

# TOPICS IN PROBABILITY, PARAMETRIC ESTIMATION AND STOCHASTIC CALCULUS

LEVI LOPES DE LIMA

ABSTRACT. We start by recalling the fundamentals of Probability Theory needed to fully grasp the basic aspects of one of its more important applications to real-world problems, namely, Parametric Estimation, which encompasses topics such as concentration inequalities, limit theorems, confidence intervals, maximum likelihood, least squares and hypothesis testing. In a somewhat long appendix we also present the rudiments of Brownian motion and the corresponding stochastic calculus (incarnated in Itô's celebrated formula) decorated with a few of its most glamorous applications (the sharp Gaussian concentration inequality, the Feynman-Kac formula and the Black-Scholes strategy in Finance).

## CONTENTS

1. Introduction	2
2. The fundamentals of Probability Theory	4
3. Conditioning	17
3.1. Conditional probability	17
3.2. Conditional expectation	19
4. Normally distributed random variables and their friends	22
4.1. Normally distributed random variables	22
4.2. Random variables related to the normal	30
5. Concentration inequalities	34
5.1. Sub-exponential random variables and the Johnson-Lindenstrauss Lemma	35
5.2. The Gaussian concentration inequality	40
5.3. Chernoff-type bounds for binomial trials and the Erdős-Rényi model	45
6. The fundamental limit theorems	48
7. Estimation	54
7.1. Parametric estimation and the mean squared error	54
7.2. Computing the mean squared error of $\hat{\sigma}_c^2$	60
7.3. Confidence intervals	65
8. Maximum likelihood	74
8.1. Maximum likelihood estimators	74
8.2. Fisher information and Cramér-Rao lower bound	77
8.3. Optimal asymptotic normality of ML estimators	84
9. The method of least squares	91

9.1. The statistical model behind MLS	92
9.2. Inference and goodness of fit for MLS	95
9.3. Regularization in high dimension, sparsity and the LASSO	106
10. Sufficiency	112
11. Hypothesis testing	115
11.1. A glimpse at the Neyman-Pearson setup	115
11.2. Testing via likelihood ratios	116
12. A brief overview of “classical” Parametric Estimation	128
13. The Bayesian pathway	130
Appendix A. Brownian motion and some of its applications	134
A.1. Brownian motion: its construction and basic regularity properties	134
A.2. Martingales	139
A.3. Itô’s integral and Itô’s formula	140
A.4. The Gaussian concentration inequality (again)	149
A.5. The Feynman-Kac formula and the path integral representation of the heat kernel	150
A.6. The Black-Scholes strategy in Finance	154
References	158

## 1. INTRODUCTION

Probability theory is a multifaceted intellectual enterprise with countless applications, both in pure and applied mathematics. Among its many usages in real-world problems, there is at least one body of knowledge which stands out by its relevance, namely, Parametric Estimation, a collection of statistical inference techniques initially put forward by K.F. Gauss and P.-S. Laplace, mainly in the setting of the method of least squares, and further transformed into a full-fledged research program by R. Fisher, J. Neyman, E. Pearson, among others, with a myriad of applications ranging from life (Biology, Evolution, Medicine, etc) to behavioral (Psychology, Sociology, Economics, etc) sciences and of course passing through the more familiar consumption in STEM disciplines (such as Physics, Chemistry and Engineering). More recently, suitable variants of this classical theory revealed themselves crucial in probing the efficiency of standard procedures in modern Data Science, as attested by the first six chapters of [JWHT13], in which (possibly sparse) methods for regression and classification are discussed in the broader context of Statistical Learning.

The purpose of these notes is to provide a modest introduction to this circle of ideas specifically designed for those with an adequate training in the required prerequisites (essentially, Linear Algebra, Multivariate Calculus and Measure Theory, including a bit of Fourier Analysis), which in particular demands a certain amount of mathematical maturity (say, at the graduate level). With this goal in mind, we start by concisely recounting those aspects of Probability Theory needed to formulate the most basic conceptual and computational accomplishments in the field. This is done in Sections 2 and 3, which is then complemented by Section 4, where the rudiments of the theory of normally distributed random vectors are presented. With these foundational issues at hand, we then proceed in Section 5 to an introduction to the study of concentration inequalities, a collection of techniques providing *non*-asymptotic bounds on the tail probabilities of a large class of distributions.

Here we focus on certain inequalities directly accessible through the elementary yet powerful *Cramér-Chernoff method*, which requires a suitable exponential control on the underlying moment generating functions. Besides indicating how this leads both to the Johnson-Lindenstrauss lemma and to certain “phase transition” phenomena associated to the Erdős-Rényi model for random graphs, we also include a discussion of the Gaussian concentration inequality with an emphasis on its connection to Poincaré’s limit theorem and, more generally, to the “concentration of measure phenomenon”. In Section 6 we return to the asymptotic realm by indicating how Fourier Analysis, incarnated in the theory of characteristic functions, naturally leads to the fundamental limit theorems (Law of Large Numbers and Central Limit Theorem). Taken together, these classical accomplishments in a sense constitute the core results in any introductory course on the subject, as they immediately yield “large sample” confidence interval estimates for the expected value of *any* random sample. With these preliminaries out of the way, we proceed in Section 7 to estimation proper with an exposition of the most elementary aspects of Parametric Estimation, where the key concept of a *statistical model*, due to R. Fisher, stands out. We thus explain how the most commonly used measure of performance of an associated estimator, the mean squared error, relates to other important notions (consistency, bias-variance trade-off, asymptotic normality, etc.). This works as a preparation for the Fisherian approach to estimation developed in Section 8, where maximum likelihood estimators are introduced and their asymptotic properties are examined in the light of the celebrated (and *non*-asymptotic) Cramér-Rao bound for the variance of estimators, an analysis ultimately relying on the concept of *Fisher information*. We next developed in Section 9 the theoretical framework behind Linear Regression, certainly the most often used statistical tool in applications, with a view toward inspecting its inferential properties, including both interpretability (via parameter recovering) and prediction. Although we work here mainly in the classical “low dimensional” setting, where the number of features is much smaller than the number of observations ( $p \ll n$ ), hence making use of the Method of Least Squares, we also make it available an introduction to “sparsity”, with a discussion of the most basic prediction properties of the LASSO estimator. Needless to say, this works as an incitement to those interested in the “high dimensional” setting, where the classical inferential methods fail to deliver reliable results and which is so prominent in modern Data Science (see, for instance, [JWHT13, Chapter 6]). The standard estimation package is finally completed in Section 11, where hypothesis testing is discussed, with an emphasis on likelihood ratio tests, which are suitably illustrated with a few notable examples. As final appetizers, we include discussions on sufficiency (Section 10), the Bayesian approach to inference (Section 13) and an overview of classical estimation theory, with a brief assessment of R. Fisher’s main contributions to the subject (Section 12). We also include in these notes a somewhat long Appendix in which the rudiments of Itô’s calculus are presented, starting with the construction of Brownian motion. Although the main motivation here is to provide a complete proof of the sharp Gauss concentration inequality, as the elementary proof in Subsection 5.2 fails to deliver the best possible standard deviation, we also add a few selected applications of this formalism, including the Feynman-Kac formula, which exhibits a path integral representation for the heat kernel associated to certain heat-type operators in terms of the Brownian bridge, and the Black-Scholes strategy in Finance (which has been worth a Nobel Prize in 1997).

Our presentation borrows a lot from the excellent monographs and papers cited throughout the text, which we highly recommend to anyone who wants to delve deeper into these fascinating topics. Also, exactly because they abound in these monographs (specially those oriented to (under)graduate studies) we refrained from including proposed exercises throughout the text. On the other hand, in order to align an otherwise theoretical manuscript with our digital era, in which Data Science is a primary motivation for studying most of the topics treated here, we make available a few labs intended to effectively reconcile theoretical recipes with practical applications [dL25]<sup>1</sup>, where updates of this text will appear as well. We insist that the labs do not reduce themselves to a mere collection of code snippets in R, as they also include concise theoretical recaps (“A glimpse at the theory”), thus connecting the general principles in the main text to the specific problems treated

<sup>1</sup>The topics included so far are: linear regression, the Central Limit Theorem, the Johnson-Lindenstrauss lemma, maximum likelihood estimation and James-Stein estimators (<https://github.com/levilopesdelima/stat-inference-labs>).

in the step-by-step simulations. Moreover, they contain insightful visualizations of the results and additional interpretations of the (numerical and visual) outputs, as we believe that this hands-on approach is crucial for building a deeper understanding of the underlying conceptual framework.

## 2. THE FUNDAMENTALS OF PROBABILITY THEORY

In this section we present a crash course in Probability Theory (or, more precisely, on those parts of the theory playing a role in the various applications we have in mind). Since this material is certainly discussed at length in any of the many available monographs on the subject, proofs are either sketched or simply omitted.

Let  $\Omega \neq \emptyset$  be a set and consider  $\mathcal{F}$  a collection of subsets of  $\Omega$ .

**Definition 2.1.** We say that  $\mathcal{F}$  is a  $\sigma$ -algebra if

- $\Omega \in \mathcal{F}$ ;
- $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$ ;
- $\{A_i\}_{i=1}^{+\infty} \subset \mathcal{F} \Rightarrow \bigcup_{i=1}^{+\infty} A_i \in \mathcal{F}$ .

Trivial examples of  $\sigma$ -algebras are  $\mathcal{F} = \{\emptyset, \Omega\}$  and  $\mathcal{F} = 2^\Omega$ , the set of all subsets of  $\Omega$ . More generally, if  $\mathcal{U} = \{U_\lambda\}_{\lambda \in \Lambda}$  is any collection of subsets of  $\Omega$ , we denote by  $\mathcal{F}_\mathcal{U} = \mathcal{F}(U_\lambda)$  the  $\sigma$ -algebra generated by  $\mathcal{U}$ . By definition, this is the smallest  $\sigma$ -algebra contained all elements of  $\mathcal{U}$ . For example, if  $\mathcal{O}^n$  is the set of open subsets in  $\mathbb{R}^n$  then  $\mathcal{B}^n := \mathcal{F}_{\mathcal{O}^n}$  is the  $\sigma$ -algebra of Borel subsets.

**Definition 2.2.** A *measure* on  $(\Omega, \mathcal{F})$  is a real valued function  $P$  on  $\mathcal{F}$  so that:

- $P(\emptyset) = 0$ ;
- $P(A) \geq 0$ , for any  $A \in \mathcal{F}$ ;
- if  $\{A_i\}_{i=1}^{+\infty} \subset \mathcal{F}$  satisfies  $A_i \cap A_j = \emptyset, i \neq j$ , then

$$P(\bigcup_i A_i) = \sum_i P(A_i).$$

We say that a triple  $(\Omega, \mathcal{F}, P)$  is a *measure space*. A classical example is  $(\mathbb{R}^n, \mathcal{L}^n, \lambda^n)$ , where  $\mathcal{L}^n$  is the standard completion of  $\mathcal{B}^n$  and  $\lambda^n$  is Lebesgue measure. If  $P(\Omega) = 1$  then we say that  $(\Omega, \mathcal{F}, P)$  is a *probability space*, a basic notion in Probability Theory. In this setting, and when no confusion arises, we will represent the corresponding Lebesgue spaces simply by  $L^p(\Omega)$ ,  $1 \leq p < +\infty$ , with no further reference to the additional data defining the associated probability space. Also, each set  $A \in \mathcal{F}$  is called an *event*.

Another key notion is that of *random vector*. If  $(\Omega, \mathcal{F}, P)$  is a probability space, this is a function  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^n, \mathcal{L}^n)$  which is measurable in the sense that  $X^{-1}(B) \in \mathcal{F}$  for any  $B \in \mathcal{B}^n$ . When  $n = 1$  we say that  $X$  is a *random variable*. If  $X$  is a random vector then we denote by  $\mathcal{F}_X$  the  $\sigma$ -algebra generated by  $X$ , i.e. the  $\sigma$ -algebra generated by  $\{X^{-1}(B); B \in \mathcal{L}^n\}$ . A similar definition holds for any collection  $\{X_\lambda\}_{\lambda \in \Lambda}$  of random vectors and the corresponding  $\sigma$ -algebra is represented by  $\mathcal{F}_{X_\lambda}$ .

A central notion in Probability is that of *independence*<sup>2</sup>. Here we define it at several levels:

- (1) A finite collection  $A_1, \dots, A_k \in \mathcal{F}$  of events is independent if

$$P(A_1 \cap \dots \cap A_k) = P(A_1) \cdots P(A_k).$$

<sup>2</sup>As already noted in [Kol18, Section I.5]: “Historically, the independence of experiments and random variables represents the very mathematical concept that has given the theory of probabilities its peculiar stamp”.

- (2) Let  $\{\mathcal{F}_\lambda\}_{\lambda \in \Lambda}$ , where each  $\mathcal{F}_\lambda \subset \mathcal{F}$  is a  $\sigma$ -algebra. Then we say that  $\{\mathcal{F}_\lambda\}$  is independent if for any *finite* collection  $\{\mathcal{F}_{\lambda_l}\}_{l=1}^k$  and events  $A_{\lambda_l} \in \mathcal{F}_{\lambda_l}$ , we have that  $\{A_{\lambda_l}\}_{l=1}^k$  is independent.
- (3) We say that a collection  $\{X_\lambda\}$  of random variables is independent if  $\{\mathcal{F}_{X_\lambda}\}$  is independent. (Notation:  $X \perp Y$  for a pair of independent random variables).

We now consider the *expectation* (or *expected value*) of a random vector  $X : \Omega \rightarrow \mathbb{R}^n$ , which is given by

$$\mathbb{E}(X) := \int_{\Omega} X dP \in \mathbb{R}^n.$$

Usually we assume that this is finite (that is,  $X$  is integrable:  $X \in L^1(\Omega)$ ). Another key notion is that of *covariance* of two random variables: if  $X, Y : \Omega \rightarrow \mathbb{R}$  then this is given by

$$\begin{aligned} \text{cov}(X, Y) &= \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) \\ &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y). \end{aligned}$$

Here, we usually require that  $X, Y \in L^2(\Omega)$  as this implies that  $XY \in L^1(\Omega)$  by Cauchy-Schwarz.

**Definition 2.3.** We say that  $X$  and  $Y$  are *uncorrelated* if  $\text{cov}(X, Y) = 0$ .

That uncorrelatedness pertains to independence is a consequence of the next fundamental result.

**Proposition 2.4.** If  $X, Y : \Omega \rightarrow \mathbb{R}$  are independent random variables then  $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ .

*Proof.* (sketch) We may assume that  $|X|, |Y| \leq M < +\infty$ . Approximating by simple functions we have<sup>3</sup>

$$X \approx \varphi = \sum_i a_i \mathbf{1}_{F_i}, \quad Y \approx \psi = \sum_j b_j \mathbf{1}_{G_j},$$

which implies

$$XY \approx \sum_{ij} a_i b_j \mathbf{1}_{F_i \cap G_j}.$$

Thus,

$$\begin{aligned} \mathbb{E}(XY) &\approx \sum_{ij} a_i b_j P(F_i \cap G_j) \\ &\stackrel{(*)}{=} \sum_{ij} a_i b_j P(F_i) P(G_j) \\ &= \sum_i a_i P(F_i) \cdot \sum_j b_j P(G_j) \\ &\approx \mathbb{E}(X)\mathbb{E}(Y), \end{aligned}$$

where we used the independence in (\*). □

<sup>3</sup>Here and in the following, if  $A \in \mathcal{F}$  we shall denote by  $\mathbf{1}_A$  the corresponding *indicator function*:

$$\mathbf{1}_A(x) = \begin{cases} 1 & \text{if } x \in A; \\ 0 & \text{otherwise} \end{cases}$$

**Corollary 2.5.** *If  $X$  and  $Y$  are independent then they are uncorrelated.*

A more conceptual proof of Proposition 2.4 may be found in Remark 2.14 below. Also, the converse to Corollary 2.5 is not always true. It holds, however, in the important case of normally distributed random variables; see Proposition 4.10.

**Definition 2.6.** If  $X : \Omega \rightarrow \mathbb{R}^n$  is a random vector then its *distribution* (or *law*) is the probability measure  $X_{\#}P$  on  $\mathbb{R}^n$  given by

$$X_{\#}P(B) = P(X^{-1}(B)), \quad B \in \mathcal{B}^n.$$

We also represent  $X_{\#}P$  by  $P_X$  and set  $P(X \geq a) := P_X([a, +\infty))$ , etc. Also, an element  $x \in \mathbb{R}^n$  in the image of a random vector  $X : \Omega \rightarrow \mathbb{R}^n$  (or equivalently, in  $\text{supp}(X)$ ) is called a *realization* (or *observed value*) of  $X$ .

The moral here is that *any* random variable  $X$  is doomed to mediate between the (rather abstract) probability measure  $P$  and its distribution  $P_X$ , an arguably more concrete probability measure on  $\mathbb{R}$ , thus linking two levels of description of randomness. Although  $P$  may be very hard to visualize or describe, as it is defined on  $\Omega$ , a purely mathematical gadget whose inner design is rarely made explicit,  $P_X$  is a probability measure on  $\mathbb{R}$ , hence amenable to be scrutinized by means of the various summarizing quantities appearing below (densities, cumulative distribution functions, tail probabilities, quantiles, etc.). Moreover, since various aspects of the structure of  $P_X$  may in principle be efficiently probed by collecting realizations of suitable copies of  $X$  (random samples), this naturally makes it the practical object statisticians, physicists, and applied scientists actually work with.

**Definition 2.7.** If  $X : \Omega \rightarrow \mathbb{R}$  is a random variable then its *cumulative distribution function* (cdf) is the function  $F_X : \mathbb{R} \rightarrow [0, 1]$  given by  $F_X(x) = X_{\#}((-\infty, x])$ .

Notice that  $F_X$  completely determines  $X_{\#}P = P_X$ . Moreover,

$$(2.1) \quad F_X(x) = \int_{-\infty}^x dP_X, \quad x \in \mathbb{R}.$$

**Definition 2.8.** We say that random variables  $X : \Omega \rightarrow \mathbb{R}^n$  and  $Y : \Omega' \rightarrow \mathbb{R}^n$  are *identically distributed* (i.d.) if  $P_X = P_Y$ .

**Proposition 2.9.** *We have*

$$\mathbb{E}(X) = \int_{\mathbb{R}^n} \mathbf{x} dP_X,$$

where  $\mathbf{x} = (x_1, \dots, x_n)$  is the position vector. More generally, if  $f : \mathbb{R}^n \rightarrow \mathbb{R}^p$  is measurable, so that  $f(X) = f \circ X : \Omega \rightarrow \mathbb{R}^p$  is a random vector, then

$$(2.2) \quad \mathbb{E}(f(X)) = \int_{\mathbb{R}^n} f(\mathbf{x}) dP_X,$$

where  $f(\mathbf{x}) = f \circ \mathbf{x}$ .

The notion of distribution may be used to single out two important classes of random vectors.

**Definition 2.10.** Let  $X : \Omega \rightarrow \mathbb{R}^n$  be a random vector. We say that

- $X$  is *discrete* if its range  $\text{Ran}(X) := X(\Omega) \subset \mathbb{R}^n$  is countable:

$$\text{Ran}(X) = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots\}.$$

In this case, the map

$$(2.3) \quad \mathbf{x}_j \in \text{Ran}(X) \mapsto p_j := P_X(\{\mathbf{x}_j\}) \in \mathbb{R}, \quad j = 1, 2, \dots,$$

is called the *mass density function (mdf)* and satisfies

$$\sum_j p_j = 1,$$

with (2.2) meaning that

$$(2.4) \quad \mathbb{E}(f(X)) = \sum_j p_j f(\mathbf{x}_j).$$

If needed, in (2.3) we may replace  $\text{Ran}(X)$  by  $\text{supp}(P_X)$ , in which case each  $p_j > 0$ .

- $X$  is *continuous* if  $P_X$  is absolutely continuous with respect to the Lebesgue measure  $dx$ . In this case, the Radon-Nykodim derivative

$$(2.5) \quad \psi_X := \frac{dP_X}{dx} : \mathbb{R}^n \rightarrow \mathbb{R},$$

is called the *probability density function (pdf)* of  $X$ , with (2.2) meaning that

$$(2.6) \quad \mathbb{E}(f(X)) = \int_{\mathbb{R}^n} f(\mathbf{x}) \psi_X(\mathbf{x}) dx.$$

We recall that the *support* of  $P_X$  is given by

$$\text{supp}(P_X) = \{x \in \mathbb{R}; F_X(x + \varepsilon) - F_X(x - \varepsilon) > 0 \text{ for all } \varepsilon > 0\}.$$

Also notice that, at least for a distribution whose support  $\text{supp}(P_X)$  is contained in some closed, bounded interval, the absolute continuity in the second item above means that the corresponding cdf  $F_X$  is absolutely continuous, or equivalently, the following assertions hold:

- (1)  $F'_X$  exists a.s. and is integrable (both with respect to Lebesgue measure);
- (2) there holds

$$F_X(b) - F_X(a) = \int_a^b F'_X(t) dt,$$

where  $[a, b] \subset \text{supp}(P_X)$ <sup>4</sup>.

In particular,

$$\int_{-\infty}^x dP_X \stackrel{(2.1)}{=} F_X(x) = \int_{-\infty}^x F'_X(t) dt, \quad x \in \text{supp}(P_X),$$

so that from (2.5) we see that

$$(2.7) \quad \psi_X = F'_X \quad \text{a.s.}$$

<sup>4</sup>For proofs of these claims we refer to [RF10, Sections 6.4, 6.5, 18.4 and 20.3].

**Convention 2.11.** Given a random vector  $X : \Omega \rightarrow \mathbb{R}^n$ , unless otherwise stated we will always assume in the sequel that it is continuous in the sense above so that it possesses a pdf  $\psi_X$ . Moreover, in order to being able to use the standard methods of calculus (including its fundamental theorem) we further assume that  $\psi_X$  is piecewise smooth with at most finitely many singularities. We insist, however, that virtually any assertion involving integration in the continuous case admits, if properly interpreted, an immediate rewording in the discrete case (and vice versa). For instance, the right-hand side of (2.4) may be written as  $\int_{\mathbb{R}^n} f(\mathbf{x}) dP_X$ , the “abstract” Lebesgue integral of  $f$  against the discrete measure  $dP_X$ , which, if we take (2.5) into account, turns it formally indistinguishable from the “continuous” integral in the right-hand side of (2.6). Of course, this only attests in favor of the benefits of working with Lebesgue integration from the very beginning in modern probability theory.

We now consider random vectors  $X_j : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^{p_j}, \mathcal{B}^{p_j})$ ,  $j = 1, \dots, n$  with distributions  $P_{X_j}$ . We may form the random vector

$$(X_1, \dots, X_n) : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}, \mathcal{B}^{p_1} \otimes \dots \otimes \mathcal{B}^{p_n})$$

given by  $(X_1, \dots, X_n)(\omega) = (X_1(\omega), \dots, X_n(\omega))$ ,  $\omega \in \Omega$ , so that the *joint distribution*  $P_{(X_1, \dots, X_n)}$  on  $\mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}$  is well defined. Moreover, each choice of  $k$  distinct indexes  $\{i_1, \dots, i_k\} \subset \{1, \dots, n\}$  determines a *marginal distribution*  $P_{(X_{i_1}, \dots, X_{i_k})}$  induced by

$$(X_{i_1}, \dots, X_{i_k}) : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^{p_{i_1}} \times \dots \times \mathbb{R}^{p_{i_k}}, \mathcal{B}^{p_{i_1}} \otimes \dots \otimes \mathcal{B}^{p_{i_k}}).$$

**Proposition 2.12.** (*Joint distribution and pdf of a marginal*) With the notation above,

$$P_{(X_{i_1}, \dots, X_{i_k})}(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}) = \int_{\mathbb{R}^{j_1 + \dots + j_{n-k}}} dP_{(X_1, \dots, X_n)}(\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_{n-k}}),$$

where  $\{1, \dots, n\} = \{i_1, \dots, i_k\} \cup \{j_1, \dots, j_{n-k}\}$ . In particular, in the continuous case,

$$\psi_{(X_{i_1}, \dots, X_{i_k})}(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}) = \int_{\mathbb{R}^{j_1 + \dots + j_{n-k}}} \psi_{(X_1, \dots, X_n)}(\mathbf{x}_1, \dots, \mathbf{x}_n) d\mathbf{x}_{j_1} \dots d\mathbf{x}_{j_{n-k}}.$$

The next result provides a way of handling independence which is quite satisfactory from an operational viewpoint.

**Proposition 2.13.**  $\{X_j\}_{j=1}^n$  is independent if and only if

$$P_{(X_1, \dots, X_n)} = P_{X_1} \otimes \dots \otimes P_{X_n},$$

the product measure. Equivalently, in terms of the corresponding pdfs,

$$\psi_{(X_1, \dots, X_n)}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \psi_{X_1}(\mathbf{x}_1) \dots \psi_{X_n}(\mathbf{x}_n), \quad (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_n}.$$

**Remark 2.14.** We may now give another proof of Proposition 2.4: since

$$\mathbb{E}(XY) = \iint_{\mathbb{R}^2} xy dP_{(X,Y)},$$



it follows that

$$\begin{aligned}
 \mathbb{E}(XY) &= \iint_{\mathbb{R}^2} xy dP_X \otimes dP_Y \\
 &= \int_{\mathbb{R}} y \left( \int_{\mathbb{R}} x dP_X \right) dP_Y \\
 &= \mathbb{E}(X) \int_{\mathbb{R}} y dP_Y \\
 &= \mathbb{E}(X)\mathbb{E}(Y),
 \end{aligned}$$

as desired.  $\square$

**Remark 2.15.** If  $X : \Omega \rightarrow \mathbb{R}$  is a random variable and  $Y = X^2$  we may compute  $\psi_Y$  in terms of  $\psi_X$  as follows. Since  $F_Y \leq y$  if and only if  $-\sqrt{y} \leq F_X \leq \sqrt{y}$ , it follows that

$$F_{X^2}(y) = (F_X(\sqrt{y}) - F_X(-\sqrt{y})) \mathbf{1}_{[0, +\infty)}(y),$$

so that from (2.7) we get

$$\psi_{X^2}(y) = F'_{X^2}(y) = \frac{1}{2\sqrt{y}} (\psi_X(\sqrt{y}) + \psi_X(\sqrt{-y})) \mathbf{1}_{(0, +\infty)}(y).$$

A similar computation shows that

$$(2.8) \quad \psi_{\sqrt{X}}(x) = 2x\psi_X(x^2) \mathbf{1}_{[0, +\infty)}(x)$$

if  $X \geq 0$ . Also, if  $Z$  and  $V$  are given with  $V > 0$  and  $Z \perp V$  then, in terms of the joint distribution  $P_{(V, Z)}$ ,

$$F_{Z/V}(x) = P(Z \leq xV) = \iint_{\{z \leq xv\}} dP_{(V, Z)}(v, z).$$

By the independence and Proposition 2.13 we may write this as an iterated integral,

$$\begin{aligned}
 F_{Z/V}(x) &= \int_0^{+\infty} \left( \int_{\{z \leq xv\}} dP_Z(z) \right) dP_V(v) \\
 &= \int_0^{+\infty} F_Z(xv) \psi_V(v) dv,
 \end{aligned}$$

and upon derivation we find that

$$(2.9) \quad \psi_{Z/V}(x) = \int_0^{+\infty} \psi_Z(xv) v \psi_V(v) dv.$$

Under the same conditions we also have that

$$(2.10) \quad F_{ZV}(x) = \int_0^{+\infty} F_Z(xv^{-1}) \psi_V(v) dv,$$

and hence we obtain

$$\psi_{ZV}(x) = \int_0^{+\infty} \psi_Z(xv^{-1}) v^{-1} \psi_V(v) dv,$$

again upon derivation  $\square$

**Remark 2.16.** If  $X_j : (\Omega_j, \mathcal{F}_j, P^{(j)}) \rightarrow (\mathbb{R}, \mathcal{B})$ ,  $j = 1, \dots, n$ , are random variables, form the product probability space

$$(\Omega^\#, \mathcal{F}^\#, P^\#) = \otimes_j (\Omega_j, \mathcal{F}_j, P^{(j)}),$$

and define the random variables  $Y_j : (\Omega^\sharp, \mathcal{F}^\sharp, P^\sharp) \rightarrow (\mathbb{R}, \mathcal{B})$ ,  $Y_j = X_j \circ \pi_j$ , where  $\pi_j : \mathbb{R}^n \rightarrow \mathbb{R}$  is the canonical projection onto the  $j^{\text{th}}$  factor. Now, if  $B_j \in \mathcal{B}$  we have

$$\begin{aligned} P_{X_j}^{(j)}(B_j) &= P^{(j)}(X_j^{-1}(B_j)) \\ &= P^\sharp(\Omega_1 \times \cdots \times X_j^{-1}(B_j) \times \cdots \times \Omega_n) \\ &= P^\sharp(\pi_j^{-1}(X_j^{-1}(B_j))) \\ &= P^\sharp(Y_j^{-1}(B_j)), \end{aligned}$$

so that  $P_{X_j}^{(j)} = P_{Y_j}^\sharp$  for any  $j$ . On the other hand, if  $P_Y^\sharp$  is the joint distribution of the random vector  $Y = (Y_1, \dots, Y_n) : (\Omega^\sharp, \mathcal{F}^\sharp, P^\sharp) \rightarrow (\mathbb{R}^n, \mathcal{B}^n)$ ,

$$\begin{aligned} P_Y^\sharp(B_1 \times \cdots \times B_n) &= P^\sharp(Y^{-1}(B_1 \times \cdots \times B_n)) \\ &= P^\sharp(X_1^{-1}(B_1) \times \cdots \times X_n^{-1}(B_n)) \\ &= \Pi_j P^{(j)}(X_j^{-1}(B_j)) \\ &= \Pi_j P_{X_j}^{(j)}(B_j) \\ &= \Pi_j P_{Y_j}^\sharp(B_j) \\ &= (\otimes_j P_{Y_j}^\sharp)(B_1 \times \cdots \times B_n), \end{aligned}$$

so that  $P_Y^\sharp = \otimes_j P_{Y_j}^\sharp$ . Thus, by Proposition 2.13,  $\{Y_j\}_{j=1}^n$  is independent. Note that if  $\{X_j\}$  is identically distributed then  $\{Y_j\}$  is identically distributed as well, so we conclude that given any random variable  $X$  there exist  $\{Y_j\}_{j=1}^n$  which is independent and identically distributed to  $X$ . It turns out to be a bit more involved to extend this construction to *countably* many random variables [FG13, Section 9.6]; this latter assertion turns out to be a rather special case of Kolmogorov's extension (Theorem A.3 below).  $\square$

**Definition 2.17.** If  $X_1, \dots, X_n : \Omega \rightarrow \mathbb{R}$  are random variables (equivalently,  $X = (X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n$  is a random vector) we define its covariance matrix by

$$\text{cov}(X, X)_{ij} = \text{cov}(X_i, X_j),$$

where

$$\text{cov}(X_i, X_j) = \mathbb{E}(X_i X_j) - \mathbb{E}(X_i) \mathbb{E}(X_j).$$

a symmetric matrix.

In case  $n = 1$ ,  $X = X_1$ , this defines the *variance* of  $X$ :

$$(2.11) \quad \text{var}(X) = \text{cov}(X, X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2.$$

Since  $\text{var}(X) \geq 0$  we usually set  $\sigma^2 := \text{var}(X)$  and  $\sigma := \sqrt{\text{var}(X)}$ , the *standard deviation*, which we also denote by  $\text{sd}(X)$ . Note also that

$$\text{var}\left(\sum_i X_i\right) = \sum_i \text{var}(X_i) + \sum_{i \neq j} \text{cov}(X_i, X_j).$$

Thus, if the  $X_i$ 's are pairwise uncorrelated (in particular, if they are independent) then

$$(2.12) \quad \text{var}\left(\sum_i X_i\right) = \sum_i \text{var}(X_i).$$

In order to simplify the notation we sometimes also set  $\text{var}(X) = \text{cov}(X, X)$  in case  $X$  is vector valued.

In order to properly compare distinct random variables it is sometimes convenient to pass to a suitable normalization. In most cases, this is accomplished by proceeding as follows.

**Definition 2.18.** If  $X : \Omega \rightarrow \mathbb{R}$  is a random variable with  $\mathbb{E}(X) = \mu$  and  $\text{var}(X) = \sigma^2$  then its *standardization* is

$$Z = \frac{X - \mu}{\sigma}.$$

Note that  $\mathbb{E}(Z) = 0$  and  $\text{var}(Z) = 1$ , hence the terminology.

In many applications it is useful to estimate from above the tail probabilities of a random variable whose expectation/variance is given. We now present a couple of elementary results in this direction, which can be regarded as examples of (quite conservative) concentration inequalities.

**Proposition 2.19.** (*Markov's inequality*) If  $X : \Omega \rightarrow \mathbb{R}$  is a non-negative random variable and  $a > 0$  then

$$(2.13) \quad P(X \geq a) \leq \frac{\mathbb{E}(X)}{a}.$$

*Proof.* Using (2.6) with  $f(x) = x$  we compute

$$\begin{aligned} \mathbb{E}(X) &= \int_0^{+\infty} x \psi_X(x) dx \\ &\geq \int_a^{+\infty} x \psi_X(x) dx \\ &= a \int_a^{+\infty} \psi_X(x) dx \\ &= a P(X \geq a), \end{aligned}$$

as desired. □

**Corollary 2.20.** (*Chebyshev's inequality*) Let  $X : \Omega \rightarrow \mathbb{R}$  be a random variable with  $0 < \sigma^2 := \text{var}(X) < +\infty$ . Then

$$(2.14) \quad P(|X - \mathbb{E}(X)| \geq a) \leq \frac{\sigma^2}{a^2}.$$

Equivalently,

$$(2.15) \quad P(|X - \mathbb{E}(X)| \geq c\sigma) \leq c^{-2}, \quad c > 0.$$

*Proof.* Note that

$$P(|X - \mathbb{E}(X)| \geq a) = P(|X - \mathbb{E}(X)|^2 \geq a^2)$$

and use (2.13) with  $X$  replaced by  $|X - \mathbb{E}(X)|^2$  and  $a$  replaced by  $a^2$ . □

We now discuss the various modes of convergence of random variables.

**Definition 2.21.** Let  $\{X_j\}_{j=1}^{+\infty}$  a sequence of random variables and let  $X$  be another random variable (all defined on the same sample space  $(\Omega, \mathcal{F}, P)$ ). We say that

- $X_j$  converges to  $X$  *almost surely* (notation:  $X_j \xrightarrow{a.s.} X$ ) if

$$P(\lim_{j \rightarrow +\infty} X_j = X) = 1.$$

- $X_j$  converges to  $X$  *in probability* (notation:  $X_j \xrightarrow{p} X$ ) if, for any  $\varepsilon > 0$ ,

$$\lim_{j \rightarrow +\infty} P(|X_j - X| < \varepsilon) = 1.$$

- $X_j$  converges to  $X$  *in distribution* (notation:  $X_j \xrightarrow{d} X$ ) if

$$\lim_{j \rightarrow +\infty} F_{X_j}(x) = F_X(x),$$

for any  $x \in \mathbb{R}$  where  $F_X$  is continuous. Equivalently,  $\mathbb{E}(\xi(X_j)) \rightarrow \mathbb{E}(\xi(X))$  for all  $\xi : \mathbb{R} \rightarrow \mathbb{R}$  uniformly bounded and continuous.

- $X_j$  converges to  $X$  *in the mean* (notation:  $X_j \xrightarrow{m} X$ ) if

$$\lim_{j \rightarrow +\infty} \mathbb{E}(|X_j - X|^2) = 0.$$

Since, as mentioned in Definition 2.6,  $F_X$  as also known as the law of  $X$ , convergence in distribution is also referred to as convergence in law (notation:  $X_j \xrightarrow{l} X$ ). Also, the equivalence between these two ways above of defining convergence in distribution is part of the Portmanteau theorem [VdV00, Lemma 2.2].

**Proposition 2.22.** One has  $(\xrightarrow{a.s.}) \Rightarrow (\xrightarrow{p}) \Rightarrow (\xrightarrow{d})$ . Also,  $(\xrightarrow{m}) \Rightarrow (\xrightarrow{p})$  and  $(\xrightarrow{d}) \Rightarrow (\xrightarrow{p})$  if the limiting variable is constant.

The following quite useful result is worth mentioning here<sup>5</sup>.

**Theorem 2.23.** (Slutsky) If  $X_j \xrightarrow{d} X$  and  $Y_j \xrightarrow{p} c$ ,  $c \in \mathbb{R}$ , then  $X_j + Y_j \xrightarrow{d} X + c$  and  $X_j Y_j \xrightarrow{d} cX$ . Also, if  $Y_j \neq 0$  and  $c \neq 0$  then  $X_j/Y_j \xrightarrow{d} X/c$ . Finally, these assertions hold true if  $(\xrightarrow{d})$  gets replaced by  $(\xrightarrow{p})$  everywhere.

We now introduce an important notion which will allow us to make use of analytical techniques in the theory<sup>6</sup>.

