} = r(x)$$

### Likelihood Ratio (LR) Test

The LR test is defined as:

$$\varphi(x) = \begin{cases} 1 & \text{if } r(x) > k, \\ \gamma & \text{if } r(x) = k, \\ 0 & \text{if } r(x) < k, \end{cases}$$

where  $k \in [0, \infty)$  and  $\gamma \in [0, 1]$ .

**Note:** LR tests are Uniformly Most Powerful (UMP) for simple hypothesis testing:

- Given a significance level  $\alpha$ , if the LR test satisfies  $E_{\theta_0}[\varphi] = \alpha$ , it controls the Type I error.
- The LR test minimizes the Type II error:

$$E_{\theta_1}[\varphi] \geq E_{\theta_1}[\varphi'] \quad \forall \varphi'$$



### Continuation of Proof (Part of UMP)

Let  $\varphi'$  be another level  $\alpha$  test such that  $E_{\theta_0}[\varphi'] \leq \alpha$ .

**Goal:** Show that  $E_{\theta_1}[\varphi] \geq E_{\theta_1}[\varphi']$ .

Let  $\mu$  be the dominating measure. Consider:

$$\int (\varphi(x) - \varphi'(x)) (p_{\theta_1}(x) - kp_{\theta_0}(x)) d\mu(x) = 0$$

**Claim:**  $p \geq 0$ .

**Observation:**

- If  $p_{\theta_1}(x) - kp_{\theta_0}(x) > 0$ , then  $\frac{p_{\theta_1}(x)}{p_{\theta_0}(x)} > k$ , implying  $\varphi(x) = 1$ .
- If  $p_{\theta_1}(x) - kp_{\theta_0}(x) < 0$ , then  $\varphi(x) = 0$ .
- If  $p_{\theta_1}(x) - kp_{\theta_0}(x) = 0$ , then the integrand is 0.

Thus,  $p = 0$ , leading to:

$$\int (\varphi - \varphi') p_{\theta_1} d\mu = \int (\varphi - \varphi') p_{\theta_0} d\mu = k [E_{\theta_0}[\varphi] - E_{\theta_0}[\varphi']] \geq 0$$

Therefore,

$$E_{\theta_1}[\varphi] \geq E_{\theta_1}[\varphi']$$

### Part (3) UMP $\Rightarrow$ LR

Assume  $\varphi^*$  is a UMP test with  $E_{\theta_0}[\varphi^*] = \alpha$ . Let  $\varphi$  be the LR test satisfying  $E_{\theta_0}[\varphi] = \alpha$ .

**Goal:** Show that  $\varphi = \varphi^*$  almost everywhere except on  $\{r(x) = k\}$ .

Define the sets:

$$x^+ = \{x : \varphi(x) > \varphi^*(x)\}, \quad x^- = \{x : \varphi(x) < \varphi^*(x)\}, \quad x^0 = \{x : \varphi(x) = \varphi^*(x)\}$$

$$\tilde{x} = (x^+ \cup x^-) \cap \{x : p_{\theta_1}(x) \neq kp_{\theta_0}(x)\}$$

It suffices to show  $\mu(\tilde{x}) = 0$ .

On  $\tilde{x}$ :

$$(\varphi - \varphi^*)(p_{\theta_1} - kp_{\theta_0}) > 0$$

If  $\mu(\tilde{x}) > 0$ , then:

$$\int_{\mathcal{X}} (\varphi - \varphi^*)(p_{\theta_1} - kp_{\theta_0}) d\mu \geq 0$$

$$\int_{\tilde{x}} (\varphi - \varphi^*)(p_{\theta_1} - kp_{\theta_0}) d\mu \geq 0$$

However,

$$E_{\theta_1}[\varphi] - E_{\theta_1}[\varphi^*] > k [E_{\theta_0}[\varphi] - E_{\theta_0}[\varphi^*]] \geq 0$$

This leads to a contradiction, implying  $\mu(\tilde{x}) = 0$  and thus  $\varphi = \varphi^*$  almost everywhere.

### Example: Gaussian Location Model

Consider:

$$X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu, \sigma^2)$$

Testing:

$$H_0 : \mu = \mu_0 \quad \text{vs.} \quad H_1 : \mu = \mu_1, \quad \mu_0 < \mu_1$$

The likelihood ratio is:

$$\frac{p_1(X_1, \dots, X_n)}{p_0(X_1, \dots, X_n)} = \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_1)^2 + \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_0)^2 \right)$$

Simplifying:

$$= \exp \left( -\frac{n}{2\sigma^2} (\mu_1^2 - \mu_0^2) - \frac{2(\mu_1 - \mu_0)}{\sigma^2} \sum_{i=1}^n X_i \right) \geq K_\alpha$$

This implies:

$$\frac{1}{n} \sum_{i=1}^n X_i \geq K_\alpha, \quad \text{for some } K_\alpha \in \mathbb{R}$$

To determine  $K_\alpha$ :

$$\bar{X}_n := \frac{1}{n} \sum X_i \stackrel{H_0}{\sim} \mathcal{N} \left( \mu_0, \frac{\sigma^2}{n} \right)$$

Thus:

$$P_{H_0} \left( \bar{X}_n \geq K_\alpha \right) = 1 - \Phi \left( \frac{\sqrt{n}}{\sigma} (K_\alpha - \mu_0) \right)$$

Solving for  $K_\alpha$ :

$$K_\alpha = \mu_0 + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha)$$

Therefore, the LR test is:

$$\varphi(X_1, \dots, X_n) = \begin{cases} 1 & \text{if } \bar{X}_n \geq \mu_0 + \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha), \\ 0 & \text{otherwise.} \end{cases}$$

## Corollary

Consider simple hypothesis testing. Let  $\varphi$  be a UMP test at level  $\alpha$ . Then:

$$\alpha = E_{H_0}[\varphi] \leq E_{\theta_1}[\varphi]$$

Suppose  $E_{\theta_1}[\varphi] = E_{\theta_1}[\varphi_0]$ . Then  $\varphi_0$  is also UMP, implying  $\varphi_0$  is an LR test:

$$\varphi_0 = \begin{cases} 1 & \text{if } \frac{p_{\theta_1}}{p_{\theta_0}} \geq K \quad \text{a.s., for some } K, \\ 0 & \text{otherwise.} \end{cases}$$

Since  $\varphi_0 \in \{\varphi, \beta\}$ , it follows that  $p_{\theta_1} = K p_{\theta_0}$  almost surely.

Moreover:

$$\int p_{\theta_0} d\mu = K \int p_{\theta_0} d\mu = 1 \quad \Rightarrow \quad K = 1$$

## Correspondence Theorem

**Statement:** There is a correspondence between tests and confidence regions.

$$\text{Tests} \quad \longleftrightarrow \quad \text{Confidence regions } C(x)$$

with

$$\Pr_{\theta}(\theta \in C(x)) \geq 1 - \alpha$$

and

$$\Pr_{\theta}(\phi_{\theta}(x) = 1) = \alpha$$

**Theorem:** Let  $\{P_{\theta} : \theta \in \Theta\}$  be a statistical model and  $\alpha \in (0, 1)$ .

(i) If  $C = C(X)$  is a level- $\alpha$  confidence set, then

$$\phi_{\theta_0}(x) = \mathbb{I}\{\theta_0 \notin C(x)\}$$

is a level- $\alpha$  test for  $H_0 : \theta = \theta_0$  vs.  $H_1 : \theta \neq \theta_0$ .