**Definition 2.24.** If  $X : \Omega \rightarrow \mathbb{R}^n$  is a random vector then its characteristic function  $\phi_X : \mathbb{R}^n \rightarrow \mathbb{C}$  is given by

$$\phi_X(\mathbf{u}) = \mathbb{E}(e^{i\langle X, \mathbf{u} \rangle}) = \int_{\mathbb{R}^n} e^{i\langle \mathbf{x}, \mathbf{u} \rangle} dP_X(\mathbf{x}).$$

<sup>5</sup>We refer [Gut06, Chapter 5] for much more on the convergence properties of random variables.

<sup>6</sup>Here and in the following, we denote the inner product of vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  either by  $\langle \mathbf{a}, \mathbf{b} \rangle$  or by  $\mathbf{a}^t \mathbf{b} = \mathbf{a} \mathbf{b}^t$ , where the superscript  $t$  indicates transpose of a vector (or of a matrix, more generally). Also, we set  $\|\mathbf{a}\|^2 = \mathbf{a}^t \mathbf{a}$  for the corresponding squared norm.

If  $X$  carries a pdf  $\psi_X$  then

$$\phi_X(\mathbf{u}) = \int_{\mathbb{R}^n} e^{i\langle \mathbf{x}, \mathbf{u} \rangle} \psi_X(\mathbf{x}) d\mathbf{x}.$$

**Remark 2.25.** It is immediate from Definition 2.21 that  $X_j \xrightarrow{d} X$  implies  $\phi_{X_j} \rightarrow \phi_X$  pointwise.  $\square$

Since  $\phi_X$  is the (inverse) Fourier transform of  $P_X$ , we expect that it completely determines the corresponding cdf  $F_X$ . A proof of this general statement, at least in case  $X$  is real, may be found in [Gne18, Section 39], where an explicit formula for  $F_X$  in terms of  $\phi_X$  is indicated; see also the discussion in [Luk70, Section 3.2]. We present here two instances where this expectation is confirmed (with explicit formulas).

**Proposition 2.26.** *The following holds:*

(1) *If  $X$  is  $\mathbb{Z}$ -valued and  $p_k := P(X = k)$ ,  $k \in \mathbb{Z}$ , then*

$$(2.16) \quad p_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-iku} \phi_X(u) du, \quad k \in \mathbb{Z}.$$

(2) *If  $X$  is real and has a characteristic function  $\phi_X : \mathbb{R} \rightarrow \mathbb{C}$  such that  $|\phi_X|$  is integrable then its distribution is absolutely continuous with respect to Lebesgue measure with the corresponding pdf being continuous and given by*

$$(2.17) \quad \psi_X(x) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} e^{-ixu} \phi_X(u) du, \quad x \in \mathbb{R}.$$

*Proof.* We only prove (2.16) here<sup>7</sup>. If  $\text{supp } P_X \subset \mathbb{Z}$  then it is immediate to check that  $\phi_X$  is  $2\pi$ -periodic. Also,

$$(2.18) \quad \phi_X(u) = \sum_{l \in \mathbb{Z}} e^{ilu} p_l,$$

where the convergence is uniform. In particular,  $\phi_X$  is continuous. Now integrate over  $[-\pi, \pi]$  the product of this series by  $e^{-iku}$  and use the well-known orthogonality relations for the basis  $\{e^{imx}\}_{m=-\infty}^{+\infty}$  in order to obtain (2.17).  $\square$

**Remark 2.27.** The inversion formula (2.17) means that  $\psi_X = \widehat{\phi_X}$ , where the hat means Fourier transform. On the other hand, (2.18) provides the Fourier series expansion of  $\phi_X$  with Fourier coefficients given by (2.16).  $\square$

We now describe a simple condition on a random variable ensuring that its characteristic function is sufficiently regular.

**Proposition 2.28.** *If a random variable  $X$  satisfies  $\mathbb{E}(|X|^r) < +\infty$  for some  $r \geq 1$  then  $\phi_X \in C^r(\mathbb{R})$  and*

$$\phi_X^{(j)}(u) = i^j \mathbb{E}(X^j e^{iXu}), \quad u \in \mathbb{R}, \quad j = 1, \dots, r.$$

*In particular, as  $u \rightarrow 0$ ,*

$$\phi_X(u) = \sum_{j=0}^r \frac{i^j}{j!} \mathbb{E}(X^j) u^j + o(|u|^r).$$

<sup>7</sup>A direct proof of the inversion formula (2.17) may be found in [FG13, Chapter 13]; see also [Luk70, Theorem 3.2.2].

*Proof.* If  $r = 1$  we have  $\mathbb{E}(|X|) < +\infty$  and since

$$\frac{e^{\mathbf{i}X(u+h)} - e^{\mathbf{i}Xu}}{h} = \mathbf{i}Xe^{\mathbf{i}Xu} + o(h)$$

we may use dominated convergence to see that

$$\phi'_X(u) = \lim_{h \rightarrow 0} \mathbb{E} \left( \frac{e^{\mathbf{i}X(u+h)} - e^{\mathbf{i}Xu}}{h} \right) = \mathbf{i} \mathbb{E} (X e^{\mathbf{i}Xu}),$$

which proves this case. The general assertion for  $r \geq 2$  follows by induction taking into account that

$$\frac{(\mathbf{i}X)^j e^{\mathbf{i}X(u+h)} - (\mathbf{i}X)^j e^{\mathbf{i}Xu}}{h} = (\mathbf{i}X)^j e^{\mathbf{i}Xu} + o(h)$$

and that  $\mathbb{E}(|X|^j) \leq \mathbb{E}(|X|^r)^{j/r}$  by Hölder inequality.  $\square$

**Proposition 2.29.** *A real random variable satisfies:*

- (1)  $\phi_{\alpha X}(u) = \phi_X(\alpha u)$ ,  $\alpha \in \mathbb{R}$ . In particular,  $\phi_{-X} = \overline{\phi_X}$ .
- (2) If  $\mu = \mathbb{E}(X)$  is finite then  $\phi_X(u) = 1 + u\mu\mathbf{i} + o(|u|)$  as  $u \rightarrow 0$ . Moreover, if  $\mu = 0$  and  $\sigma^2 = \mathbb{E}(X^2)$  is finite then

$$\phi_X(u) = 1 - \frac{1}{2}\sigma^2 u^2 + o(|u|^2)$$

*Proof.* (1) is obvious and (2) is an immediate consequence of Proposition 2.28 (after Taylor expanding  $\phi_X$  around  $u = 0$ ).  $\square$

We now examine how the characteristic functions of independent random variables contribute to the characteristic and density functions of their sum or difference.

**Proposition 2.30.** *The following properties hold for independent real random variables  $X$  and  $Y$ :*

- (1)  $\phi_{X+Y} = \phi_X \phi_Y$ .
- (2)  $\psi_{X+Y} = \psi_X \star \psi_Y$ , where  $\star$  means convolution.
- (3) moreover, if  $X$  and  $Y$  are identically distributed then  $\phi_{X-Y} = |\phi_X|^2$ .

*Proof.* For (1) note that, in terms of the joint distribution  $P_{(X,Y)}$ ,

$$\begin{aligned} \phi_{X+Y}(u) &= \iint_{\mathbb{R}^2} e^{\mathbf{i}(x+y)u} dP_{(X,Y)}(x, y) \\ &\stackrel{(*)}{=} \iint_{\mathbb{R}^2} e^{\mathbf{i}xu} e^{\mathbf{i}yu} dP_X(x) \otimes dP_Y(y) \\ &= \phi_X(u) \phi_Y(u), \end{aligned}$$

where we used Proposition 2.13 in  $(*)$  and Fubini in the last step. Also, by Remark 2.27,

$$\begin{aligned} \psi_{X+Y} &= \widehat{\phi_{X+Y}} \\ &\stackrel{(3)}{=} \widehat{\phi_X \phi_Y} \\ &\stackrel{(**)}{=} \widehat{\phi_X} \star \widehat{\phi_Y} \\ &= \psi_X \star \psi_Y, \end{aligned}$$

where we used a well-known property of the Fourier transform in (\*\*). Finally, (3) follows from (1) and Proposition 2.29 (1).  $\square$

**Definition 2.31.** A random variable  $X$  is *symmetric* (about 0) if  $X$  and  $-X$  are identically distributed.

**Proposition 2.32.**  $X$  is symmetric if and only if  $\phi_X$  is  $\mathbb{R}$ -valued, in which case there holds

$$(2.19) \quad \phi_X(u) = \mathbb{E}(\cos(Xu)).$$

*Proof.* Immediate from the previous results.  $\square$

**Definition 2.33.** We say that  $\epsilon$  is a *Rademacher variable* if  $\text{supp } P_\epsilon = \{-1, 1\}$  with  $P(\epsilon = -1) = P(\epsilon = 1) = 1/2$ .

**Proposition 2.34.** If  $\{\epsilon, X\}$  is independent with  $X$  symmetric then  $X$  and  $\epsilon X$  are identically distributed.

*Proof.* The cdf of  $\epsilon X$  is

$$\begin{aligned} F_{\epsilon X}(x) &= P(\epsilon X \leq x) \\ &= P(\{X \leq x\} \cap \{\epsilon = 1\}) + P(\{-X \leq x\} \cap \{\epsilon = -1\}), \end{aligned}$$

so independence gives

$$\begin{aligned} F_{\epsilon X}(x) &= P(X \leq x) P(\epsilon = 1) + P(-X \leq x) P(\epsilon = -1) \\ &= \frac{1}{2} (F_X(x) + F_{-X}(x)) \\ &= F_X(x), \end{aligned}$$

where in the last step we used that  $F_X = F_{-X}$ .  $\square$

We now introduce another important notion which is closely related to characteristic functions.

**Definition 2.35.** The *moment generating function* (mgf) of a random vector  $X : \Omega \rightarrow \mathbb{R}^n$  is given by

$$(2.20) \quad \varphi_X(\mathbf{u}) = \mathbb{E}(e^{\langle X, \mathbf{u} \rangle}), \quad \mathbf{u} \in \mathbb{R}^n.$$

**Remark 2.36.** Here we always assume that  $\varphi_X$  is defined at least in a neighborhood  $V \subset \mathbb{R}^n$  of the origin, which happens if  $\phi_X$  is analytic there [Luk70, Section 7.2]. In this case we have  $\varphi_X(\mathbf{u}) = \phi_X(-i\mathbf{u})$ ,  $\mathbf{u} \in V$ , a replacement we shall use in the sequel without further notice. Also, if  $X \in \mathbb{R}$  then all the *moments* of  $X$ ,

$$\alpha_k(X) := \int_{-\infty}^{+\infty} x^k dP_X(s), \quad k = 0, 1, 2, \dots,$$

are finite with

$$\beta^{-1} := \limsup_k \left( \frac{\alpha_k(X)}{k!} \right)^{1/k} < +\infty,$$

so there holds

$$(2.21) \quad \varphi_X(u) = \sum_k \frac{\alpha_k(X)}{k!} u^k, \quad u \in (-\beta, \beta).$$

Thus,  $\alpha_k(X) = \varphi_X^{(k)}(0)$ , which justifies the mgf terminology. Note that the expectation and variance of  $X$  are given by

$$(2.22) \quad \mathbb{E}(X) = \varphi'_X(0), \quad \text{var}(X) = \varphi''_X(0) - (\varphi'_X(0))^2,$$

with similar formulae holding for higher order centered moments.  $\square$

**Example 2.37.** (Binomial trials as the sum of independent Bernoulli trials) Set  $\mathbb{N}^{(n)} := \{0, 1, \dots, n\}$ ,  $n \geq 1$ , and consider a discrete random variable  $X$  whose probability distribution is supported in  $\mathbb{N}^{(n)}$  and satisfies, for some  $0 < p < 1$ ,

$$(2.23) \quad P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}, \quad k \in \mathbb{N}^{(n)}.$$

We then say that  $X \sim \text{Bin}(p; n)$ , the *binomial distribution* determined by the pair  $(p, n)$ . Using that the characteristic function of  $X$  is

$$\begin{aligned} \phi_X(u) &= \sum_{k=0}^n e^{iku} \binom{n}{k} p^k (1-p)^{n-k} \\ &= \sum_{k=0}^n \binom{n}{k} (pe^{iu})^k (1-p)^{n-k}, \end{aligned}$$

which gives

$$(2.24) \quad \phi_X(u) = (1-p + pe^{iu})^n,$$

together with Proposition 2.13, Proposition 2.30 (1) and Proposition 2.26 (1), we see that, by eventually changing the underlying sample space, we may assume that  $X = X_1 + \dots + X_n$ , where  $\{X_j\}_{j=1}^n$  is independent and each  $X_j \sim \text{Bin}(p; 1) =: \text{Ber}(p)$ , the *Bernoulli distribution*, so that  $\mathbb{E}(X_j) = p$  and  $\text{var}(X_j) = p(1-p)$ . Finally,

$$(2.25) \quad \varphi_X(u) = (1-p + pe^u)^n$$

follows immediately from (2.24).  $\square$

**Example 2.38.** (Poisson trials). For each  $n \geq 1$  consider the discrete random variable  $Y$  supported on  $\mathbb{N}_0 = \{0\} \cup \mathbb{N}$  with

$$P(Y = k) = \frac{\lambda^k e^{-\lambda}}{k!}.$$

We represent this as  $Y \sim \mathcal{P}(\lambda)$ , the *Poisson distribution* with parameter  $\lambda$ . We compute:

$$\begin{aligned} \phi_Y(u) &= \sum_{k \geq 0} e^{iku} \frac{\lambda^k e^{-\lambda}}{k!} \\ &= e^{-\lambda} \sum_{k \geq 0} \frac{(\lambda e^{iu})^k}{k!} \\ &= e^{-\lambda} e^{\lambda e^{iu}}, \end{aligned}$$

which gives

$$(2.26) \quad \phi_Y(u) = e^{\lambda(e^{iu} - 1)},$$



and hence

$$(2.27) \quad \varphi_Y(u) = e^{\lambda(e^u - 1)}.$$

In particular,  $\mathbb{E}(Y) = \text{var}(Y) = \lambda$ . Also, if  $Y \sim \mathcal{P}(n)$ ,  $n \in \mathbb{N}$ , then it follows from (2.26) with  $\lambda = n$  that we may decompose  $Y = Y_1 + \cdots + Y_n$  with  $\{Y_j\}_{j=1}^n$  independent and each  $Y_j \sim \mathcal{P}(1)$ , so that  $\mathbb{E}(Y_j) = \text{var}(Y_j) = 1$ .  $\square$

### 3. CONDITIONING

We now discuss the various ways of conditioning a given random variable.

**3.1. Conditional probability.** Let us consider random vectors  $X : \Omega \rightarrow \mathbb{R}^m$  and  $Y : \Omega \rightarrow \mathbb{R}^p$  with distributions  $P_X$  and  $P_Y$ , respectively. Here,  $(\Omega, \mathcal{F}, P)$  is the underlying probability space. We also consider the joint distribution  $P_{(X,Y)}$  associated to  $(X, Y) : \Omega \rightarrow \mathbb{R}^m \times \mathbb{R}^p$ . As usual, we assume that all these distributions are absolutely continuous with respect to the underlying Lebesgue measure and hence admit pdfs. Given  $B \in \mathcal{B}^p$ , our aim here is to define the *conditional* probability that  $Y \in B$  given a realization  $\mathbf{x} \in \mathbb{R}^m$  of  $X$ .

**Definition 3.1.** A *Markov kernel* is a map  $\kappa : \mathbb{R}^m \times \mathcal{B}^p \rightarrow [0, 1]$  such that:

- $x \mapsto \kappa(\mathbf{x}, B)$  is  $\mathcal{B}^m$ -measurable for any  $B \in \mathcal{B}^p$ ;
- $B \mapsto \kappa(\mathbf{x}, B)$  is a probability measure on  $(\mathbb{R}^p, \mathcal{B}^p)$  for any  $\mathbf{x} \in \mathbb{R}^m$ .

Given a Markov kernel  $\kappa$  and a probability measure  $\mu$  on  $(\mathbb{R}^m, \mathcal{B}^m)$ , the rule

$$(A, B) \mapsto (\mu \star \kappa)(A, B) := \int_A \kappa(\mathbf{x}, B) d\mu(\mathbf{x}), \quad (A, B) \in \mathcal{B}^m \times \mathcal{B}^p,$$

defines a probability measure in  $(\mathbb{R}^m \times \mathbb{R}^p, \mathcal{B}^m \times \mathcal{B}^p)$ .

**Proposition 3.2.** There exists a unique Markov kernel  $P_{Y|X}$  such that

$$P_{(X,Y)} = P_X \star P_{Y|X}.$$

In other words,

$$(3.1) \quad P_{(X,Y)}(A, B) = \int_A P_{Y|X}(\mathbf{x}, B) dP_X(\mathbf{x}), \quad (A, B) \in \mathcal{B}^m \times \mathcal{B}^p.$$

*Proof.* See [Kle13, Chapter 8].  $\square$

**Definition 3.3.** If  $\mathbf{x} \in \mathbb{R}^m$  we define the *conditional probability*

$$P_{Y|X=\mathbf{x}} = P_{Y|X}(\mathbf{x}, \cdot),$$

which is a probability measure in  $(\mathbb{R}^p, \mathcal{B}^p)$ .

Thus,

$$(3.2) \quad P(Y \in B | X=\mathbf{x}) := P_{Y|X=\mathbf{x}}(B) = P_{Y|X}(\mathbf{x}, B), \quad B \in \mathcal{B}^p,$$

should be interpreted as the *conditional probability* that  $Y \in B$  given that  $X = \mathbf{x}$ . It is immediate from (3.1) that the corresponding pdf's satisfy

$$(3.3) \quad \psi_{Y|X=\mathbf{x}}(\mathbf{y}) = \frac{\psi_{(X,Y)}(\mathbf{x}, \mathbf{y})}{\psi_X(\mathbf{x})}, \quad \mathbf{y} \in \mathbb{R}^p,$$

whenever  $\psi_X(\mathbf{x}) > 0$ , so that the corresponding *conditional expectation function* and *conditional covariance function* are

$$(3.4) \quad \mathbb{E}(Y|X=\mathbf{x}) = \int_{\mathbb{R}^p} \mathbf{y} \psi_{Y|X=\mathbf{x}}(\mathbf{y}) d\mathbf{y}$$

and

$$(3.5) \quad \text{cov}(Y|X=\mathbf{x}) = \int_{\mathbb{R}^p} (\mathbf{y} - \mathbb{E}(Y|X=\mathbf{x}))^2 \psi_{Y|X=\mathbf{x}}(\mathbf{y}) d\mathbf{y},$$

respectively.

**Remark 3.4.** Whenever possible, we may simply dispense with the existence theory sketched above and adopt (3.3) as the definition of the *conditional pdf* of  $Y$  given  $X = \mathbf{x}$ .  $\square$

We now present a few elementary (but useful!) consequences of the theory above. To simplify matters, in the rest of this subsection we shall assume that all random variables are real, continuous and have positive pdfs everywhere. In particular, with a little abuse we may view (3.3) as the pdf of the “conditioned random variable”  $Y|_{X=\mathbf{x}}$ , so that all the concepts pertaining random variables introduced so far may be easily extended to this broader setting.

**Proposition 3.5.** *The following hold:*

- (1) *If  $\{X, Y\}$  is independent then  $\mathbb{E}(Y) = \mathbb{E}(Y|X=\mathbf{x})$  for any  $\mathbf{x}$ ;*
- (2) *If  $\{X, Y, Z\}$  is independent then  $\{X|Z=\mathbf{z}, Y|Z=\mathbf{z}\}$  is independent for any  $\mathbf{z}$ .*

*Proof.* For (1) note that by (3.3) and Proposition 2.13,

$$(3.6) \quad \psi_{Y|X=\mathbf{x}}(\mathbf{y}) = \frac{\psi_{(X,Y)}(\mathbf{x}, \mathbf{y})}{\psi_X(\mathbf{x})} = \frac{\psi_X(\mathbf{x})\psi_Y(\mathbf{y})}{\psi_X(\mathbf{x})} = \psi_Y(\mathbf{y}).$$

As for (2), again by Proposition 2.13,

$$\begin{aligned} \psi_{(X|Z=\mathbf{z}, Y|Z=\mathbf{z})}(\mathbf{x}, \mathbf{y}) &= \frac{\psi_{(X,Y,Z)}(\mathbf{x}, \mathbf{y}, \mathbf{z})}{\psi_Z(\mathbf{z})} \\ &= \frac{\psi_X(\mathbf{x})\psi_Y(\mathbf{y})\psi_Z(\mathbf{z})}{\psi_Z(\mathbf{z})} \\ &\stackrel{(3.6)}{=} \psi_{X|Z=\mathbf{z}}(\mathbf{x})\psi_{Y|Z=\mathbf{z}}(\mathbf{y}), \end{aligned}$$

and the result follows from the “conditioned” version of Proposition 2.13.  $\square$

**Remark 3.6.** In general the converses to both items in Proposition 3.5 fail to hold true if the independence assumptions are removed.  $\square$

Finally, we present another consequence of (3.3) with notable applications to the so-called Bayesian approach to Statistical Inference; see Subsection 13.

**Theorem 3.7.** (Bayes rule) If both  $\psi_X$  and  $\psi_Y$  are everywhere positive then

$$\psi_{X|Y=\mathbf{y}}(\mathbf{x}) = \frac{\psi_{Y|X=\mathbf{x}}(\mathbf{y})\psi_X(\mathbf{x})}{\psi_Y(\mathbf{y})}, \quad (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^m \times \mathbb{R}^q,$$

with

$$\psi_Y(\mathbf{y}) = \int_{\mathbb{R}^m} \psi_{Y|X=\mathbf{x}}(\mathbf{y})\psi_X(\mathbf{x})d\mathbf{x}.$$

*Proof.* Just exchange the roles of  $X$  and  $Y$  in (3.3) and eliminate the common term  $\psi_{(X,Y)}(\mathbf{x}, \mathbf{y}) = \psi_{(Y,X)}(\mathbf{y}, \mathbf{x})$  in the resulting formulas.  $\square$

**3.2. Conditional expectation.** Here we discuss how conditioning the expectation of a random variable with respect to a  $\sigma$ -subalgebra and then relate this to the discussion in the previous subsection (Proposition 3.14).

**Proposition 3.8.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space and let  $\mathcal{G} \subset \mathcal{F}$  a  $\sigma$ -subalgebra. Given a random vector  $X : \Omega \rightarrow \mathbb{R}^n$  there exists a unique random vector  $Y : \Omega \rightarrow \mathbb{R}^n$  which is  $\mathcal{G}$ -measurable and satisfies

$$\int_G Y dP = \int_G X dP, \quad G \in \mathcal{G}.$$

*Proof.* Define a measure  $Q$  on  $\mathcal{P}$  by

$$Q(G) = \int_G X dP, \quad G \in \mathcal{G}.$$

Clearly,  $Q$  is absolutely continuous with respect to  $P|_{\mathcal{G}}$ . Now take  $Y = dQ/dP|_{\mathcal{G}}$ .  $\square$

**Remark 3.9.** Given  $X \in L^2(\Omega, \mathcal{F}, P)$  consider the closed subspace  $L^2(\Omega, \mathcal{G}, P|_{\mathcal{G}}) \subset L^2(\Omega, \mathcal{F}, P)$  and let  $\pi : L^2(\Omega, \mathcal{F}, P) \rightarrow L^2(\Omega, \mathcal{G}, P|_{\mathcal{G}})$  be the standard orthogonal projection. Hence,  $Y = \pi X$ .  $\square$

**Definition 3.10.** We call  $Y = \mathbb{E}(X|\mathcal{G})$  the conditional expectation of  $X$  given  $\mathcal{G}$ . If  $\mathcal{G} = \mathcal{F}_Z$  for some other  $Z$  then we set  $\mathbb{E}(X|Z) := \mathbb{E}(X|\mathcal{F}_Z)$ .

Note that  $\mathbb{E}(X|\mathcal{G})$  is characterized by

$$(3.7) \quad \int_G \langle Z, \mathbb{E}(X|\mathcal{G}) \rangle dP = \int_G \langle Z, X \rangle dP, \quad G \in \mathcal{G},$$

for any  $Z : \Omega \rightarrow \mathbb{R}^n$   $\mathcal{G}$ -measurable.

**Proposition 3.11.** Conditional expectation satisfies the following properties:

- (1)  $\mathbb{E}(aX + bX'|\mathcal{G}) = a\mathbb{E}(X|\mathcal{G}) + b\mathbb{E}(X'|\mathcal{G})$ ;
- (2)  $\mathbb{E}(\mathbb{E}(X|\mathcal{G})) = \mathbb{E}(X)$ ;
- (3) if  $X$  is  $\mathcal{G}$ -measurable then  $\mathbb{E}(X|\mathcal{G}) = X$ ;
- (4) if  $X \perp \mathcal{G}$  then  $\mathbb{E}(X|\mathcal{G}) = \mathbb{E}(X)$ ;
- (5) if  $\mathcal{G} \subset \mathcal{H}$  then  $\mathbb{E}(X|\mathcal{G}) = \mathbb{E}(\mathbb{E}(X|\mathcal{H})|\mathcal{G})$ ;
- (6) if  $X \perp \mathcal{G}$  then  $\mathbb{E}(XY|\mathcal{G}) = \mathbb{E}(X)\mathbb{E}(Y|\mathcal{G})$ . In particular, if  $Y$  is  $\mathcal{G}$ -measurable then  $\mathbb{E}(XY|\mathcal{G}) = Y\mathbb{E}(X)$ ;
- (7) if  $X$  is  $\mathcal{G}$ -measurable then  $\mathbb{E}(XY|\mathcal{G}) = X\mathbb{E}(Y|\mathcal{G})$ .

*Proof.* (1) and (2) are obvious. For (3), just think of  $X : (\Omega, \mathcal{G}) \rightarrow \mathbb{R}^n$  as a random vector. For (4),

$$\begin{aligned} \int_G X dP &= \int_\Omega X \mathbf{1}_G dP \\ &\stackrel{(*)}{=} \int_\Omega X dP \int_\Omega \chi_G dP \\ &= \mathbb{E}(X) P(G) \\ &= \int_G \mathbb{E}(X) dP, \end{aligned}$$

where the assumption was used in (\*). The result then follows by uniqueness. For (5), note that  $G \in \mathcal{G}$  implies  $G \in \mathcal{H}$  and hence

$$\int_G \mathbb{E}(X|\mathcal{H}) dP = \int_G X dP = \int_G \mathbb{E}(X|\mathcal{G}) dP.$$

Also, (6) is the obvious generalization of (4): using that  $X \perp Y \mathbf{1}_G$ ,  $G \in \mathcal{G}$ , we show that  $\mathbb{E}(XY|\mathcal{G}) = \mathbb{E}(X|\mathcal{G})\mathbb{E}(Y|\mathcal{G}) = \mathbb{E}(X)\mathbb{E}(Y|\mathcal{G})$ . Finally, if  $G \in \mathcal{G}$ ,

$$\begin{aligned} \int_G X \mathbb{E}(Y|\mathcal{G}) dP &\stackrel{(3)}{=} \int_G \mathbb{E}(X|\mathcal{G}) \mathbb{E}(Y|\mathcal{G}) dP \\ &\stackrel{(3.7)}{=} \int_G \mathbb{E}(X|\mathcal{G}) Y dP \\ &\stackrel{(3)}{=} \int_G XY dP, \end{aligned}$$

which proves (7). □

**Example 3.12.** (Birkhoff ergodic theorem) If  $(\Omega, \mathcal{F}, P)$  is a probability space then  $T : \Omega \rightarrow \Omega$  is *measure preserving* if  $T$  is  $\mathcal{F}$ -measurable and satisfies  $P(T^{-1}(A)) = P(A)$  for any event  $A \in \mathcal{F}$ . Given a random variable  $X : \Omega \rightarrow \mathbb{R}$  we then define, for  $n \in \mathbb{N}$ ,  $X_n^{(T)} : \Omega \rightarrow \mathbb{R}$  by

$$X_n^{(T)}(\omega) = \frac{1}{n} \sum_{j=0}^{n-1} X(T^j \omega).$$

A version of *Birkhoff's ergodic theorem* [Kre11, Section 1.2] says that there exists a random variable  $X^{(T)}$  such that:

$$(3.8) \quad P\left(\omega \in \Omega; X_n^{(T)}(\omega) \rightarrow_{n \rightarrow +\infty} X^{(T)}(\omega)\right) = 1,$$

from which it follows that  $X^{(T)} \circ T = X^{(T)}$  and  $\mathbb{E}(X^{(T)}) = \mathbb{E}(X)$ . In order to identify  $X^{(T)}$  let us consider

$$\mathcal{G}_T = \{A \in \mathcal{F}; T^{-1}(A) = A\},$$

the  $\sigma$ -subalgebra of  $T$ -invariant events. If  $G \in \mathcal{G}_T$  define  $X_G := \mathbf{1}_G X$ . Thus,

$$X_G(T^j \omega) = \mathbf{1}_G(T^j \omega) X(T^j \omega) = \mathbf{1}_G(\omega) X(T^j \omega),$$

so if we use (3.8) with  $X_G$  replacing  $X$  we see that  $X_G^{(T)} = \mathbf{1}_G X^{(T)}$  and hence  $\mathbb{E}(\mathbf{1}_G X^{(T)}) = \mathbb{E}(X_G) = \mathbb{E}(\mathbf{1}_G X)$ , which means that  $X^{(T)} = \mathbb{E}(X|\mathcal{G}_T)$ . Also, if  $T$  is *ergodic* in the sense that

$$\mathcal{G}_T = \{A \in \mathcal{F}; P(A) = 0 \text{ or } P(A) = 1\}$$

then it is immediate to check that  $X \perp \mathcal{G}_T$  and it follows from Proposition 3.11 (4) that

$$P\left(\omega \in \Omega; X_n^{(T)}(\omega) \rightarrow_{n \rightarrow +\infty} \mathbb{E}(X)\right) = 1,$$

which is an improvement of (3.8). □

**Example 3.13.** (von Neumann ergodic theorem) The  $L^2$  version of Example 3.12 goes as follows. For any measure preserving  $T$  as above let us consider

$$L^2(\Omega, \mathcal{G}_T, P|_{\mathcal{G}_T}) = \left\{ \tilde{X} \in L^2(\Omega, \mathcal{F}, P); \tilde{X} \circ T = \tilde{X} \right\}.$$

Then *von Neumann's mean ergodic theorem* [Kre11, Section 1.1] assures that for any random variable  $X \in L^2(\Omega, \mathcal{F}, P)$  there exists a unique  $X^{[T]} \in L^2(\Omega, \mathcal{G}_T, P|_{\mathcal{G}_T})$  such that

$$(3.9) \quad \lim_{n \rightarrow +\infty} \mathbb{E}(|X_n^{(T)} - X^{[T]}|^2) = 0.$$

Using that

$$\iota_T : L^2(\Omega, \mathcal{F}, P) \rightarrow L^2(\Omega, \mathcal{F}, P), \quad \iota_T(X) = X \circ T,$$

is an isometry we compute, for any  $\tilde{X} \in L^2(X, \mathcal{G}_T, P|_{\mathcal{G}_T})$ ,

$$\begin{aligned} \mathbb{E}(X^{[T]} \tilde{X}) &= \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{j=0}^{n-1} \mathbb{E}((X \circ T^j) \tilde{X}) \\ &= \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{j=0}^{n-1} \mathbb{E}((X \circ T^j)(\tilde{X} \circ T^j)) \\ &= \lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{j=0}^{n-1} \mathbb{E}(X \tilde{X}) \\ &= \mathbb{E}(X \tilde{X}), \end{aligned}$$

which means that  $X^{[T]}$  is the  $L^2$  projection of  $X$  over  $L^2(\Omega, \mathcal{G}_T, P|_{\mathcal{G}_T})$ . It follows from Remark 3.9 that  $\bar{X}^{[T]} = \mathbb{E}(X|\mathcal{G}_T)$ , so that (3.9) is the “mean squared” version of (3.8) above.  $\square$

We now give an useful rewording of the random variable induced by the conditional expectation function (3.4) in terms of the notion of conditional expectation appearing in Definition 3.10.

**Proposition 3.14.** *If, as in (3.4), we set*

$$(3.10) \quad g(\mathbf{x}) = \mathbb{E}(Y|_{X=\mathbf{x}}), \quad \mathbf{x} \in \mathbb{R}^m,$$

*then*

$$(3.11) \quad g(X) = \mathbb{E}(Y|X).$$

*In particular,*

$$(3.12) \quad \mathbb{E}(g(X)) = \mathbb{E}(Y).$$

*Proof.* We need to check that

$$\int_C g(X) dP = \int_C Y dP, \quad C \in \mathcal{F}_X,$$

so we write  $C = X^{-1}(A)$ ,  $A \in \mathcal{B}^m$ , in order to have

$$(3.13) \quad \mathbf{1}_C(\omega) = \mathbf{1}_A(X(\omega)), \quad \omega \in \Omega.$$

We first note that

$$\begin{aligned}
 \int_C g(X) dP &= \int_{\Omega} \mathbf{1}_C g(X) dP \\
 &\stackrel{(3.13)}{=} \int_{\Omega} \mathbf{1}_A(X) g(X) dP \\
 &= \int_{\mathbb{R}^m} \mathbf{1}_A(\mathbf{x}) g(\mathbf{x}) \psi_X(\mathbf{x}) d\mathbf{x},
 \end{aligned}$$

and using both (3.10) and (3.4),

$$\int_C g(X) dP = \int_{\mathbb{R}^m} \mathbf{1}_A(\mathbf{x}) \left( \int_{\mathbb{R}^p} \mathbf{y} \psi_{Y|X=\mathbf{x}}(\mathbf{y}) d\mathbf{y} \right) \psi_X(\mathbf{x}) d\mathbf{x}.$$

From Fubini and (3.3) we thus get

$$\begin{aligned}
 \int_C g(X) dP &= \int \int_{\mathbb{R}^m \times \mathbb{R}^p} \mathbf{1}_A(\mathbf{x}) \mathbf{y} \psi_{(X,Y)}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \\
 &= \int \int_{\mathbb{R}^m \times \mathbb{R}^p} \mathbf{1}_A(\mathbf{x}) \mathbf{y} dP_{(X,Y)} \\
 &= \mathbb{E}(\mathbf{1}_A(X) Y) \\
 &\stackrel{(3.13)}{=} \mathbb{E}(\mathbf{1}_C Y) \\
 &= \int_C Y dP,
 \end{aligned}$$

which proves (3.11). Finally, (3.14) follows from Proposition 3.11, (2). □

**Corollary 3.15.** (*Law of total expectation and variance*) *There hold*

$$(3.14) \quad \mathbb{E}(X) = \mathbb{E}(\mathbb{E}(X|Y))$$

and

$$(3.15) \quad \text{var}(X) = \mathbb{E}(\text{var}(X|Y)) + \text{var}(\mathbb{E}(X|Y)).$$

Note that (3.14) corresponds to item (2) in Proposition 3.11 with  $\mathcal{G} = \mathcal{F}_Y$ .

#### 4. NORMALLY DISTRIBUTED RANDOM VARIABLES AND THEIR FRIENDS

Here we isolate and study important families of random variables which somehow relate to the normally distributed random variables. A full account of this topic appears in [Ton90].

**4.1. Normally distributed random variables.** We start with the family of random variables which certainly is the most pervasive in Probability Theory and its applications.

**Definition 4.1.** We say that a random vector  $X : \Omega \rightarrow \mathbb{R}^n$  is normally distributed if its pdf  $\psi_X : \mathbb{R}^n \rightarrow \mathbb{R}$  is given by

$$(4.1) \quad \psi_X(\mathbf{x}) = \frac{\sqrt{\det A}}{(2\pi)^{n/2}} e^{-\frac{1}{2} \langle A(\mathbf{x}-\boldsymbol{\mu}), \mathbf{x}-\boldsymbol{\mu} \rangle},$$

where  $A$  is a positive definite, symmetric matrix and  $\boldsymbol{\mu} \in \mathbb{R}^n$ . We then write  $X \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ , or simply  $X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma} = A^{-1}$ .

The next proposition shows that this is well defined.

**Proposition 4.2.** *One has  $\int_{\mathbb{R}^n} \psi_X(\mathbf{x}) d\mathbf{x} = 1$ .*

*Proof.* Take  $O$  an orthogonal matrix so that

$$OAO^{-1} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n),$$

and define  $\mathbf{y} = O(\mathbf{x} - \boldsymbol{\mu})$ . It follows that

$$\psi_X(\mathbf{x}) = \frac{\sqrt{\det A}}{(2\pi)^{n/2}} e^{-\frac{1}{2} \langle \Lambda \mathbf{y}, \mathbf{y} \rangle},$$

so that

$$\int_{\mathbb{R}^n} \psi_X(\mathbf{x}) d\mathbf{x} = \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \Pi_i \int_{\mathbb{R}} e^{-\frac{1}{2} \lambda_i y_i^2} dy_i.$$

Thus, if  $z_i = \sqrt{\lambda_i/2} y_i$  then

$$\int_{\mathbb{R}} e^{-\frac{1}{2} \lambda_i y_i^2} dy_i = \sqrt{\frac{2}{\lambda_i}} \int_{\mathbb{R}} e^{-z_i^2} dz_i = \sqrt{\frac{2\pi}{\lambda_i}},$$

so that

$$\int_{\mathbb{R}^n} \psi_X(\mathbf{x}) d\mathbf{x} = \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \frac{(2\pi)^{n/2}}{\Pi_i \sqrt{\lambda_i}} = 1,$$

where we used that  $\det A = \Pi_i \lambda_i$  in the last step. □

In general, if  $X : \Omega \rightarrow \mathbb{R}^n$  is a random vector we have defined its *expectation vector*

$$\boldsymbol{\mu}(X) = \mathbb{E}(X),$$

and its *covariance matrix*

$$\text{cov}(X, X)_{ij} = \text{cov}(X_i, X_j).$$

Sometimes we write this simply as  $\text{cov}(X)$ . We now compute these invariants assuming that  $X$  is normally distributed.

**Proposition 4.3.** *If  $X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $\boldsymbol{\Sigma} = A^{-1}$ , then  $\boldsymbol{\mu}(X) = \boldsymbol{\mu}$  and  $\text{cov}(X) = \boldsymbol{\Sigma}$ .*

*Proof.* In terms of the substitution above we have  $\mathbf{x} = \boldsymbol{\mu} + Q\mathbf{y}$ ,  $Q = O^{-1}$ . Thus,

$$\begin{aligned} \mathbb{E}(X_i) &= \int_{\mathbb{R}^n} x_i \psi_X(\mathbf{x}) d\mathbf{x} \\ &= \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} \left( \mu_i + \sum_j Q_{ij} y_j \right) \Pi_k e^{-\frac{1}{2} \lambda_k y_k^2} dy_1 \cdots dy_n. \end{aligned}$$

But

$$(4.2) \quad \int_{\mathbb{R}} y_j e^{-\frac{1}{2} \lambda_j y_j^2} dy_j = 0,$$

so we get  $\mu(X)_i = \mu_i$ . Also,

$$\begin{aligned}\mathbb{E}(X_i X_j) &= \int_{\mathbb{R}^n} x_i x_j \psi_X(\mathbf{x}) d\mathbf{x} \\ &= \int_{\mathbb{R}^n} \left( \mu_i + \sum_k Q_{ik} y_k \right) \left( \mu_j + \sum_l Q_{jl} y_l \right) \psi_X(\mathbf{x}) d\mathbf{x},\end{aligned}$$

and using (4.2) again we get

$$\begin{aligned}\mathbb{E}(X_i X_j) &= \int_{\mathbb{R}^n} \mu_i \mu_j \psi_X(\mathbf{x}) d\mathbf{x} + \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \sum_{kl} Q_{ik} Q_{jl} \int_{\mathbb{R}^n} y_k y_l \Pi_p(e^{-\frac{1}{2}\lambda_p y_p^2}) dy_p \\ &= \mu_i \mu_j + \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \sum_k Q_{ik} Q_{jk} \int_{\mathbb{R}^n} y_k^2 \Pi_p e^{-\frac{1}{2}\lambda_p y_p^2} dy_p \\ &= \mu_i \mu_j + \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \sum_k Q_{ik} Q_{jk} \int_{\mathbb{R}} y_k^2 e^{-\frac{1}{2}y_k^2} dy_k \times \Pi_{p \neq k} \int_{\mathbb{R}} e^{-\frac{1}{2}y_p^2} dy_p \\ &= \mu_i \mu_j + \frac{\sqrt{\det A}}{(2\pi)^{n/2}} (2\pi)^{\frac{n-1}{2}} \sum_k Q_{ik} Q_{jk} \int_{\mathbb{R}} y_k^2 e^{-\frac{1}{2}y_k^2} dy_k \times \Pi_{p \neq k} \frac{1}{\lambda_p^{1/2}}.\end{aligned}$$

But

$$\begin{aligned}\int_{\mathbb{R}} y_k^2 e^{-\frac{1}{2}\lambda_k y_k^2} dy_k &= -\frac{y_k}{\lambda_k} e^{-\frac{1}{2}\lambda_k y_k^2} \Big|_{-\infty}^{+\infty} + \frac{1}{\lambda_k} \int_{\mathbb{R}^n} e^{-\frac{1}{2}\lambda_k y_k^2} dy_k \\ &= \frac{1}{\lambda_k} \frac{(2\pi)^{1/2}}{\lambda_k^{1/2}},\end{aligned}$$

so that

$$\mathbb{E}(X_i X_j) = \mu_i \mu_j + \sum_k \frac{Q_{ik} Q_{jk}}{\lambda_k} = \mu_i \mu_j + \Sigma_{ij},$$

where we used that  $\Sigma = A^{-1} = Q\Lambda^{-1}Q^{-1}$ . This completes the proof.  $\square$

We now compute the characteristic function of a normally distributed random variable.

**Proposition 4.4.** *If  $X \sim \mathcal{N}(\mu, \Sigma)$  then*

$$(4.3) \quad \phi_X(\mathbf{u}) = e^{\langle \mu, \mathbf{u} \rangle - \frac{1}{2} \langle \Sigma \mathbf{u}, \mathbf{u} \rangle}.$$

*Proof.* Recalling that  $Q^{-1}AQ = \Lambda$  and  $\mathbf{x} = \mu + Q\mathbf{y}$ , we have

$$\begin{aligned}\phi_X(\mathbf{u}) &= \int_{\mathbb{R}^n} e^{\langle \mathbf{x}, \mathbf{u} \rangle} \frac{\sqrt{\det A}}{(2\pi)^{n/2}} e^{-\frac{1}{2} \langle A(\mathbf{x} - \mu), \mathbf{x} - \mu \rangle} d\mathbf{x} \\ &= \frac{\sqrt{\det A}}{(2\pi)^{n/2}} \int_{\mathbb{R}^n} e^{\langle \mu, \mathbf{u} \rangle + \langle \mathbf{y}, \mathbf{v} \rangle - \frac{1}{2} \langle \Lambda \mathbf{y}, \mathbf{y} \rangle} d\mathbf{y},\end{aligned}$$

where  $\mathbf{v} = Q^t \mathbf{u}$ . Now observe that if  $(\cdot, \cdot)$  is the sesquilinear product in  $\mathbb{C}^n$  then

$$\begin{aligned}-\frac{1}{2} \left( \Lambda^{1/2} \mathbf{y} - \mathbf{i} \Lambda^{-1/2} \mathbf{v}, \Lambda^{1/2} \mathbf{y} - \mathbf{i} \Lambda^{-1/2} \mathbf{v} \right) &= -\frac{1}{2} \langle \Lambda \mathbf{y}, \mathbf{y} \rangle \\ &\quad + \frac{1}{2} \langle \Lambda^{-1} \mathbf{v}, \mathbf{v} \rangle + \mathbf{i} \langle \mathbf{v}, \mathbf{y} \rangle,\end{aligned}$$



which gives

$$\begin{aligned}
 \phi_X(\mathbf{u}) &= \frac{\sqrt{\det A}}{(2\pi)^{n/2}} e^{\langle \mu, \mathbf{u} \rangle - \frac{1}{2} \langle \Lambda^{-1} \mathbf{v}, \mathbf{v} \rangle} \int_{\mathbb{R}^n} e^{-\frac{1}{2} (\Lambda^{1/2} \mathbf{y} - i \Lambda^{-1/2} \mathbf{v}, \Lambda^{1/2} \mathbf{y} - i \Lambda^{-1/2} \mathbf{v})} d\mathbf{y} \\
 &\stackrel{(*)}{=} \frac{\sqrt{\det A}}{(2\pi)^{n/2}} e^{\langle \mu, \mathbf{u} \rangle - \frac{1}{2} \langle Q \Lambda^{-1} Q^{-1} \mathbf{u}, \mathbf{u} \rangle} \int_{\mathbb{R}^n} e^{-\frac{1}{2} (\Lambda^{1/2} \mathbf{y}, \Lambda^{1/2} \mathbf{y})} d\mathbf{y} \\
 &= e^{\langle \mu, \mathbf{u} \rangle - \frac{1}{2} \langle \Sigma \mathbf{u}, \mathbf{u} \rangle},
 \end{aligned}$$

where in  $(*)$  we changed the contour of integration (and used the appropriate multi-dimensional version of Cauchy's theorem).  $\square$

**Corollary 4.5.** *If  $X \sim \mathcal{N}(\mu, \Sigma)$  then its mgf is*

$$(4.4) \quad \varphi_X(\mathbf{u}) = e^{\langle \mu, \mathbf{u} \rangle + \frac{1}{2} \langle \Sigma \mathbf{u}, \mathbf{u} \rangle}.$$

*In particular, if  $n = 1$  and  $X \sim \mathcal{N}(\mu, \sigma^2)$  then*

$$(4.5) \quad \varphi_X(u) = e^{u\mu + \frac{1}{2} \sigma^2 u^2}.$$

**Corollary 4.6.** *If  $X$  is a normally distributed random vector,  $X \sim \mathcal{N}(\mu, \Sigma)$ , then its distribution  $P_X$  is completely determined by its characteristic function.*

*Proof.* Note that

$$\int_{\mathbb{R}^n} |\phi_X(\mathbf{u})| d\mathbf{u} = \int_{\mathbb{R}^n} e^{-\frac{1}{2} \langle \Sigma \mathbf{u}, \mathbf{u} \rangle} d\mathbf{u} < +\infty$$

and apply (the appropriate multi-variate version of) Proposition 2.26 (2).  $\square$

These latter results have a number of consequences that we now explore.

**Proposition 4.7.** *If  $X : \Omega \rightarrow \mathbb{R}$  is normally distributed,  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then*

- (1)  $\phi_X(u) = e^{\mu u - \frac{1}{2} \sigma^2 u^2}$ ;
- (2) if  $(r, s) \in \mathbb{R}^2$ ,  $r \neq 0$ , then  $rX + s \sim \mathcal{N}(r\mu + s, r^2 \sigma^2)$ .
- (3) If  $Y \sim \mathcal{N}(\bar{\mu}, \bar{\sigma}^2)$  and  $Y \perp X$  then  $X + Y \sim \mathcal{N}(\mu + \bar{\mu}, \sigma^2 + \bar{\sigma}^2)$ . As a consequence, if  $\{X_j\}_{j=1}^n$  is independent with  $X_j \sim \mathcal{N}(\mu_j, \sigma_j^2)$  then

$$(4.6) \quad \sum_j a_j X_j \sim \mathcal{N} \left( \sum_j a_j \mu_j, \sum_j a_j^2 \sigma_j^2 \right), \quad a_j \in \mathbb{R}.$$

*In particular, if  $\mu_j = 0$  and  $\sigma_j = \sigma$  then*

$$(4.7) \quad \langle X, \vec{a} \rangle \sim \mathcal{N}(0, \|\vec{a}\|^2 \sigma^2),$$

*where  $\vec{a} = (a_1, \dots, a_n)$ .*

*Proof.* (1) is a special case of (4.3). As for (2), note that  $\phi_s(u) = e^{sui}$ ,  $s \perp rX$  and use Proposition 2.30 (1) and Proposition 2.29 (1) to check that

$$\phi_{rX+s}(u) = e^{(r\mu+s)ui - \frac{1}{2} r^2 \sigma^2 u^2},$$

and finally use Corollary 4.6. Clearly, (3) is proved with the same kind of argument.  $\square$

**Example 4.8.** (Moments of a normal) If  $X \sim \mathcal{N}(0, \sigma^2)$  then (4.5) gives

$$\varphi_X(u) = e^{\frac{1}{2}\sigma^2 u^2} = \sum_{k \geq 0} \frac{\sigma^{2k}}{2^k k!} u^{2k},$$

so if we compare with (2.21) we conclude that

$$(4.8) \quad \alpha_l(X) = \begin{cases} \frac{l!}{2^{l/2}(l/2)!} \sigma^{2k}, & l \text{ even} \\ 0, & l \text{ odd} \end{cases}$$

which provides explicit expressions for all the moments of  $X$ .  $\square$

**Example 4.9.** (The log-normal distribution) If  $n = 1$  and  $X \sim \mathcal{N}(\mu, \sigma^2)$  then (4.5) implies that  $Y = e^X$  satisfies

$$(4.9) \quad \mathbb{E}(Y) = \mathbb{E}(e^X) = e^{\mu + \frac{1}{2}\sigma^2}$$

and

$$\mathbb{E}(Y^2) = \mathbb{E}(e^{2X}) = e^{2\mu + 2\sigma^2},$$

so that

$$(4.10) \quad \text{var}(Y) = \mathbb{E}(Y^2) - \mathbb{E}(Y)^2 = (e^{\sigma^2} - 1)e^{2\mu + \sigma^2}.$$

Hence, we may summarize (4.9) and (4.10) by writing

$$(4.11) \quad Y = e^X \sim \mathcal{LN}(e^{\mu + \frac{1}{2}\sigma^2}, (e^{\sigma^2} - 1)e^{2\mu + \sigma^2}),$$

where  $\mathcal{LN}$  stands for “log-normal” (which means that  $X = \ln Y$  follows a normal). Alternatively, we may write

$$(4.12) \quad Y \sim \Lambda(\mu, \sigma^2),$$

if emphasis on the parameters of the underlying normal distribution is needed [AB69]. In this notation, it is immediate from Proposition 4.7-(2) that (4.12) implies

$$(4.13) \quad e^a Y \sim \Lambda(\mu + a, \sigma^2), \quad a \in \mathbb{R}.$$

Now, an (obvious) generalization of (4.11) is

$$Y^u \sim \mathcal{LN}(e^{\mu u + \frac{1}{2}\sigma^2 u^2}, (e^{\sigma^2 u^2} - 1)e^{2\mu u + \sigma^2 u^2}), \quad u \in \mathbb{R},$$

so that the corresponding *coefficient of variation*,

$$(4.14) \quad \text{cv}(Y^u) := \frac{\text{sd}(Y^u)}{\mathbb{E}(Y^u)},$$

is given by

$$\text{cv}(Y^u) = \sqrt{e^{\sigma^2 u^2} - 1}.$$

In particular, it does not depend on  $\mu = \mathbb{E}(\ln Y)$  and satisfies the “scaling-plus-inversion invariance property”

$$(4.15) \quad \text{cv}(\alpha Y^u) = \text{cv}(Y^u) = \text{cv}(Y^{-u}), \quad \alpha > 0.$$

For later reference, we also note that the pdf  $\psi_Y$  of  $Y = e^X$  is

$$(4.16) \quad \psi_Y(x) = \frac{1}{\sigma x \sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\ln x - \mu}{\sigma}\right)^2},$$

so that the corresponding cdf is

$$(4.17) \quad F_Y(x) = \Phi\left(\frac{\ln x - \mu}{\sigma}\right),$$

where

$$(4.18) \quad \Phi(x) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}y^2} dy$$

is the cdf of  $\mathcal{N}(0, 1)$ . □

We proceed with another nice application of the formalism of characteristic functions to normally distributed random variables. We have seen that if  $X$  and  $Y$  are independent random variables then they are uncorrelated (Corollary 2.5). We now check that the converse statement holds for the entries of a normally distributed random vector.

**Proposition 4.10.** *If  $X = (X_1, \dots, X_k) : \Omega \rightarrow \mathbb{R}^k$  is a normally distributed random vector, say  $X \sim \mathcal{N}(\mu, \Sigma)$ , with  $\Sigma$  diagonal then:*

- $\{X_j\}_{j=1}^k$  is independent;
- $X_j \sim \mathcal{N}(\mu_j, \sigma_{X_j}^2)$ ,  $\sigma_{X_j}^2 = \text{var}(X_j)$ .

*Proof.* By assumption,  $\Sigma$  is diagonal, that is,  $\Sigma = \text{diag}(\sigma_{X_1}^2, \dots, \sigma_{X_k}^2)$ . If  $\mu = \mathbb{E}(X)$  we have from (4.3) that

$$\begin{aligned} \phi_X(\mathbf{u}) &= e^{\langle \mu, \mathbf{u} \rangle - \frac{1}{2} \langle \Sigma \mathbf{u}, \mathbf{u} \rangle} \\ &= e^{(\sum_j \mu_j u_j) - \frac{1}{2} \sum_j \Sigma_{jj} u_j^2} \\ &= \prod_j e^{\mu_j u_j - \frac{1}{2} \Sigma_{jj} u_j^2} \\ &= \prod_j \phi_{Y_j}(u_j), \end{aligned}$$

where  $Y_j \sim \mathcal{N}(\mu_j, \sigma_{X_j}^2)$  by (4.5) and Proposition 2.26, (2). Using (the multi-variate version of) (2.17) we then have

$$\begin{aligned} \psi_X(\mathbf{x}) &= \frac{1}{(2\pi)^n} \int_{\mathbb{R}^n} e^{-i\langle \mathbf{x}, \mathbf{u} \rangle} \phi_X(\mathbf{u}) d\mathbf{u} \\ &= \prod_j \frac{1}{2\pi} \int_{\mathbb{R}} e^{-ix_j u_j} \phi_{Y_j}(u_j) du_j \\ &= \prod_j \psi_{Y_j}(x_j), \end{aligned}$$

which not only proves that  $\{X_j\}_{j=1}^k$  is independent (by Proposition 2.13) but also that  $\psi_{X_j} = \psi_{Y_j}$  (by Proposition 2.12), which concludes the proof. □

**Corollary 4.11.** *The following assertions are equivalent:*

- $X = (X_1, \dots, X_k) \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_k)$ .
- $\{X_j\}_{j=1}^k$  is independent and  $X_j \sim \mathcal{N}(0, \sigma^2)$ .

**Corollary 4.12.** *(Rotational invariance) Let  $\{X_j\}_{j=1}^k$  be independent with  $X_j \sim \mathcal{N}(0, \sigma^2)$  and consider  $Y_l = \sum_{j=1}^k C_{lj} X_j$ , where  $C = \{C_{lj}\}$  is orthogonal. Then  $\{Y_l\}_{l=1}^k$  is independent with  $Y_l \sim \mathcal{N}(0, \sigma^2)$ .*

*Proof.* Write  $Y = CX$  with  $X \sim \mathcal{N}(\vec{0}, \sigma^2 I)$ . It follows that

$$\begin{aligned}\phi_Y(\mathbf{u}) &= \phi_{CX}(\mathbf{u}) \\ &\stackrel{(*)}{=} \phi_X(C\mathbf{u}) \\ &= e^{-\frac{1}{2}\langle \sigma^2 I C\mathbf{u}, C\mathbf{u} \rangle} \\ &= e^{-\frac{1}{2}\sigma^2 \|\mathbf{u}\|^2},\end{aligned}$$

where we used Proposition 2.29 (1) in (\*). Hence,  $Y \sim \mathcal{N}(\vec{0}, \sigma^2 I)$  as well (by Corollary 4.6) and the independence of  $\{Y_j\}$  now follows from the proposition.  $\square$

**Definition 4.13.** If  $X = (X_1, \dots, X_k)$  satisfies any of the conditions in Corollary 4.11 with  $\sigma = 1$  (so that  $X \sim \mathcal{N}(\vec{0}, \text{Id}_k)$ ) then we say that  $X$  is a *standard* normal random vector.

**Remark 4.14.** The projection property in (4.7) is an easy consequence of rotational invariance: set  $Y = OX$ , where  $O$  is an orthogonal matrix whose first line is  $\|\vec{a}\|^{-1}\vec{a}$  and note that  $\|\vec{a}\|^{-1}\langle X, \vec{a} \rangle = (OX)_1 \sim \mathcal{N}(0, \sigma^2)$ .  $\square$

**Remark 4.15.** The same computation as above shows that if  $X \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  then  $Y := CX + \mathbf{d}$ , where  $C$  is an invertible matrix and  $\mathbf{d}$  is a vector, then

$$Y \sim \mathcal{N}(C\boldsymbol{\mu} + \mathbf{d}, C\boldsymbol{\Sigma}C^t),$$

so that normality is preserved under an affine transformation.  $\square$

**Remark 4.16.** In Proposition 4.10 it is essential to assume that the normal random variables  $X_j$ ,  $j = 1, \dots, n$ , are *jointly* normally distributed in the sense that  $X = (X_1, \dots, X_n)$  is normally distributed. In fact, there exist random vectors  $X = (X_1, X_2)$  with  $\text{cov}(X_1, X_2) = 0$ ,  $X_1, X_2 \sim \mathcal{N}(0, 1)$  but with  $\{X_1, X_2\}$  *not* being independent. The classical example is obtained by taking  $X_1 \sim \mathcal{N}(0, 1)$ ,  $\epsilon$  a Rademacher random variable as in Definition 2.33 which is independent from  $X_1$  and  $X_2 = \epsilon X_1$ . To check the claims above, we first compute

$$\begin{aligned}\text{cov}(X_1, X_2) &= \mathbb{E}(X_1 X_2) - \mathbb{E}(X_1)\mathbb{E}(X_2) \\ &= \mathbb{E}(X_1 X_2) \\ &= \mathbb{E}(X_1^2 \epsilon) \\ &= \mathbb{E}(X_1^2)\mathbb{E}(\epsilon) \\ &= 0,\end{aligned}$$

where we used that  $\mathbb{E}(\epsilon) = 0$  in the last step. Also, the fact that  $X_2 \sim \mathcal{N}(0, 1)$  follows from Proposition 2.34. Finally, if  $\{X_1, X_2\}$  were independent then  $\{|X_1|, |X_2|\}$  would be independent as well, which is a contradiction because  $|X_1| = |X_2|$ .  $\square$

**Remark 4.17.** (The effectiveness of the characteristic function) The simplest case  $n = 2$  already illustrates the difficulty in trying to prove Proposition 4.10 by means of pdfs (thus directly relying on Proposition 2.13). Let us assume that

$$(4.19) \quad (X_1, X_2) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where, with self-explanatory notation,

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_{X_1} \\ \mu_{X_2} \end{pmatrix}$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{X_1}^2 & \sigma_{X_1 X_2} \\ \sigma_{X_1 X_2} & \sigma_{X_2}^2 \end{pmatrix} = \begin{pmatrix} \sigma_{X_1}^2 & \rho \sigma_{X_1} \sigma_{X_2} \\ \rho \sigma_{X_1} \sigma_{X_2} & \sigma_{X_2}^2 \end{pmatrix},$$

where

$$\rho = \frac{\sigma_{X_1 X_2}}{\sigma_{X_1} \sigma_{X_2}}$$

is the *correlation coefficient*. Hence, one must check that  $\rho = 0$  implies that  $\{X_1, X_2\}$  is independent, with each marginal following the appropriate normal distribution. To proceed, note that

$$\det \Sigma = (1 - \rho^2) \sigma_{X_1}^2 \sigma_{X_2}^2,$$

so that

$$\begin{aligned} \Sigma^{-1} &= \frac{1}{(1 - \rho^2) \sigma_{X_1}^2 \sigma_{X_2}^2} \begin{pmatrix} \sigma_{X_2}^2 & -\rho \sigma_{X_1} \sigma_{X_2} \\ -\rho \sigma_{X_1} \sigma_{X_2} & \sigma_{X_1}^2 \end{pmatrix} \\ &= \frac{1}{1 - \rho^2} \begin{pmatrix} 1/\sigma_{X_1}^2 & -\rho/\sigma_{X_1} \sigma_{X_2} \\ -\rho/\sigma_{X_1} \sigma_{X_2} & 1/\sigma_{X_2}^2 \end{pmatrix}, \end{aligned}$$

and leading this to (4.1), with  $A = \Sigma^{-1}$ , we see that the joint density of  $(X_1, X_2)$  is

$$\begin{aligned} \psi_{(X_1, X_2)}(x_1, x_2) &= \frac{1}{2\pi \sigma_{X_1} \sigma_{X_2} \sqrt{1 - \rho^2}} \times \\ (4.20) \quad &\times e^{-\frac{1}{2(1 - \rho^2)} \left( \frac{(x_1 - \mu_{X_1})^2}{\sigma_{X_1}^2} - \frac{2\rho(x_1 - \mu_{X_1})(x_2 - \mu_{X_2})}{\sigma_{X_1} \sigma_{X_2}} + \frac{(x_2 - \mu_{X_2})^2}{\sigma_{X_2}^2} \right)}. \end{aligned}$$

Thus, if  $\rho = 0$  this decomposes as

$$\psi_{(X_1, X_2)}(x_1, x_2) = \frac{1}{\sqrt{2\pi} \sigma_{X_1}} e^{-\frac{(x_1 - \mu_{X_1})^2}{2\sigma_{X_1}^2}} \times \frac{1}{\sqrt{2\pi} \sigma_{X_2}} e^{-\frac{(x_2 - \mu_{X_2})^2}{2\sigma_{X_2}^2}},$$

from which the claim follows immediately. However, it is not clear how this argument, which involves explicitly inverting the covariance matrix  $\Sigma$ , carries over as  $n$  gets indefinitely large. This should be compared with the general proof displayed above, which relies on the inversion formula (2.17) combined with the fact that  $\Sigma$  appears *linearly* in the exponent of (4.3). As yet another nice application of characteristic functions, let us note that, in general, if  $X = (X_1, X_2)$  is given then the characteristic function of the marginal  $X_1$  is

$$\phi_{X_1}(u_1) = \mathbb{E}(e^{iu_1 X_1}) = \mathbb{E}(e^{i(u_1 X_1 + 0 X_2)}),$$

that is,

$$\phi_{X_1}(u_1) = \phi_{(X_1, X_2)}(u_1, 0),$$

which tells us how to calculate the characteristic function of a marginal in terms of the characteristic function of the joint distribution. In particular, when applied to a bi-variate normal as in (4.19), and *not* necessarily assuming that  $\{X_1, X_2\}$  is independent, this clearly implies that the marginals are normally distributed in the expected way:  $X_j \sim \mathcal{N}(\mu_{X_j}, \sigma_{X_j}^2)$ . Needless to say, a similar result holds for the the marginals of a multivariate, normally distributed random vector, with essentially the same proof.  $\square$

**Remark 4.18.** (Regression to the mean) If  $X = (X_1, X_2)$  is a bi-variate normal as in (4.19) then we know from Remark 4.17 that  $X_j \sim \mathcal{N}(\mu_{X_j}, \sigma_{X_j}^2)$ ,  $j = 1, 2$ . Using this, (4.20), (3.3) and a little algebra we get

$$\begin{aligned} \psi_{X_2|X_1=x_1}(x_2) &= \frac{1}{\sqrt{2\pi} \sqrt{1 - \rho^2} \sigma_{X_2}} \times \\ &\times e^{-\frac{1}{2(1 - \rho^2) \sigma_{X_2}^2} \left( x_2 - \mu_{X_2} - \rho \frac{\sigma_{X_2}}{\sigma_{X_1}} (x_1 - \mu_{X_1}) \right)^2}, \end{aligned}$$

so that

$$X_2|_{X_1=x_1} \sim \mathcal{N}\left(\mu_{X_2} + \rho \frac{\sigma_{X_2}}{\sigma_{X_1}} (x_1 - \mu_{X_1}), (1 - \rho^2) \sigma_{X_2}^2\right)$$

or equivalently,

$$\frac{X_2|_{X_1=x_1} - \mu_{X_2}}{\sigma_{X_2}} \sim \mathcal{N}\left(\rho \frac{x_1 - \mu_{X_1}}{\sigma_{X_1}}, 1 - \rho^2\right).$$

In particular,

$$(4.21) \quad \frac{\mathbb{E}(X_2|_{X_1=x_1}) - \mu_{X_2}}{\sigma_{X_2}} = \rho \frac{x_1 - \mu_{X_1}}{\sigma_{X_1}},$$

which says that, on average, the proper standardization of  $X_2|_{X_1=x_1}$  is proportional to the observed standardization of  $X_1$  by a factor which is strictly less than 1 (in absolute value) unless  $X_1$  and  $X_2$  are perfectly correlated ( $|\rho| = 1$ ). More specifically, let us suppose that the variables model random measurements of some hereditary trait (stature, for instance) that passes from parents ( $X_1$ ) to offspring ( $X_2$ ) and happens to be “stable” in the sense that both variables follow the same normal  $\mathcal{N}(\mu, \sigma^2)$  (and of course are jointly normally distributed as well). We then obtain from (4.21) that

$$\mathbb{E}(X_2|_{X_1=x_1}) - x_1 = -(1 - \rho)(x_1 - \mu),$$

which means that, on average,  $X_2|_{X_1=x_1}$  lies somewhere between  $x_1$  and  $\mu$ . This “regression to the mean”, first (empirically) discovered by F. Galton, has played a fundamental role in the conceptual development of Multivariate Analysis [Sti90, Sti97, Gor16]. In order to relate this to the simple linear regression model as discussed in Remark 9.17, note from (4.20) that the ellipses of “equal frequency” for the joint distribution are given by

$$(x_1 - \mu)^2 - 2\rho(x_1 - \mu)(x_2 - \mu) + (x_2 - \mu)^2 = \text{const.},$$

so the contact points of the corresponding vertical tangent lines satisfy

$$x_2 - \mu = \rho(x_1 - \mu),$$

which, upon comparison with (9.47) and (9.48), identifies  $\rho$  to the slope of the associated regression line<sup>8</sup>.  $\square$

**4.2. Random variables related to the normal.** We now present a few distributions closely related to the normal.

**Definition 4.19.** A random variable  $Y : \Omega \rightarrow \mathbb{R}$  is Gamma( $\alpha, \lambda$ )-distributed, where  $\alpha, \lambda > 0$ , if its pdf is

$$(4.22) \quad \Gamma_{\alpha, \lambda}(x) = \frac{\alpha^\lambda}{\Gamma(\lambda)} x^{\lambda-1} e^{-\alpha x} \mathbf{1}_{(0, +\infty)}(x),$$

where

$$\Gamma(\lambda) = \int_0^{+\infty} y^{\lambda-1} e^{-y} dy,$$

is the Gamma function. We then say that  $\alpha$  and  $\lambda$  are the *inverse scale* and *shape* parameters of  $X$ , respectively. In particular,  $Y$  is *chi-squared distributed* with  $k \geq 1$  degrees of freedom if its pdf is  $\chi_k^2 := \Gamma_{1/2, k/2}$ . Explicitly,

$$(4.23) \quad \chi_k^2(x) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2} \mathbf{1}_{(0, +\infty)}(x).$$

**Proposition 4.20.** If  $Y \sim \text{Gamma}(\alpha, \lambda)$  then its mgf is  $\varphi_Y(u) = (1 - \alpha^{-1}u)^{-\lambda}$ ,  $|u| < \alpha$ . In particular,  $\mathbb{E}(Y) = \lambda/\alpha$  and  $\text{var}(Y) = \lambda/\alpha^2$ .

<sup>8</sup>For assessments of the social and intellectual contexts of his time and the nasty ideology behind Galton’s pursuit of this statistical result, we refer to [Cow72, Hil73, Mac81, Bul03].

*Proof.* We have

$$\begin{aligned}\varphi_Y(u) &= \frac{\alpha^\gamma}{\Gamma(\lambda)} \int_0^{+\infty} x^{\lambda-1} e^{-(\alpha-u)x} dx \\ &\stackrel{y=(\alpha-u)x}{=} \frac{\alpha^\gamma}{\Gamma(\lambda)} (\alpha-u)^{-\lambda} \int_0^{+\infty} y^{\lambda-1} e^{-y} dy \\ &= \alpha^\gamma (\alpha-u)^{-\lambda}.\end{aligned}$$

The last assertion follows from Remark 2.36.  $\square$

**Corollary 4.21.** *If  $Y \sim \chi_k^2$  then  $\varphi_Y(u) = (1 - 2u)^{-k/2}$ ,  $|u| < 1/2$ . In particular,  $\mathbb{E}(Y) = k$  and  $\text{var}(Y) = 2k$ .*

**Corollary 4.22.** *If  $Y \sim \text{Gamma}(\alpha, \lambda)$  then its characteristic function is given by  $\phi_Y(u) = (1 - \alpha^{-1}u\mathbf{i})^{-\lambda}$ . In particular, if  $Y \sim \chi_k^2$  then  $\phi_Y(u) = (1 - 2u\mathbf{i})^{-k/2}$ .*

**Corollary 4.23.** *If  $a > 0$  and  $Y \sim \text{Gamma}(\alpha, \lambda)$  then  $aY \sim \text{Gamma}(\alpha/a, \lambda)$ . In particular, if  $Y \sim \chi_k^2$  then  $aY \sim \Gamma_{1/2a, k/2}$ .*

*Proof.* Recall from Proposition 2.29 (1) that

$$\phi_{aY}(u) = \phi_Y(au) = (1 - \alpha^{-1}au\mathbf{i})^{-k/2}.$$

$\square$

Note that this justifies the adopted terminology for  $\alpha$ .

**Corollary 4.24.** *If  $\{Y_j\}_{j=1}^k$  is independent and  $Y_j \sim \text{Gamma}(\alpha, \lambda_j)$  then*

$$\sum_j Y_j \sim \Gamma_{\alpha, \sum_j \lambda_j}.$$

*In particular, if  $Y_j \sim \chi_{k_j}^2$  then*

$$\sum_j Y_j \sim \chi_{\sum_j k_j}^2.$$

**Corollary 4.25.** *If  $\{Z_j\}_{j=1}^k$  is independent with  $Z_j \sim \mathcal{N}(0, 1)$  then*

$$\sum_j Z_j^2 \sim \chi_k^2.$$

*Proof.* By Remark 2.15 and Proposition 4.7 (1),

$$\psi_{Z_j^2}(x) = \frac{1}{\sqrt{2\pi}} x^{-1/2} e^{-x/2} \mathbf{1}_{[0, +\infty)}(x),$$

so that  $Z_j^2 \sim \chi_1^2$  (recall that  $\Gamma(1/2) = \sqrt{\pi}$ ). The result now follows from Corollary 4.24.  $\square$

**Remark 4.26.** (The geometric way to  $\chi_k^2$ ) Corollary 4.25 can be elegantly retrieved as an application of the “ $n$ -space computations” introduced by R. Fisher [Fis15, Fis21, Fis25]. Indeed, from Proposition 2.13 we know that the amount of probability density spanned by a standard normal vector  $Z = (Z_1, \dots, Z_k) \in \mathbb{R}^k$  in an infinitesimal region of volume  $dz = dz_1 \cdots dz_k$  is

$$\begin{aligned} \frac{1}{(2\pi)^{k/2}} e^{-\|z\|^2/2} dz &= \frac{1}{(2\pi)^{k/2}} e^{-\|z\|^2/2} \|z\|^{k-1} d\|z\| d\theta \\ &= \frac{1}{2} \frac{1}{(2\pi)^{k/2}} e^{-\|z\|^2/2} (\|z\|^2)^{\frac{k}{2}-1} d\|z\|^2 d\theta, \end{aligned}$$

where  $z = (\|z\|, \theta) \in (0, +\infty) \times \mathbb{S}^{k-1}$  is the polar decomposition of  $Z$ <sup>9</sup>. Again by Proposition 2.13, if we view this latter expression as the joint distribution of  $(\|Z\|^2, \Theta)$  then  $\{\|Z\|^2, \Theta\}$  is independent with  $\Theta = X/\|X\|$  being *uniformly* distributed over  $\mathbb{S}^{n-1}$ . Hence, by Proposition 2.12 the infinitesimal density of  $\|Z\|^2$  is

$$\psi_{\|Z\|^2}(\|z\|^2) d\|z\|^2 = \frac{\omega_{k-1}}{2} \frac{1}{(2\pi)^{k/2}} e^{-\|z\|^2/2} (\|z\|^2)^{\frac{k}{2}-1} d\|z\|^2,$$

where  $\omega_{k-1}$  is the volume of  $\mathbb{S}^{k-1}$ . Since

$$(4.24) \quad \omega_{k-1} = \frac{2\pi^{k/2}}{\Gamma(k/2)}$$

it suffices to set  $x = \|z\|^2$  in order to recover (4.23). Note that the same computation gives that  $Y = (Y_1, \dots, Y_k) \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_k)$  implies  $\|Y\|^2 \sim \Gamma_{1/2\sigma^2, k/2}$ ; cf. Corollary 4.24. For  $k = 3$  and  $\sigma^2 = \kappa T/2$ , where  $\kappa$  is the Boltzmann constant and  $T$  is the temperature, this gives

$$\psi_E(\epsilon) d\epsilon = \frac{2\sqrt{\epsilon}}{\sqrt{\pi}(\kappa T)^{3/2}} e^{-\frac{\epsilon}{\kappa T}} d\epsilon,$$

the *energy distribution* of a Maxwellian gas [Kit04, Section 13]. In particular, by Proposition 4.20,  $\mathbb{E}(E) = 3\kappa T/2$ , which confirms the *principle of equipartition of energy*.  $\square$

By Corollary 4.11 we may rephrase Corollary 4.25 as saying that  $Z \sim \mathcal{N}(\vec{0}, \text{Id}_k)$  implies  $\|Z\|^2 \sim \chi_k^2$ . It turns that this is just a special case of a more general result which makes it clear the geometric meaning of the notion of degree of freedom for a chi-square distribution.