(ii) If  $\{\phi_{\theta} : \theta \in \Theta\}$  is a family of level- $\alpha$  tests, then

$$C(X) = \{\theta \in \Theta : \phi_{\theta}(X) = 0\}$$

is a  $1 - \alpha$  confidence set.

**Proof:**

(i)

$$\Pr_{\theta_0}(\phi_{\theta_0} = 1) = \Pr_{\theta_0}(\theta_0 \notin C(X)) \leq \alpha$$

(ii)

$$\Pr_{\theta}(\theta \notin C(X)) = \Pr_{\theta}(\phi_{\theta}(X) = 1) \leq \alpha$$

## UMPT Tests in Models with Monotone Likelihoods

**Proposition:** Let  $\Theta \subseteq \mathbb{R}$ . Consider testing:

$$H_0 : \theta \leq \theta_0 \quad \text{vs.} \quad H_1 : \theta > \theta_0,$$

for some  $\theta_0 \in \mathbb{R}$ .

Assume there exists a test statistic  $T : X \rightarrow \mathbb{R}$  and a function  $h : \mathbb{R} \times \Theta \times \Theta \rightarrow \mathbb{R}$  such that:

$$\frac{P_{\theta}(X)}{P_{\tilde{\theta}}(X)} = h(T(X), \theta, \tilde{\theta})$$

and for all  $\theta \geq \tilde{\theta}$ , the function  $t \mapsto h(t, \theta, \tilde{\theta})$  is monotone increasing.

**Conclusion:** LR tests are also UMP for level  $\alpha$ . Specifically, the LR test of  $H_0 : \theta = \theta_0$  against  $H_1 : \theta = \theta_1$  for any  $\theta_1 > \theta_0$  will be UMP.

# 7 Linear Models

## 7.1 Introduction to Linear Regression Models

Linear regression models are fundamental tools in statistical analysis, enabling us to understand and quantify the relationship between a dependent variable and one or more independent variables. In this lecture, we will explore the simplest form of linear regression, discuss the statistical framework underpinning it, and delve into key estimation techniques and their optimality properties.

## 7.2 Simple Linear Regression Model

Consider the simplest scenario where we model the relationship between a dependent variable  $Y_i$  and an independent variable  $X_i$  using a linear relationship:

$$Y_i = aX_i + b + \varepsilon_i \quad (7.1)$$

for  $i = 1, \dots, n$ , where:

- $a$  and  $b$  are unknown parameters representing the slope and intercept, respectively.
- $\varepsilon_i$  is the error term, assumed to be centered, i.e.,  $E(\varepsilon_i) = 0$ , and having constant variance  $\text{Var}(\varepsilon_i) = \sigma^2$ .

We further assume that the error terms are normally distributed,  $\varepsilon_i \sim N(0, \sigma^2)$ , with  $\sigma$  known.

### 7.2.1 Statistical Model

The statistical model can be formally defined as:

$$(\mathbb{R}, \mathcal{B}(\mathbb{R}), \left( \bigotimes_{i=1}^n N(aX_i + b, \sigma^2) \right)_{(a,b) \in \mathbb{R}^2})$$

Here:

- $\mathbb{R}$  denotes the real line.
- $\mathcal{B}(\mathbb{R})$  is the Borel sigma-algebra on  $\mathbb{R}$ .
- $\bigotimes_{i=1}^n N(aX_i + b, \sigma^2)$  represents the product measure of the normal distributions for each observation.

### 7.2.2 Likelihood Function

Within this statistical model, the likelihood function is given by:

$$L((a, b) \mid y) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - aX_i - b)^2}{2\sigma^2}\right)$$

Simplifying, we obtain:

$$L((a, b) \mid y) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - aX_i - b)^2\right)$$

### 7.2.3 Maximum Likelihood Estimation (MLE)

The maximum likelihood estimators (MLE) for the parameters  $a$  and  $b$  are obtained by maximizing the likelihood function  $L((a, b) \mid y)$ . Equivalently, since the logarithm is a monotonically increasing function, we can maximize the log-likelihood:

$$\log L((a, b) \mid y) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - aX_i - b)^2$$

Maximizing the log-likelihood is equivalent to minimizing the sum of squared residuals:

$$(\hat{a}, \hat{b}) = \arg \min_{(a, b) \in \mathbb{R}^2} \sum_{i=1}^n (y_i - aX_i - b)^2$$

Provided that  $X_i \neq X_j$  for  $i \neq j$ , the least squares problem has a unique solution, historically attributed to Gauss (1801), given by:

$$(\hat{a}, \hat{b}) = \left( \frac{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}, \bar{Y} - \hat{a}\bar{X} \right)$$

where  $\bar{X}$  and  $\bar{Y}$  denote the sample means of  $X$  and  $Y$ , respectively.

## 7.3 General Linear Models

To generalize the simple linear regression model, we introduce the framework of linear models in multiple dimensions.

**Definition 18 (Linear Model).** A random vector  $Y = (Y_1, \dots, Y_n)^T \in \mathbb{R}^n$  stems from a linear model if there exists a parameter vector  $\beta \in \mathbb{R}^p$ , a matrix  $X \in \mathbb{R}^{n \times p}$ , and a random vector  $\varepsilon \in \mathbb{R}^n$  such that

$$Y = X\beta + \varepsilon$$

1. **Regular Linear Model:** A linear model is called *regular* if
  - a) The number of parameters does not exceed the sample size, i.e.,  $p \leq n$ .
  - b) The design matrix  $X$  has full rank,  $\text{rank}(X) = p \leq n$ , ensuring a unique solution.
  - c) The error vector satisfies  $E(\varepsilon) = 0$ , meaning the noise is centered.
  - d) The covariance matrix of the errors is positive definite,  $\Sigma = \text{Cov}(\varepsilon_i, \varepsilon_j)_{i,j \in [n]}$ .
2. **Ordinary Linear Model:** A linear model is called *ordinary* if  $\Sigma = \sigma^2 I_n$ , where  $I_n$  is the  $n \times n$  identity matrix. Typically, the noise is assumed to be Gaussian in this case.

*Remark 11.* 1. **Terminology:** Various synonyms are used in the literature:

- $Y$ : Dependent variable, response, regressand.
  - $X$ : Independent variable, predictor, design matrix, regressor.
  - $\varepsilon$ : Error, perturbation, regression function.
2. **Covariance Matrix  $\Sigma$ :** The matrix  $\Sigma$  is symmetric and diagonalizable, i.e.,  $\Sigma = UDU^T$  for some orthogonal matrix  $U$  and diagonal matrix  $D = \text{diag}(\lambda_1, \dots, \lambda_n) \in \mathbb{R}^{n \times n}$ .
  3. **Positive Semi-definiteness:** The covariance matrix  $\Sigma$  is positive semi-definite, meaning  $\lambda_i \geq 0$  for all  $i$ , which can be shown as:

$$\langle \Sigma u, u \rangle = \langle E[(\varepsilon - E[\varepsilon])(\varepsilon - E[\varepsilon])^T]u, u \rangle = E[(\varepsilon - E[\varepsilon])^2] \geq 0, \quad \forall u \in \mathbb{R}^n$$

4. **Positive Definiteness and Inverses:** If  $\Sigma$  is positive definite (i.e.,  $\lambda_i > 0$  for all  $i$ ), then its inverse  $\Sigma^{-1}$  and square root  $\Sigma^{-1/2}$  exist and can be expressed as:

$$\Sigma^{-1} = UD^{-1}U^T \quad \text{and} \quad \Sigma^{-1/2} = UD^{-1/2}U^T$$

5. **Random Design:** If the matrix  $X$  is not deterministic but random, the model is referred to as having a *random design*.