**Proposition 4.27.** *If  $Z \sim \mathcal{N}(\vec{0}, \text{Id}_k)$  and  $W = \langle Y, QY \rangle$ , where  $Q$  is a  $n \times n$  symmetric and idempotent matrix with  $\text{rank } Q = r \leq k$ , then  $W \sim \chi_r^2$ .*

*Proof.* Since  $Q : \mathbb{R}^n \rightarrow \mathbb{R}^n$  defines a projection onto its range  $\text{Im } Q$ , a linear subspace of dimension  $r$ , we may use the projection property in (4.7), with  $\sigma = 1$  and  $\vec{a}$  running over an orthonormal basis of  $\text{Im } Q$ , to conclude that  $QZ \sim \mathcal{N}(\vec{0}, \text{Id}_r)$ . Thus,  $W = \langle Z, QZ \rangle = |QZ|^2 \sim \chi_r^2$  by Corollary 4.25.  $\square$

This kind of geometric argument has many useful applications, including the next one, whose proof we omit.

**Proposition 4.28.** *Let  $Z \sim \mathcal{N}(\vec{0}, \text{Id}_n)$ ,  $c \in \mathbb{R}^n$  and  $A$  a symmetric  $n \times n$  matrix. Then  $\langle c, Z \rangle$  and  $\langle Z, AZ \rangle$  are independent if and only if  $Ac = \vec{0}$ .*

<sup>9</sup>In this and similar computations, as in Remark 7.31 and Example 7.39, we represent a realization of a random variable, say  $Z_j$  or  $\Theta$ , by the corresponding lower-case symbol (in this case,  $z_j$  or  $\theta$ ).



We now discuss some more random variables related to the normal distribution.

**Definition 4.29.** A random variable  $X$  is *t-Student distributed* with  $k \geq 1$  degrees of freedom if

$$(4.25) \quad \psi_X(x) = t_k(x) := \frac{\Gamma(\frac{k+1}{2})}{\sqrt{k\pi}\Gamma(k/2)} (1 + k^{-1}x^2)^{-(k+1)/2}.$$

**Proposition 4.30.** If  $Z \sim \mathcal{N}(0, 1)$  and  $W \sim \chi_k^2$  with  $Z \perp W$  then  $Z/\sqrt{W/k} \sim t_k$ .

*Proof.* Note that  $Z/\sqrt{W/k} = \sqrt{k}Z/V$ , where  $\sqrt{k}Z \sim \mathcal{N}(0, k)$  and  $V := \sqrt{W}$  so that  $\psi_V(v) = 2v\chi_k^2(v^2)$  by (2.8). It follows from (2.9) that

$$\psi_{Z/\sqrt{W/k}}(x) = \frac{1}{\sqrt{k\pi}2^{(k-1)/2}\Gamma(k/2)} \int_0^{+\infty} e^{-\frac{1}{2}(1+k^{-1}x^2)v^2} v^k dv.$$

The substitution  $w = \frac{1}{2}(1 + k^{-1}x^2)v^2$  then finishes the job.  $\square$

**Remark 4.31.** Although its derivation along these lines is a bit more involved, Proposition 4.30 can also be accessed by means of Fisher's geometric method illustrated in Remark 4.26. This kind of argument appeared in [Fis25] and it is reproduced here in Remark 7.31, where the method is employed to obtain the pdf of Student's sampling distribution defined in (7.26) below.  $\square$

**Definition 4.32.** Given  $k_1, k_2 \in \mathbb{N}$  we say that a random variable  $X$  is  $F_{k_1, k_2}$ -distributed if

$$\psi_X(x) = F_{k_1, k_2}(x) := c_{k_1, k_2} x^{k_1/2-1} \left(1 + \frac{k_1}{k_2}x\right)^{-\frac{k_1+k_2}{2}} \mathbf{1}_{(0, +\infty)}(x),$$

where

$$c_{k_1, k_2} = \frac{\Gamma(\frac{k_1+k_2}{2})}{\Gamma(k_1/2)\Gamma(k_2/2)} \left(\frac{k_1}{k_2}\right)^{k_1/2}.$$

**Proposition 4.33.** If  $W_1 \sim \chi_{k_1}^2$  and  $W_2 \sim \chi_{k_2}^2$  with  $W_1 \perp W_2$  then

$$\frac{W_1/k_1}{W_2/k_2} \sim F_{k_1, k_2}.$$

*Proof.* From (2.9) and (4.23) we find that

$$\psi_{W_1/W_2}(x) = \frac{x^{\frac{k_1}{2}-1}}{2^{\frac{k_1+k_2}{2}}\Gamma(k_1/2)\Gamma(k_2/2)} \int_0^{+\infty} v^{\frac{k_1+k_2}{2}-1} e^{-(1+x)v/2} dv,$$

so that the substitution  $w = (1+x)v/2$  transforms this into

$$(4.26) \quad \psi_{W_1/W_2}(x) = \frac{\Gamma(\frac{k_1+k_2}{2})}{\Gamma(k_1/2)\Gamma(k_2/2)} x^{\frac{k_1}{2}-1} (1+x)^{-\frac{k_1+k_2}{2}}.$$

The result now follows because

$$\psi_Y(x) = \frac{k_1}{k_2} \psi_{W_1/W_2} \left( \frac{k_1}{k_2} x \right).$$

□

**Corollary 4.34.** *If  $Y \sim F_{k_1, k_2}$  then  $Y^{-1} \sim F_{k_2, k_1}$ . Also, if  $T \sim t_k$  then  $T^2 \sim F_{1, k}$ .*

**Definition 4.35.** A random variable  $X$  is *Beta-distributed* with shape parameters  $\alpha, \beta > 0$  if

$$(4.27) \quad \psi_X(x) = \text{Beta}(\alpha, \beta)(x) := \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \mathbf{1}_{(0,1)}(x), \quad x \in \mathbb{R}.$$

**Proposition 4.36.** *If  $W_1 \sim \chi_{k_1}^2$  and  $W_2 \sim \chi_{k_2}^2$  with  $W_1 \perp W_2$  then*

$$\widehat{W} := \frac{W_1}{W_1 + W_2} \sim \text{Beta} \left( \frac{k_1}{2}, \frac{k_2}{2} \right).$$

*Proof.* We have

$$\widehat{W} = \frac{W_1/W_2}{1 + W_1/W_2},$$

so that, for  $x \in (0, 1)$ ,

$$\begin{aligned} F_{\widehat{W}}(x) &= P(\widehat{W} \leq x) \\ &= P\left(\frac{W_1}{W_2} \leq \frac{x}{1-x}\right) \\ &= F_{W_1/W_2}\left(\frac{x}{1-x}\right). \end{aligned}$$

By taking derivative with respect to  $x$ ,

$$\psi_{\widehat{W}}(x) = (1-x)^{-2} \psi_{W_1/W_2} \left( \frac{x}{1-x} \right),$$

and the result follows from (4.26). □

## 5. CONCENTRATION INEQUALITIES

Here we elaborate a bit on the suggestive idea that an exponential control on the mgf of a random variable yields corresponding bounds for its tail probabilities which are usually referred to as *concentration inequalities*. As we shall illustrate below with a few examples, this *Cramér-Chernoff method* of obtaining such tail inequalities has many applications, both in pure and applied mathematics; see [Ver18, Wai19] for accounts on these “concentration inequalities” and their applications to many problems in Data Science, including the estimation theory of the “high dimensional” regression schemes mentioned in Subsection 9.3. Appetizers to this circle of ideas, as applied to the celebrated Johnson-Lindenstrauss Lemma and the Erdős-Rényi model for random graphs, appear below.

**5.1. Sub-exponential random variables and the Johnson-Lindenstrauss Lemma.** If  $X \sim \mathcal{N}(\mu, \sigma^2)$  and  $t > 0$  then the tail probability

$$P(|X - \mu| \geq t) = \frac{2}{\sqrt{2\pi}\sigma} \int_t^{+\infty} e^{-x^2/2\sigma^2} dx$$

may be easily estimated by observing that  $x/t \geq 1$  implies

$$P(|X - \mu| \geq t) \leq \frac{2}{\sqrt{2\pi}\sigma} \int_t^{+\infty} \frac{x}{t} e^{-x^2/2\sigma^2} dx,$$

thus yielding the exponential tail bound

$$(5.1) \quad P(|X - \mu| \geq t) \leq \sqrt{\frac{2}{\pi}} \frac{\sigma}{t} e^{-t^2/2\sigma^2},$$

which happens to blow up as  $t \rightarrow 0$ . We may remedy this by considering the function

$$\begin{aligned} \xi(t) &= P(X - \mu \geq t) - \frac{1}{2} e^{-t^2/2\sigma^2} \\ &= 1 - F_{X-\mu}(t) - \frac{1}{2} e^{-t^2/2\sigma^2}. \end{aligned}$$

Since

$$\xi'(t) = \left( \frac{t}{2\sigma^2} - \frac{1}{\sqrt{2\pi}\sigma} \right) e^{-t^2/2\sigma^2},$$

we see that  $\xi$  decreases in the interval  $(0, \sqrt{2/\pi}\sigma)$  and increases in the interval  $(\sqrt{2/\pi}\sigma, +\infty)$ . Since  $\xi(0) = 0$  and  $\xi(t) \rightarrow 0$  as  $t \rightarrow +\infty$  we conclude that

$$(5.2) \quad P(|X - \mu| \geq t) \leq e^{-t^2/2\sigma^2},$$

which provides a sharp bound as  $t \rightarrow 0$ . As yet another way to obtain an exponential tail probability bound, observe that

$$\begin{aligned} P(|X - \mu| \geq t) &= 2P(X - \mu \geq t) \\ &\stackrel{u \geq 0}{=} 2P(e^{u(X-\mu)} \geq e^{ut}) \\ (2.13) \quad &\leq 2 \frac{\mathbb{E}(e^{u(X-\mu)})}{e^{ut}} \quad (\text{Markov's inequality}) \\ (4.5) \quad &\stackrel{=}{=} 2e^{\frac{1}{2}\sigma^2 u^2 - ut}, \end{aligned}$$

and since the right-hand side reaches its minimal value at  $u = t/\sigma^2$  we conclude that

$$(5.3) \quad P(|X - \mu| \geq t) \leq 2e^{-t^2/2\sigma^2}.$$

Although in general (5.1) and (5.2) give sharper bounds than (5.3), this latter argument seems to be more promising as it suggests that a suitable exponential control on the mgf of a random variable might yield corresponding bounds for its tail probabilities, a trick usually referred to as the *Cramér-Chernoff method*. Its flexibility is illustrated by considering the following class of random variables.

**Definition 5.1.** We say that  $X$  is *Sub-Gaussian* if

$$(5.4) \quad \mathbb{E} \left( e^{(X-\mu)u} \right) \leq e^{\frac{\sigma^2 u^2}{2}}, \quad u \in \mathbb{R},$$

which we represent as  $X \in \text{SubG}(\sigma)$ .

Clearly, any  $X \in \text{SubG}(\sigma)$  satisfies (5.3). As a simple example distinct from a normal to which this applies, note that the Rademacher variable  $\epsilon$  in Definition 2.33 satisfies

$$(5.5) \quad \mathbb{E}_\epsilon(e^{\epsilon u}) = \cosh u \leq e^{\frac{u^2}{2}}, \quad u \in \mathbb{R},$$

so that  $\epsilon \in \text{SubG}(1)$ . Also, if  $\{X_j\}_{j=1}^N$  is independent with  $X_j \in \text{SubG}(\sigma_j)$  then

$$(5.6) \quad \sum_j X_j \in \text{SubG} \left( \sqrt{\sum_j \sigma_j^2} \right).$$

Thus, if we apply this to  $X^{(N)} = \epsilon_1 + \dots + \epsilon_N$ , where  $\{\epsilon_j\}_{j=1}^N$  is a collection of independent Rademacher variables, we see that  $X^{(N)} \in \text{SubG}(\sqrt{N})$  and hence

$$(5.7) \quad P(|X^{(N)}| \geq t) \leq 2e^{-t^2/2N},$$

or equivalently,

$$(5.8) \quad P\left(\frac{|X^{(N)}|}{N} \geq t\right) \leq 2e^{-Nt^2/2}.$$

The next result substantially enriches the class of sub-Gaussian random variables.

**Proposition 5.2.** *If  $X$  is bounded, say  $a \leq X \leq b$ , then  $X \in \text{SubG}(b - a)$ .*

*Proof.* We may assume that  $\mathbb{E}(X) = 0$ . Let  $Y$  be an independent copy of  $X$ , so that  $X - Y$  is symmetric by Proposition 2.30 (3). Hence, if  $\epsilon$  is a Rademacher variable independent from both  $X$  and  $Y$  then  $\epsilon(X - Y)$  and  $X - Y$  are identically distributed by Proposition 2.34. It follows that

$$\begin{aligned} \mathbb{E}_{(X,Y)}(e^{(X-Y)u}) &= \mathbb{E}_{(X,Y)} \left( \mathbb{E}_\epsilon \left( e^{(X-Y)u} \right) \right) \\ &= \mathbb{E}_{(X,Y)} \left( \mathbb{E}_\epsilon \left( e^{\epsilon(X-Y)u} \right) \right) \\ &\stackrel{(5.5)}{\leq} \mathbb{E}_{(X,Y)} \left( e^{\frac{(X-Y)^2 u^2}{2}} \right), \end{aligned}$$

and since  $|X - Y| \leq b - a$ , we see that

$$\mathbb{E}_{(X,Y)}(e^{(X-Y)u}) \leq e^{\frac{(b-a)^2 u^2}{2}}.$$

On the other hand,

$$\begin{aligned} \mathbb{E}_{(X,Y)}(e^{(X-Y)u}) &= \mathbb{E}_X (e^{Xu} \mathbb{E}_Y (e^{-Yu})) \\ &\geq \mathbb{E}_X (e^{Xu} e^{-\mathbb{E}_Y(Y)u}) \\ &= \mathbb{E}_X (e^{Xu}), \end{aligned}$$

where we used Jensen inequality and  $\mathbb{E}_Y(Y) = 0$ . Putting all the pieces of this computation together we conclude that

$$\mathbb{E}_X (e^{Xu}) \leq e^{\frac{(b-a)^2 u^2}{2}},$$

as desired. □

**Remark 5.3.** Under the conditions of Proposition 5.2, it is possible to show that  $X \in \text{SubG}((b - a)/2)$ , which is known as the *Hoeffding lemma*. □

The Hoeffding-type concentration inequalities stemming from Proposition 5.2 (or Remark 5.3) completely neglect the dispersion of a random variable as measured by its variance. For instance, if we take  $X \sim \text{Bin}(p; n)$  and use Example 2.37 to express it as a sum of independent Bernoulli variables so that (5.6) applies, we see that

$$(5.9) \quad P(X - np > t) \leq e^{-\frac{t^2}{2n}}, \quad t > 0,$$

an estimate not involving the sampling probability  $p \in (0, 1)$  (for more on this, see Subsection 5.3 below). Besides, the class of sub-Gaussian variables fails to accommodate certain concentration inequalities commonly occurring in applications. This justifies the investigation of other classes of random variables for which concentration inequalities are available.

**Definition 5.4.** A random variable  $Y$  with  $\mathbb{E}(Y) = \mu$  is *sub-exponential* if there exist positive parameters  $(\nu, \beta)$  such that

$$(5.10) \quad \mathbb{E}\left(e^{u(Y-\mu)}\right) \leq e^{\frac{\nu^2 u^2}{2}}, \quad |u| < \frac{1}{\beta}.$$

We represent this as  $Y \in \text{SubE}(\nu, \beta)$ .

**Remark 5.5.** If  $\mu = 0$  then  $Y \in \text{SubE}(\nu, \beta)$  implies  $-Y \in \text{SubE}(\nu, \beta)$ . Also, if  $\{Y_j\}_{j=1}^n$  is independent and  $Y_j \in \text{SubE}(\nu_j, \beta_j)$  with

$$\sum_j Y_j \in \text{SubE}\left(\sqrt{\sum_j \nu_j^2}, \max_j \{\beta_j\}\right),$$

where we assume that  $\mu_j = \mathbb{E}(Y_j) = 0$ . □

**Example 5.6.** From Corollary 4.21 we have that  $Y \sim \chi_k^2$  implies

$$\mathbb{E}\left(e^{u(Y-k)}\right) = e^{-ku} (1 - 2u)^{-k/2}, \quad |u| < \frac{1}{2}.$$

Note that  $\mathbb{E}(e^{\frac{1}{2}(Y-k)}) = +\infty$  and hence  $Y$  is not sub-Gaussian. However, by Taylor expanding around  $u = 0$  we find that

$$f(u) := (1 - 2u)^{-k/2} = 1 + ku + k\left(\frac{k}{2} + 1\right)u^2 + \dots$$

and

$$g(u) := e^{k(2u^2+u)} = 1 + ku + k\left(\frac{k}{2} + 2\right)u^2 + \dots,$$

from which we easily see that  $g(u) \geq f(u)$  for  $|u| < 1/4$  since both functions remain convex in this interval. In other words,

$$\mathbb{E}\left(e^{u(Y-k)}\right) \leq e^{2ku^2}, \quad |u| < \frac{1}{4},$$

and we conclude that  $Y \in \text{SubE}(2\sqrt{k}, 4)$ . □

It turns out that sub-exponential random variables satisfy a concentration inequality exhibiting a clear-cut threshold between the sub-Gaussian and the purely sub-exponential regimes.

**Proposition 5.7.** *If  $Y \in \text{SubE}(\nu, \beta)$  with  $\mathbb{E}(Y) = \mu$  then*

$$P(|Y - \mu| \geq t) \leq \begin{cases} 2e^{-\frac{t^2}{2\nu^2}}, & 0 \leq t < \frac{\nu^2}{\beta} \\ 2e^{-\frac{t}{2\beta}}, & t \geq \frac{\nu^2}{\beta} \end{cases}$$

*Proof.* Clearly, we may assume that  $\mu = 0$ , so Markov inequality (2.13) gives

$$P(Y \geq t) \leq e^{h_t(u)}, \quad 0 \leq u < \frac{1}{\beta},$$

with the graph of  $h_t(u) = -ut + \nu^2 u^2/2$  being an upward pointing parabola passing through  $(0, 0)$  and with vertex located at  $(t/\nu^2, -t^2/2\nu^2)$ . Thus, in the first case, when  $t/\nu^2 < 1/\beta$ , the tail probability is bounded by

$$e^{h_t(t/\nu^2)} = e^{-\frac{t^2}{2\nu^2}},$$

whereas in the second case, when  $t/\nu^2 \geq 1/\beta$ , it is bounded by

$$e^{h_t(1/\beta)} \leq e^{-\frac{t}{\beta} + \frac{\nu^2}{2\beta^2}} \leq e^{-\frac{t}{2\beta}}.$$

By Remark 5.5, an identical estimate holds for  $P(Y \leq -t)$ , which concludes the proof.  $\square$

**Corollary 5.8.** *If  $Y \sim \chi_k^2$  then*

$$(5.11) \quad P(|k^{-1}Y - 1| \geq t) \leq \begin{cases} 2e^{-\frac{kt^2}{8}}, & 0 < t < 1 \\ 2e^{-\frac{kt}{8}}, & t \geq 1 \end{cases}$$

*Proof.* From Example 5.6 and the proposition (with  $\mu = k$  and  $(\nu, \beta) = (2\sqrt{k}, 4)$ ) we know that

$$(5.12) \quad P(|Y - k| \geq \tau) \leq \begin{cases} 2e^{-\frac{\tau^2}{8k}}, & 0 < \tau < k \\ 2e^{-\frac{\tau}{8}}, & \tau \geq k \end{cases}$$

Since

$$P(|k^{-1}Y - 1| \geq k^{-1}\tau) = P(|Y - k| \geq \tau),$$

the substitution  $\tau = kt$  finishes the proof.  $\square$

We now use the concentration inequality (5.11) to establish a celebrated result with a number of applications both in pure and applied mathematics.

**Theorem 5.9.** (Johnson-Lindenstrauss). *If  $\mathcal{C} = \{x_1, \dots, x_n\} \subset \mathbb{R}^p$  is a collection of  $n$  points then, given  $\epsilon, \delta \in (0, 1)$ , there exist  $m = O(\epsilon^{-2} \ln(n/\delta))$  and a map  $F : \mathbb{R}^p \rightarrow \mathbb{R}^m$  such that*

$$(5.13) \quad 1 - \epsilon \leq \frac{\|F(x_i) - F(x_j)\|^2}{\|x_i - x_j\|^2} \leq 1 + \epsilon,$$

for all  $x_i \neq x_j$  in  $\mathcal{C}$ .

*Proof.* Use Remark 2.16 to construct a  $m \times p$  random matrix  $\mathbf{A}$  whose entries are independent and  $\mathcal{N}(0, 1)$ -distributed random variables and define the linear map

$$(5.14) \quad F : \mathbb{R}^p \rightarrow \mathbb{R}^m, \quad F(x) = \frac{\mathbf{A}x}{\sqrt{m}},$$

with  $m$  to be chosen later on. If  $\mathbf{a}_i$  is the  $i^{\text{th}}$  line of  $\mathbf{A}$  and  $x \neq \vec{0}$  then  $\langle \mathbf{a}_i, x/\|x\| \rangle \sim \mathcal{N}(0, 1)$  by Proposition 4.7 (3), so Corollary 4.25 applies to ensure that

$$\frac{\|\mathbf{A}x\|^2}{\|x\|^2} = \sum_{i=1}^m \left\langle \mathbf{a}_i, \frac{x}{\|x\|} \right\rangle^2 \sim \chi_m^2.$$

Thus, by Corollary 5.8,

$$P\left(\left|m^{-1} \frac{\|\mathbf{A}x\|^2}{\|x\|^2} - 1\right| \geq \epsilon\right) \leq 2e^{-m\epsilon^2/8}, \quad 0 < \epsilon < 1,$$

or equivalently,

$$P\left(\frac{\|F(x)\|^2}{\|x\|^2} \notin (1 - \epsilon, 1 + \epsilon)\right) \leq 2e^{-m\epsilon^2/8}, \quad 0 < \epsilon < 1.$$

From this we easily deduce that

$$P\left(\frac{\|F(x_i) - F(x_j)\|^2}{\|x_i - x_j\|^2} \notin (1 - \epsilon, 1 + \epsilon) \text{ for some } x_i \neq x_j\right) \leq 2 \binom{n}{2} e^{-m\epsilon^2/8},$$

and the result follows if we impose that the right-hand side equals  $\delta$  (which determines  $m$  as in the statement of the theorem), since this means that the probability that (5.13) holds true is  $\geq 1 - \delta > 0$ .  $\square$

**Remark 5.10.** The elegant argument above illustrates the celebrated *probabilistic method*, which roughly consists of upgrading the assertion we intend to prove (in our case, the purely deterministic statement in (5.13)) to a random event (in our case, through the choice of the random matrix  $\mathbf{A}$  in (5.14)) and then checking that this event occurs with *positive* probability, from which the validity of the statement follows immediately. We refer to [AS16] for many other illustrations of the method, notably in Combinatorics.  $\square$

**Remark 5.11.** In applications to Data Science [Ver18, BHK20], the number  $n$  of “samples” is much smaller than the number  $p$  of “features”. The remarkable aspect of Theorem 5.9 is that if we allow for a controlled distortion  $\epsilon > 0$  on the “approximate projection”  $F$ ,  $\mathcal{C}$  has a sort of intrinsic dimension  $m^{10}$  which scales as  $\ln n$  and happens to be completely insensitive to the ambient dimension  $p$ , which might even be infinite.  $\square$

**Remark 5.12.** If  $Y \in \text{SubE}(\nu, \beta)$  then Proposition 5.7 says that

$$P(|Y - \mathbb{E}(Y)| \geq t) \leq 2e^{-\min\left\{\frac{t^2}{2\nu^2}, \frac{t}{2\beta}\right\}},$$

thus confirming that as  $t \rightarrow +\infty$  this tail bound is much heavier than the one for  $X \in \text{SubG}(\sigma)$  in (5.3) since in this regime the upper bound here is  $2e^{-t/2\beta}$ . In a sense this reflects the fact that  $X \in \text{SubG}(\sigma)$  with  $\mathbb{E}(X) = 0$  implies  $Y = X^2 \in \text{SubE}(\nu, \beta)$  for some  $(\nu, \beta)$  to be determined below. To check this claim, take  $v \in (0, 1)$ , multiply both sides of the defining condition for  $X$  in (5.4) (with  $\mu = 0$ ) by  $e^{-\sigma^2 u^2/2v}$  and integrate to obtain

$$\begin{aligned} \int_{-\infty}^{+\infty} e^{\frac{\sigma^2 u^2(v-1)}{2v}} du &\geq \int_{-\infty}^{+\infty} \mathbb{E}\left(e^{uX - \frac{\sigma^2 u^2}{2v}}\right) du \\ &= \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} e^{ux - \frac{\sigma^2 u^2}{2v}} dP_X(x)\right) du \\ &= \int_{-\infty}^{+\infty} \left(\int_{-\infty}^{+\infty} e^{ux - \frac{\sigma^2 u^2}{2v}} du\right) dP_X(x). \end{aligned}$$

<sup>10</sup>This means that there exists a subset  $\mathcal{I} \subset \{1, \dots, p\}$  of indexes with  $\#\mathcal{I} = p - m$  such that  $y_j = 0$  for all  $y \in F(\mathcal{C})$  and  $j \in \mathcal{I}$ . In other words, up to an overall distortion measured by  $\epsilon$ ,  $m$  turns out to be the number of relevant (nonzero) features of the elements of  $\mathcal{C}$ , hence the relevance of the result for data reduction.

Using that

$$\int_{-\infty}^{+\infty} e^{-au^2+bu+c} du = e^{\frac{b^2+4ac}{4a}} \sqrt{\frac{\pi}{a}}, \quad a > 0, \quad b, c \in \mathbb{R},$$

we may compute the Gaussian integrals above to conclude that

$$\mathbb{E} \left( e^{\frac{X^2 v}{2\sigma^2}} \right) \leq \frac{1}{\sqrt{1-v}}, \quad 0 \leq v < 1,$$

or equivalently,

$$\mathbb{E} \left( e^{X^2 u} \right) \leq \frac{1}{\sqrt{1-2\sigma^2 u}}, \quad 0 \leq u < \frac{1}{2\sigma^2}.$$

In particular, there holds

$$\mathbb{E} \left( e^{(X^2 - \mathbb{E}(X^2))u} \right) < +\infty$$

for  $u \in [0, \varepsilon)$ ,  $\varepsilon > 0$ . By Remark 2.36,  $\mathbb{E}(|X|^4) < +\infty$  and we may Taylor expand up to second order:

$$\mathbb{E} \left( e^{(X^2 - \mathbb{E}(X^2))u} \right) = 1 + \frac{\mathbb{E}((X^2 - \mathbb{E}(X^2))^2)}{2} u^2 + o(u^2).$$

Comparing this with

$$e^{\frac{\nu^2 u^2}{2}} = 1 + \frac{\nu^2 u^2}{2} + o(u^2),$$

we easily see that  $Y = X^2 \in \text{SubE}(\nu, \beta)$  if we take

$$\nu^2 > \mathbb{E}((X^2 - \mathbb{E}(X^2))^2)$$

and  $\beta > 0$  large enough. Finally, note that if  $\sigma = 1$  then we can take  $\beta = 1$ . □

**Remark 5.13.** If  $\{Y_j\}_{j=1}^k$  is independent with  $Y_j \in \text{SubE}(\nu_j, \beta_j)$  and  $\mathbb{E}(Y_j) = 0$  then we easily see that  $Y = \sum_j Y_j \in \text{SubE}(\nu, \beta)$ , where

$$\nu^2 = \sum_j \nu_j^2, \quad \beta = \max_j \beta_j.$$

You may apply this to  $Y_j = X_j^2$ ,  $X_j \sim \mathcal{N}(0, 1)$ . Using the well-known fact that, for a standard normal,  $\mathbb{E}((X_j^2 - \mathbb{E}(X_j^2))^2) = 3$ , and the computation in Example 5.12 above, we conclude that  $\sum_j X_j^2 \in \text{SubE}(\sqrt{2}k, \beta)$ , for some  $\beta > 0$ , which provides essentially the same result as in Example 5.6. Clearly, this suffices for the applications culminating in Theorem 5.9 above. Notice that in this derivation no appeal to the explicit expression for the pdf of a chi-squared distribution has been needed. □

**5.2. The Gaussian concentration inequality.** A suitable rewording of the inequalities in Corollary 5.8 provides valuable insights on the “high dimensional” behavior of standard normal random vectors, including the precise formulation of a remarkable dimension-free concentration inequality for Lipschitz functions of such vectors; see (5.18) below. Indeed, let  $X = (X_1, \dots, X_k)$  be such a vector, which means by definition that  $\{X_j\}_{j=1}^k$  is independent with  $X_j \sim \mathcal{N}(0, 1)$ . Recalling that

$$|a - 1| \geq \delta \implies |a^2 - 1| \geq \max \{ \delta, \delta^2 \}, \quad a \geq 0, \quad \delta > 0,$$



we have

$$\begin{aligned}
P\left(\left|\frac{\|X\|}{\sqrt{k}} - 1\right| \geq \frac{\tau}{k}\right) &\leq P\left(\left|\frac{\|X\|^2}{k} - 1\right| \geq \max\left\{\frac{\tau}{k}, \frac{\tau^2}{k^2}\right\}\right) \\
&= P\left(\left|\|X\|^2 - k\right| \geq \max\left\{\tau, \frac{\tau^2}{k}\right\}\right) \\
&= P\left(\left|\|X\|^2 - k\right| \geq \begin{cases} \tau, & \tau < k \\ \frac{\tau^2}{k}, & \tau \geq k \end{cases}\right) \\
&\leq 2e^{-\frac{\tau^2}{8k}}, \quad \tau > 0,
\end{aligned}$$

where we used (5.12) with  $Y = \|X\|^2$  in the last step. By setting  $t = \tau/k$  we then obtain

$$(5.15) \quad P\left(\left|\frac{\|X\|}{\sqrt{k}} - 1\right| \geq t\right) \leq 2e^{-\frac{kt^2}{8}}, \quad t > 0,$$

or equivalently,

$$(5.16) \quad P\left(\left|\|X\| - \sqrt{k}\right| \geq t\right) \leq 2e^{-\frac{t^2}{8}}, \quad t > 0,$$

These concentration inequalities say that, with a very high<sup>11</sup> probability,  $X/\sqrt{k}$  remains at an *arbitrarily small* distance from the unit sphere  $\mathbb{S}^{k-1} \subset \mathbb{R}^k$  or, equivalently,  $X$  remains at a *bounded* distance from the round sphere  $\mathbb{S}_{\sqrt{k}}^{k-1} \subset \mathbb{R}^k$  of radius  $\sqrt{k}$  as  $k \rightarrow +\infty$ . We note the striking similarity between (5.15) and (5.8): in both cases the relevant random variable, which turns out to be a function of a large collection of independent variables, becomes *almost constant* when properly re-scaled; for more on this perspective, which in a sense underlies the modern applications of the “concentration of measure phenomenon” to (high dimensional) probability, see [Tal96].

**Remark 5.14.** (Poincaré’s limit theorem) From Remark 4.26 we know that a random vector  $Z^{[k]}$  uniformly distributed over  $\mathbb{S}_{\sqrt{k}}^{k-1}$  may be expressed as

$$(5.17) \quad Z^{[k]} = \sqrt{k}\Theta^{[k]},$$

where  $\Theta^{[k]} = X^{[k]}/\|X^{[k]}\|$ ,  $X^{[k]} \sim \mathcal{N}(\vec{0}, \text{Id}_k)$ . Now, it follows from (5.15) that  $\|X^{[k]}\|/\sqrt{k} \xrightarrow{P} 1$  as  $k \rightarrow +\infty$ <sup>12</sup>. On the other hand, if  $x \in \mathbb{R} = \mathbb{R}^1 \subset \mathbb{R}^2 \subset \cdots \subset \mathbb{R}^k \subset \cdots$ ,  $\|x\| = 1$ , we know that  $\langle X^{[k]}, x \rangle \sim \mathcal{N}(0, 1)$  by Proposition 4.7 (3). Thus, from the identity

$$\langle Z^{[k]}, x \rangle = \frac{\sqrt{k}}{\|X^{[k]}\|} \langle X^{[k]}, x \rangle$$

and Theorem 2.23 we conclude that  $\langle Z^{[k]}, x \rangle \xrightarrow{d} \mathcal{N}(0, 1)$ . Put in another way, as  $k \rightarrow +\infty$  the marginals of  $Z^{[k]}$  corresponding to a given set of  $l \geq 1$  coordinates converge in distribution to a standard normal vector  $Z^{[\infty]} \sim \mathcal{N}(\vec{0}, \text{Id}_l)$ , a statement usually referred to as “Poincaré’s limit theorem” [HN64, McK73, DF87].  $\square$

<sup>11</sup>That is, as close to 1 as we wish!

<sup>12</sup>This assertion also follows from the Law of Large Numbers (Theorem 6.2 below). Indeed, if  $X \sim \mathcal{N}(\vec{0}, \text{Id}_k)$  then

$$\|X\|^2 = \sum_{j=1}^k X_j^2,$$

where  $X_j \sim \mathcal{N}(0, 1)$  and hence  $X_j^2 \sim \chi_1^2$  with  $\mathbb{E}(X_j^2) = 1$  by Corollary 4.25. Thus, by making  $X = X^{[k]}$ , Theorem 6.2 applies to ensure that  $\|X^{[k]}\|^2/k \xrightarrow{P} 1$ , as desired.

As yet another instance of an insight coming from the concentration inequalities above, if we compare the bounds in (5.15) and (5.16) with the normal bound in (5.3), we see that the corresponding fluctuations, as measured by the standard deviation, are  $O(1/\sqrt{k})$  and  $O(1)$ , respectively. Noticing that

$$\frac{1}{\sqrt{k}} = \text{Lip} \left( x \mapsto \frac{\|x\|}{\sqrt{k}} \right)$$

and

$$1 = \text{Lip} (x \mapsto \|x\|),$$

where  $\text{Lip}$  denotes the Lipschitz constant of a function on  $\mathbb{R}^k$  (with respect to the euclidean norm), we are thus led to suspect that the dimension-free inequality

$$(5.18) \quad P(|F(X) - \mathbb{E}(F(X))| > t) \leq 2e^{-\frac{Ct^2}{\text{Lip}(F)^2}}, \quad t > 0,$$

should hold true for some universal constant  $C > 0$  not depending on  $k$ , where  $F : \mathbb{R}^k \rightarrow \mathbb{R}$  is assumed to be Lipschitz. Of course, the optimal possibility is  $C = 1/2$ , in which case (5.18) says that if  $\text{Lip}(F) = 1$  then  $F(X)$  is at least as much concentrated around its mean as each  $X_j$ , regardless of the size  $k$  of the sample. For discussions on this *Gaussian Concentration Inequality* which explore the connection with several other mathematical topics, including the geometric notion of *isoperimetry* and the analytical concept of *hypercontractivity*, we refer to [Led01, Led06, LT13, BLM13]; see also Remark 5.15 below. An elegant approach to the sharpest version of (5.18), due to Maurey and Pisier [Pis06, Chapter 2] and relying heavily on Itô's Stochastic Calculus, is presented in Subsection A.4 below. We provide here a more pedestrian (but no less elegant!) argument, also available in [Pis06, Chapter 2], which delivers a constant slightly smaller than the optimal one ( $1/2$  gets replaced by  $2/\pi^2$ ).

Without loss of generality, we may assume that  $F$  is smooth (so that  $\|\nabla F\| \leq \text{Lip}(f)$  a.s.) and  $\mathbb{E}(F(X)) = 0$ . This latter assumption implies, via Jensen's inequality, that  $\mathbb{E}(e^{vF(X)}) \geq e^{v\mathbb{E}(F(X))} = 1$  for any  $v \in \mathbb{R}$ , which gives

$$(5.19) \quad \mathbb{E} \left( e^{uF(X)} \right) \leq \mathbb{E}_{(X, X')} \left( e^{u(F(X) - F(X'))} \right), \quad u \geq 0,$$

where  $X'$  is an independent copy of  $X$  (so that  $(X, X') \sim \mathcal{N}(\vec{0}, \text{Id}_{2k})$ ). For each  $\theta \in [0, \pi/2]$  define

$$X_\theta = \cos \theta X + \sin \theta X'$$

and

$$X'_\theta = \frac{dX_\theta}{d\theta} = -\sin \theta X + \cos \theta X'.$$

Since

$$\begin{pmatrix} X_\theta \\ X'_\theta \end{pmatrix} = \begin{pmatrix} \cos \theta \text{Id}_k & \sin \theta \text{Id}_k \\ -\sin \theta \text{Id}_k & \cos \theta \text{Id}_k \end{pmatrix} \begin{pmatrix} X \\ X' \end{pmatrix},$$

it follows from the rotational invariance in Corollary 4.12 that  $(X_\theta, X'_\theta)$  is identically distributed to  $(X, X')$  with  $X'_\theta$  being an independent copy of  $X_\theta$  as well. Now,

$$F(X) - F(X') = \int_0^{\pi/2} \langle (\nabla F)(X_\theta), X'_\theta \rangle d\theta,$$

so that Jensen's inequality gives

$$e^{u(F(X) - F(X'))} \leq \frac{2}{\pi} \int_0^{\pi/2} e^{\frac{\pi}{2} u \langle (\nabla F)(X_\theta), X'_\theta \rangle} d\theta,$$

and hence,

$$\begin{aligned} \mathbb{E}_{(X, X')} \left( e^{u(F(X) - F(X'))} \right) &\leq \frac{2}{\pi} \int_0^{\pi/2} \mathbb{E}_{(X, X')} \left( e^{\frac{\pi}{2} u \langle (\nabla F)(X_\theta), X'_\theta \rangle} \right) d\theta \\ &= \mathbb{E}_{(X, X')} \left( e^{\frac{\pi}{2} u \langle (\nabla F)(X), X' \rangle} \right) \\ &= \int_{\mathbb{R}^k} \left( \mathbb{E}_{X'} \left( e^{\frac{\pi}{2} u \langle (\nabla F)(x), X' \rangle} \right) \right) dP_X(x), \end{aligned}$$

where in the second step we used the consequences of the rotational invariance mentioned above to ensure that the expectation integrand does not depend on  $\theta$  and in the last step we used the independence of  $\{X, X'\}$ . Now, for each  $x$  such that  $\nabla F(x) \neq \vec{0}$  we know from Proposition 4.7 (3) that

$$\frac{\pi}{2} \langle (\nabla F)(x), X' \rangle \sim \mathcal{N} \left( 0, \frac{\pi^2}{4} \|\nabla F(x)\|^2 \right),$$

so that (4.5) gives

$$\mathbb{E}_{X'} \left( e^{\frac{\pi}{2} u \langle (\nabla F)(x), X' \rangle} \right) \leq e^{\frac{\pi^2}{8} L(f)^2 u^2}.$$

Note that the right-hand side does not depend on  $x$  and the estimate remains true if  $(\nabla F)(x) = \vec{0}$ . Thus, if we put all the pieces of our calculation together we obtain

$$\mathbb{E} \left( e^{uF(X)} \right) \leq e^{\frac{\pi^2}{8} L(f)^2 u^2},$$

which amounts to saying that  $F(X) \in \text{SubG}(\pi \text{Lip}(f)/2)$ , a quite good but not entirely satisfactory estimate due to the  $\pi/2$  factor appearing in the sub-Gaussian parameter<sup>13</sup>. In any case, we may now appeal to the Cramér-Chernoff method introduced above: for each  $t > 0$ ,

$$P(|F(X)| > t) \leq 2e^{\frac{\pi^2}{8} L(f)^2 u^2 - ut}, \quad u \geq 0,$$

and minimizing the right-hand side we finally get

$$P(|F(X)| > t) \leq 2e^{-\frac{2t^2}{\pi^2 L(f)^2}},$$

as desired.

**Remark 5.15.** (Poncaré's limit revisited) Let  $\Pi_{k,l} : \mathbb{R}^k \rightarrow \mathbb{R}^l$  be the orthogonal projection associated to the natural embedding  $\mathbb{R}^l \hookrightarrow \mathbb{R}^k$  and let  $P_k$  be the uniform probability measure on  $\mathbb{S}_{\sqrt{k}}^{k-1}$ . With this terminology, the "Poncaré's limit theorem" in Remark 5.14 says that the random vectors  $\tilde{\Pi}_{k,l} = \Pi_{k,l}|_{\mathbb{S}_{\sqrt{k}}^{k-1}} : (\mathbb{S}_{\sqrt{k}}^{k-1}, P_k) \rightarrow \mathbb{R}^l$  converge in distribution to  $Z^{[\infty]} \sim \mathcal{N}(\vec{0}, \text{Id}_l)$ , which means that  $\mathbb{E}(\xi(\tilde{\Pi}_{k,l})) \rightarrow \mathbb{E}(\xi(Z^{[\infty]}))$  for any  $\xi : \mathbb{R}^l \rightarrow \mathbb{R}$  uniformly bounded and continuous; cf. Definition 2.21. It turns out that with a bit more of effort it may be checked that this statement actually holds true with  $\xi = \mathbf{1}_A$ , the indicator function of an arbitrary Borel set  $A \in \mathcal{B}^l$ . Precisely,

$$(5.20) \quad \lim_{k \rightarrow +\infty} P_{\tilde{\Pi}_{k,l}}(A) = \frac{1}{(2\pi)^{l/2}} \int_A e^{-\|y\|^2/2} dy.$$

To prove this claim, let us first observe that, with the notation of Remark 5.14,

$$\begin{aligned} P_{\tilde{\Pi}_{k,l}}(A) &= P_k(\tilde{\Pi}_{k,l}^{-1}(A)) \\ &= P_k \left( \Pi_{k,l}^{-1}(A) \cap \mathbb{S}_{\sqrt{k}}^{k-1} \right) \\ &= P \left( \Pi_{k,l}(Z^{[k]}) \in A \right), \end{aligned}$$

<sup>13</sup>This should be compared to the sharp estimate in (A.23) obtained by means of the full machinery of the Stochastic Calculus.

so if  $R_m^2 := X_1^2 + \dots + X_m^2$ ,  $1 \leq m \leq k$ , (5.17) gives

$$\begin{aligned} P_{\tilde{\Pi}_{k,l}}(A) &= P\left(\frac{\sqrt{k}}{R_k}(X_1, \dots, X_l) \in A\right) \\ &= P\left(\left(k \frac{R_l^2}{R_k^2}\right)^{1/2} \frac{1}{R_l}(X_1, \dots, X_l) \in A\right). \end{aligned}$$

Now, by Remark 4.26,

$$\left\{R_l^2, R_k^2 - R_l^2, \frac{1}{R_l}(X_1, \dots, X_l)\right\}$$

is independent and therefore

$$\left\{\frac{R_l^2}{R_k^2}, \frac{1}{R_l}(X_1, \dots, X_l)\right\}$$

is independent as well. On the other hand, by Corollary 4.25 and Proposition 4.36,

$$\frac{R_l^2}{R_k^2} = \frac{R_l^2}{R_l^2 + (R_k^2 - R_l^2)} \sim \text{Beta}\left(\frac{l}{2}, \frac{k-l}{2}\right),$$

so if we put together these facts we get, by Proposition 2.13,

$$\begin{aligned} P_{\tilde{\Pi}_{k,l}}(A) &= \frac{\Gamma(k/2)}{\Gamma(l/2)\Gamma((k-l)/2)} \times \\ &\quad \omega_{l-1}^{-1} \int_{\mathbb{S}_1^{l-1}} \int_0^1 \mathbf{1}_A(\sqrt{kx}\theta) x^{\frac{l}{2}-1} (1-x)^{\frac{k-l}{2}-1} d\mathbb{S}_1^{l-1}(\theta) dx, \end{aligned}$$

where  $d\mathbb{S}_1^{l-1}(\theta)$  is the (unnormalized) volume element of the unit sphere  $\mathbb{S}_1^{l-1}$  (induced by the embedding  $\mathbb{S}_1^{l-1} \hookrightarrow \mathbb{R}^l$ ) and  $\omega_{l-1} = \text{vol}_{l-1}(\mathbb{S}_1^{l-1})$ . Thus, if we set  $u = \sqrt{kx}$  and use (4.24) we find that

$$\begin{aligned} P_{\tilde{\Pi}_{k,l}}(A) &= \frac{\Gamma(k/2)}{\Gamma((k-l)/2)} \frac{1}{\pi^{l/2} k^{l/2}} \times \\ &\quad \int_{\mathbb{S}_1^{l-1}} \int_0^{\sqrt{k}} \mathbf{1}_A(u\theta) u^{l-1} \left(1 - \frac{u^2}{k}\right)^{\frac{k-l}{2}-1} d\mathbb{S}_1^{l-1}(\theta) du. \end{aligned}$$

Since  $l$  is held fixed, the Stirling approximation in (6.1) below gives

$$\frac{\Gamma(k/2)}{\Gamma((k-l)/2)} \approx_{k \rightarrow +\infty} 2^{-l/2} k^{l/2} \left(\frac{k-l}{k}\right)^{-l/2},$$

so we end up with

$$\lim_{k \rightarrow +\infty} P_{\tilde{\Pi}_{k,l}}(A) = \frac{1}{(2\pi)^{l/2}} \int_{\mathbb{S}_1^{l-1}} \int_0^{+\infty} \mathbf{1}_A(u\theta) u^{l-1} e^{-u^2/2} d\mathbb{S}_1^{l-1}(\theta) du,$$

which proves the claim because this double integral clearly equals the right-hand side of (5.20) under the substitution  $y = u\theta$ . The limit theorem in (5.20) is the key ingredient in explicitly solving the isoperimetric problem for the Gaussian space  $(\mathbb{R}^l, \delta, (2\pi)^{-l/2} e^{-|y|^2/2} dy)$  by essentially viewing it as the limit of the corresponding problem for large, high-dimensional spheres  $\mathbb{S}_{\sqrt{k}}^{k-1}$  as  $k \rightarrow +\infty$  [Bor75, ST78], a celebrated result which by its turn may be used to establish a version of (5.18) with the optimal constant  $C = 1/2$ , but this time with the mean replaced by the median [Led06, Chapter 2]. Needless to say, this is a prominent instance of the “concentration of measure phenomenon” extensively studied elsewhere [GKPS99, Led01, BLM13, Shi16].  $\square$

**5.3. Chernoff-type bounds for binomial trials and the Erdős-Rényi model.** As illustrated in (5.9), the sub-Gaussian version of the Cramér-Chernoff method yields an estimate which fails to account for the real dispersion of a binomial trial (as measured by its variance). However, we may remedy this by directly applying the method to  $X \sim \text{Bin}(p; n)$  as in Example 2.37 in order to get, for  $t > 0$ ,

$$\begin{aligned} P(X \geq t) &= P(e^{Xu} \geq e^{tu}) \\ &= e^{-tu} \mathbb{E}(e^{Xu}) \quad (\text{Markov}) \\ &\stackrel{(2.25)}{=} e^{-tu} (1 - p + pe^u)^n. \end{aligned}$$

Using that  $1 + x \leq e^x$ ,  $x \geq 0$ , we obtain

$$(5.21) \quad P(X \geq t) \leq e^{-tu + np(e^u - 1)},$$

and since the function on the exponent is minimized at  $u = \ln(t/\lambda)$ , where  $\lambda = np = \mathbb{E}(X)$  is the expectation, we end up with the *Chernoff-type inequality*

$$(5.22) \quad P(X \geq t) \leq e^{-\lambda} \left( \frac{e\lambda}{t} \right)^t, \quad t > \lambda.$$

Regarding this analysis, the following comments are worth mentioning:

- From (5.22) we have

$$(5.23) \quad P(X \geq t) \leq C_1 e^{C_2 t - t \ln t},$$

where  $C_1 = e^{-\lambda}$  and  $C_2 = 1 + \ln \lambda$ , which for  $t$  large gives a tail behavior somehow interpolating between the sub-Gaussian and sub-exponential regimes.

- Past experience with the sub-exponential case in Proposition 5.7 suggests that we should be able to recover a sub-Gaussian tail for *small* deviations around the mean  $\lambda$  (which is the only critical point of the exponential function in the right-hand side of (5.23)). This is the case indeed: if we insert  $t = (1 + \varepsilon)\lambda$ ,  $|\varepsilon| < 1$ , in (5.22) we see that

$$\begin{aligned} P(X \geq (1 + \varepsilon)\lambda) &\leq \left( e^{\varepsilon - (1 + \varepsilon) \ln(1 + \varepsilon)} \right)^\lambda \\ &= \left( e^{-\frac{\varepsilon^2}{2} + \frac{\varepsilon^3}{6} + o(|\varepsilon|^4)} \right)^\lambda, \end{aligned}$$

which easily leads to the estimate

$$(5.24) \quad P(|X - \lambda| \geq \varepsilon\lambda) \leq 2e^{-\frac{\varepsilon^2}{3}\lambda}, \quad 0 < \varepsilon < 1.$$

- From (5.21) and (2.27) with  $\lambda = np$  we see that the Cramér-Chernoff method delivers the same estimate as in (5.22) had we started with  $X \sim \mathcal{P}(\lambda)$ , the Poisson variable with the same expectation as our original binomial variable; see Example 2.38. This clearly suggests that, at least in the asymptotic regime  $n \rightarrow +\infty$ , the classes  $\text{Bin}(\lambda/n; n)$  and  $\mathcal{P}(\lambda)$  are closely related, a claim substantiated by the Law of Rare Events (Theorem 6.13 below). This deep relationship between binomial and Poisson distributions finds many applications in the theory of random graphs, notably in connection with the Erdős-Rényi model studied in the sequel; see the proof of Proposition 5.18 below for a simple manifestation of this connection and [VDH24] for the general theory.

We now illustrate the of the Chernoff-type bounds for binomial trials developed so far in the art of precisely determining the exact threshold for the emergence of certain “phase transitions” in the most commonly studied class of random graphs.

For each  $N \geq 2$  let us define  $[N] = \{1, \dots, N\}$ , which we call a *set of vertices*. We represent by  $[ij]$  the *unordered pair* derived from  $\{i, j\} \in [N] \times [N]$  with  $i \neq j$  (so that  $[ij] = [ji]$ ). The union of all such objects is

the set of potential edges, denoted  $E$ . Note that  $\sharp(E) = N(N-1)/2$ . We now inject a probabilistic ingredient in the construction of graphs starting with  $E$ . The most obvious possibility, which we adopt here, is simply to flip a (possibly biased) coin for each potential edge in order to decide whether it effectively occurs as a link between two vertices, with the provision that the flippings should comprise independent events. In formal terms, for each  $e \in E$  we consider a random variable  $X_e \sim \text{Ber}(p)$  so that  $\text{supp } P_{X_e} = \{0, 1\}$  with  $P(X_e = 1) = 1 - P(X_e = 0) = p$ . Now define  $\Omega = \{0, 1\}^{\binom{N}{2}}$ , the cartesian product of  $N(N-1)/2$  copies of  $\{0, 1\}$ , one for each  $e \in E$ , and  $P = \otimes_{e \in E} P_{X_e}$ , the product probability on  $\Omega$ . Note that each element  $\omega \in \Omega$  may be viewed as a function  $\omega : \Omega \rightarrow \{0, 1\}$  and hence defines a graph whose edge set is

$$E_\omega = \{e \in E; \omega(e) = 1\}.$$

For this reason, the sample space  $(\Omega, 2^\Omega, P)$  is called the *Erdős-Rényi model* for a random graph, usually denoted by  $\mathbb{G}(N; p)$ .

**Proposition 5.16.** *Let  $\pi_e : \Omega \rightarrow \{0, 1\}$  be the canonical projection onto the factor corresponding to  $e$ . Then each  $\pi_e$  is identically distributed to  $X_e$  (in particular,  $\pi_e \sim \text{Ber}(p)$ ) with  $\{\pi_e\}_{e \in E}$  being independent.*

*Proof.* This is a special case of the general procedure in Remark 2.16. □

By construction of  $\mathbb{G}(N; p)$ , any event (a subset of  $\Omega$ ) defines a specific collection of graphs. For instance, for each  $e \in E$  we may consider

$$\Omega_e = \{\omega \in \Omega; \omega(e) = 1\},$$

the set of all graphs having  $e$  as a vertex. The next result confirms that a random graph in the Erdős-Rényi model is obtained by flipping a coin for each potential vertex with the flippings being independent moves.

**Proposition 5.17.**  *$\{\Omega_e\}_{e \in E}$  is a set of independent events.*

*Proof.* Note that  $\omega(e) = \pi_e(\omega)$  so that  $\Omega_e = \pi_e^{-1}(1)$  and then apply Proposition 5.16. □

Since  $\pi_e = \mathbf{1}_{\Omega_e}$  for each  $e$ , the total number of edges in a random graph  $\omega$  is given by  $\mathcal{E}_N(\omega)$ , where

$$\mathcal{E}_N = \sum_{e \in E} \pi_e,$$

so that, from Proposition 5.16 and Example 2.37,

$$\mathcal{E}_N \sim \text{Bin}\left(p; \binom{N}{2}\right) \implies \mathbb{E}(\mathcal{E}_N) = \binom{N}{2}p.$$

It follows from (5.24) that

$$P(|\mathcal{E}_N - \mathbb{E}(\mathcal{E}_N)| < \varepsilon \mathbb{E}(\mathcal{E}_N)) \geq 1 - 2e^{-\frac{\varepsilon^2}{3} \mathbb{E}(\mathcal{E}_N)}, \quad 0 < \varepsilon < 1.$$

Thus, as  $N \rightarrow +\infty$ ,

$$\frac{\mathcal{E}_N}{\mathbb{E}(\mathcal{E}_N)} \xrightarrow{p} 1,$$

so that the total number of edges *asymptotically* approaches its expected value. More generally, if we assume that  $p = p_N$  (that is, the biased coin possibly changes with  $N$ ) then the same conclusion holds as long as

$$(5.25) \quad \mathbb{E}(\mathcal{E}_N) = \binom{N}{2}p_N \rightarrow +\infty,$$

with the expectation now being computed with respect to  $\text{Bin}\left(p_N; \binom{N}{2}\right)$ . At this point, a slightly more ambitious task would be to make sure that, with very high probability, a minimal amount of edge emerges in the regime determined by (5.25). That this is the case indeed follows from the next result, which actually shows that the *asymptotic* emergence of a fixed number of edges in the Erdős-Rényi model is explicitly determined by the limiting value of  $\mathbb{E}(\mathcal{E}_N)$ .

**Proposition 5.18.** *Under the conditions above, if  $m \in \mathbb{N}$ ,*

$$\lim_{N \rightarrow +\infty} P(\mathcal{E}_N > m) = \begin{cases} 0 & \mathbb{E}(\mathcal{E}_N) \rightarrow 0 \\ 1 - e^{-\lambda} \sum_{k=0}^m \frac{\lambda^k}{k!} & \mathbb{E}(\mathcal{E}_N) \rightarrow \lambda \in \mathbb{R}_+ \\ 1 & \mathbb{E}(\mathcal{E}_N) \rightarrow +\infty \end{cases}$$

*Proof.* We only prove the convergence in the middle since the remaining items, at least formally, follow from this case. By the Law of Rare Events (Theorem 6.13 below) there exists a Poisson variable  $Z \sim \mathcal{P}(\lambda)$  such that  $\mathcal{E}_N \xrightarrow{d} Z$  as  $n \rightarrow +\infty$ . Hence,

$$\begin{aligned} \lim_{N \rightarrow +\infty} P(\mathcal{E}_N > m) &= P(Z > m) \\ &= 1 - P(Z \leq m) \\ &= 1 - e^{-\lambda} \sum_{k=0}^m \frac{\lambda^k}{k!}, \end{aligned}$$

as desired. □

We now turn to the incidence properties of  $\mathbb{G}(N; p)$ . For each vertex  $i \in [N]$  consider the random variable

$$d_i = \sum_{j: j \neq i} \pi_{[ij]}.$$

Clearly, for each graph  $\omega \in \Omega$ ,  $d_i(\omega)$  measures the number of edges of  $\omega$  having  $i$  as a vertex. We call  $d_i$  the *degree*.

**Proposition 5.19.** *For each  $i$ ,  $d_i \sim \text{Bin}(p; N - 1)$ . In particular,  $d := \mathbb{E}(d_i) = (N - 1)p$ .*

*Proof.* Immediate from Proposition 5.16 and Example 2.37. □

Recall that a random graph is *almost regular* if the degree of each vertex equals its expected value with very high probability. The next result identifies the threshold on the degree function beyond which almost regularity holds in the Erdős-Rényi model.

**Proposition 5.20.** *For any  $\varepsilon, \delta \in (0, 1)$  there exists  $C = C_{\varepsilon, \delta} > 0$  such that  $d \geq C \ln N$  implies*

$$P(|d_i - d| \leq \varepsilon d \text{ for all } i) \geq 1 - \delta.$$

*Proof.* For each  $i \in [N]$  we have from Proposition 5.19 and (5.24) that

$$P(|d_i - d| > \varepsilon d) \leq 2e^{-\frac{\varepsilon^2}{3}d},$$

so that

$$P(|d_i - d| > \varepsilon d \text{ for some } i) \leq 2Ne^{-\frac{\varepsilon^2}{3}d},$$

and hence

$$P(|d_i - d| \leq \varepsilon d \text{ for all } i) \geq 1 - 2Ne^{-\frac{\varepsilon^2}{3}d}.$$

Thus, we must find  $C$  such that

$$2Ne^{-\frac{\varepsilon^2}{3}C \ln N} \leq \delta,$$

or equivalently,

$$C \geq \frac{3}{\varepsilon^2} h_\delta(N), \quad h_\delta(N) = \frac{\ln(2/\delta) + \ln N}{\ln N}.$$

Now, as  $N$  varies  $h_\delta(N)$  is uniformly bounded by  $M_\delta = \ln(4/\delta)/\ln 2$ , so it suffices to take  $C \geq 3M_\delta/\varepsilon^2$ .  $\square$

Note that the almost regularity in Proposition 5.20 implies a sort of homogeneous behavior of the random graph around each of its vertices<sup>14</sup>. In particular, the event that no vertex is isolated occurs with high probability. Now, it turns out that in the regime where  $p_N \approx \ln N/N$  with  $N \rightarrow +\infty$ , this event is essentially equiprobable to the event defining connectedness of a random graph, which suggests that  $\ln N/N$  should be a sharp threshold for the *asymptotic* occurrence of this topological property. Indeed, arguing along these lines it may be shown that if  $p_N = c_N \ln N/N$  and

$$K := \lim_{N \rightarrow +\infty} (c_N - 1) \ln N = \lim_{N \rightarrow +\infty} (Np_N - \ln N)$$

exists as an extended real number then

$$\lim_{N \rightarrow +\infty} P(\{\omega \in \mathbb{G}(N; p_N) : \omega \text{ is connected}\}) = e^{-e^{-K}}.$$

In particular,

$$\lim_{N \rightarrow +\infty} P(\{\omega \in \mathbb{G}(N; p_N) : \omega \text{ is connected}\}) = \begin{cases} 0 & c_N \rightarrow c < 1 \\ 1 & c_N \rightarrow c > 1 \end{cases}$$

For full discussions on this and similar “phase transition” phenomena exhibiting a sharp threshold in the Erdős-Rényi model, see [JLR11, FK16, VDH24].

## 6. THE FUNDAMENTAL LIMIT THEOREMS

We present here a couple of asymptotic results which are central in the theory. We emphasize, however that, differently from the concentration estimates in Section 5, which are quite effective due its manifestly *non-asymptotic* character, eventual applications of the limits theorems only become reliable in the asymptotic regime (when the number of random variables gets larger and larger); see Remark 6.7. The proofs we describe below depend on a rather special case of a deep convergence result due to Lévy [Wil91, Theorem 18.1]. The version we present here is the appropriate converse to Remark 2.25 and may be approached via Fourier Analysis.

**Theorem 6.1.** (*Lévy’s convergence*) Let  $\{Z_j\}_{j=1}^\infty$  be a random variable such that  $\phi_{Z_j}$  converges pointwise to  $\phi_Z$ , where  $Z$  is another random variable. Then  $Z_j \rightarrow Z$  in distribution.

<sup>14</sup>Incidentally, this homogeneity confirms that the Erdős-Rényi random graph fails to reliably model real-world complex networks, where a sizable amount of variability of the incidence pattern of the vertices is observed.



*Proof.* By a simple approximation we may assume that  $\xi$  in Definition 2.21 is an arbitrary Schwartz function. Hence,

$$\begin{aligned}\mathbb{E}(\xi(Z_j)) &= \int_{-\infty}^{+\infty} \xi(z_j) dP_{Z_j}(z_j) \\ &= \int_{-\infty}^{+\infty} \left( \int_{-\infty}^{+\infty} \widehat{\xi}(u) e^{iz_j u} du \right) dP_{Z_j}(z_j),\end{aligned}$$

where we used Fourier inversion in order to recover  $\xi$  from its Fourier transform  $\widehat{\xi}$ , which is Schwarz as well and hence uniformly bounded. Using Fubini and dominated convergence we get

$$\begin{aligned}\mathbb{E}(\xi(Z_j)) &= \int_{-\infty}^{+\infty} \widehat{\xi}(u) \left( \int_{-\infty}^{+\infty} e^{iz_j u} dP_{Z_j}(z_j) \right) du \\ &= \int_{-\infty}^{+\infty} \widehat{\xi}(u) \phi_{Z_j}(u) du \\ &\xrightarrow{j \rightarrow +\infty} \int_{-\infty}^{+\infty} \widehat{\xi}(u) \phi_Z(u) du \\ &\quad \vdots \quad (\text{the same computation as above in reverse}) \\ &= \mathbb{E}(\xi(Z)),\end{aligned}$$

and the result follows.  $\square$

We may now present the first fundamental limit theorem.

**Theorem 6.2.** (Law of large numbers, LLN) If  $\{X_j\}_{j \geq 1}$  is a sequence of i.i.d. (that is, independent and identically distributed) real random variables with  $\mathbb{E}(X_j) = \mu$  then the sequence of random variables

$$\overline{X}_n := \frac{1}{n}(X_1 + \cdots + X_n)$$

converges in probability to  $\mu$  as  $n \rightarrow +\infty$ .

*Proof.* By Proposition 2.22, it suffices to prove that  $\overline{X}_n \rightarrow \mu$  in distribution. By Propositions 2.29 and 2.30, if  $|u|/n$  is small,

$$\phi_{\overline{X}_n}(u) = \Pi_{j=1}^n \phi_{X_j}(u/n) = \left[ 1 + \mu \frac{u}{n} \mathbf{i} + o\left(\frac{|u|}{n}\right) \right]^n, \quad n \rightarrow +\infty,$$

so that, for any  $u \in \mathbb{R}$ ,

$$\lim_{n \rightarrow +\infty} \phi_{\overline{X}_n}(u) = e^{u\mu \mathbf{i}} = \phi_\mu(u),$$

Now apply Theorem 6.1.  $\square$

**Remark 6.3.** We append two complements to this result:

- (1) If we further assume that  $\mathbb{E}(|X_j|^2) < +\infty$  then it also follows from the argument based on (7.32) below, which relies on Chebyshev's inequality and hence provides a quite effective (i.e. *non-asymptotic*) estimate;
- (2) For obvious reasons, Theorem 6.2 is usually referred as the *weak* LLN. With some more effort we may show that the convergence holds in a rather strong sense:  $\overline{X}_n \xrightarrow{a.s.} \mu$ . This latter result is usually known as *Kolmogorov's LLN* (in [Kre11, Section 1.4] it is shown how it follows from Birkhoff's ergodic theorem discussed in Example 3.12).  $\square$

**Remark 6.4.** The limiting behavior of the Student's  $t$ -distribution  $t_k$  in Definition 4.29 as the number of degrees of freedom  $k$  grows indefinitely may be determined if one makes use of the Stirling asymptotics for the gamma function:

$$(6.1) \quad \Gamma(k) \approx_{k \rightarrow +\infty} \sqrt{2\pi} k^{k-\frac{1}{2}} e^{-k};$$

see Remark 6.8 below for a probabilistic proof of this result. Using this, a little computation starting with (4.25) then shows that

$$\lim_{k \rightarrow +\infty} t_k(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad x \in \mathbb{R},$$

so that

$$(6.2) \quad t_k \xrightarrow{d} \mathcal{N}(0, 1)$$

by Scheffé's lemma [Sch47]. We point out that this may also be justified with a simple application of Theorem 6.2. Indeed, by Remark 2.16 we may pick an independent sample  $Z, X_1, \dots, X_k \sim \mathcal{N}(0, 1)$ , so that  $X_j^2 \sim \chi_1^2$  by Proposition 4.25. By LLN, as  $k \rightarrow +\infty$  we have that  $W_k := \sum_{j=1}^k X_j^2$  satisfies

$$\sqrt{\frac{W_k}{k}} \rightarrow \sqrt{\mathbb{E}(\chi_1^2)} = 1$$

in probability. Hence,  $Z/\sqrt{W_k/k} \rightarrow \mathcal{N}(0, 1)$  in distribution so that (6.2) may be verified using that  $Z/\sqrt{W_k/k}$  is  $t_k$ -distributed by Proposition 4.30. As another instance of this kind of argument, let us check that if  $X \sim F_{k_1, k_2}$  then

$$k_1 X \xrightarrow{d} \chi_{k_1}^2 \quad \text{as } k_2 \rightarrow +\infty.$$

Indeed, from Proposition 4.33 we may write

$$X = \frac{W_1/k_1}{W_2/k_2},$$

where  $W_1 \perp W_2$  and  $W_j \sim \chi_{k_j}^2$ ,  $j = 1, 2$ . Thus,

$$k_1 X = \frac{W_1}{W_2/k_2}$$

and since  $W_2/k_2 \xrightarrow{p} 1$  the claim follows.  $\square$

The next result provides an accurate description of a re-scaled version of  $\overline{X}_n$  and illustrates the ubiquitous character of the normal distribution in Probability Theory. A possible rationale behind it goes as follows. From Theorem 6.2 we suspect that there exists a (possibly monotone) function  $n \in \mathbb{N} \mapsto \nu(n) \in \mathbb{R}$  satisfying  $\nu(n) \rightarrow +\infty$  as  $n \rightarrow +\infty$  and such that, rather informally,

$$\overline{X}_n \approx \mu + O(\nu(n)^{-1}),$$

so that

$$\nu(n) (\overline{X}_n - \mu) \approx O(1)$$

has a chance to converge to something finite. To guess who  $\nu$  might be, let us assume for a moment that the sample is normally distributed:  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ . It follows from Proposition 4.7 that  $\sqrt{n}(\overline{X}_n - \mu) \sim \mathcal{N}(0, \sigma^2)$ , which suggests that  $\nu(n) = \sqrt{n}$ . The striking feature of the next result is that, in alignment with Theorem 6.2, this exact relation holding for normal samples becomes in general an asymptotic convergence (in distribution) no matter how the original sample is distributed (as long as it has a finite variance).

**Theorem 6.5.** (Central Limit Theorem, CLT) Let  $\{X_j\}_{j \geq 1}$  be a sequence of i.i.d. real random variables with  $\mathbb{E}(X_j) = \mu$  and  $\text{var}(X_j) = \sigma^2 > 0$ . Then the sequence

$$(6.3) \quad Z_n := \frac{\sum_{j=1}^n X_j - n\mu}{\sqrt{n}\sigma} = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

converges in distribution to a random variable whose pdf is the standard normal distribution  $\mathcal{N}(0, 1)$ . Equivalently,  $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$ .

*Proof.* Define  $Y_j = (X_j - \mu)/\sigma$  so that  $\mathbb{E}(Y_j) = 0$  and  $\text{cov}(Y_j) = 1$ . Since  $Z_n = \sum_j Y_j/\sqrt{n}$ , by Propositions 2.29 and 2.30 we get, for  $|u|/\sqrt{n}$  small,

$$\phi_{Z_n}(u) = \Pi_{j=1}^n \phi_{Y_j} \left( \frac{u}{\sqrt{n}} \right) = \left[ 1 - \frac{u^2}{2n} + o \left( \frac{|u|^2}{n} \right) \right]^n,$$

so that

$$(6.4) \quad \lim_{n \rightarrow +\infty} \phi_{Z_n}(u) = e^{-\frac{1}{2}u^2}, \quad u \in \mathbb{R}.$$

By Corollary 4.6 and Proposition 4.7 (1), the right-hand side is the characteristic function of a random variable  $Z \sim \mathcal{N}(0, 1)$ , so we may apply Theorem 6.1 to conclude the proof.  $\square$

**Remark 6.6.** An enlightening discussion of several proofs of Theorem 6.5, including the one above, may be found in [Tao12, Chapter 2].  $\square$

**Remark 6.7.** A common misconception in practical applications of Theorem 6.5 is to take it for granted that the sample mean  $\bar{X}_n$  converges (in distribution) to a normal. To make things even worse, this is often extrapolated to the declaration that there is a threshold on the sample size (usually taken around  $n = 30$ ) beyond which  $\bar{X}_n$  is *exactly* distributed as a normal. Clearly, there is nothing in the statement of the theorem ensuring that this guess makes any sense and in fact both assertions above contradict Theorem 6.2, which says that  $\bar{X}_n$  converges (even almost surely by Remark 6.3 (2)) to the (constant) population mean  $\mu$ . On the other hand, it is immediate that Theorem 6.5 justifies *approximating*  $\bar{X}_n$  in distribution by  $\mathcal{N}(\mu, \sigma^2/n)$  as  $n \rightarrow +\infty$ . Thus, for any  $-\infty \leq a < b \leq +\infty$  we may assume for all practical purposes that

$$(6.5) \quad P(a \leq \bar{X}_n \leq b) \approx_{n \rightarrow +\infty} \frac{\sqrt{n}}{\sqrt{2\pi}\sigma} \int_a^b e^{-\frac{n(x-\mu)^2}{2\sigma^2}} dx$$

where  $\approx_{n \rightarrow +\infty}$  means that the equality only holds *asymptotically* in the regime of “very large” samples. Using this notation, CLT may be rephrased either as

$$(6.6) \quad \bar{X}_n \approx_{n \rightarrow +\infty} \mathcal{N}(\mu, \sigma^2/n)$$

or as

$$(6.7) \quad X^{(n)} := X_1 + \cdots + X_n \approx_{n \rightarrow +\infty} \mathcal{N}(n\mu, n\sigma^2).$$

In any case, those more inclined to applications should be aware that the convergence rates in the approximations above are typically quite slow. For instance, under the conditions of Theorem 6.5, and assuming further that  $\rho := \mathbb{E}(|X - \mu|^3) < +\infty$ , a version of the famous *Berry-Esseen theorem* ensures that

$$\sup_{x \in \mathbb{R}} |F_{\bar{X}_n}(x) - \Phi(x)| \leq \frac{C\rho}{\sqrt{n}\sigma^3}, \quad C > 0.$$

Although one could argue that a much better convergence rate than  $O(n^{-1/2})$  might work for a given  $\{X_j\}$ , certainly there exist examples of distributions for which this worst case upper bound is asymptotically achieved

even if moments of much higher order than three are required to be finite, the simplest example being the Rademacher variable in Remark 4.16 [Saz81, Chapter 1].  $\square$

**Remark 6.8.** We insist that the proof presented above *does* cover the case in which the initial i.i.d. sequence  $\{X_j\}$  is *discrete*. In fact, this is how the CLT first appeared, incarnated in the famous De Moivre-Laplace formulas (6.9)-(6.10) below [Fis11]. Let  $\{X_j\}_{j=1}^n$  be independent with  $X_j \sim \text{Ber}(p)$ , the Bernoulli distribution. From Example 2.37, we know that  $X^{(n)} = X_1 + \dots + X_n \sim \text{Bin}(p; n)$ , the binomial distribution. Since  $\mathbb{E}(X_j) = p$  and  $\text{var}(X_j) = p(1-p)$ , CLT applies<sup>15</sup> to give

$$(6.8) \quad Z_n = \sqrt{n} \frac{n^{-1} X^{(n)} - p}{\sqrt{p(1-p)}} \xrightarrow{d} \mathcal{N}(0, 1), \quad n \rightarrow +\infty,$$

or equivalently, if we combine (6.7) and (2.23),

$$(6.9) \quad \sum_{a \leq k \leq b} \binom{n}{k} p^k (1-p)^{n-k} \approx_{n \rightarrow +\infty} \frac{1}{\sqrt{2\pi np(1-p)}} \int_a^b e^{-\frac{(x-np)^2}{2np(1-p)}} dx, \quad a < b.$$

It is not hard to check that this is the same as having

$$(6.10) \quad \binom{n}{k} p^k (1-p)^{n-k} \approx_{n \rightarrow +\infty} \frac{1}{\sqrt{2\pi np(1-p)}} e^{-\frac{(k-np)^2}{2np(1-p)}}$$

uniformly in  $k$  satisfying

$$(6.11) \quad k = np + \sqrt{np(1-p)} O(1),$$

which may be proved by using Stirling's formula in (6.13) below and the fact that (6.11) implies that  $k/n \rightarrow p$  as  $n \rightarrow +\infty$ ; see [CA06, Section 7.3]. As yet another application of CLT in the discrete setting, let us assume that  $\{Y_j\}_{j=1}^n$  is independent with  $Y_j \sim \mathcal{P}(1)$ , the Poisson distribution as in Example 2.38. Thus,  $Y^{(n)} = Y_1 + \dots + Y_n \sim \mathcal{P}(n)$ , and CLT applies to yield

$$(6.12) \quad \sqrt{n}(n^{-1} Y^{(n)} - 1) \xrightarrow{d} \mathcal{N}(0, 1), \quad n \rightarrow +\infty,$$

that is,

$$\sum_{a \leq k \leq b} \frac{n^k e^{-n}}{k!} \approx_{n \rightarrow +\infty} \frac{1}{\sqrt{2\pi n}} \int_a^b e^{-\frac{(x-n)^2}{2n}} dx,$$

which is the same as having

$$\frac{n^k e^{-n}}{k!} \approx_{n \rightarrow +\infty} \frac{1}{\sqrt{2\pi n}} e^{-\frac{(k-n)^2}{2n}}$$

uniformly in  $k$  such that

$$k = n + \sqrt{n} O(1).$$

Taking  $k = n$  gives

$$(6.13) \quad n! \approx_{n \rightarrow +\infty} \sqrt{2\pi n} n^{n+\frac{1}{2}} e^{-n},$$

which is *Stirling's asymptotic formula*. If we take into account that  $\Gamma(n) = (n-1)!$ , this clearly implies (6.1).  $\square$

**Remark 6.9.** We may directly justify (6.8), the CLT for a Bernoulli population, as follows. Propositions 2.29 and 2.30 applied to (6.8) give

$$\phi_{Z_n}(u) = \phi_{X^{(n)}}(u) e^{-i \frac{np}{\sqrt{npq}} u},$$

<sup>15</sup>See Remark 6.9 below for a direct justification of this step along the lines of the proof of Theorem 6.5.

so that (2.24) leads to

$$\begin{aligned}\phi_{Z_n}(u) &= \left(q + pe^{\frac{u}{\sqrt{npq}}}\right)^n e^{-\frac{np}{\sqrt{npq}}u} \\ &= \left(qe^{-\frac{u}{\sqrt{npq}}} + pe^{\frac{u}{\sqrt{npq}}}\right)^n.\end{aligned}$$

Expanding the exponential terms in parentheses and performing some cancellations we find that

$$(6.14) \quad \phi_{Z_n}(u) = \left(1 - \frac{u^2}{2n} + o\left(\frac{u^2}{n}\right)\right)^n \xrightarrow{n \rightarrow +\infty} e^{-u^2/2},$$

which reproduces (6.4) in this case. We may also obtain a proof of (6.12), the CLT for a Poisson population, along the same lines. Indeed, this time the left-hand side of (6.12) is

$$Z_n = \frac{Y^{(n)}}{\sqrt{n}} - \sqrt{n},$$

so that

$$\begin{aligned}\phi_{Z_n}(u) &= \phi_{Y^{(n)}}\left(\frac{u}{\sqrt{n}}\right) e^{-i\sqrt{n}u} \\ (2.26) \quad &\stackrel{=}{=} e^{n\left(e^{\frac{i}{\sqrt{n}}u} - 1\right)} e^{-i\sqrt{n}u} \\ &= \left(e^{\frac{i}{\sqrt{n}}u} - 1 - \frac{i}{\sqrt{n}}u\right)^n \\ &= \left(e^{-u^2/2n + o(u^2/n)}\right)^n \\ &= \left(1 - \frac{u^2}{2n} + o\left(\frac{u^2}{n}\right)\right)^n,\end{aligned}$$

so we may proceed as in (6.14), as desired.  $\square$

**Example 6.10.** If  $\{X_j\}$  is i.i.d. with a common cdf  $F$  then its *empirical distribution function* is the random variable

$$\mathbb{F}_n(x) = \frac{1}{n} \sum_{j=1}^n \mathbf{1}_{\{X_j \leq x\}}, \quad n \in \mathbb{N}, \quad x \in \mathbb{R}.$$

Since, for each fixed  $x$ ,  $\mathbf{1}_{\{X_j \leq x\}} \sim \text{Ber}(F(x))$ , CLT applies:

$$\sqrt{n}(\mathbb{F}_n(x) - F(x)) \xrightarrow{d} \mathcal{N}(0, F(x)(1 - F(x))).$$

In particular,  $\mathbb{F}_n(x) \xrightarrow{p} F(x)$  and in fact a.s. convergence takes place; cf. Remark 6.3. With a bit more of work we can prove that the convergence is actually uniform in  $x$ :  $\|\mathbb{F}_n - F\|_\infty \rightarrow 0$  a.s. and assertion known as the *Glivenko-Cantelli theorem* [VdV00].  $\square$

The following immediate consequence of CLT, which uses the notation of Example 4.9, is also worth mentioning here.