## 7.4 Least Squares Estimation

In the context of a regular linear model, the least squares estimator (LSE) seeks to minimize the weighted sum of squared residuals. Specifically, the LSE  $\hat{\beta}$  satisfies:

$$\|\sigma^{-1/2}(Y - X\hat{\beta})\|^2 = \inf_{\beta \in \mathbb{R}^p} \|\sigma^{-1/2}(Y - X\beta)\|^2 = \inf_{\beta \in \mathbb{R}^p} \|\Sigma^{-1/2}Y - X_{\Sigma}\beta\|^2$$

where  $X_{\Sigma} = \Sigma^{-1/2}X$ .

### 7.4.1 Geometric Interpretation

The estimator  $X_{\Sigma}\hat{\beta}$  represents the point within the subspace:

$$U = \{X_{\Sigma}\beta \mid \beta \in \mathbb{R}^p\} \subseteq \mathbb{R}^n$$

that is closest to the vector  $\Sigma^{-1/2}Y$  in terms of Euclidean distance. Formally, this can be expressed using the orthogonal projection  $\Pi_U$  onto  $U$ :

$$X_{\Sigma}\hat{\beta} = \Pi_U(\Sigma^{-1/2}Y)$$

The orthogonal projection satisfies:

- $\Pi_U u = u$  for all  $u \in U$ .
- $\langle \Pi_U v - v, u \rangle = 0$  for all  $u \in U$  and  $v \in \mathbb{R}^n$ .

Provided that  $(X_{\Sigma}^T X_{\Sigma})^{-1}$  exists, we can confirm by direct computation that the projection operator  $\Pi_U$  is given by:

$$\Pi_U = X_{\Sigma}(X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T$$

For any  $u = X_{\Sigma}\beta$ , we have:

$$X_{\Sigma}(X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T X_{\Sigma}\beta = X_{\Sigma}\beta = u$$

By symmetry, for any  $u \in U$  and  $v \in \mathbb{R}^n$ :

$$\langle \Pi_U v - v, u \rangle = \langle v, \Pi_U u \rangle - \langle v, u \rangle = \langle v, u \rangle - \langle v, u \rangle = 0$$



## 7.5 Representation for the Least Squares Estimator

**Lemma 1 (Representation for the LSE).** *Consider a regular linear model. Then the least squares estimator (LSE) exists uniquely and is given by:*

$$\hat{\beta} = (X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T \Sigma^{-1/2} Y = X_{\Sigma}^+ \Sigma^{-1/2} Y$$

where  $X_{\Sigma}^+$  denotes the Moore-Penrose pseudoinverse of  $X_{\Sigma}$ .

*Proof.* Since  $X$  has full rank and  $p \leq n$ , the matrix  $X_{\Sigma}^T X_{\Sigma}$  is invertible. Suppose  $v \in \ker(X_{\Sigma}^T X_{\Sigma})$ , then:

$$0 = v^T X_{\Sigma}^T X_{\Sigma} v = (X_{\Sigma} v)^T (X_{\Sigma} v) = \|X_{\Sigma} v\|^2 = \|\Sigma^{-1/2} X v\|^2$$

This implies  $\|X v\|^2 = 0$ , hence  $X v = 0$ . Given that  $X$  has full rank, the only solution is  $v = 0$ . Therefore,  $X_{\Sigma}^T X_{\Sigma}$  is invertible.

The projection of  $\Sigma^{-1/2} Y$  onto  $U$  is:

$$X_{\Sigma} \hat{\beta} = \Pi_U(\Sigma^{-1/2} Y) = X_{\Sigma} (X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T \Sigma^{-1/2} Y$$

Multiplying both sides by  $X_{\Sigma}^T$ :

$$X_{\Sigma}^T X_{\Sigma} \hat{\beta} = X_{\Sigma}^T X_{\Sigma} (X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T \Sigma^{-1/2} Y = X_{\Sigma}^T \Sigma^{-1/2} Y$$

Hence, solving for  $\hat{\beta}$ :

$$\hat{\beta} = (X_{\Sigma}^T X_{\Sigma})^{-1} X_{\Sigma}^T \Sigma^{-1/2} Y \quad \square$$

*Remark 12.* 1. If  $p > n$ , the matrix  $X_{\Sigma}^T X_{\Sigma}$  is not invertible, and the LSE is not unique. In this case, the set of solutions is a  $(p - n)$ -dimensional subspace, and each solution interpolates the data perfectly, i.e.,

$$\{\beta \in \mathbb{R}^p \mid \|\Sigma^{-1/2} Y - X_{\Sigma} \beta\|^2 = 0\}$$

## 7.6 Optimality of the Least Squares Estimator

The least squares estimator possesses several optimality properties under the ordinary linear model. This is formalized in the Gauss-Markov Theorem.

**Theorem 2 (Gauss-Markov Theorem).** Consider an ordinary linear model with  $\sigma > 0$ . Then:

1. The least squares estimator  $\hat{\beta} = (X^T X)^{-1} X^T Y$  is a linear and unbiased estimator for the parameter  $\beta$ .
2. For any desired linear combination of parameters  $\alpha = \langle \beta, v \rangle$  where  $v \in \mathbb{R}^p$ , the estimator  $\hat{\alpha} = \langle \hat{\beta}, v \rangle$  is the best linear unbiased estimator (BLUE) of  $\alpha$ . This means that  $\hat{\alpha}$  has the smallest variance among all linear unbiased estimators of  $\alpha$ .
3. The estimator  $\hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|^2}{n-p}$  is an unbiased estimator of  $\sigma^2$ .

*Proof.* 1. **Linearity and Unbiasedness of  $\hat{\beta}$ :**

The estimator  $\hat{\beta}$  is linear because it can be expressed as a linear transformation of  $Y$ :

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

For any two vectors  $y$  and  $\tilde{y}$  in  $\mathbb{R}^n$ ,

$$\hat{\beta}(y + \tilde{y}) = (X^T X)^{-1} X^T (y + \tilde{y}) = \hat{\beta}(y) + \hat{\beta}(\tilde{y})$$

demonstrating linearity.

To show unbiasedness, compute the expectation of  $\hat{\beta}$ :

$$\begin{aligned} E[\hat{\beta}] &= (X^T X)^{-1} X^T E[Y] \\ &= (X^T X)^{-1} X^T E[X\beta + \varepsilon] \\ &= (X^T X)^{-1} X^T (X\beta + E[\varepsilon]) \\ &= (X^T X)^{-1} X^T X\beta \quad (\text{since } E[\varepsilon] = 0) \\ &= \beta \end{aligned}$$

Hence,  $\hat{\beta}$  is unbiased.