**Theorem 6.11.** (Multiplicative CLT) If  $\{Y_j\}_{j \geq 1}$  is a i.i.d. sequence of positive random variables satisfying  $\mathbb{E}(\ln Y_j) = \mu$  and  $\text{var}(\ln Y_j) = \sigma^2$  then

$$\sqrt[n]{\prod_{j=1}^n Y_j} \approx_{n \rightarrow +\infty} \Lambda(\mu, \sigma^2/n) = \mathcal{LN}(e^{\mu + \frac{\sigma^2}{2n}}, (e^{\frac{\sigma^2}{n}} - 1)e^{2\mu + \frac{\sigma^2}{n}}).$$

**Example 6.12.** We say that a sequence of positive random variables  $\{X_j\}_{j=0}^{+\infty}$  satisfy *Gibrat's law of proportionate effect* if there exist random variables  $\{Y_j\}_{j=1}^{+\infty}$  such that  $Y_j \perp X_{j-1}$  and the corresponding cdfs satisfy

$$F_{X_j}(z) = \int_0^{+\infty} F_{Y_j}(xu^{-1}) dF_{X_{j-1}}(u), \quad j \geq 1.$$

It then follows from (2.10) that  $X_j = Y_j X_{j-1}$  and hence

$$X_0^{-1} X_n = \prod_{j=1}^n Y_j, \quad n \geq 1.$$

Thus, if  $\{Y_j\}$  is as in Theorem 6.11 we see that

$$X_0^{-1} X_n \approx_{n \rightarrow +\infty} \Lambda(n\mu, n\sigma^2) = \mathcal{LN}(e^{n\mu + n\frac{\sigma^2}{2}}, (e^{n\sigma^2} - 1)e^{2n\mu + n\sigma^2}).$$

Variations of this simple argument go a long way toward explaining the occurrence of lognormal distributions in a large class of natural and social phenomena [AB69].  $\square$

As a final illustration of the usefulness of Theorem 6.1, we now present a result describing the limiting distribution of a sequence of binominal distributions  $\text{Bin}(p_n; n)$ , with  $np_n$  approaching a positive constant, as a Poisson distribution.

**Theorem 6.13.** (*Law of Rare Events*) If  $X_n \sim \text{Bin}(p_n; n)$  and  $np_n \rightarrow \lambda > 0$  as  $n \rightarrow +\infty$  then there exists  $Z \sim \mathcal{P}(\lambda)$  such that

$$X_n \xrightarrow{d} Z.$$

*Proof.* We compute for any  $u \in \mathbb{R}$ :

$$\begin{aligned} \phi_{X_n}(u) &\stackrel{(2.24)}{=} (1 - p_n + p_n e^{iu})^n \\ &= \left(1 - \frac{\lambda}{n} + \frac{\lambda}{n} e^{iu} + o(n^{-1})\right)^n \\ &= \left(1 + \frac{\lambda}{n} (e^{iu} - 1) + o(n^{-1})\right)^n \\ &\rightarrow e^{\lambda(e^{iu} - 1)} \\ &\stackrel{(2.26)}{=} \phi_Z(u). \end{aligned}$$

Now apply Theorem 6.1.  $\square$

## 7. ESTIMATION

Here we shall use the theory developed so far to provide an introduction to Estimation Theory, an important topic in Statistics with countless applications.

**7.1. Parametric estimation and the mean squared error.** With the preliminary “large sample” results of Section 6 at hand, we now turn our attention to a *non-asymptotic* problem that appears very often in real world applications, where we only have access to *finitely* many measurements.

**Definition 7.1.** A *random sample* is a finite family  $\{X_j\}_{j=1}^n$  of i.i.d. random variables.

We usually represent a random sample by

$$X_1, \dots, X_n \sim \psi,$$

or simply by  $X_j \sim \psi$ , where  $\psi$  is the common pdf. Also, in the following we set  $\mathbb{E}(X_j) = \mu$  and  $\text{var}(X_j) = \sigma^2$ ,  $j = 1, \dots, n$ .

**Definition 7.2.** A (parametric) statistical model is a random sample

$$X_1, \dots, X_n \sim \psi_\theta,$$

where the associated pdf is allowed to depend on the unknown parameter  $\theta$  running in a given subset  $\Theta \subset \mathbb{R}^q$ .

**Remark 7.3.** Implicit in this definition is the existence of an underlying family of probability spaces, say  $(\Omega, \mathcal{F}, \{\mathcal{P}_\theta\}_{\theta \in \Theta})$ , so that  $\{X_j\}$  is i.i.d. with respect to each element in this family. Also, by Proposition 2.13 the joint pdf of  $(X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n$  is

$$(7.1) \quad \mathbf{x} = (x_1, \dots, x_n) \mapsto \psi_\theta(\mathbf{x}) := \prod_{j=1}^n \psi(x_j; \theta),$$

where  $\psi(x_j; \theta) = \psi_\theta(x_j)$ . □

**Remark 7.4.** Sometimes it is convenient to enlarge the scope of Definition 7.2 above in order to include samples  $X_j \sim \psi_\theta$  for which the “identically distributed” assumption no longer holds, so that only independence is retained. In the following, whenever we make use of this extended version of a statistical model, we will make explicit reference to this remark. □

Given the statistical model  $X_j \sim \psi_\theta$  as above, we will always assume that it is *identifiable* in the sense that the map  $\theta \mapsto \psi_\theta$  is injective. In any case, the corresponding *point estimator problem* consists of finding an *estimator*

$$(7.2) \quad \hat{\theta} = h(X_1, \dots, X_n)$$

for some *statistic*<sup>16</sup>  $h : \mathbb{R}^n \rightarrow \Theta \subset \mathbb{R}^q$ , which is supposed to yield an “efficient” guess of the true (and unknown) parameter  $\theta \in \Theta$ . The evaluation  $\hat{\theta}(\mathbf{x})$  of an estimator at a realization  $\mathbf{x} \in \mathbb{R}^n$  of a given random sample  $X = (X_1, \dots, X_n)$  is called an *estimate*.

The statistical analysis of point estimators falls naturally into two parts:

- First, one has to appraise the performance of a given estimator in comparison with competing ones so as to be able to select the “best” choice for the problem at hand. In this regard, it is usually assumed that the contending estimators are allowed to vary in a previously chosen family, although there is no guarantee that the “best” estimator in the family retains this property in general (see the performance analysis, under mean squared error, of the variance estimators  $\hat{\sigma}_c^2$ ,  $c > 0$ , in Subsection 7.2).
- Second, after selecting an estimator to work with, we should be aware that the true value of the relevant parameter never equals the corresponding point estimate, so we should develop tools to measure the “dispersion” of the random estimate around the true value in addition to just reporting a point estimate (this gives rise to the notion of a confidence interval discussed in Subsection 7.3 below).

We thus start here by considering the first issue.

<sup>16</sup>Just to get things right here, a statistic is any measurable function  $h = h(X_1, \dots, X_n)$ , whereas an estimator for  $\theta$ , denoted here by  $\hat{\theta}$ , is any statistic that does not depend on the unknown parameter  $\theta$ .

**Definition 7.5.** The *bias* of an estimator  $\hat{\theta}$  is given by

$$\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \theta).$$

An estimator  $\hat{\theta}$  is said to be *unbiased* if  $\text{bias}(\hat{\theta}) = 0$  (equivalently,  $\mathbb{E}(\hat{\theta}) = \theta$  for any  $\theta$ ). Also, the *mean squared error* (mse) of  $\hat{\theta}$  is

$$(7.3) \quad \text{mse}(\hat{\theta}) = \mathbb{E}(\|\hat{\theta} - \theta\|^2).$$

**Remark 7.6.** Strictly speaking, the dependence of the invariants above on  $\theta$  should be emphasized. For instance, the unbiasedness condition actually means that  $\mathbb{E}_{\mathcal{P}_\theta}(\hat{\theta}) = \theta$ , where

$$\mathbb{E}_{\mathcal{P}_\theta}(X) = \int_{\mathbb{R}^n} \mathbf{x} dP_\theta(\mathbf{x}) \stackrel{(7.1)}{=} \int_{\mathbb{R}^n} \mathbf{x} \psi_\theta(\mathbf{x}) d\mathbf{x},$$

where  $P_\theta = X_\# \mathcal{P}_\theta$  is the distribution of  $X$  coming from  $\mathcal{P}_\theta$  (see Remark 7.3). However, in order to keep the notation light, we usually refrain from doing so. Notice also that our notation ignores the dependence of  $\hat{\theta}$  on the size  $n$  of the random sample. Whenever emphasizing this is needed, we write  $\hat{\theta} = \hat{\theta}_n$ .  $\square$

**Proposition 7.7.** (*bias-variance trade-off*) There holds

$$(7.4) \quad \text{mse}(\hat{\theta}) = \text{tr cov}(\hat{\theta}) + \|\text{bias}(\hat{\theta})\|^2.$$

*Proof.* If  $\theta \in \mathbb{R}$  expand

$$\text{mse}(\hat{\theta}) = \mathbb{E}((\hat{\theta} - \mathbb{E}(\hat{\theta}) + \mathbb{E}(\hat{\theta}) - \theta)^2)$$

and check that the crossed terms cancel, thus yielding (7.5) below. The vector case then follows because

$$\begin{aligned} \text{mse}(\hat{\theta}) &= \sum_j \mathbb{E}((\hat{\theta}_j - \theta_j)^2) \\ &\stackrel{(7.5)}{=} \sum_j \text{var}(\hat{\theta}_j) + \sum_j |\text{bias}(\hat{\theta}_j)|^2 \\ &= \text{tr cov}(\hat{\theta}) + \|\text{bias}(\hat{\theta})\|^2. \end{aligned}$$

$\square$

**Convention 7.8.** Unless otherwise explicitly stated, we always assume in the sequel that  $\theta \in \Theta \subset \mathbb{R}$ , the *uni-dimensional* case, so that (7.4) reduces to

$$(7.5) \quad \text{mse}(\hat{\theta}) = \text{var}(\hat{\theta}) + |\text{bias}(\hat{\theta})|^2.$$

Here we adopt the viewpoint that the measure of the “performance” of an estimator is encoded in the “smallness” of the corresponding mse. In particular, a bound of the type  $\text{mse}(\hat{\theta}) \leq Cn^{-\alpha}$ ,  $\alpha > 0$ , immediately provides an  $O(n^{-\alpha/2})$  convergence rate estimate (in the mean) on how  $\hat{\theta}$  approaches  $\theta$  as  $n \rightarrow +\infty$ . Another kind of convergence of estimators appears in the next definition.

**Definition 7.9.** We say that an estimator  $\hat{\theta}_n = \hat{\theta}$  as above is *consistent* if  $\hat{\theta}_n \rightarrow \theta$  in probability (with respect to  $\theta$ ).



**Proposition 7.10.** *If  $\hat{\theta}_n$  is consistent with a uniformly bounded variance then  $\text{bias}(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow +\infty$  (thus,  $\hat{\theta}_n$  is asymptotically unbiased).*

*Proof.* By Proposition 2.22,  $\hat{\theta}_n \rightarrow \theta$  in distribution so that  $\mathbb{E}(\hat{\theta}_n) \rightarrow \mathbb{E}(\theta) = \theta$  and hence  $\mathbb{E}(\hat{\theta}_n)$  is uniformly bounded (for each  $\theta$ ). Combining this with the bound on the variance we see that  $\mathbb{E}(|\hat{\theta}_n|^2) \leq M_\theta$  for some  $M_\theta > 0$ . Now, for any  $\varepsilon > 0$  we have

$$\begin{aligned} |\mathbb{E}(\hat{\theta}_n - \theta)| &\leq |\mathbb{E}((\hat{\theta}_n - \theta)\mathbf{1}_{|\hat{\theta}_n - \theta| < \varepsilon})| + |\mathbb{E}((\hat{\theta}_n - \theta)\mathbf{1}_{|\hat{\theta}_n - \theta| \geq \varepsilon})| \\ &< \varepsilon + \mathbb{E}(|\hat{\theta}_n|\mathbf{1}_{|\hat{\theta}_n - \theta| \geq \varepsilon}) + \mathbb{E}(|\theta|\mathbf{1}_{|\hat{\theta}_n - \theta| \geq \varepsilon}) \\ &\leq \varepsilon + \sqrt{\mathbb{E}(|\hat{\theta}_n|^2)}\sqrt{\mathbb{E}(\mathbf{1}_{|\hat{\theta}_n - \theta| \geq \varepsilon})} + |\theta|\mathbb{E}(\mathbf{1}_{|\hat{\theta}_n - \theta| \geq \varepsilon}), \end{aligned}$$

where we used Cauchy-Schwarz in the last step. Thus,

$$|\text{bias}(\hat{\theta}_n)| < \varepsilon + \sqrt{M_\theta}\sqrt{P(|\hat{\theta}_n - \theta| \geq \varepsilon)} + |\theta|P(|\hat{\theta}_n - \theta| \geq \varepsilon)$$

and since  $P(|\hat{\theta}_n - \theta| \geq \varepsilon) \rightarrow 0$  the result follows.  $\square$

**Proposition 7.11.** *If  $\text{mse}(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow +\infty$  then  $\hat{\theta}_n$  is consistent.*

*Proof.* By Chebychev's inequality (2.15), for any  $\varepsilon > 0$ ,

$$P_\theta(|\hat{\theta}_n - \theta - \text{bias}(\hat{\theta}_n)| \geq \varepsilon) \leq \frac{\text{var}(\hat{\theta}_n - \theta)}{\varepsilon^2} \rightarrow 0,$$

which means that  $\hat{\theta}_n - \text{bias}(\hat{\theta}_n) \rightarrow \theta$  in probability (with respect to  $\theta$ ). Since  $\text{bias}(\hat{\theta}_n) \rightarrow 0$  as well, Theorem 2.23 applies to ensure that  $\hat{\theta}_n \rightarrow \theta$  in probability.  $\square$

**Definition 7.12.** An estimator  $\hat{\theta}_n$  as above is *asymptotically normal* with *asymptotic variance*  $\sigma_\theta^2 > 0$ ,  $\theta \in \Theta$ , if there exists  $Z_\theta \sim \mathcal{N}(0, \sigma_\theta^2)$  such that  $\sqrt{n}(\hat{\theta}_n - \theta) \rightarrow Z_\theta$  in distribution (with respect to  $\theta$ ).

**Proposition 7.13.** *If  $\hat{\theta}_n$  is asymptotically normal then it is consistent.*

*Proof.* If  $Z_{\theta,n} = Z_\theta/\sqrt{n} \sim \mathcal{N}(0, \sigma_\theta^2/n)$  then

$$\hat{\theta}_n - \theta - Z_{\theta,n} = \frac{1}{\sqrt{n}} \left( \sqrt{n}(\hat{\theta}_n - \theta - Z_{\theta,n}) \right) \xrightarrow{P} 0.$$

But Chebychev's inequality gives, for any  $\varepsilon > 0$ ,

$$P(|Z_{\theta,n}| \geq \varepsilon) \leq \frac{\sigma_\theta^2}{n\varepsilon^2} \rightarrow 0,$$

that is,  $Z_{\theta,n} \xrightarrow{P} 0$  and the result follows by Theorem 2.23.  $\square$

**Remark 7.14.** The true nature of the asymptotic variance  $\sigma_\theta^2$  has not been explored in the previous argument, which makes sense because consistence only pertains to position (not to dispersion). But notice that asymptotically normality clearly implies that  $\text{var}(\hat{\theta}_n)$  is uniformly bounded (for each  $\theta$ ) so  $\hat{\theta}_n$  is asymptotically unbiased by Proposition 7.10. Of course, this also follows directly from the definition: one has  $\sqrt{n} \mathbb{E}(\hat{\theta}_n - \theta) \rightarrow \mathbb{E}(Z_\theta) = 0$ .  $\square$

We include here an useful consequence of asymptotic normality.

**Proposition 7.15.** (*the delta method*) If  $\hat{\theta}_n$  is asymptotically normal (as in Definition 7.12) and  $g : \Theta \subset \mathbb{R} \rightarrow \mathbb{R}$  is a  $C^1$  function whose derivative vanishes nowhere then  $\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) \rightarrow Z_{\theta,g}$  in distribution, where  $Z_{\theta,g} \sim \mathcal{N}(0, |g'(\theta)|^2 \sigma_\theta^2)$ .

*Proof.* By Taylor,

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) = g'(\tilde{\theta}_n) \sqrt{n}(\hat{\theta}_n - \theta),$$

for some  $\tilde{\theta}_n$  between  $\hat{\theta}_n$  and  $\theta$ . Since  $\hat{\theta}_n \xrightarrow{P} \theta$  by Proposition 7.13, it is not hard to check that  $g'(\tilde{\theta}_n) \xrightarrow{P} g'(\theta)$ , so the result follows from Theorem 2.23.  $\square$

**Example 7.16.** For any random sample  $\{X_j\}$  as above it is immediate to check that the *sample mean*

$$(7.6) \quad \bar{X}_n := \frac{1}{n} (X_1 + \cdots + X_n)$$

is an unbiased estimator for the expected value of the underlying distribution. In other words,  $\mathbb{E}(\bar{X}_n) = \mu$ , where  $\mu = \mathbb{E}(X_j)$  is the common expectation. Also, if  $\sigma^2 = \text{var}(X_j)$  is the common variance of the sample (the population variance) then it follows from (2.12) that

$$\text{var}(\bar{X}_n) = \frac{1}{n^2} \sum_{j=1}^n \text{var}(X_j) = \frac{1}{n^2} \sum_{j=1}^n \sigma^2,$$

that is,

$$(7.7) \quad \text{var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

and hence  $\text{mse}(\bar{X}_n) = \sigma^2/n$ . This is the reason why we call

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

the *standardization* of the sample mean; compare with (6.3). Note that  $\mathbb{E}(Z_n) = 0$  and  $\text{var}(Z_n) = 1$ . Finally, note that  $\bar{X}_n$  is consistent (as an estimator for  $\mu$ ) either by LLN or by Proposition 7.11 and that for any  $g$  as in Proposition 7.15, CLT applies to ensure that  $g(\bar{X}_n)$ , as an estimator of  $g(\mu)$ , satisfies

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

hence being asymptotically normal as well.  $\square$

**Example 7.17.** (weighted estimators for the population mean) If  $w = (w_1, \dots, w_n)$  is a *weight vector* (which means that  $\sum_j w_j = 1$ ), then we may consider the corresponding *weighted estimator* for  $\mu$  given by

$$\bar{X}_n^w = \sum_{j=1}^n w_j X_j,$$

which includes (7.6) as a rather special case. One easily verifies that  $\text{bias}(\bar{X}_n^w) = 0$  and  $\text{mse}(\bar{X}_n^w) = \text{var}(\bar{X}_n^w) = \sigma^2 \sum_j w_j^2$ , so the estimator with the least mse in this class is obtained by minimizing  $w \mapsto |w|^2$  under the constraint  $\sum_j w_j = 1$ , which gives  $w = (1/n, \dots, 1/n)$ , corresponding to the sample mean  $\bar{X}_n$ .  $\square$

**Example 7.18.** (Monte Carlo estimator) Let  $X : \Omega \rightarrow \mathbb{R}^m$  be a random vector with a pdf  $\psi$  whose support is contained in the unit cube  $[0, 1]^m$  and let  $f : [0, 1]^m \rightarrow \mathbb{R}$  be such that  $f\psi$  is Lebesgue integrable. If  $X_j \sim \psi$  is a i.i.d. sample it is immediate from Example 7.16 that the *Monte Carlo estimator*

$$\hat{\mu}_{(n)}^f := \frac{1}{n} \sum_{j=1}^n f(X_j),$$

is an unbiased estimator for the unknown parameter

$$(7.8) \quad \mu^f := \mathbb{E}(f(X_j)) = \int_{[0,1]^m} f(x)\psi(x)dx$$

which is consistent because

$$(7.9) \quad \lim_{n \rightarrow +\infty} \hat{\mu}_{(n)}^f = \mu^f$$

in probability by LLN. As usual, this also follows from Proposition 7.11, given that

$$(7.10) \quad \text{mse}(\hat{\mu}_{(n)}^f) = \text{var}(\hat{\mu}_{(n)}^f) = \frac{\sigma_f^2}{n},$$

where  $\sigma_f^2$  is the (common) variance of  $f(X_j)$ . Notice that this conveys a  $O(n^{-1/2})$  convergence rate for (7.9) which can be made explicit if we apply Chebyshev's inequality (2.15) with  $X = \hat{\mu}_{(n)}^f$ ,  $\sigma = \sigma_f/\sqrt{n}$  and  $c = 1/\sqrt{\delta}$ ,  $\delta > 0$ , so that

$$P\left(\left|\hat{\mu}_{(n)}^f - \mu^f\right| \leq \frac{1}{\sqrt{\delta}} \frac{\sigma_f}{\sqrt{n}}\right) \geq 1 - \delta.$$

Thus, one needs at least

$$(7.11) \quad n \approx \frac{1}{\delta} \frac{\sigma_f^2}{\varepsilon^2}$$

samples in order to obtain a dispersion of at most  $\varepsilon$  of the estimator around the expected value  $\mu^f$  with probability at least  $1 - \delta$ . This may be substantially improved if we appeal to CLT (Theorem 6.5) to obtain

$$\lim_{n \rightarrow +\infty} P\left(\left|\hat{\mu}_{(n)}^f - \mu^f\right| \leq \eta \frac{\sigma_f}{\sqrt{n}}\right) = \frac{1}{\sqrt{2\pi}} \int_{-\eta}^{\eta} e^{-x^2/2} dx \stackrel{(5.2)}{\approx} 1 - e^{-\eta^2/2}, \quad \eta \rightarrow 0,$$

which allows us to replace the previous estimate by

$$(7.12) \quad n \approx 2 \ln(1/\delta) \frac{\sigma_f^2}{\varepsilon^2}.$$

One should be aware, however, that whereas the estimate (7.11) holds non-asymptotically, as it comes from Chebyshev's inequality, the estimate (7.12) becomes reliable only in the asymptotic regime ( $n \rightarrow +\infty$ ); see Remarks 6.3 and 6.7. Regardless of the shape of their dependence on  $\delta$ , the estimates above share the nice property of not depending on  $m$ , so that the “dimensionality curse” is not present here. Of course this is one of the reasons why Monte Carlo methods, based on (7.9), are quite versatile in approximating multiple integrals as those in the right-hand side of (7.8)<sup>17</sup>. Besides its slow  $O(n^{-1/2})$  convergence rate, an obvious drawback of this method is its explicit dependence on the standard deviation  $\sigma_f$ , which is at least as hard to compute as  $\mu^f$  itself. The simplest choices avoiding this latter problem (by explicitly bounding the variance) corresponds

<sup>17</sup>This method should be compared with the usual numerical approach which requires evaluation of  $f\psi$  on a  $\varepsilon$ -net and hence has a complexity that grows like  $\varepsilon^{-m}$ .

to taking  $\{X_j\}$  *uniformly* distributed in  $[0, 1]^m$  (so that  $\psi = \mathbf{1}_{[0,1]^m}$ ), and  $f = \mathbf{1}_B$ , the indicator of a Borel set  $B \subset [0, 1]^m$ , so that  $\mu^f = \text{vol}_m(B)$ , the  $m$ -volume of  $M$ , and  $\sigma_f^2 = \text{vol}_m(B)(1 - \text{vol}_m(B)) \leq 1/4$ . Thus, at least the volumes of (well-behaved) Borel subsets can be efficiently calculated if we are able to provide low cost simulations of independent, uniformly distributed random variables on  $[0, 1]^m$  [RCC10].  $\square$

**7.2. Computing the mean squared error of  $\hat{\sigma}_c^2$ .** We further illustrate the concepts introduced in Definition 7.5 by determining the “best” estimator for the variance  $\sigma^2$  in the family

$$(7.13) \quad \hat{\sigma}_c^2 = h_c(X_1, \dots, X_n), \quad c > 0,$$

where

$$h_c(X_1, \dots, X_n) = c \sum_{j=1}^n (X_j - \bar{X}_n)^2$$

and

$$\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n)$$

is the sample mean; see Example 7.16. Thus,  $\theta = \sigma^2 > 0$  and  $\Theta = \mathbb{R}_+$ . This involves minimizing the corresponding mean squared error  $\text{mse}(\hat{\sigma}_c^2)$  viewed as a function of  $c$ .

**Proposition 7.19.** *One has*

$$(7.14) \quad \text{bias}(\hat{\sigma}_c^2) = (c(n-1) - 1) \sigma^2.$$

*Proof.* Using that

$$(7.15) \quad X_j - \bar{X}_n = \frac{n-1}{n} X_j - \frac{1}{n} \sum_{k \neq j} X_k$$

we first note that  $\mathbb{E}(X_j - \bar{X}_n) = 0$  and hence

$$\mathbb{E}(\hat{\sigma}_c^2) = c \sum_{i=1}^n \mathbb{E}((X_j - \bar{X}_n)^2) = c \sum_{j=1}^n \text{var}(X_j - \bar{X}_n).$$

From the independence assumption and (2.12) we get, again using (7.15),

$$\text{var}(X_j - \bar{X}_n) = \frac{(n-1)^2}{n^2} \text{var}(X_j) + \frac{1}{n^2} \sum_{k \neq j} \text{var}(X_k) = \frac{n-1}{n} \sigma^2,$$

so that

$$(7.16) \quad \mathbb{E}(\hat{\sigma}_c^2) = cn \frac{n-1}{n} \sigma^2 = c(n-1) \sigma^2,$$

and the result follows.  $\square$

**Corollary 7.20.**  $\hat{\sigma}_c^2$  is unbiased only if  $c = (n-1)^{-1}$ .

**Definition 7.21.** The *sample variance* of the sample  $\{X_j\}$  as above is

$$S_n^2 := \hat{\sigma}_{(n-1)^{-1}}^2 = \frac{1}{n-1} \sum_{j=1}^n (X_j - \bar{X}_n)^2.$$

Also,  $S_n = \sqrt{S_n^2}$  is the *sample standard deviation*.

Recall that in general an unbiased estimator  $\hat{\theta}$  satisfies  $\mathbb{E}(\hat{\theta}) = \theta$ , so that the target parameter  $\theta$  is the expected value of the corresponding sample distribution. Intuitively, the unbiasedness property says that on average the estimator hits the right target. This is the main reason why unbiased estimators are often used in applications and we provide below two classical results (Theorems 9.7 and 8.17) ensuring that unbiased estimators minimize their variance (and hence their mse) within certain classes of competing unbiased estimators. However, we point out that the family  $\hat{\sigma}_c^2$  above may be used to illustrate that in general the best variance estimator might not be unbiased (that is, the function  $c \mapsto \text{mse}(\hat{\sigma}_c^2)$  is minimized for some  $c \neq (n-1)^{-1}$ ), which turns out to be a manifestation of the variance-bias trade-off in (7.5). For this we need to compute  $\text{var}(\hat{\sigma}_c^2)$ , which we do by assuming in the rest of the calculation that each  $X_j$  is normally distributed:  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ . Let us set

$$U^2 = \sum_{j=1}^n \left( \frac{X_j - \mu}{\sigma} \right)^2, \quad V^2 = n \left( \frac{\bar{X}_n - \mu}{\sigma} \right)^2.$$

**Proposition 7.22.** *If  $\{X_j\}$  is independent with  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  then  $U^2 \sim \chi_n^2$  and  $V^2 \sim \chi_1^2$ .*

*Proof.* Note that  $\sigma^{-1}(X_j - \mu) \sim \mathcal{N}(0, 1)$  by Proposition 4.7, which can also be used to check that  $V^2 = W^2$ , where

$$W = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1).$$

The results then follow from Proposition 4.25. □

**Proposition 7.23.** *If  $\{X_j\}$  is independent and  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  then  $\bar{X}_n$  and*

$$S^2 := \sum_j (X_j - \bar{X}_n)^2$$

*are independent. In particular,  $\bar{X}_n \perp \hat{\sigma}_c^2$ ,  $c > 0$ .*

*Proof.* We may assume that  $\mu = 0$ . Note that

$$S^2 = \sum_j X_j^2 - Y_1^2,$$

where

$$(7.17) \quad Y_1 = \sqrt{n} \bar{X}_n = \sum_j \frac{X_j}{\sqrt{n}}$$

is normal. By Gramm-Schmidt there exists an orthogonal  $n \times n$  matrix, say  $O$ , whose first line is the vector  $(1/\sqrt{n}, \dots, 1/\sqrt{n})$  i.e.  $O_{1j} = 1/\sqrt{n}$ . If  $Y = OX$  then  $Y_1$  is indeed given by (7.17) so that

$$S^2 = \|X\|^2 - Y_1^2 = \|Y\|^2 - Y_1^2 = \sum_{l=2}^n Y_l^2.$$

From Corollary 4.12,  $Y' = (Y_2, \dots, Y_n)$  is normally distributed and  $\{Y_l\}_{l=2}^n$  is independent as well so its covariance matrix is diagonal. Moreover, using again the independence of  $\{X_j\}$ , we compute for  $l \geq 2$  that

$$\begin{aligned} \text{cov}(Y_1, Y_l) &= \text{cov} \left( \sum_j O_{1j} X_j, \sum_k O_{lk} X_k \right) \\ &= \sigma^2 \sum_j O_{1j} O_{lj} \\ &= 0, \end{aligned}$$

so that  $\{Y_j\}_{j=1}^n$  is independent by Proposition 4.10. In particular,  $\mathcal{S}^2 = \|Y'\|^2 \perp Y_1/\sqrt{n} = \bar{X}_n$ , as desired.  $\square$

**Proposition 7.24.** *If*

$$(7.18) \quad \Sigma^2 := \sigma^{-2} \sum_j (X_j - \bar{X}_n)^2$$

*then*

$$(7.19) \quad \Sigma^2 \sim \chi_{n-1}^2.$$

*In particular,*

$$(7.20) \quad \text{var}(\Sigma^2) = 2(n-1).$$

*Proof.* Upon multiplication by  $\sigma^{-2}$ , the elementary algebraic identity

$$(7.21) \quad \sum_{j=1}^n (X_j - \mu)^2 = \sum_{j=1}^n (X_j - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2$$

becomes

$$(7.22) \quad U^2 = \Sigma^2 + V^2.$$

Since  $\{\Sigma^2, V^2\}$  is independent (by Proposition 7.23) we have  $\phi_{U^2} = \phi_{\Sigma^2} \phi_{V^2}$  so that Corollary 4.22 applies to give

$$\phi_{\Sigma^2}(u) = (1 - 2u\mathbf{i})^{-(n-1)/2},$$

which yields (7.19). Finally, (7.20) follows from Corollary 4.21.  $\square$

**Remark 7.25.** The parameter  $\mu$  plays no essential role in the validity of the identity (7.21), the only relevant point being that  $\bar{X}_n$  is the arithmetic mean of  $\{X_j\}_{j=1}^n$ . Thus, (7.21) remains true if  $\mu$  gets replaced by any real number:

$$(7.23) \quad \sum_{j=1}^n (\eta_j - c)^2 = \sum_{j=1}^n (\eta_j - \bar{\eta}_n)^2 + n(\bar{\eta}_n - c)^2, \quad c \in \mathbb{R},$$

where

$$\bar{\eta}_n = \frac{\eta_1 + \dots + \eta_n}{n}.$$

In this more general form, this important identity resurfaces at many points below (see Remark 7.31, and Examples 7.39 and 7.40).  $\square$

We may record the result of our computation as follows.

**Proposition 7.26.** If  $\{X_j\}$  is independent with  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  then

$$(7.24) \quad \text{var}(\hat{\sigma}_c^2) = 2(n-1)c^2\sigma^4.$$

As a consequence,

$$(7.25) \quad \text{mse}(\hat{\sigma}_c^2) = ((n-1)(n+1)c^2 - 2(n-1)c + 1)\sigma^4.$$

*Proof.* Combine (7.5), (7.14) and (7.20).  $\square$

**Corollary 7.27.** Under the conditions above,  $\text{mse}(\hat{\sigma}_c^2)$  is minimized for  $c = (n+1)^{-1}$ .

**Remark 7.28.** Since

$$\text{bias}\hat{\sigma}_{(n+1)^{-1}}^2 = -2(n+1)^{-1}\sigma^2,$$

which only vanishes in the asymptotic limit  $n \rightarrow +\infty$ , as already advertised Corollary 7.27 illustrates that an unbiased estimator may fail to be the most efficient one (if the “performance” is measured by mse); a quite similar phenomenon, involving the so-called James-Stein estimator for the mean of certain normal random vectors, appears in Example 8.21.  $\square$

**Remark 7.29.** (Studentized mean) If  $Z \sim \mathcal{N}(0, 1)$  and  $W \sim \chi_k^2$  then Proposition 4.30 says that

$$\frac{Z}{\sqrt{W/k}} \sim t_k,$$

the Student’s  $t$ -distribution with  $k \geq 1$  degrees of freedom [Stu08b, Fis25]. In the setting of Proposition 7.26 (that is, under sample normality) we may apply this to  $Z = \sqrt{n}(\bar{X}_n - \mu)/\sigma$  (which is  $\mathcal{N}(0, 1)$  by Proposition 4.7) and  $W = \hat{\sigma}_{\sigma^{-2}}^2$  (just use (7.19)) to conclude that

$$(7.26) \quad T_{n-1} := \frac{\bar{X}_n - \mu}{S_n/\sqrt{n}} \sim t_{n-1},$$

where  $S_n$  is the sample standard deviation (Definition 7.21). Notice that  $Z \perp W$  here by Proposition 7.23<sup>18</sup>. We say that  $T_{n-1}$  is the *studentized mean* of the normal sample  $\{X_j\}$ . It plays a key role in finding “small sample” estimates for the population mean of a normally distributed sample with no prior knowledge of the population variance; see Subsection 7.3 below.  $\square$

**Remark 7.30.** If the random sample  $\{X_j\}_{j=1}^n$  is not necessarily normal then a somewhat tedious computation gives

$$(7.27) \quad \text{var}(\hat{\sigma}_c^2) = \frac{(n-1)^2}{n}c^2\sigma^4 \left( \kappa(X_j) - \frac{n-3}{n-1} \right),$$

where

$$\kappa(X) = \frac{\mathbb{E}((X - \mathbb{E}(X))^4)}{\text{var}(X)^2}$$

is the *kurtosis* of  $X$ , which is finite if we require that  $\mathbb{E}(|X|^4) < +\infty$ ; see [ONe14] for this and many other moment computations. In the normal case we have  $\kappa(X_j) = 3$  and (7.27) reduces to (7.24). Of course, this general computation suffices if we are merely interested in the conclusions of Proposition 7.26, but we stress that the elegant argument above based on sample normality has the added bonus of yielding an explicit expression for the sampling distribution of the studentized mean  $T_{n-1}$  considered in Remark 7.29.  $\square$

<sup>18</sup>This independence between the sample mean  $\bar{X}_n$  and the sample standard deviation  $S_n$ , which is crucial in precisely determining the shape of the sampling distribution of  $T_{n-1}$ , turns out to be a characteristic feature of normal samples; see [Luk42] for a proof which is a clever application of Proposition 2.30 (3).

**Remark 7.31.** (The geometric way to Student) The calculation leading to Proposition 7.24, in particular the independence between the sample mean  $\bar{X}_n$  and the sample variance  $S_n^2$  in Proposition 7.23 which plays a central role in accessing Student's distribution in (7.26) above, may be retrieved by means of the “ $n$ -space computations” due to R. Fisher already mentioned in Remarks 4.26 and 4.31<sup>19</sup>. Indeed, the probability density spanned by the independent normal random vector  $X = (X_1, \dots, X_n)$ ,  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ , in an infinitesimal region of volume  $dx = dx_1 \cdots dx_n$  is

$$(7.28) \quad \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-\frac{\sum_j (x_j - \mu)^2}{2\sigma^2}} dx = \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2}} e^{-\frac{(n-1)s_n^2}{2\sigma^2}} dx,$$

where Remark 7.25 has been used. Following [Fis25] we now observe that  $\bar{x}_n$  is proportional to the height of  $x$  with respect to the hyperplane  $H^{n-1}$  defined by  $\sum_j x_j = 0$ , whereas  $s_n$  is proportional to the distance of  $x$  to the line  $l^1$  given by  $x_1 = \dots = x_n$ , with  $\mathcal{X}_n := (\bar{x}_n, \dots, \bar{x}_n) \in l^1$  realizing this distance. Since  $H^{n-1}$  and  $l^1$  are perpendicular to each other, we may use the corresponding “cylindrical” coordinate system to check that  $dx$  is proportional to  $s_n^{n-2} d\bar{x}_n ds_n d\theta$ , where  $d\theta$  is the volume element of the unit sphere  $\mathbb{S}^{n-2} \subset H^{n-1}$  with center located at  $H^{n-1} \cap l^1$ , the origin of  $H^{n-1}$ . Leading this to (7.28), integrating with respect to  $\theta$  and using Proposition 2.12 shows that the infinitesimal joint probability density of the random vector  $(\bar{X}_n, S_n^2)$  is

$$(7.29) \quad \psi_{(\bar{X}_n, S_n^2)}(\bar{x}_n, s_n^2) d\bar{x}_n ds_n^2 \approx e^{-\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2}} d\bar{x}_n \times e^{-\frac{(n-1)s_n^2}{2\sigma^2}} (s_n^2)^{\frac{n-3}{2}} ds_n^2,$$

where  $\approx$  here means that we are neglecting certain normalizing constants which will take care of themselves. Incidentally, this geometric argument makes it obvious the connection to the previous computational proof of Proposition 7.23: the orthogonal map  $Y = OX$  used there carries  $H^{n-1}$  onto the coordinate hyperplane  $Y_1 = 0$ , which has the net effect of reducing the size of the sample data by one, thus allowing for an induction argument based on Corollary 4.12. Moreover, it has a number of consequences which we now describe.

- Clearly, (7.29) implies that  $\{\bar{X}_n, S_n^2\}$  is independent.
- Also, it follows from (7.29) that

$$\psi_{S_n^2}(s_n^2) ds_n^2 \approx e^{-\frac{(n-1)s_n^2}{2\sigma^2}} (s_n^2)^{\frac{n-3}{2}} ds_n^2,$$

so if we combine this with (7.18) we see that

$$(7.30) \quad \psi_{\Sigma^2}(s_\sigma^2) ds_\sigma^2 \approx e^{-s_\sigma^2/2} (s_\sigma^2)^{\frac{n-3}{2}} ds_\sigma^2, \quad s_\sigma^2 = \frac{(n-1)s_n^2}{\sigma^2},$$

from which we easily deduce (7.19); compare with the computation in Remark 4.26.

- Finally, the geometric argument also provides another way of explicitly computing the probability density of the studentized mean in (7.26). Indeed, (7.29) implies that

$$\psi_{\bar{X}_n}(\bar{x}_n) d\bar{x}_n \approx e^{-\frac{n(\bar{x}_n - \mu)^2}{2\sigma^2}} d\bar{x}_n \approx (s_n^2)^{1/2} e^{-\frac{s_n^2 t_{n-1}^2}{2\sigma^2}} dt_{n-1},$$

where for  $s_n^2$  fixed we set

$$t_{n-1} = \frac{\bar{x}_n - \mu}{s_n/\sqrt{n}}.$$

Transplanting this to (7.29) yields an explicit expression for the joint density  $\psi_{(T_{n-1}, S_n^2)}(t_{n-1}, s_n^2) dt_{n-1} ds_n^2$ , so that integration with respect to  $s_n^2$  gives

$$\begin{aligned} \psi_{T_{n-1}}(t_{n-1}) dt_{n-1} &\approx \left( \int_0^{+\infty} (s_n^2)^{\frac{n-2}{2}} e^{-\frac{s_n^2(n-1+t_{n-1}^2)}{2\sigma^2}} ds_n^2 \right) dt_{n-1} \\ &\approx (n-1+t_{n-1}^2)^{-n/2} dt_{n-1} \\ &\approx t_{n-1}(t_{n-1}) dt_{n-1}, \end{aligned}$$

<sup>19</sup>Recall that we always represent a realization of a random variable, say  $\bar{X}_n$ , by the corresponding lower-case symbol, in this case  $\bar{x}_n$ .



as desired.  $\square$

**7.3. Confidence intervals.** If  $\hat{\theta}$  is an *unbiased* estimator for the parameter  $\theta$  whose standard deviation  $\sigma_{\hat{\theta}}$  is known then Chebyshev's inequality (2.15) gives

$$P(|\hat{\theta} - \theta| \leq c\sigma_{\hat{\theta}}) \geq 1 - c^{-2}, \quad c > 1,$$

which translates into a “confidence interval” estimate for the unknown parameter:

$$(7.31) \quad \theta \in [\hat{\theta} \mp c\sigma_{\hat{\theta}}] \text{ with prob. at least } 1 - c^{-2},$$

where here and in the following we denote the interval  $[a - b, a + b]$  simply by  $[a \mp b]$  whenever convenient. In particular, this applies to  $\hat{\theta} = \bar{X}_n$ , the mean sample estimator, which is an unbiased estimator for the population mean  $\mu$  (by Example 7.16). Since  $\sigma_{\bar{X}_n} = \sigma/\sqrt{n}$ , if we take  $c^{-2} \approx \delta$  and  $c\sigma/\sqrt{n} \approx \epsilon$ , where  $\delta$  and  $\epsilon$  are arbitrarily small positive real numbers, then we see that

$$(7.32) \quad n > \frac{\sigma^2}{\epsilon^2 \delta} \implies \mu \in [\bar{X}_n \mp \epsilon] \text{ with prob. } \approx 1 - \delta.$$

Notice that this retrieves the “convergence in probability” version of LLN (under the additional assumption that  $\sigma$  is finite); see Remarks 6.3 and 6.7. We may also turn this into a confidence interval estimate as in (7.31):

$$(7.33) \quad \mu \in \left[ \bar{X}_n \mp \frac{1}{\sqrt{\delta}} \frac{\sigma}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta,$$

If we further require that the sample is normally distributed ( $X_j \sim \mathcal{N}(\mu, \sigma^2)$ ) then we can employ

$$(7.34) \quad Z_n := \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1),$$

where we used Proposition 4.7, to find a “small sample” confidence interval for the unknown expected value:

$$(7.35) \quad \mu \in \left[ \bar{X}_n \mp z_{\delta/2} \frac{\sigma}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta,$$

where, for a given  $\beta > 0$ , the *quantile*  $z_{\beta} > 0$  is determined by  $P(Z \geq z_{\beta}) = P(Z \leq -z_{\beta}) = \beta$ , the “tail” probability associated to  $Z \sim \mathcal{N}(0, 1)$ <sup>20</sup>. Notice that if we (more realistically!) relax the normality assumption then (7.35) becomes a “large sample” estimate since (7.34) holds asymptotically as  $n \rightarrow +\infty$  due to CLT. Upon comparison with (7.35) we see that this amounts to replacing  $1/\sqrt{\delta}$  by  $z_{\delta/2}$  in the estimate for the dispersion around the sample mean<sup>21</sup>. In any case, the estimates (7.33) and (7.35) remain ineffective as long as  $\sigma$  is unknown, in which case (and coming back to a normal sample  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ ), Remark 7.29 suggests replacing (7.34) by (7.26) so as to obtain the “small sample” estimate

$$(7.36) \quad \mu \in \left[ \bar{X}_n \mp t_{n-1, \delta/2} \frac{S_n}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta,$$

where  $P(T_{n-1} \geq t_{n-1, \beta}) = \beta$  is the “tail” probability associated to the t-distribution  $t_{n-1}$  which defines the corresponding quantile  $t_{n-1, \beta}$ . The obvious advantage of (7.36) over (7.35) is that no previous knowledge of  $\sigma$  is used. Finally, note that in Remark 8.26 below it is shown that  $\hat{\sigma}_{n-1}^2 \xrightarrow{p} \sigma^2$  under this normality assumption. Since

$$\hat{\sigma}_{(n-1)-1}^2 = \frac{n}{n-1} \hat{\sigma}_{n-1}^2,$$

<sup>20</sup>For instance, if  $\delta = 0.05$  then  $z_{\delta/2} \approx 1.96$ , which shows that roughly two standard deviations around the normal mean suffice to ensure the customary 95% confidence statement.

<sup>21</sup>Thus, if  $\delta = 0.05$  we are replacing  $1/\sqrt{0.05} \approx 4.47$  by 1.96, which shrinks the dispersion by a factor of  $4.47/1.96 \approx 2.28$  while still retaining the same confidence level. But recall that this reduction only becomes reliable in the asymptotic regime (Remarks 6.3 and 6.7).

we see from Theorem 2.23 that

$$(7.37) \quad S_n = \sqrt{\hat{\sigma}_{(n-1)^{-1}}^2} \xrightarrow{P} \sigma.$$

It then follows from

$$T_{n-1} = \frac{\sigma}{S_n} Z_n$$

and Theorem 2.23 that (7.35) and (7.36) provide essentially the same information in this asymptotic regime (in the sense that  $T_{n-1} - Z_n \xrightarrow{P} 0$ ). We stress, however, the usefulness of (7.36) when dealing with small samples, which attests in favor of Student's fundamental contribution coming from Remark 7.29.

**Remark 7.32.** The convergence in (7.37) holds more generally (that is, with no normality assumption) if we assume that the random sample satisfies  $\mathbb{E}(|X_j|^4) < +\infty$ . Indeed, we already know from Corollary 7.20 that  $\text{bias}(S_n^2) = 0$ . Also, from (7.27) with  $c = (n-1)^{-1}$  we see that  $\text{var}(S_n^2) \rightarrow 0$ . Thus,  $\text{mse}(S_n^2) \rightarrow 0$  and Proposition 7.11 applies to ensure that  $S_n^2 \xrightarrow{P} \sigma^2$ , from which (7.37) follows.  $\square$

**Remark 7.33.** It is important to have in mind that if we evaluate the sample mean  $\bar{X}_n$  through a measurement so as to obtain a numerical value, say  $\mu_n$ , then the corresponding realization of (7.35), namely,

$$(7.38) \quad \mu \in \left[ \mu_n \mp z_{\delta/2} \frac{\sigma}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta,$$

is completely devoid of sense. Indeed, since any trace of randomness has been removed from the interval in (7.38) (it has now become deterministic!) then either  $\mu$  definitely belongs to it or not, with probability 0 or 1. Thus, the proper way to interpret (7.35) is to regard the corresponding interval as stochastic in nature and to adopt the “frequentist” perspective according to which the “relative frequency” that

$$\mu \in \left[ \mu_n - z_{\delta/2} \frac{\sigma}{\sqrt{n}}, \mu_n + z_{\delta/2} \frac{\sigma}{\sqrt{n}} \right]$$

approaches  $1 - \delta$  as the number of successive realizations  $\mu_n$  of  $\bar{X}_n$  becomes larger and larger. Needless to say, similar remarks hold for (7.36) under realizations of  $\bar{X}_n$  and  $S_n$ .  $\square$

**Remark 7.34.** Strictly speaking,  $Z_n$  and  $T_{n-1}$  do not qualify as estimators as they are statistics which depend on the underlying parameters. Instead, they are referred to as *pivotal quantities*, a terminology incorporating the appreciated property that their distributions do not depend on these parameters.  $\square$

**Example 7.35.** (Sampling from a Bernoulli population) If  $X_j \sim \text{Ber}(p)$  as in Remark 6.8 then the computation leading to (6.8) also gives

$$(7.39) \quad \sqrt{n} \frac{\bar{X}_n - p}{\sqrt{p(1-p)}} \xrightarrow{d} \mathcal{N}(0, 1),$$

where  $\bar{X}_n = X^{(n)}/n$  is the corresponding sample mean. This translates into the “large sample” estimate

$$p \in \left[ \bar{X}_n \mp z_{\delta/2} \frac{\sqrt{p(1-p)}}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta,$$

which displays the usual drawback, namely, the size of the confidence interval depends on  $p$ , the parameter we want to estimate. The conservative way of remedying this is to implement the rather crude estimate  $p(1-p) \leq 1/4$  to eliminate the dependence on  $p$ . Another, certainly more effective, route consists of combining Theorem 2.23 and LLN to replace (7.39) by

$$\sqrt{n} \frac{\bar{X}_n - p}{\sqrt{\bar{X}_n(1-\bar{X}_n)}} \xrightarrow{d} \mathcal{N}(0, 1),$$

which gives

$$p \in \left[ \bar{X}_n \mp z_{\delta/2} \frac{\sqrt{\bar{X}_n(1 - \bar{X}_n)}}{\sqrt{n}} \right] \text{ with prob. } \approx 1 - \delta.$$

Since  $\sqrt{\bar{X}_n(1 - \bar{X}_n)} \leq 1/2$ , we see that an increase of order  $\gamma > 0$  on the sample size only allows for a decrease of order  $\gamma^{-1/2}$  on the dispersion around  $\bar{X}_n$ , the center of the confidence interval. Applications of this elementary observation abound since Bernoulli populations are quite useful in modeling the possible outcome of any elementary experiment that asks a yes-no question (coin flipping, election poll with two contenders, male-female birth rates, among others.)  $\square$

**Example 7.36.** (The difference of means of two normal samples) Let  $\{X_j\}_{j=1}^m$  and  $\{Y_k\}_{k=1}^n$  be normally distributed random samples, say with  $X_j \sim \mathcal{N}(\mu_X, \sigma_X^2)$  and  $Y_k \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$ , which we assume to be independent to one another. In general, we also assume that the true parameter  $\theta = (\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2)$  is unknown. In order to estimate the difference of means  $\mu := \mu_X - \mu_Y$  we first note that

$$\mathbb{E}(\bar{X}_m - \bar{Y}_n) = \mu, \quad \text{var}(\bar{X}_m - \bar{Y}_n) = \frac{\sigma_X^2}{m} + \frac{\sigma_Y^2}{n},$$

and hence

$$Z_{X,Y} := \frac{\bar{D}_{mn} - \mu}{\sqrt{\frac{\sigma_X^2}{m} + \frac{\sigma_Y^2}{n}}} \sim \mathcal{N}(0, 1),$$

where  $\bar{D}_{mn} = \bar{X}_m - \bar{Y}_n$  is an unbiased estimator for  $\mu$ . Thus, if both  $\sigma_X$  and  $\sigma_Y$  are known we get the “small sample” confidence interval

$$\mu \in \left[ \bar{D}_{mn} \mp z_{\delta/2} \sqrt{\frac{\sigma_X^2}{m} + \frac{\sigma_Y^2}{n}} \right] \text{ with prob. } \approx 1 - \delta.$$

In the general case we may proceed as follows. If

$$S_X^2 := \frac{1}{m-1} \sum_j (X_j - \bar{X}_m)^2, \quad S_Y^2 := \frac{1}{n-1} \sum_k (Y_k - \bar{Y}_n)^2,$$

are the unbiased estimators for  $\sigma_X^2$  and  $\sigma_Y^2$  respectively (by Corollary 7.20) then Corollary 7.24 implies that

$$(7.40) \quad \Sigma_X^2 := (m-1) \frac{S_X^2}{\sigma_X^2} \sim \chi_{m-1}^2 \text{ and } \Sigma_Y^2 := (n-1) \frac{S_Y^2}{\sigma_Y^2} \sim \chi_{n-1}^2$$

are independent and hence, by Corollary 4.24,

$$W_{X,Y} := \Sigma_X^2 + \Sigma_Y^2 \sim \chi_{m+n-2}^2.$$

If we set

$$S_{X,Y}^2(\eta) = \frac{(m-1)S_X^2 + (n-1)\eta S_Y^2}{m+n-2}, \quad \eta := \frac{\sigma_X^2}{\sigma_Y^2},$$

then it follows from Proposition 4.30 that

$$\begin{aligned} T_{X,Y} &:= \frac{Z_{X,Y}}{\sqrt{W_{X,Y}/(m+n-2)}} \\ &= \frac{\bar{D}_{mn} - \mu}{\sqrt{c_{mn}(\eta) S_{X,Y}^2(\eta)}}, \quad c_{mn}(\eta) = \frac{1}{m} + \frac{1}{n\eta}, \end{aligned}$$

satisfies

$$(7.41) \quad T_{X,Y} \sim t_{m+n-2},$$

thus being a pivotal quantity with respect to the parameter  $\eta$ ; see Remark 7.34. Hence, if the population variances, though unknown, are such that their ratio  $\eta$  is known, then there holds

$$(7.42) \quad \mu \in \left[ \bar{D}_{mn} \mp t_{m+n-2, \delta/2} \sqrt{c_{mn}(\eta) S_{XY}^2(\eta)} \right] \text{ with prob. } \approx 1 - \delta,$$

a “small sample” confidence interval for  $\mu$  quite similar in spirit to (7.36), which handles the case of a single normal sample. In particular, if the population variances are assumed to be equal then (7.42) holds with  $\eta = 1$ , in which case one has

$$(7.43) \quad T_{X,Y} = \frac{\bar{D}_{mn} - \mu}{\sqrt{c_{mn}(1) S_{XY}^2}},$$

where

$$S_{XY}^2 = \frac{(m-1)S_X^2 + (n-1)S_Y^2}{m+n-2},$$

is the *pooled variance*, a weighted sum of the sample variances. Otherwise, (7.42) has no practical usefulness as it provides a dispersion around  $\bar{D}_{mn}$  depending on the unknown parameter  $\eta$ . Proceeding as before, we are thus led to consider the *Behrens-Fisher-Welch statistics*

$$Z_{X,Y} = \frac{\bar{D}_{mn} - \mu}{\sqrt{\frac{S_X^2}{m} + \frac{S_Y^2}{n}}},$$

the obvious counterpart of  $Z_{X,Y}$ . Unfortunately, no simple, closed expression for the pdf of  $Z_{X,Y}$  does seem to exist. To appreciate the difficulties involved, note that  $Z_{X,Y} = Z_{X,Y} / \sqrt{W_{X,Y}}$ , where if

$$(7.44) \quad \beta_X = \frac{1}{m-1} \left( 1 + \frac{m}{n} \eta^{-1} \right)^{-1}, \quad \beta_Y = \frac{1}{n-1} \left( 1 + \frac{n}{m} \eta \right)^{-1}$$

then

$$W_{X,Y} = \beta_X \Sigma_X^2 + \beta_Y \Sigma_Y^2,$$

so that Corollary 4.23 and (7.40) may be used to ensure that

$$(7.45) \quad W_{X,Y} \sim \text{Gamma} \left( \frac{1}{2\beta_X}, \frac{m-1}{2} \right) + \text{Gamma} \left( \frac{1}{2\beta_Y}, \frac{n-1}{2} \right),$$

where  $(m-1)\beta_X + (n-1)\beta_Y = 1$  by (7.44). Since  $Z_{X,Y} \perp W_{X,Y}$  and  $Z_{X,Y} \sim \mathcal{N}(0, 1)$ , it follows from Remark 2.15 that computing  $\psi_{Z_{X,Y}}$  essentially reduces to figuring out  $\psi_{W_{X,Y}}$ , the pdf of a sum of independent Gamma-distributed random variables whose inverse scale parameters are distinct except when

$$\eta = \frac{m(m-1)}{n(n-1)}.$$

Thus, in most cases this sum fails to be Gamma-distributed; cf. Corollary 4.24. In fact, it is known that the exact formula for  $\psi_{W_{X,Y}}$  (and more generally for the pdf of a linear combination of chi-squared distributions) involves an infinite series representation in terms of certain transcendental functions [RP61, Mos85, HCP22], which demands the use of suitable approximations for  $\psi_{Z_{X,Y}}$ , a state of affairs that has inspired numerous studies on the computational performance of such methods<sup>22</sup>. To make things even worse, as it is apparent from (7.44) and (7.45),  $\psi_{Z_{X,Y}}$  is expected to explicitly depend on the “nuisance” parameter  $\eta$ , which means that  $Z_{X,Y}$  should fail to qualify as a pivotal quantity; see Remark 7.34. Thus, despite considerable progress in its

<sup>22</sup>See [Bau13] for a recent critical appraisal of recent contributions to this problem.

practical implementation [KC98], in a sense the general task of finding the most efficient estimate for  $\mu$  when  $\eta$  is unknown, usually referred to as the *Behrens-Fisher problem* [Wel96, Section 3.8], remains elusive.  $\square$

**Remark 7.37.** (F-test for the equality of variances) The reliability of the assumption on the equality of the population variances which led to (7.43) may be statistically justified (or not!) by running an F-test; for more on this see Section 11 below. We start by noticing that under the corresponding *null hypothesis*

$$H_0 : \sigma_X^2 = \sigma_Y^2,$$

(7.40) and Proposition 4.33 ensure that the appropriate test statistics

$$(7.46) \quad U := \frac{S_X^2}{S_Y^2} = \frac{\Sigma_X^2/(m-1)}{\Sigma_Y^2/(n-1)} \sim F_{m-1, n-1}.$$

Now, a very rough analysis of the departure from  $H_0$  goes as follows. If this hypothesis is not satisfied (so that  $\sigma_X^2 = \eta\sigma_Y^2$ ,  $\eta \neq 1$ ) then (7.46) becomes

$$U = \frac{\frac{\eta}{m-1}\Sigma_X^2}{\frac{1}{n-1}\Sigma_Y^2},$$

and we see from Corollary 4.23 and Proposition 4.20 that:

- The denominator satisfies

$$\frac{1}{n-1}\Sigma_Y^2 \sim \text{Gamma}\left(\frac{n-1}{2}, \frac{n-1}{2}\right),$$

so its distribution remains the same regardless of the validity of  $H_0$ ;

- On the other hand, the numerator satisfies

$$\frac{\eta}{m-1}\Sigma_X^2 \sim \text{Gamma}\left(\frac{m-1}{2\eta}, \frac{m-1}{2}\right),$$

a Gamma distribution whose scale factor is proportional to

$$\eta = \mathbb{E}\left(\frac{\eta}{m-1}\Sigma_X^2\right),$$

the parameter quantifying the departure from  $H_0$ .

Thus, at least on average, very small or very large realizations for  $U$  (substantially departing from  $u = 1$ ) provide statistical evidence for *rejecting*  $H_0$ . Precisely, if we fix  $0 < \alpha < 1$  and consider the corresponding quantiles  $f_{m-1, n-1, \alpha/2}^\pm$  determined by

$$(7.47) \quad F_{F_{m-1, n-1}}(f_{m-1, n-1, \alpha/2}^-) = \frac{\alpha}{2}, \quad F_{F_{m-1, n-1}}(f_{m-1, n-1, \alpha/2}^+) = 1 - \frac{\alpha}{2},$$

where  $F_{F_{m-1, n-1}}$  is the cdf of  $F_{m-1, n-1}$ , then  $H_0$  should be rejected “at significance level  $\alpha$ ” if the realization  $u$  of the statistics in (7.46) satisfies

$$(7.48) \quad u \in \left(0, f_{m-1, n-1, \alpha/2}^-\right] \cup \left[f_{m-1, n-1, \alpha/2}^+, +\infty\right).$$

A more convincing justification for this rather informal argument may be found in Section 11 below.  $\square$

**Remark 7.38.** (Reciprocity of the f-quantiles) Regarding the f-quantiles defined in (7.47), let us take  $X \sim F_{m, n}$  so that  $X^{-1} \sim F_{n, m}$  by Corollary 4.34. For any  $\alpha > 0$  we then have

$$\frac{\alpha}{2} = P\left(X \leq f_{m, n, \alpha/2}^-\right) = P\left(X^{-1} \geq \frac{1}{f_{m, n, \alpha/2}^-}\right),$$

so that

$$P\left(X^{-1} \leq \frac{1}{f_{m,n,\alpha/2}^-}\right) = 1 - \frac{\alpha}{2} = P\left(X^{-1} \leq f_{n,m,\alpha/2}^-\right),$$

from which the identity

$$f_{m,n,\alpha/2}^- f_{n,m,\alpha/2}^+ = 1$$

follows. □

**Example 7.39.** (The sampling distribution of the correlation coefficient) Let us retain the notation of Example 7.36, but this times with  $m = n$  and assuming that the *independent* random sample

$$(X, Y) := \{(X_1, Y_1), \dots, (X_m, Y_m)\}$$

has been drawn from a bi-variate normal population whose marginals are not necessarily independent. Thus,

$$(7.49) \quad (X_j, Y_j) \sim \mathcal{N}\left(\begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{XY} & \sigma_Y^2 \end{pmatrix}\right), \quad j = 1, \dots, n,$$

where  $\sigma_{XY} = \text{cov}(X_j, Y_j)$  is the *population covariance*<sup>23</sup>, so by (4.20) the joint distribution of  $(X, Y)$  is

$$(7.50) \quad \psi_{(X,Y)}(x, y) dx dy = \frac{1}{(2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2})^m} \times \\ \times e^{-\frac{1}{2(1-\rho^2)} \sum_j \left( \frac{(x_j - \mu_X)^2}{\sigma_X^2} - \frac{2\rho(x_j - \mu_X)(y_j - \mu_Y)}{\sigma_X\sigma_Y} + \frac{(y_j - \mu_Y)^2}{\sigma_Y^2} \right)} dx dy,$$

where

$$\rho = \frac{\sigma_{XY}}{\sigma_X\sigma_Y}$$

is the *correlation coefficient*, a population parameter whose estimation is a central theme in Multivariate Statistical Analysis [And03]<sup>24</sup>. It turns out that Fisher's geometric approach in Remark 7.31 can be successfully employed to this end [Fis15]. Indeed, using Remark 7.25 we have the identities

$$\sum_j (x_j - \mu_X)^2 = m \left( (\bar{x}_m - \mu_X)^2 + \hat{\sigma}_{m-1}^2(x) \right)$$

and

$$\sum_j (y_j - \mu_Y)^2 = m \left( (\bar{y}_m - \mu_Y)^2 + \hat{\sigma}_{m-1}^2(y) \right),$$

where  $\hat{\sigma}_{m-1}^2(x)$  and  $\hat{\sigma}_{m-1}^2(y)$  are the realizations of the variance estimators appearing in (7.13) with  $c = m^{-1}$ . Also, we will need their polarized version

$$\sum_j (x_j - \mu_X)(y_j - \mu_Y) = m \left( (\bar{x}_m - \mu_X)(\bar{y}_m - \mu_Y) + \hat{\sigma}_{m-1}^2(x, y) \right),$$

where

$$\hat{\sigma}_{m-1}^2(X, Y) := \frac{1}{m} \sum_j (X_j - \bar{X}_m)(Y_j - \bar{Y}_m).$$

<sup>23</sup>By Proposition 4.10,  $\{X_j, Y_j\}$  is independent if and only if  $\sigma_{XY} = 0$ .

<sup>24</sup>As usual, we assume that  $|\rho| < 1$ , thus avoiding the degenerate cases  $\rho = \pm 1$ .

Leading these identities to (7.50) we get

$$\begin{aligned} \psi_{(X,Y)} dx dy &= \frac{1}{(2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2})^m} \times \\ &\times e^{-\frac{m}{2(1-\rho^2)}\left(\frac{(\bar{x}_m-\mu_X)^2}{\sigma_X^2} - \frac{2\rho(\bar{x}_m-\mu_X)(\bar{y}_m-\mu_Y)}{\sigma_X\sigma_Y} + \frac{(\bar{y}_m-\mu_Y)^2}{\sigma_Y^2}\right)} \times \\ &\times e^{-\frac{m}{2(1-\rho^2)}\left(\frac{\hat{\sigma}_{m-1}^2(x)}{\sigma_X^2} - \frac{2\rho\hat{\sigma}_{m-1}(x)\hat{\sigma}_{m-1}(y)}{\sigma_X\sigma_Y} + \frac{\hat{\sigma}_{m-1}^2(y)}{\sigma_Y^2}\right)} dx dy, \end{aligned}$$

where

$$\hat{\rho} = \frac{\hat{\sigma}_{m-1}^2(X, Y)}{\hat{\sigma}_{m-1}(X)\hat{\sigma}_{m-1}(Y)}$$

is the *sample correlation coefficient*, the natural estimator for  $\rho$ ; see Example 8.5 for a justification of this latter claim. If we could view realizations of the samples  $X$  and  $Y$  as *independent* elements of  $\mathbb{R}_X^m$  and  $\mathbb{R}_Y^m$ , respectively, then the geometric reasoning in Remark 7.31 would ensure that

$$(7.51) \quad dx dy \approx \hat{\sigma}_{m-1}^{m-2}(x) d\hat{\sigma}_{m-1}(x) d\bar{x}_m d\theta_X \times \hat{\sigma}_{m-1}^{m-2}(y) d\hat{\sigma}_{m-1}(y) d\bar{y}_m d\theta_Y,$$

from which we would compute  $\psi_{(\hat{\sigma}_{m-1}(X), \hat{\sigma}_{m-1}(Y), \hat{\rho})}$  after integrating  $\psi_{(X,Y)} dx dy$  above against  $d\bar{x}_m d\theta_X d\bar{y}_m d\theta_Y$ . However, and this is the key point here,  $x$  and  $y$  are *not* allowed to vary freely as they are constrained to move in such a way that  $x \in \mathbb{S}_{\sqrt{m}\hat{\sigma}_{m-1}(x)}^{m-2}(\mathcal{X}_m)$  and  $y \in \mathbb{S}_{\sqrt{m}\hat{\sigma}_{m-1}(y)}^{m-2}(\mathcal{Y}_m)$  with

$$\hat{\rho} = \cos \theta, \quad \theta = \angle(\mathcal{X}_m x, \mathcal{Y}_m y).$$

Thus,

$$\begin{aligned} \psi_{(\hat{\sigma}_{m-1}(X), \hat{\sigma}_{m-1}(Y), \hat{\rho})} dv &\approx e^{-\frac{m}{2(1-\rho^2)}\left(\frac{\hat{\sigma}_{m-1}^2(s)}{\sigma_X^2} - \frac{2\rho\hat{\sigma}_{m-1}(x)\hat{\sigma}_{m-1}(y)}{\sigma_X\sigma_Y} + \frac{\hat{\sigma}_{m-1}^2(y)}{\sigma_Y^2}\right)} \times \\ &\times \hat{\sigma}_{m-1}^{m-2}(x) \hat{\sigma}_{m-1}^{m-2}(y) f(\hat{\rho}) dv, \end{aligned}$$

where  $dv = d\hat{\sigma}_{m-1}(x) d\hat{\sigma}_{m-1}(y) d\hat{\rho}$  and the extra factor  $f(\hat{\rho})$  comes from the constraint referred to above. Incidentally, this already shows that  $\{\bar{X}_m, \bar{Y}_m\}$  is independent from  $\{\hat{\sigma}_{m-1}^2(X), \hat{\sigma}_{m-1}^2(Y), \hat{\sigma}_{m-1}(XY)\}$ , as in the uni-variate case; cf. Proposition 7.23. Now, in order to determine  $f(\hat{\rho})$  note that if  $x$  is fixed then, corresponding to an infinitesimal displacement  $d\theta$ , the segment  $\overline{\mathcal{Y}_m y}$  describes an infinitesimal spherical slab in  $\mathbb{S}_{\sqrt{m}\hat{\sigma}_{m-1}(y)}^{m-2}(\mathcal{Y}_m)$  with radius

$$\sqrt{m}\hat{\sigma}_{m-1}(y) \sin \theta = \sqrt{m}\hat{\sigma}_{m-1}(y) \sqrt{1-\hat{\rho}^2}$$

and height

$$\sqrt{m}\hat{\sigma}_{m-1}(y) |d\theta| = \sqrt{m}\hat{\sigma}_{m-1}(y) \frac{d\hat{\rho}}{\sqrt{1-\hat{\rho}^2}},$$

thus tracing a volume proportionate to

$$(\sqrt{m}\hat{\sigma}_{m-1}(y) \sin \theta)^{m-3} \sqrt{m}\hat{\sigma}_{m-1}(y) |d\theta| = \hat{\sigma}_{m-1}^{m-2}(Y) \underbrace{(1-\hat{\rho}^2)^{\frac{m-4}{2}}}_{=f(\hat{\rho})} d\hat{\rho},$$

which finally gives

$$\begin{aligned} \psi_{(\hat{\sigma}_{m-1}(X), \hat{\sigma}_{m-1}(Y), \hat{\rho})} dv &\approx e^{-\frac{m}{2(1-\rho^2)}\left(\frac{\hat{\sigma}_{m-1}^2(x)}{\sigma_X^2} - \frac{2\rho\hat{\sigma}_{m-1}(x)\hat{\sigma}_{m-1}(y)}{\sigma_X\sigma_Y} + \frac{\hat{\sigma}_{m-1}^2(y)}{\sigma_Y^2}\right)} \times \\ &\times \hat{\sigma}_{m-1}^{m-2}(x) \hat{\sigma}_{m-1}^{m-2}(y) d\hat{\sigma}_{m-1}(x) d\hat{\sigma}_{m-1}(y) \times \\ &\times (1-\hat{\rho}^2)^{\frac{m-4}{2}} d\hat{\rho}. \end{aligned}$$

As usual, explicit, albeit quite complicated, expressions for the desired pdf  $\psi_{\hat{\rho}}$ , which may even be chosen so as to only involve elementary functions, are obtained by integrating this against the area element  $d\hat{\sigma}_{m-1}(x)d\hat{\sigma}_{m-1}(y)$ , with a further integration against  $d\hat{\rho}$  being needed to restore the normalizing constant. Of course, the computational difficulty here comes from the mixed term in the exponential which prevents  $\{\hat{\sigma}_{m-1}^2(X), \hat{\sigma}_{m-1}^2(Y), \hat{\rho}\}$  from being independent (except when  $\rho = 0$ ). In any case, the resulting expressions are found to depend on the parameters of the underlying normal bi-variate population only through  $\rho$  (and not on any other combination of the entries of the variance matrix in (7.49)), and in fact they all reduce to

$$\psi_{\hat{\rho}}(r) \approx (1 - r^2)^{\frac{m-4}{2}} \mathbf{1}_{(-1,1)}(r), \quad r \in \mathbb{R},$$

when  $\rho = 0$ , which suffices to efficiently testing the mutual independence of  $\{X_j, Y_j\}$  for any value of  $m$  along the lines of the general theory developed in Section 11. Otherwise, one has to appeal to asymptotic methods in order to construct “large sample” confidence intervals for  $\rho$ ; see Example 8.28 below. We refer to the original sources [Stu08a, Fis15], as well as to [Ken46, Chapter 14] and [And03, Chapter 4], for discussions of the basic properties of  $\psi_{\hat{\rho}}$  and their applications.  $\square$