## 2. Optimality of $\hat{\alpha}$ as BLUE:

Consider a linear unbiased estimator  $\tilde{\alpha}$  for  $\alpha = \langle \beta, v \rangle$ . Since  $\tilde{\alpha}$  is linear, there exists a vector  $w \in \mathbb{R}^n$  such that:

$$\tilde{\alpha} = \langle Y, w \rangle$$

For  $\tilde{\alpha}$  to be unbiased, we require:

$$E[\tilde{\alpha}] = \langle E[Y], w \rangle = \langle X\beta, w \rangle = \langle \beta, X^T w \rangle = \alpha = \langle \beta, v \rangle$$

This implies that:

$$v = X^T w$$

The variance of  $\tilde{\alpha}$  is:

$$\text{Var}(\tilde{\alpha}) = \text{Var}(\langle Y, w \rangle) = \text{Var}(\langle X\beta + \varepsilon, w \rangle) = \text{Var}(\langle \varepsilon, w \rangle) = \sigma^2 \|w\|^2$$

Now, consider the variance of  $\hat{\alpha} = \langle \hat{\beta}, v \rangle$ :

$$\begin{aligned} \text{Var}(\hat{\alpha}) &= \text{Var}(\langle \hat{\beta} - \beta, v \rangle) \\ &= \text{Var} \left( \langle (X^T X)^{-1} X^T \varepsilon, v \rangle \right) \\ &= \text{Var} \left( \langle \varepsilon, X (X^T X)^{-1} v \rangle \right) \\ &= \sigma^2 \|X (X^T X)^{-1} v\|^2 \\ &= \sigma^2 \|X (X^T X)^{-1} X^T w\|^2 \quad (\text{since } v = X^T w) \\ &= \sigma^2 \|\Pi_U w\|^2 \\ &\leq \sigma^2 \|w\|^2 \quad (\text{since projection does not increase the norm}) \end{aligned}$$

Therefore,  $\text{Var}(\hat{\alpha}) \leq \text{Var}(\tilde{\alpha})$ , showing that  $\hat{\alpha}$  has the smallest variance among all linear unbiased estimators of  $\alpha$ .

### 3. Unbiasedness of $\hat{\sigma}^2$ :

The estimator  $\hat{\sigma}^2$  is given by:

$$\hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}\|^2}{n - p}$$

To show that  $\hat{\sigma}^2$  is unbiased, compute its expectation:

$$\begin{aligned} E[\hat{\sigma}^2] &= \frac{E[\|Y - X\hat{\beta}\|^2]}{n - p} \\ &= \frac{E[\|\varepsilon\|^2 - \|X\hat{\beta}\|^2 + 2\langle \varepsilon, X\hat{\beta} \rangle]}{n - p} \\ &= \frac{E[\|\varepsilon\|^2]}{n - p} \quad (\text{since } E[\langle \varepsilon, X\hat{\beta} \rangle] = 0) \\ &= \frac{n\sigma^2}{n - p} \quad (\text{since } \|\varepsilon\|^2 \text{ is chi-squared with } n \text{ degrees of freedom}) \\ &= \sigma^2 \quad (\text{since } E[\|\varepsilon\|^2] = n\sigma^2) \end{aligned}$$

Hence,  $\hat{\sigma}^2$  is an unbiased estimator of  $\sigma^2$ . □

*Remark 13.* 1. The Gauss-Markov theorem establishes that under the ordinary linear

model assumptions, the least squares estimator is the best (in the sense of having the smallest variance) among all linear unbiased estimators. This optimality holds without requiring the error terms to be normally distributed.

## **7.7 Conclusion**

In this lecture, we have introduced the fundamental concepts of linear regression models, both in their simplest form and in a more general framework. We discussed the estimation of model parameters using the method of least squares, explored the geometric interpretation of the estimator, and established its optimality through the Gauss-Markov theorem. Understanding these foundational elements is crucial for further studies in statistical modeling and inference.

## 8 Lecture 9: Ridge Regression and Its Properties

In this lecture, we delve into Ridge Regression, a technique used to address multicollinearity in linear models. We will explore its derivation, properties, and how it compares to Ordinary Least Squares (OLS). Additionally, we will discuss the bias-variance trade-off inherent in Ridge Regression and present propositions regarding its Mean Squared Error (MSE).

### Proof of Ridge Regression Estimator

We begin by deriving the Ridge Regression estimator through optimization of the penalized loss function.

**Objective:** Minimize the cost function

$$\mathcal{J}(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \frac{\sigma^2}{\tau^2} \|\beta\|^2$$

with respect to the coefficient vector  $\beta$ .

**Steps:**

1. **\*\*Compute the Gradient:\*\***

To find the minimum, we take the gradient of  $\mathcal{J}(\beta)$  with respect to  $\beta$ :

$$\nabla_{\beta} \mathcal{J}(\beta) = 2\mathbf{X}^{\top}(\mathbf{Y} - \mathbf{X}\beta) + \frac{2\sigma^2}{\tau^2} \beta$$

2. **\*\*Set Gradient to Zero:\*\***

Setting the gradient equal to zero to find the critical point:

$$2\mathbf{X}^{\top}(\mathbf{Y} - \mathbf{X}\beta) + \frac{2\sigma^2}{\tau^2} \beta = 0$$

Simplifying:

$$(\mathbf{X}^{\top} \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I}) \beta = \mathbf{X}^{\top} \mathbf{Y}$$

3. **Solve for  $\beta$ :**

Finally, we solve for  $\beta$ :

$$\beta = (\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}$$

This is the Ridge Regression estimator, which introduces a penalty term  $\lambda = \frac{\sigma^2}{\tau^2}$  to shrink the coefficients.

## Posterior Mean

In a Bayesian framework, the posterior mean of  $\beta$  given the data can be expressed as:

$$\mu_{\text{post}} = \Sigma_{\text{post}}^{-1} (\mathbf{X}^\top \mathbf{Y} + \mathbf{M}_0^{-1} \mu_0)$$

Upon substituting  $\mathbf{M}_0^{-1} = \frac{\sigma^2}{\tau^2} \mathbf{I}_p$ , we obtain:

$$\mu_{\text{post}} = (\sigma^{-2} \mathbf{X}^\top \mathbf{X} + \tau^{-2} \mathbf{I}_p)^{-1} \sigma^{-2} \mathbf{X}^\top \mathbf{Y}$$

Simplifying further:

$$\mu_{\text{post}} = (\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}$$

This aligns with the Ridge Regression estimator derived earlier.

## Remark on Ridge Regression's Existence

An important property of Ridge Regression is that the estimator  $\beta$  is always defined, even when the design matrix  $\mathbf{X}$  does not have full rank. Specifically, Ridge Regression provides a unique solution even in cases where  $n < p$  (i.e., when there are more predictors than observations), addressing issues of multicollinearity and overfitting inherent in OLS.

## Definition of Ridge Regression Estimator

$$\hat{\beta}_{\text{ridge}} = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2,$$

is termed the **Ridge Regression** estimator. Here,  $\lambda > 0$  serves as the regularization parameter that controls the strength of the penalty on the size of the coefficients. A key feature of Ridge Regression is that  $\hat{\beta}_{\text{ridge}}$  is always uniquely defined, making it a robust alternative to OLS, especially in high-dimensional settings.

When the model is specified as  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , where  $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$ , the Ridge estimator takes the form:

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Notably, this estimator is independent of  $\sigma^2$ , highlighting its focus on minimizing the penalized residual sum of squares without direct dependence on the noise variance.

**Proposition: Mean Squared Error (MSE) of  $\hat{\boldsymbol{\beta}}_{\text{ridge}}$**

We now analyze the MSE of the Ridge estimator under specific conditions.

**Assumptions:**

- The linear model is given by  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ .
- The error term  $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$ , with  $\sigma^2 > 0$  known.
- The design matrix satisfies  $\mathbf{X}^\top \mathbf{X} = n \mathbf{I}_p$ , indicating an orthonormal design.

**Definition:** Let  $\mathcal{J} := \langle \boldsymbol{\beta}, \mathbf{v} \rangle$  for some  $\mathbf{v} \in \mathbb{R}^p$ , and define:

$$\delta_{\text{ridge}} = \langle \hat{\boldsymbol{\beta}}_{\text{ridge}}, \mathbf{v} \rangle.$$

**Statements:**

1. The expected squared error of  $\delta_{\text{ridge}}$  is:

$$\mathbb{E}_{\boldsymbol{\beta}}[(\delta_{\text{ridge}} - \mathcal{J})^2] = (1 + \lambda)^{-2} \langle \boldsymbol{\beta}_v, \mathbf{v} \rangle^2 + \frac{\sigma^2}{n} \|\mathbf{v}\|^2 (1 + \lambda)^{-2}.$$

2. The expected squared error of the Ridge estimator vector is:

$$\mathbb{E}_{\boldsymbol{\beta}}[\|\hat{\boldsymbol{\beta}}_{\text{ridge}} - \boldsymbol{\beta}\|^2] = (1 + \lambda)^{-2} \|\boldsymbol{\beta}\|^2 + \frac{p\sigma^2}{n} \frac{1}{(1 + \lambda)^2}.$$

**Derivation:**

Starting from the Ridge estimator:

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y}.$$

Given the orthonormal design  $\mathbf{X}^\top \mathbf{X} = n \mathbf{I}_p$ , this simplifies to:

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (n \mathbf{I}_p + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{Y} = \frac{1}{1 + \frac{\lambda}{n}} (\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} + \mathbf{X}^\top \boldsymbol{\varepsilon}).$$

Simplifying further:

$$\hat{\beta}_{\text{ridge}} = \frac{1}{1 + \frac{\lambda}{n}} \beta + \frac{1}{1 + \frac{\lambda}{n}} \mathbf{X}^\top \boldsymbol{\varepsilon}.$$

This expression separates the estimator into a bias term and a variance term, which will be instrumental in understanding the bias-variance decomposition.

### Bias-Variance Decomposition

The Mean Squared Error (MSE) of an estimator can be decomposed into the sum of the squared bias and the variance of the estimator:

$$\mathbb{E} \left[ (\hat{\beta}_{\text{ridge}} - \mathcal{J})^2 \right] = (\mathbb{E}[\hat{\beta}_{\text{ridge}}] - \mathcal{J})^2 + \text{Var}(\hat{\beta}_{\text{ridge}}).$$

Substituting from our Ridge estimator:

$$= \left( \left(1 + \frac{\lambda}{n}\right)^{-1} \langle \beta, \mathbf{v} \rangle \right)^2 + \frac{\lambda^2}{(1 + \lambda)^2} \text{Var}(\mathbf{X}^\top \boldsymbol{\varepsilon}, \nu).$$

**Observations:**

$$\left(1 + \frac{\lambda}{n}\right)^{-1} = \frac{1}{\left(1 + \frac{\lambda}{n}\right)}.$$

Furthermore, the variance term can be expressed as:

$$\text{Var}(\mathbf{X}^\top \boldsymbol{\varepsilon}, \nu) = \nu^\top \mathbf{X} \text{Cov}(\boldsymbol{\varepsilon}) \mathbf{X}^\top \nu = \sigma^2 \|\nu\|^2.$$

Thus, the bias-variance decomposition clearly illustrates the trade-off controlled by the regularization parameter  $\lambda$ .

### Corollary: MSE Under Orthonormal Design

Under the same assumptions, specifically the orthonormal design  $\mathbf{X}^\top \mathbf{X} = n \mathbf{I}_p$ , the MSE of the Ridge estimator vector is:

$$\mathbb{E}[\|\hat{\beta}_{\text{ridge}} - \beta\|^2] = \frac{1}{\left(1 + \frac{\lambda}{n}\right)^2} \|\beta\|^2 + \frac{p\sigma^2}{n\left(1 + \frac{\lambda}{n}\right)^2}.$$

This expression highlights how both the bias and variance components of the estimator's error are influenced by  $\lambda$ .



### Remark on Ridge Regression Behavior

An insightful observation about Ridge Regression is its behavior in relation to the norm of  $\beta$ :

For small  $\|\beta\|$ , Ridge Regression converges to OLS.

This implies that the optimal choice of  $\lambda$  depends on the true underlying parameter  $\|\beta\|$ . When  $\|\beta\|$  is small, a smaller  $\lambda$  is preferred to minimize bias, whereas larger  $\lambda$  values are beneficial when  $\|\beta\|$  is large to control variance.

### Confidence Sets & Tests in Linear Model

While Ridge and OLS estimators provide point estimates for  $\beta$ , it's crucial to quantify the uncertainty associated with these estimates. This involves constructing confidence sets and performing hypothesis tests, which inherently depend on the variance  $\sigma^2$ .

**Assumption:** Throughout this discussion, we assume  $\varepsilon \sim N(0, \sigma^2 \mathbf{I}_n)$ , ensuring the errors are normally distributed with known variance.

#### Easy Case: Known $\sigma^2$

When the variance  $\sigma^2$  is known, the distribution of the OLS estimator is straightforward:

$$\hat{\beta}_{\text{OLS}} \sim N(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}).$$

Given an orthonormal design ( $\mathbf{X}^\top \mathbf{X} = n \mathbf{I}_p$ ), this simplifies to:

$$\hat{\beta}_{\text{OLS}} \sim N\left(\beta, \frac{\sigma^2}{n} \mathbf{I}_p\right).$$

For any linear combination  $\mathcal{J} = \langle \beta, \nu \rangle$ , the estimator is:

$$\hat{\mathcal{J}} = \langle \hat{\beta}_{\text{OLS}}, \nu \rangle \sim N\left(\mathcal{J}, \sigma^2 \nu^\top (\mathbf{X}^\top \mathbf{X})^{-1} \nu\right).$$

This allows us to construct a 95% confidence interval for  $\mathcal{J}$ :

$$I_{95\%}(\mathcal{J}) = \left[ \hat{\mathcal{J}} \pm 1.96 \sqrt{\nu^\top (\mathbf{X}^\top \mathbf{X})^{-1} \nu} \right].$$

### Notes on $t$ - and $F$ -distributions

In practice, the variance  $\sigma^2$  is often unknown and must be estimated from the data. This necessitates the use of the  $t$ -distribution and the  $F$ -distribution for constructing

confidence intervals and conducting hypothesis tests.

## Definitions

**t-distribution:** The  $t$ -distribution with  $n \geq 1$  degrees of freedom on  $\mathbb{R}$  has the probability density function:

$$f_n(x) = C_n \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}},$$

where  $C_n$  is the normalizing constant ensuring that the total probability integrates to 1.

**Special Case:** For  $n = 1$ , the  $t$ -distribution simplifies to:

$$f_1(x) = C_1 \frac{1}{1 + x^2},$$

which corresponds to the **Cauchy distribution**, a distribution with heavy tails.

**F-distribution:** The  $F$ -distribution with  $(m, n) \in \mathbb{N}^2$  degrees of freedom has the density:

$$f_{m,n}(x) = C_{m,n} \frac{x^{\frac{m}{2}-1}}{(mx + n)^{\frac{m+n}{2}}}, \quad x \in (0, \infty),$$

where  $C_{m,n}$  is the normalizing constant. The  $F$ -distribution is pivotal in variance ratio tests and regression analysis.

## Utility of $t$ - and $F$ -distributions

These distributions are instrumental in constructing confidence intervals and performing hypothesis tests when dealing with normally distributed errors, especially when the variance is unknown and must be estimated from the data.

**Lemma 2.** *Let  $X_1, \dots, X_n, Y_1, \dots, Y_n$  be independent and identically distributed (i.i.d.)  $N(0, \Delta)$  random variables. Then:*

1.

$$T_n := \frac{X_n}{\sqrt{\frac{1}{n} \sum_{i=1}^n Y_i^2}} \sim t_n.$$

2.

$$F_{m,n} := \frac{\frac{1}{m} \sum_{i=1}^m X_i^2}{\frac{1}{n} \sum_{j=1}^n Y_j^2} \sim F_{m,n}.$$

## Remarks on the Lemma

1. The  $t$ -distribution emerges when evaluating the ratio of a normally distributed variable to the square root of a scaled chi-squared variable, reflecting the relationship between the sample mean and sample variance.
2. As the degrees of freedom  $n$  increase, the  $t$ -distribution  $T_n$  converges in distribution to the standard normal distribution  $N(0, 1)$ , illustrating the asymptotic behavior of the estimator.

## Proof of the Lemma

*Proof.* 1. To show that  $T_n$  follows a  $t$ -distribution with  $n$  degrees of freedom, observe that:

$$T_n^2 = F_{1,n}.$$

Considering the symmetry of the  $t$ -distribution around 0, we utilize a change of variables:

$$f_{F_{m,n}}(x) = f_{F_{m,n}}(x^2) \cdot 2x, \quad x > 0.$$

Given that  $T_n$  is symmetric, we derive:

$$f_{T_n}(x) = f_{F_{m,n}}(x^2)|x| = C_n \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}.$$

This aligns with the known density of the  $t$ -distribution.

2. To establish the distribution of  $F_{m,n}$ , define:

$$X = \sum_{i=1}^m X_i^2, \quad Y = \sum_{j=1}^n Y_j^2.$$

Since each  $X_i$  and  $Y_j$  is normally distributed, it follows that:

$$X \sim \chi_m^2, \quad Y \sim \chi_n^2,$$

where  $\chi_m^2$  and  $\chi_n^2$  denote chi-squared distributions with  $m$  and  $n$  degrees of freedom, respectively. The density of  $\chi_m^2$  is:

$$f(x) \propto x^{m/2-1} e^{-x/2}, \quad x > 0. \quad \square$$

## Derivation of the $F$ -distribution

To derive the distribution of  $F_{m,n}$ , consider the ratio:

$$W = \frac{X}{Y}.$$

The cumulative distribution function (CDF) of  $W$  is:

$$\mathbb{P}\left(\frac{X}{Y} < z\right) = \int_0^\infty \int_0^{zy} 1 \cdot f_X(x) f_Y(y) dx dy.$$

By substituting  $x = wy$ , we obtain:

$$= \int_0^\infty \int_0^z f_X(wy) f_Y(y) \cdot y dw dy.$$

This simplifies to:

$$= \int_0^\infty f_X(zy) f_Y(y) \cdot y dy.$$

Substituting the chi-squared densities:

$$\propto \int_0^\infty (zy)^{\frac{m}{2}-1} y^{\frac{n}{2}-1} e^{-(z+y)/2} dy.$$

Introducing the change of variable  $a = \frac{z}{z+1}y$ , the integral becomes:

$$\propto z^{\frac{m}{2}-1} (z+1)^{-\frac{m+n}{2}} \int_0^\infty a^{\frac{m}{2}-1} e^{-a} da.$$

Recognizing that the integral is the gamma function  $\Gamma\left(\frac{m}{2}\right)$ , we conclude:

$$f_F(z) \propto z^{\frac{m}{2}-1} (z+1)^{-\frac{m+n}{2}}.$$

Finally, by appropriate scaling, we arrive at the density function of the  $F$ -distribution:

$$f_{m,n}(z) = \frac{m^{m/2} n^{n/2}}{B\left(\frac{m}{2}, \frac{n}{2}\right)} \frac{z^{m/2-1}}{(mz+n)^{(m+n)/2}},$$

where  $B(\cdot, \cdot)$  is the beta function.

## 9 Lecture 10: Confidence Sets and Hypothesis Testing in Linear Models

Building upon the foundations laid in the previous lecture, we now focus on constructing confidence sets and performing hypothesis tests within the context of linear models. We will explore both  $t$ -tests and  $F$ -tests, essential tools for statistical inference in regression analysis.

### Ridge and OLS Estimators

Recall that in the linear model  $Y = X\beta + \varepsilon$ , with  $\varepsilon \sim N(0, \sigma^2 I_n)$ , the OLS estimator is given by:

$$\hat{\beta} = (X^T X)^{-1} X^T Y \sim N(\beta, \sigma^2 (X^T X)^{-1}).$$

When  $\sigma^2$  is unknown, it is typically estimated by:

$$\hat{\sigma}^2 = \frac{\|Y - X\hat{\beta}_{OLS}\|^2}{n - p} \sim \sigma^2 \frac{\chi^2(n - p)}{n - p}.$$

This estimator is unbiased and follows a scaled chi-squared distribution.

### Distributions of Estimators

$$t_n = \frac{N(0, 1)}{\sqrt{\chi^2(n)/n}} \quad \text{and} \quad F_{m,n} = \frac{\chi^2(m)/m}{\chi^2(n)/n}$$

These ratios form the basis of the  $t$ - and  $F$ -distributions, respectively.

**Lemma 3.** Let  $\xi \sim N(0, I_n)$ , a random variable in  $\mathbb{R}^n$ , and let  $R \in \mathbb{R}^{n \times n}$  be an orthogonal projection matrix ( $R = R^2$ ,  $R = R^T$ ), with  $\text{rank}(R) = r \leq n$ . Then:

1.  $\xi^T R \xi = \|R\xi\|^2 \sim \chi^2(r)$ .
2. If  $B \in \mathbb{R}^{p \times n}$  is such that  $BR = 0$ , then  $B\xi$  is independent of  $R\xi$ .
3. If  $S \in \mathbb{R}^{n \times n}$  is another orthogonal projection with  $\text{rank}(S) = s \leq n$  and  $RS = 0$ ,

then

$$\frac{s}{r} \frac{\xi^T R \xi}{\xi^T S \xi} \sim F(r, s).$$

*Proof.* 1. Since  $R$  is an orthogonal projection, there exists an orthogonal matrix  $T$  such that:

$$R = T \begin{pmatrix} I_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} T^T = T D_r T^T.$$

Applying the orthogonal transformation  $T^T$  to  $\xi$ , we maintain the distribution:

$$T^T \xi \sim N(0, I_n).$$

Therefore,

$$\xi^T R \xi = \xi^T (T D_r T^T) \xi = (T^T \xi)^T D_r (T^T \xi) = \sum_{i=1}^r (T^T \xi)_i^2 \sim \chi^2(r).$$

2. Let  $A_1 = B\xi$  and  $A_2 = R\xi$ . The covariance between  $A_1$  and  $A_2$  is:

$$\text{Cov}(A_1, A_2) = \text{Cov}(B\xi, R\xi) = B \text{Cov}(\xi, \xi) R^T = B R^T = B R = 0,$$

since  $B R = 0$  by assumption. Given that both  $A_1$  and  $A_2$  are linear combinations of jointly normal variables, zero covariance implies independence.

3. Given that  $R S = 0$ , by part (2),  $S\xi$  and  $R\xi$  are independent. From part (1), we have  $S\xi \sim \chi^2(s)$  and  $R\xi \sim \chi^2(r)$ . Thus, the ratio:

$$\frac{s}{r} \frac{\xi^T R \xi}{\xi^T S \xi} \sim F(r, s)$$

follows directly from the definition of the  $F$ -distribution as the ratio of two scaled chi-squared variables.  $\square$

### Theorem: Confidence Sets in Linear Models with Unknown $\sigma^2$

**Theorem 3 (Linear Model Confidence Sets - Unknown  $\sigma^2$ ).** Assume the regular linear model  $Y = X\beta + \varepsilon$ , with  $\text{rank}(X) = p \leq n$  and  $\varepsilon \sim N(0, \sigma^2 I_n)$ . Let  $\alpha \in (0, 1)$ . Then:

1. Let  $q_{F_{p, n-p}, 1-\alpha}$  be the  $1 - \alpha$  quantile of the  $F_{p, n-p}$  distribution. Then the set

$$C(Y, X) = \left\{ \beta \in \mathbb{R}^p \mid \frac{\|X(\beta - \hat{\beta}_{OLS})\|^2}{p\hat{\sigma}^2} \leq q_{F_{p, n-p}, 1-\alpha} \right\}$$

is a  $1 - \alpha$  confidence set for  $\beta$ .

2. For a specific linear combination  $\alpha = \langle \beta, v \rangle$  for some  $v \in \mathbb{R}^p$ , a  $1 - \alpha$  confidence interval is given by:

$$C = C(Y, X) = \left\{ \alpha \in \mathbb{R} \mid \left| \frac{\alpha - \hat{\alpha}}{\hat{\sigma} \sqrt{v^T (X^T X)^{-1} v}} \right| < q \right\},$$

where  $\hat{\alpha} = \langle \hat{\beta}_{OLS}, v \rangle$  and  $q$  is the  $1 - \frac{\alpha}{2}$  quantile of the  $t_{n-p}$  distribution.

*Proof.* 1. Consider the residual sum of squares under the null hypothesis  $H_0 : \beta = \hat{\beta}_{OLS}$ :

$$RSS = \|Y - X\hat{\beta}_{OLS}\|^2.$$

The residual sum of squares under the alternative hypothesis is:

$$RSS_{H_0} = \|Y - X\beta\|^2.$$

Given the orthonormal design and the properties of the Ridge estimator, we have:

$$RSS_{H_0} - RSS = (K\hat{\beta}_{OLS} - d)^2 (K(X^T X)^{-1} K^T)^{-1},$$

where  $K$  is the contrast matrix. Under  $H_0$ , this difference follows a chi-squared distribution with  $r$  degrees of freedom:

$$\frac{RSS_{H_0} - RSS}{\sigma^2} \sim \chi^2(r).$$

Consequently, the ratio:

$$\frac{RSS_{H_0} - RSS}{r\hat{\sigma}^2} \sim F(r, n - p).$$

This establishes that:

$$\frac{\|X(\beta - \hat{\beta}_{OLS})\|^2}{p\hat{\sigma}^2} \sim F(p, n - p),$$

leading to the construction of the confidence set  $C(Y, X)$  as stated.

2. For the specific linear combination  $\alpha = \langle \beta, v \rangle$ , the OLS estimator satisfies:

$$\hat{\alpha} = \langle \hat{\beta}_{OLS}, v \rangle \sim N\left(\alpha, \sigma^2 v^T (X^T X)^{-1} v\right).$$

The standardized statistic is:

$$\frac{\alpha - \hat{\alpha}}{\sigma \sqrt{v^T (X^T X)^{-1} v}} \sim t_{n-p}.$$

However, since  $\sigma^2$  is unknown and estimated by  $\hat{\sigma}^2$ , the statistic becomes:

$$\frac{\alpha - \hat{\alpha}}{\hat{\sigma} \sqrt{v^T (X^T X)^{-1} v}} \sim t_{n-p}.$$

Thus, the confidence interval  $C$  is constructed by:

$$\left| \frac{\alpha - \hat{\alpha}}{\hat{\sigma} \sqrt{v^T (X^T X)^{-1} v}} \right| < q,$$

where  $q$  is the critical value from the  $t_{n-p}$  distribution corresponding to the desired confidence level.  $\square$

### 9.0.1 The $t$ - and $F$ -Tests

Confidence sets facilitate hypothesis testing. Here, we outline the procedures for conducting  $t$ -tests and  $F$ -tests in the linear model context.

#### Remark: Method (t-test)

In a regular linear model with  $\varepsilon \sim N(0, \sigma^2 I_n)$ , consider the null hypothesis:

$$H_0 : \gamma = \gamma_0 \quad \text{vs} \quad H_1 : \gamma \neq \gamma_0,$$

where  $\gamma = \langle \beta, v \rangle$  for some  $v \in \mathbb{R}^p$ .

The two-sided  $t$ -test statistic is defined as:

$$T_{\gamma_0, n-p} = \frac{\gamma_0 - \hat{\gamma}}{\hat{\sigma} \sqrt{v^T (X^T X)^{-1} v}},$$

where  $\hat{\gamma} = \langle \hat{\beta}_{OLS}, v \rangle$ .

The corresponding test function is:

$$\varphi_{\alpha_0}(Y, X) = \mathbf{1}(|T_{\gamma_0, n-p}(Y, X)| > q),$$

where  $q$  is the critical value from the  $t_{n-p}$  distribution corresponding to the significance



level  $\alpha_0$ .

**Remark: Method (F-test)**

For the same linear model setting, consider the null hypothesis:

$$H_0 : \beta = \beta_0 \quad \text{vs} \quad H_1 : \beta \neq \beta_0,$$

where  $\beta_0 \in \mathbb{R}^p$ .

The  $F$ -test statistic is defined as:

$$F_{\beta_0, n-p}(Y, X) = \frac{\|X(\beta - \hat{\beta}_{OLS})\|^2}{p\hat{\sigma}^2},$$

which follows an  $F_{p, n-p}$  distribution under  $H_0$ .

The corresponding test function is:

$$\varphi_{\beta_0}(Y, X) = \mathbb{I}(|F_{\beta_0, n-p}(Y, X)| > q),$$

where  $q$  is the critical value from the  $F_{p, n-p}$  distribution corresponding to the desired significance level.

## 9.0.2 General Linear Hypothesis Testing Problems

Beyond testing individual parameters, we often encounter more complex hypotheses involving multiple linear constraints on  $\beta$ . This section addresses such scenarios.

**Definition 19.** A **linear hypothesis testing problem** is of the form:

$$H_0 : K\beta = d \quad \text{vs} \quad H_1 : K\beta \neq d,$$

where  $K \in \mathbb{R}^{r \times p}$  is a matrix with  $\text{rank}(K) = r \leq p$  and  $d \in \mathbb{R}^r$ . Essentially, this represents  $r$  linear constraints on the parameter vector  $\beta$ . The matrix  $K$  is referred to as the **contrast matrix**.

**Theorem 4.** Assume the regular linear model  $Y = X\beta + \varepsilon$ , with  $\text{rank}(X) = p \leq n$  and  $\varepsilon \sim N(0, \sigma^2 I_n)$ . Consider the hypothesis:

$$H_0 : K\beta = d \quad \text{vs} \quad H_1 : K\beta \neq d.$$

Define the residual sum of squares (RSS) as:

$$RSS = \|Y - X\hat{\beta}_{OLS}\|^2 \quad \text{and} \quad RSS_{H_0} = \|Y - X\hat{\beta}_{H_0}\|^2,$$

where  $\hat{\beta}_{H_0}$  is the estimator of  $\beta$  under the constraint  $K\beta = d$ .

Then:

1. The constrained estimator is:

$$\hat{\beta}_{H_0} = \hat{\beta}_{OLS} - (X^T X)^{-1} K^T (K (X^T X)^{-1} K^T)^{-1} (K \hat{\beta}_{OLS} - d).$$

2. The difference in RSS satisfies:

$$RSS_{H_0} - RSS = (K \hat{\beta}_{OLS} - d)^T (K (X^T X)^{-1} K^T)^{-1} (K \hat{\beta}_{OLS} - d) \sim \sigma^2 \chi^2(r).$$

3. The test statistic:

$$F = \frac{(RSS_{H_0} - RSS)/r}{RSS/(n-p)} \sim F_{r, n-p}$$

under the null hypothesis  $H_0$ .

*Proof.* 1. The constrained estimator  $\hat{\beta}_{H_0}$  minimizes the RSS subject to  $K\beta = d$ . By applying the method of Lagrange multipliers or directly solving the constrained optimization problem, we obtain:

$$\hat{\beta}_{H_0} = \hat{\beta}_{OLS} - (X^T X)^{-1} K^T (K (X^T X)^{-1} K^T)^{-1} (K \hat{\beta}_{OLS} - d).$$

2. The difference in RSS between the unconstrained and constrained models is:

$$RSS_{H_0} - RSS = (K \hat{\beta}_{OLS} - d)^T (K (X^T X)^{-1} K^T)^{-1} (K \hat{\beta}_{OLS} - d).$$

Under  $H_0$ ,  $K \hat{\beta}_{OLS} \sim N(d, \sigma^2 K (X^T X)^{-1} K^T)$ . Thus, the quadratic form above follows a chi-squared distribution with  $r$  degrees of freedom:

$$\frac{RSS_{H_0} - RSS}{\sigma^2} \sim \chi^2(r).$$

3. Defining the  $F$ -statistic as:

$$F = \frac{(RSS_{H_0} - RSS)/r}{RSS/(n-p)} = \frac{RSS_{H_0} - RSS}{r \hat{\sigma}^2} \sim F_{r, n-p},$$

where  $\hat{\sigma}^2 = \frac{RSS}{n-p}$  is the unbiased estimator of  $\sigma^2$ . Under  $H_0$ , this statistic follows the  $F$ -distribution with  $r$  and  $n-p$  degrees of freedom.  $\square$

### 9.0.3 The $t$ - and $F$ -Tests

#### Remark: $t$ -Test Method

Consider a regular linear model where  $\varepsilon \sim N(0, \sigma^2 I_n)$ . To test the null hypothesis  $H_0 : \gamma = \gamma_0$  against the alternative  $H_1 : \gamma \neq \gamma_0$ , where  $\gamma = \langle \beta, v \rangle$  for some  $v \in \mathbb{R}^p$ , we employ the two-sided  $t$ -test.

The  $t$ -test statistic is defined as:

$$T_{\gamma_0, n-p} = \frac{\gamma_0 - \hat{\gamma}}{\hat{\sigma} \sqrt{v^T (X^T X)^{-1} v}},$$

where  $\hat{\gamma} = \langle \hat{\beta}_{OLS}, v \rangle$ .

The corresponding test function is:

$$\varphi_{\alpha_0}(Y, X) = \mathbf{1}(|T_{\gamma_0, n-p}(Y, X)| > q),$$

where  $q$  is the critical value from the  $t_{n-p}$  distribution corresponding to the significance level  $\alpha_0$ .

#### Remark: $F$ -Test Method

In the same linear model setting, to test the null hypothesis  $H_0 : \beta = \beta_0$  against  $H_1 : \beta \neq \beta_0$ , we use the  $F$ -test.

The  $F$ -test statistic is defined as:

$$F_{\beta_0, n-p}(Y, X) = \frac{\|X(\beta - \hat{\beta}_{OLS})\|^2}{p\hat{\sigma}^2},$$

which follows an  $F_{p, n-p}$  distribution under  $H_0$ .

The corresponding test function is:

$$\varphi_{\beta_0}(Y, X) = \mathbf{1}(|F_{\beta_0, n-p}(Y, X)| > q),$$

where  $q$  is the critical value from the  $F_{p, n-p}$  distribution at the desired significance level.

### 9.0.4 General Linear Hypothesis Testing Problems

Beyond testing individual coefficients or simple contrasts, we often need to test more general linear hypotheses involving multiple parameters. This section formalizes such testing problems and provides the necessary statistical framework.

**Definition 20.** A linear hypothesis testing problem involves testing  $H_0 : K\beta = d$  against  $H_1 : K\beta \neq d$ , where:

- $K \in \mathbb{R}^{r \times p}$  is the contrast matrix with  $\text{rank}(K) = r \leq p$ .
- $d \in \mathbb{R}^r$  is the vector of constants.

**Theorem 5.** Assume the regular linear model  $Y = X\beta + \varepsilon$ , with  $\text{rank}(X) = p \leq n$  and  $\varepsilon \sim N(0, \sigma^2 I_n)$ . Consider the linear hypothesis:

$$H_0 : K\beta = d \quad \text{vs} \quad H_1 : K\beta \neq d.$$

Define the residual sum of squares (RSS) under  $H_0$  as:

$$RSS_{H_0} = \|Y - X\hat{\beta}_{H_0}\|^2,$$

where  $\hat{\beta}_{H_0}$  is the constrained estimator minimizing RSS subject to  $K\beta = d$ .

Then:

1. The constrained estimator is:

$$\hat{\beta}_{H_0} = \hat{\beta}_{OLS} - (X^T X)^{-1} K^T (K(X^T X)^{-1} K^T)^{-1} (K\hat{\beta}_{OLS} - d).$$

2. The difference in RSS is given by:

$$RSS_{H_0} - RSS = (K\hat{\beta}_{OLS} - d)^T (K(X^T X)^{-1} K^T)^{-1} (K\hat{\beta}_{OLS} - d) \sim \sigma^2 \chi^2(r).$$

3. The  $F$ -statistic is defined as:

$$F = \frac{RSS_{H_0} - RSS}{r\hat{\sigma}^2} \sim F_{r, n-p},$$

under the null hypothesis  $H_0$ .

*Proof.* 1. To derive the constrained estimator  $\hat{\beta}_{H_0}$ , we minimize the RSS subject to  $K\beta = d$ . Utilizing the method of Lagrange multipliers or projecting the OLS estimator onto the constraint set, we obtain:

$$\hat{\beta}_{H_0} = \hat{\beta}_{OLS} - (X^T X)^{-1} K^T (K(X^T X)^{-1} K^T)^{-1} (K\hat{\beta}_{OLS} - d).$$

2. The difference in RSS between the constrained and unconstrained models is a quadratic form:

$$RSS_{H_0} - RSS = (K\hat{\beta}_{OLS} - d)^T (K(X^T X)^{-1} K^T)^{-1} (K\hat{\beta}_{OLS} - d).$$

Under  $H_0$ ,  $K\hat{\beta}_{OLS}$  follows a normal distribution centered at  $d$ , and the quadratic form follows a scaled chi-squared distribution:

$$\frac{RSS_{H_0} - RSS}{\sigma^2} \sim \chi^2(r).$$

3. Defining the  $F$ -statistic as:

$$F = \frac{RSS_{H_0} - RSS}{r\hat{\sigma}^2},$$

and noting that  $\hat{\sigma}^2 = \frac{RSS}{n-p}$  is an unbiased estimator of  $\sigma^2$ , we conclude:

$$F \sim F_{r,n-p}.$$

This completes the proof. □