**Example 7.40.** (One way ANOVA) Fix  $p \in \mathbb{N}$ ,  $p \geq 3$ , a finite sequence  $\{n_j\}_{j=1}^p \subset \mathbb{N}$  and for each  $j$  consider a random sample  $\{X_{jk}\}_{k=1}^{n_j}$  with  $X_{jk} \sim \mathcal{N}(\mu_j, \sigma^2)$  such that all these  $n := \sum_j n_j$  samples  $X_{ij}$  form an independent set. In other words, we are dealing here with  $p$  independent normal random samples with varied sizes and expectations but sharing the *same* variance, with the parameters  $\{\mu_1, \dots, \mu_p, \sigma^2\}$  being regarded as unknown. Within each sample we have the decomposition coming from Remark 7.25,

$$(7.52) \quad \sum_{k=1}^{n_j} (X_{jk} - \bar{X}_{\bullet\bullet})^2 = \sum_{k=1}^{n_j} (X_{jk} - \bar{X}_{j\bullet})^2 + n_j (\bar{X}_{j\bullet} - \bar{X}_{\bullet\bullet})^2,$$

where

$$\bar{X}_{j\bullet} = \frac{1}{n_j} \sum_{k=1}^{n_j} X_{jk}$$

and

$$\bar{X}_{\bullet\bullet} = \frac{1}{n} \sum_{j=1}^p \sum_{k=1}^{n_j} X_{ij} = \frac{1}{n} \sum_{j=1}^p n_j \bar{X}_{j\bullet}.$$

Note that

$$(7.53) \quad \mathbb{E}(\bar{X}_{j\bullet}^2) = \mu_j^2 + \frac{\sigma^2}{n_j},$$

and

$$(7.54) \quad \mathbb{E}(\bar{X}_{\bullet\bullet}^2) = \frac{1}{n^2} \left( \sum_{j=1}^p n_j \mu_j \right)^2 + \frac{\sigma^2}{n}.$$

Now, the same argument leading to the proof of Proposition 7.24 implies that

$$(7.55) \quad \sigma^{-2} \sum_{k=1}^{n_j} (X_{jk} - \bar{X}_{j\bullet})^2 \sim \chi_{n_j-1}^2$$

is independent of  $\bar{X}_{j\bullet}$  and hence of  $(\bar{X}_{j\bullet} - \bar{X}_{\bullet\bullet})^2$ . Thus, if we set

$$S_{\text{Total}}^2 = \sum_{j=1}^p \sum_{k=1}^{n_j} (X_{jk} - \bar{X}_{\bullet\bullet})^2,$$



the *total* sum of squares,

$$S_{\text{Within}}^2 = \sum_{j=1}^p \sum_{k=1}^{n_j} (X_{jk} - \bar{X}_{j\bullet})^2,$$

the sum of squares *within* the samples, and

$$S_{\text{Between}}^2 = \sum_{j=1}^p n_j (\bar{X}_{j\bullet} - \bar{X}_{\bullet\bullet})^2,$$

the sum of squares *between* the samples, then

$$(7.56) \quad S_{\text{Total}}^2 = S_{\text{Within}}^2 + S_{\text{Between}}^2,$$

with

$$\sigma^{-2} S_{\text{Within}}^2 \sim \chi_{n-p}^2$$

being independent of  $S_{\text{Between}}^2$ . On the other hand, again by the argument leading to (7.55), but this time under the *null hypothesis*

$$(7.57) \quad H_0 : \mu_1 = \cdots = \mu_p,$$

we have

$$\sigma^{-2} S_{\text{Total}}^2 \sim \chi_{n-1}^2$$

and hence

$$\sigma^{-2} S_{\text{Between}}^2 \sim \chi_{p-1}^2,$$

which gives

$$(7.58) \quad \mathbb{E}(S_{\text{Between}}^2) = (p-1)\sigma^2$$

by Corollary 4.21. It then follows from Proposition 4.33 that

$$(7.59) \quad V := \frac{S_{\text{Between}}^2/(p-1)}{S_{\text{Within}}^2/(n-p)} \sim F_{p-1, n-p} \text{ under } H_0.$$

In order to proceed we now observe that:

- The decomposition (7.56), which is an easy consequence of the fundamental algebraic identity in Remark 7.25, plays a central role in our analysis as it displays  $S_{\text{Total}}^2$ , the total sum of squares, as resulting from the contribution of two terms of rather distinct types:  $S_{\text{Within}}^2$  collects together the variations *within* the various samples whereas  $S_{\text{Between}}^2$  measures the variation *between* the samples;
- In consonance with the previous item, the computation leading to the statistics in (7.59) shows that the distribution of its numerator is conditioned to the validity of  $H_0$  whereas the distribution of its denominator remains the same regardless of the validity of this hypothesis;
- If  $H_0$  is not necessarily satisfied then starting from the fact that

$$S_{\text{Between}}^2 = \sum_{j=1}^p n_j \bar{X}_{j\bullet}^2 - n \bar{X}_{\bullet\bullet}^2,$$

we easily deduce by means of (7.53) and (7.54) that

$$\mathbb{E}(S_{\text{Between}}^2) = (p-1)\sigma^2 + \sum_{j=1}^p n_j (\mu_j - \bar{\mu})^2, \quad \bar{\mu} = \frac{1}{n} \sum_{j=1}^p n_j \mu_j,$$

which assumes its minimal value, given by (7.58), exactly when  $H_0$  holds true.

Thus, at least on average, a sufficiently large value of  $V$  provides statistical evidence for *rejecting*  $H_0$ . Precisely, if we fix  $0 < \alpha < 1$  and consider the corresponding quantile  $f_{p-1, n-p, \alpha}$  determined by

$$F_{F_{p-1, n-p}}(f_{p-1, n-p, \alpha}) = 1 - \alpha,$$

where  $F_{F_{p-1, n-p}}$  is the cdf of  $F_{p-1, n-p}$ , then  $H_0$  should be rejected “at significance level  $\alpha$ ” if the realization  $v$  of  $V$  in (7.59) satisfies

$$(7.60) \quad v \in [f_{p-1, n-p, \alpha}, +\infty).$$

Again, we refer to Section 11 for a more theoretically inclined justification of this procedure, in particular for the proper understanding of why the “rejection subsets” in the right-hand side of the F-tests in (7.48) and (7.60) differ in their “connectedness”.  $\square$

## 8. MAXIMUM LIKELIHOOD

We now present a remarkable class of estimators, introduced by R. Fisher, which displays, under suitable regularity assumptions, many desirable asymptotic properties, including asymptotic normality (Theorem 8.23).

**8.1. Maximum likelihood estimators.** We start with an *independent* family  $\{X_j\}_{j=1}^n$  of random variables with  $X_j \sim \psi_j(x_j; \theta) > 0$ ,  $\theta \in \Theta$ .

**Definition 8.1.** The *likelihood function* of the random vector  $X = (X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n$  is

$$(8.1) \quad L(\mathbf{x}; \theta) = \prod_{j=1}^n \psi_j(x_j; \theta),$$

where  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$  is viewed as a realization of  $X$ . In particular, if  $\{X_j\}$  is i.i.d. ( $X_j \sim \psi(x_j; \theta)$ ) then

$$(8.2) \quad L(\mathbf{x}; \theta) = \prod_{j=1}^n \psi_\theta(x_j), \quad \psi_\theta(x_j) = \psi(x_j; \theta).$$

We begin with a motivation to the effect that MLE estimators in the i.i.d. setting (which is actually the only case considered here) arise as approximate solutions to a natural variational problem.

**Definition 8.2.** If  $\theta_0 \in \Theta$  is the sought-for parameter we intend to estimate, we define the *Kullback-Leibler divergence* (centered at  $\theta_0$ ) by

$$(8.3) \quad \theta \in \Theta \mapsto D_{\theta_0}^{KL}(\theta) := \int_{\mathbb{R}^n} \psi_{\theta_0}(\mathbf{x}) \ln \left( \frac{\psi_{\theta_0}(\mathbf{x})}{\psi_\theta(\mathbf{x})} \right) d\mathbf{x}.$$

Using Jensen’s inequality and assuming as always that our statistic model  $X_j \sim \psi_\theta$  is identifiable<sup>25</sup>, we easily see that  $D_{\theta_0}^{KL}(\theta) \geq 0$  for any  $\theta$ , with the equality holding only if  $\theta = \theta_0$ , which justifies seeking for an estimator  $\theta^*$  satisfying

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} D_{\theta_0}^{KL}(\theta).$$

Since

$$D_{\theta_0}^{KL}(\theta) = \operatorname{const}_{\theta_0} - \mathbb{E}_{\theta_0}(\ln(\psi_\theta)),$$

this amounts to finding

$$\theta^* = \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_{\theta_0}(\ln(\psi_\theta)).$$

<sup>25</sup>Recall that this means that the map  $\theta \in \Theta \mapsto \psi_\theta$  is injective.

By LLN we may approximate this for  $n$  large enough by

$$(8.4) \quad \theta^* \approx \frac{1}{n} \operatorname{argmax}_{\theta \in \Theta} \sum_j \ln \psi_\theta(x_j) = \operatorname{argmax}_{\theta \in \Theta} \ln \Pi_j \psi_\theta(x_j),$$

which may be interpreted as saying that  $\theta^*$  represents the best choice for the parameter estimator based on the observed value  $\mathbf{x}$ . This justifies the following important construction due to R. Fisher, which is by far the most popular technique for deriving estimators.

**Definition 8.3.** (Maximum Likelihood Estimation, MLE) Under the conditions above, a maximum likelihood (ML) estimator  $\hat{\theta}$  is a solution of the maximization problem

$$(8.5) \quad \hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} L(\mathbf{x}; \theta).$$

Equivalently,

$$(8.6) \quad \hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} l(\mathbf{x}; \theta).$$

where

$$(8.7) \quad l(\mathbf{x}; \theta) = \ln L(\mathbf{x}; \theta)$$

is the *log-likelihood function*.

As usual, we assume that  $l$  is strictly concave in  $\theta$  so a solution to (8.5) is unique (whenever it exists). Thus, at least in the i.i.d. case, the ML estimator asymptotically minimizes the Kullback-Leibler “distance” to  $\theta_0$  in (8.3).

**Example 8.4.** (MLE from a normal population) If  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  then

$$(8.8) \quad L(\mathbf{x}; \theta) = (2\pi\theta_2)^{-n/2} e^{-\frac{1}{2\theta_2} \sum_{j=1}^n (x_j - \theta_1)^2},$$

where  $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+$ , so that

$$(8.9) \quad l(\mathbf{x}; \theta) = \ln L(\mathbf{x}; \theta) = -\frac{n}{2} \ln(2\pi\theta_2) - \frac{1}{2\theta_2} \sum_j (x_j - \theta_1)^2.$$

The usual first derivative test shows that the ML estimator  $\hat{\theta}_1$  and  $\hat{\theta}_2$  corresponding to (8.9) should satisfy

$$0 = \frac{\partial l}{\partial \theta_1}(\hat{\theta}) = \frac{1}{\theta_2} \sum_j (X_j - \hat{\theta}_1), \quad 0 = \frac{\partial l}{\partial \theta_2}(\hat{\theta}) = -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \hat{\theta}_1)^2,$$

which gives

$$\hat{\theta}_1 = \bar{X}_n = \frac{1}{n} \sum_j X_j, \quad \hat{\theta}_2 = \hat{\sigma}_{n-1}^2 = \frac{1}{n} \sum_j (X_j - \bar{X}_n)^2.$$

We thus see that the ML estimator  $\hat{\theta}_2 = \hat{\sigma}_{n-1}^2$  for the variance coming from (8.9) not only fails to be unbiased but also satisfies

$$\operatorname{mse}(\hat{\sigma}_{(n-1)-1}^2) > \operatorname{mse}(\hat{\sigma}_{n-1}^2) > \operatorname{mse}(\hat{\sigma}_{(n+1)-1}^2),$$

so its performance, as measured by mse, lies somewhere between those of the variance estimators considered so far. We point out that the asymptotic performance of a (sufficiently regular and consistent) ML estimator is examined in Theorem 8.23 below. In particular, asymptotic normality is established there, which confirms that  $\hat{\sigma}_{n-1}^2$  stands out as the most efficient estimator from this viewpoint.  $\square$

**Example 8.5.** (MLE from a jointly normal population) Using the notation of Example 7.39, we see that the right-hand side of (7.50) allows us to write down the likelihood function of the jointly normal random sample  $\{X, Y\}$  as

$$L(\mathbf{x}, \mathbf{y}; \theta) = \frac{1}{(2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2})^m} \times \\ \times e^{-\frac{1}{2(1-\rho^2)} \sum_j \left( \frac{A(x_j)}{\sigma_X^2} - \frac{2\rho B(x_j, y_j)}{\sigma_X\sigma_Y} + \frac{C(y_j)}{\sigma_Y^2} \right)},$$

where  $\theta = (\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2, \rho)$  and

$$(8.10) \quad A(x_j) = (x_j - \mu_X)^2, \quad B(x_j, y_j) = (x_j - \mu_X)(y_j - \mu_Y), \quad C(y_j) = (y_j - \mu_Y)^2,$$

so that

$$(8.11) \quad l(\mathbf{x}, \mathbf{y}; \theta) = -m \ln(2\pi\sigma_X\sigma_Y) - \frac{m}{2} \ln(1-\rho^2) - \\ - \frac{1}{2(1-\rho^2)} \sum_j \left( \frac{A(x_j)}{\sigma_X^2} - \frac{2\rho B(x_j, y_j)}{\sigma_X\sigma_Y} + \frac{C(y_j)}{\sigma_Y^2} \right),$$

Starting from this, a straightforward analysis involving the first derivative test  $\nabla_\theta l = 0$  confirms that  $\hat{\theta} := (\bar{X}_m, \bar{Y}_m, \hat{\sigma}_{m-1}^2(X), \hat{\sigma}_{m-1}^2(Y), \hat{\rho})$  is the ML estimator of  $\theta$ ; see [Ken46, Section 14.11] or [And03, Corollary 3.2.2] for the details.  $\square$

**Example 8.6.** (MLE from an exponential population) If  $X_j \sim \text{Exp}(\lambda)$ , the exponential distribution with parameter  $\lambda > 0$ , where

$$\text{Exp}(\lambda)(x) = \lambda e^{-\lambda x} \mathbf{1}_{[0, +\infty)},$$

then

$$L(\mathbf{x}; \lambda) = \lambda^n e^{-\lambda \sum_j x_j} \implies l(\mathbf{x}; \lambda) = n \ln \lambda - \lambda \sum_j x_j,$$

so that

$$\hat{\lambda} = \frac{n}{\sum_j X_j},$$

which matches the fact that  $\mathbb{E}(X_j) = 1/\lambda$ .  $\square$

**Example 8.7.** The MLE method also works in the discrete case. For instance, if  $X_j \sim \text{Ber}(p)$ , the Bernoulli distribution in Remark 6.8, then for  $x \in \{0, 1\}$  we have  $P(X_j = x) = p^x(1-p)^{1-x}$ , where  $p \in (0, 1)$  is the unknown parameter. If we assume as always that  $\{X_j\}$  is independent the associated likelihood function is

$$L(\mathbf{x}; p) = p^{\sum_j x_j} (1-p)^{n-\sum_j x_j},$$

and hence,

$$(8.12) \quad l(\mathbf{x}; p) = \left( \sum_j x_j \right) \ln p + \left( n - \sum_j x_j \right) \ln(1-p).$$

The usual first derivative test for a minimum at  $p = \hat{p}$  is

$$0 = \frac{\partial l}{\partial p}(\hat{p}) = \frac{\sum_j x_j}{\hat{p}} - \frac{n - \sum_j x_j}{1 - \hat{p}},$$

so that

$$\hat{p} = \frac{1}{n} \sum_j X_j,$$

the sample mean. Also, if  $X_j \sim \mathcal{P}(\rho)$ , the Poisson distribution with parameter  $\rho > 0$ , so that  $P(X_j = x) = \rho^x e^{-\rho} / x!$ ,  $x \in \{0, 1, 2, \dots\}$ , then a simple computation shows that

$$\hat{\rho} = \frac{1}{n} \sum_j X_j,$$

which confirms that the corresponding ML estimator is also the sample mean.  $\square$

**Remark 8.8.** The analysis in the examples above should be complemented with the usual second derivative test to check in each case that the ML estimator attains the (unique) global maximum of the corresponding likelihood function.  $\square$

**8.2. Fisher information and Cramér-Rao lower bound.** We now present a universal lower bound for the covariance matrix in each class of estimators with a prescribed expectation (in particular, for unbiased estimators) in terms of an invariant (Fisher information) depending on the likelihood function of the given statistical model. Instead of restricting ourselves to statistical models, let us assume for the moment only that  $X_j : \Omega \rightarrow \mathbb{R}$  are independent random variables,  $i = 1, \dots, n$ , giving rise to a random vector  $X = (X_1, \dots, X_n)^{26}$ . Let  $L = L(\mathbf{x}; \theta) > 0$  be the corresponding likelihood function and  $l(\mathbf{x}; \theta) = \ln L(\mathbf{x}; \theta)$  the log-likelihood function, where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\theta \in \Theta \subset \mathbb{R}^q$ , the space of parameters. Clearly, the notions of statistics and estimators can be easily adapted to this more general setting. In the following, we assume that  $l$  is regular enough so that all the differential/integral manipulations hold true.

**Definition 8.9.** Under the conditions above, we define the *score vector* (of the given sample  $X$ ) as

$$s(X; \theta) = \nabla_{\theta} l(X; \theta).$$

Also, the corresponding *Fisher information matrix* is

$$(8.13) \quad \mathcal{F}^X(\theta) = \text{cov}(s(X; \theta)).$$

We henceforth assume that the symmetric matrix  $\mathcal{F} = \mathcal{F}^X$  in (8.9) is positive definite, so the inverse matrix  $\mathcal{F}^{-1}$  exists. Since the random effects present in the sample have been averaged out after taking covariance of the score,  $\mathcal{F} = \mathcal{F}(\theta)$  is an invariant of the given model, in particular not being attached to any potential estimator.

**Remark 8.10.** In order to illustrate the importance of requiring that  $\{X_j\}$  is independent, let us assume that we are in the unidimensional case,  $\Theta \subset \mathbb{R}$ , so we call the scalar  $\mathcal{F}$  simply the *Fisher information*. One has

$$s(X; \theta) = \sum_j \frac{\frac{d}{d\theta} \psi_j(X_j; \theta)}{\psi_j(X_j; \theta)},$$

a sum of *independent* random variables, so that by (2.12),

$$\text{var}(s(X; \theta)) = \sum_j \text{var} \left( \frac{\frac{d}{d\theta} \psi_j(X_j; \theta)}{\psi_j(X_j; \theta)} \right),$$

which means that

$$\mathcal{F}^X = \sum_j \mathcal{F}^{X_j}.$$

<sup>26</sup>In other words, we consider here a statistical model in the extended sense of Remark 7.4.

Thus, independence leads to a simple additive formula describing how the Fisher information of the whole sample decomposes as a sum of contributions coming from its parts. If we additionally require that  $X_j \sim \psi_\theta$  is i.d.d., which is the only case treated in all examples below, then this becomes

$$\mathcal{F}_{(n)} = n\mathcal{F}_{(1)},$$

with  $\mathcal{F}_{(n)} = \mathcal{F}^X$  referring as before to the whole sample whereas  $\mathcal{F}_{(1)} = \mathcal{F}^{X_j}$  refers to *any* single observation.  $\square$

**Example 8.11.** If  $X_j \sim \text{Ber}(p)$  then (8.12) gives

$$s(X; p) = \frac{\partial}{\partial p} l(X; p) = \frac{\sum_j X_j}{p} - \frac{n - \sum_j X_j}{1 - p} = \frac{n}{p(1 - p)} \bar{X} - \frac{n}{1 - p},$$

and since  $\text{cov}(\bar{X}) = \text{cov}(X_j)/n = p(1 - p)/n$ , we conclude that

$$(8.14) \quad \mathcal{F}_{(n)}(p) = \frac{n}{p(1 - p)}.$$

Thus, the Fisher information increases with the sample size according to a rate which is inversely proportional to the “fluctuation” (as measured by the population variance). In a sense, this simple example justifies the qualification of “information” for this concept; for more on this point see Remark 8.24 below.  $\square$

**Example 8.12.** If  $X_j \sim \mathcal{N}(\theta_1, \theta_2)$  is drawn from a normal population, where  $\theta_1 = \mu$  and  $\theta_2 = \sigma^2$ , then (8.9) leads to

$$(8.15) \quad s(X; \theta) = \nabla_\theta l(X; \theta) = \begin{pmatrix} -\frac{1}{\theta_2} \sum_j (X_j - \theta_1) \\ -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \theta_1)^2 \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}.$$

Since  $X_j - \theta_1 \sim \mathcal{N}(0, \theta_2)$ , independence implies that

$$(8.16) \quad \text{var} \left( \frac{1}{\theta_2} \sum_j (X_j - \theta_1) \right) = \frac{n}{\theta_2}.$$

On the other hand,

$$-\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \theta_1)^2 = -\frac{n}{2\theta_2} + \frac{1}{2\theta_2} \sum_j \left( \frac{X_j - \theta_1}{\sqrt{\theta_2}} \right)^2,$$

and since  $(X_j - \theta_1)/\sqrt{\theta_2} \sim \mathcal{N}(0, 1)$ , we see that

$$\sum_j \left( \frac{X_j - \theta_1}{\sqrt{\theta_2}} \right)^2 \sim \chi_n^2$$

by Corollary 4.25, so that Corollary 4.21 applies to give

$$(8.17) \quad \text{var} \left( \sum_j \left( \frac{X_j - \theta_1}{\sqrt{\theta_2}} \right)^2 \right) = 2n,$$

and hence,

$$(8.18) \quad \text{var} \left( -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \theta_1)^2 \right) = \frac{n}{2\theta_2^2}.$$

Note that (8.16) and (8.18) provide the diagonal terms of the corresponding Fisher information matrix. In order to compute the off-diagonal terms we observe that

$$\mathbb{E} \left( \frac{1}{\theta_2} \sum_j (X_j - \theta_1) \right) = 0$$

and

$$\begin{aligned} \mathbb{E} \left( -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \theta_1)^2 \right) &= -\frac{n}{2\theta_2} + \frac{1}{2\theta_2} \mathbb{E} \left( \sum_j \left( \frac{X_j - \theta_1}{\sqrt{\theta_2}} \right)^2 \right) \\ &= -\frac{n}{2\theta_2} + \frac{n}{2\theta_2} \\ &= 0, \end{aligned}$$

where we used (8.17) and Corollary 4.21 in the next to the last step; this should be compared with the general result in Corollary 8.14 below. It follows that

$$\text{cov} \left( \frac{1}{\theta_2} \sum_j (X_j - \theta_1), -\frac{n}{2\theta_2} + \frac{1}{2\theta_2^2} \sum_j (X_j - \theta_1)^2 \right) = \frac{1}{2\theta_2^3} \sum_j \mathbb{E} ((X_j - \theta_1)^3),$$

which clearly vanishes. We thus conclude that

$$(8.19) \quad \mathcal{F}_{(n)}(\theta) = \begin{pmatrix} n/\theta_2 & 0 \\ 0 & n/2\theta_2^2 \end{pmatrix}.$$

In particular, if we consider  $\theta_2$  as known,

$$(8.20) \quad \mathcal{F}_{(n)}(\theta_1) = \frac{n}{\theta_2}.$$

We will see in Remark 8.16 below a much simpler route to retrieve (8.19). □

In both (8.20) and (8.14) we see that the Fisher information  $\mathcal{F} = \mathcal{F}_{(1)}$  of the parameter equals the reciprocal of the variance of the corresponding estimator (in both cases, the sample mean, which is unbiased). This is just a manifestation of the quite remarkable fact that in general the information matrix appears in a universal lower bound for the covariance of each class of estimators with a prescribed bias (in particular, for unbiased estimators); see Theorem 8.17 below. The next result is a first step toward this goal.

**Proposition 8.13.** *For any (sufficiently regular) vector  $t = t(\mathbf{x}; \theta)$  there holds*

$$\mathbb{E}(s \otimes t) = \nabla_\theta \mathbb{E}(t) - \mathbb{E}(\nabla_\theta t).$$

*Proof.* We compute:

$$\begin{aligned} \mathbb{E}(s \otimes t) &= \int_{\mathbb{R}^n} L(\mathbf{x}; \theta)^{-1} \nabla_\theta L(\mathbf{x}; \theta) \otimes t(\mathbf{x}; \theta) L(\mathbf{x}; \theta) d\mathbf{x} \\ &= \int_{\mathbb{R}^n} \nabla_\theta (L(\mathbf{x}; \theta) t(\mathbf{x}; \theta)) d\mathbf{x} - \int_{\mathbb{R}^n} L(\mathbf{x}; \theta) \nabla_\theta t(\mathbf{x}; \theta) d\mathbf{x} \\ &= \nabla_\theta \int_{\mathbb{R}^n} L(\mathbf{x}; \theta) t(\mathbf{x}; \theta) d\mathbf{x} - \int_{\mathbb{R}^n} L(\mathbf{x}; \theta) \nabla_\theta t(\mathbf{x}; \theta) d\mathbf{x}, \end{aligned}$$

as desired. □

**Corollary 8.14.** *There holds  $\mathbb{E}(s) = \vec{0}$ . In particular,*

$$(8.21) \quad \mathcal{F} = \mathbb{E}(s \otimes s) = -\mathbb{E}(\nabla_{\theta\theta}^2 l).$$

*Proof.* Take  $t = (1, \dots, 1)$ . □

**Corollary 8.15.** *If  $t = t(\mathbf{x})$  then  $\mathbb{E}(s \otimes t) = \nabla_{\theta} \mathbb{E}(t)$ . In particular, if  $t = \hat{\theta}$  is an estimator with  $g(\theta) := \mathbb{E}(\hat{\theta})$  then  $\mathbb{E}(s \otimes \hat{\theta}) = \nabla_{\theta} g$ .*

*Proof.* The first assertion is immediate and the second one follows from the fact that  $\hat{\theta}$ , as an estimator, does not depend on  $\theta$ . □

**Remark 8.16.** As a checking we may use (8.21) to recalculate the Fisher information matrix of a normal sample  $X_j \sim \mathcal{N}(\theta_1, \theta_2)$  as in Example 8.12. From (8.15) we have

$$\nabla_{\theta\theta} l(X; \theta) = \begin{pmatrix} -\frac{n}{\theta_1} & -\frac{1}{\theta_2^2} \sum_j (X_j - \theta_1) \\ -\frac{1}{\theta_2^2} \sum_j (X_j - \theta_1) & \frac{n}{2\theta_2^2} - \frac{1}{\theta_2^3} \sum_j (X_j - \theta_1)^2 \end{pmatrix},$$

so that

$$\mathcal{F}_{(n)}(\theta) = \begin{pmatrix} \frac{n}{\theta_1} & 0 \\ 0 & -\frac{n}{2\theta_2^2} + \frac{1}{\theta_2^3} \sum_j \mathbb{E}((X_j - \theta_1)^2) \end{pmatrix},$$

and since  $\mathbb{E}((X_j - \theta_1)^2) = \text{var}(X_j) = \theta_2$ , we recover (8.19). Note that this computation is much simpler because it is based on computing expectations, at the cost of taking one more derivative of log-likelihood function but with no need to compute covariances, and hence bypasses any appeal to the connection between sums of squares of normals and chi-squares. □

Let  $\mathcal{E}$  be the set of all estimators (for  $\theta$ ). Given  $g : \Theta \rightarrow \mathbb{R}^q$  define

$$\mathcal{E}_g = \left\{ \hat{\theta} \in \mathcal{E}; \mathbb{E}(\hat{\theta}) = g(\theta) \right\}.$$

Equivalently,  $\mathcal{E}_g$  is the set of all *unbiased* estimators for  $g(\theta)$ . Note that each  $\hat{\theta} \in \mathcal{E}_g$  satisfies

$$(8.22) \quad \text{bias}(\hat{\theta}) = g(\theta) - \theta$$

and hence

$$\text{mse}(\hat{\theta}) = \|g(\theta) - \theta\|^2 + \text{tr cov}(\hat{\theta}).$$

The next result provides a uniform lower bound for the covariance (and hence for the mse) of estimators in each class  $\mathcal{E}_g$  (provided it is not empty).

**Theorem 8.17.** (Cramér-Rao) *There holds*

$$(8.23) \quad \text{cov}(\hat{\theta}) \geq \nabla_{\theta} g \mathcal{F}(\theta)^{-1} \nabla_{\theta} g^t$$

*for any  $\hat{\theta} \in \mathcal{E}_g$ . In particular,*

$$(8.24) \quad \text{cov}(\hat{\theta}) \geq \mathcal{F}(\theta)^{-1}$$

*if  $\hat{\theta}$  is unbiased ( $g(\theta) = \theta$ ).*



*Proof.* We first consider the uni-dimensional case  $\Theta \subset \mathbb{R}$ . From Corollaries 8.15 and 8.14 we have

$$\begin{aligned} g'(\theta) &= \int_{\mathbb{R}^n} s(\mathbf{x}; \theta) \hat{\theta}(\mathbf{x}) L(\mathbf{x}; \theta) d\mathbf{x} \\ &= \int_{\mathbb{R}^n} s(\mathbf{x}; \theta) \left( \hat{\theta}(\mathbf{x}) - g(\theta) \right) L(\mathbf{x}; \theta) d\mathbf{x}, \end{aligned}$$

so that Cauchy-Schwartz inequality gives

$$|g'(\theta)|^2 \leq \text{var}(s(X; \theta)) \text{var}(\hat{\theta}(X)),$$

as desired. The proof of the multi-dimensional case is quite similar and makes use of a well-known algebraic inequality: for any random vectors  $Z, W \in \mathbb{R}^q$  with  $\text{cov}(W) > 0$  there holds

$$\text{cov}(Z) \geq \text{cov}(Z, W) \text{cov}(W)^{-1} \text{cov}(W, Z).$$

Taking  $Z = \hat{\theta}$  and  $W = s$  we get

$$\begin{aligned} \text{cov}(\hat{\theta}) &\geq \text{cov}(\hat{\theta}, s) \text{cov}(s)^{-1} \text{cov}(s, \hat{\theta}) \\ &= \mathbb{E}(\hat{\theta} \otimes s) \mathbb{E}(s \otimes s)^{-1} \mathbb{E}(s \otimes \hat{\theta}) \\ &= \nabla_{\theta} g \mathcal{F}(\theta)^{-1} \nabla_{\theta} g^t, \end{aligned}$$

as desired. □

**Corollary 8.18.** *The best estimator in  $\mathcal{E}_g$  (if it exists) is the one whose covariance matrix attains the lower bound in (8.23). In particular, an unbiased estimator whose covariance matrix attains the lower bound in (8.24) has the best performance (as measured by the mse).*

**Example 8.19.** It follows from (8.14) that, for a Bernoulli population,

$$\mathcal{F}_{(1)}(p) = \frac{1}{p(1-p)} = \frac{1}{\text{var}(\hat{p})},$$

so Corollary 8.18 applies and the sample mean is the best unbiased estimator for the expected value. A similar reasoning, based on the explicit computation of the corresponding Fisher information, confirms that the sample mean is the best unbiased estimator for the expected value of a Poisson random sample as in Example 8.7.

**Example 8.20.** (The sample mean as the best estimator of the expected value of a normal population) As observed in [HL51], the Cramér-Rao inequality in Theorem 8.17 may be used to prove that the sample mean  $\hat{\theta}_1 = \bar{X}$  is the *best* estimator for the mean  $\theta_1 = \mu$  of a normal population (in the sense that it has the least possible mse among *all* such estimators<sup>27</sup>); see also [LC06, Example 5.2.8] for another approach to this result. Here we use the notation of Example 8.12 so that our sample satisfies  $X_j \sim \mathcal{N}(\theta_1, \theta_2)$ ,  $j = 1, \dots, n$ . Now, let us take an estimator  $\hat{\theta}_{\bullet}$  of  $\theta_1$  satisfying  $\text{mse}(\hat{\theta}_{\bullet}) \leq \text{mse}(\hat{\theta}_1) = \theta_2/n$ . Setting  $b(\theta_1) = \text{bias}_{\theta_1}(\hat{\theta}_{\bullet})$  and using that  $\mathcal{F}(\theta_1) = n/\theta_2$  by (8.19) we then see from (8.23) and (8.22) that

$$b(\theta_1)^2 + \frac{\theta_2}{n} (1 + b'(\theta_1))^2 \leq \frac{\theta_2}{n},$$

from which we easily deduce that  $b \equiv 0$ . Hence,  $\hat{\theta}_{\bullet}$  is unbiased and satisfies  $\text{mse}(\hat{\theta}_{\bullet}) = \text{mse}(\hat{\theta}_1)$ , as desired. □

<sup>27</sup>We then say that  $\hat{\theta}_1$  is *admissible*, which means that no other estimator  $\hat{\theta}_{\bullet}$  satisfies  $\text{mse}(\hat{\theta}_{\bullet}) < \text{mse}(\hat{\theta}_1)$ .

**Example 8.21.** (The James-Stein estimator [JS61]) Starting with a single sample  $X \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \text{Id}_p)$ , where  $\sigma^2$  is known, the log-likelihood function

$$(8.25) \quad l(\mathbf{x}; \boldsymbol{\mu}) = -\frac{p}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{x} - \boldsymbol{\mu}\|^2$$

tells us that the MLE for  $\boldsymbol{\mu}$  is  $\hat{\boldsymbol{\mu}} = X$ . Since  $\sigma^{-2} \|X - \boldsymbol{\mu}\|^2 \sim \chi_p^2$  we know that

$$(8.26) \quad \text{mse}(\hat{\boldsymbol{\mu}}) = \mathbb{E}(\|X - \boldsymbol{\mu}\|^2) = p\sigma^2,$$

with only the variance contributing (since  $\hat{\boldsymbol{\mu}}$  is unbiased). Now let us compare  $\hat{\boldsymbol{\mu}}$  with the *James-Stein estimator*

$$(8.27) \quad \hat{\boldsymbol{\mu}}_{JS} = \left(1 - \frac{(p-2)\sigma^2}{\|\hat{\boldsymbol{\mu}}\|^2}\right) \hat{\boldsymbol{\mu}} = X - (p-2)\sigma^2 \frac{X}{\|X\|^2},$$

whose mean squared error is

$$\begin{aligned} \text{mse}(\hat{\boldsymbol{\mu}}_{JS}) &= \mathbb{E}(\|\hat{\boldsymbol{\mu}}_{JS} - \boldsymbol{\mu}\|^2) \\ &= \mathbb{E}\left(\left\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu} - (p-2)\sigma^2 \frac{\hat{\boldsymbol{\mu}}}{\|\hat{\boldsymbol{\mu}}\|^2}\right\|^2\right) \\ &= \text{mse}(\hat{\boldsymbol{\mu}}) - 2(p-2)\sigma^2 \mathbb{E}\left(\frac{\langle X, X - \boldsymbol{\mu} \rangle}{\|X\|^2}\right) + (p-2)^2 \sigma^4 \mathbb{E}(\|X\|^{-2}) \\ &= p\sigma^2 - 2(p-2)\sigma^2 \mathbb{E}\left(\frac{\langle X, X - \boldsymbol{\mu} \rangle}{\|X\|^2}\right) + (p-2)^2 \sigma^4 \mathbb{E}(\|X\|^{-2}), \end{aligned}$$

where we used (8.26) in the last step. In order to handle the mixed term in the right-hand side we first note from (8.25) that the score vector is

$$s(\mathbf{x}; \boldsymbol{\mu}) = \nabla_{\boldsymbol{\mu}} l(\mathbf{x}; \boldsymbol{\mu}) = \sigma^{-2} (\mathbf{x} - \boldsymbol{\mu}),$$

so if we make  $t = t(X) \in \mathbb{R}$  in Corollary 8.15 we obtain

$$(8.28) \quad \mathbb{E}\left(\frac{\partial}{\partial x_j} t(X)\right) = \frac{\partial}{\partial \mu_j} \mathbb{E}(t(X)) = \sigma^{-2} \mathbb{E}(t(X)(X_j - \mu_j)), \quad j = 1, \dots, p,$$

a result usually known as *Stein's equation*<sup>28</sup>. By taking  $t(X) = X_j/\|X\|^2$  and summing over  $j$  we realize that

$$(8.29) \quad \mathbb{E}\left(\frac{\langle X, X - \boldsymbol{\mu} \rangle}{\|X\|^2}\right) = (p-2)\sigma^2 \mathbb{E}(\|X\|^{-2}),$$

which gives

$$\text{mse}(\hat{\boldsymbol{\mu}}_{JS}) = p\sigma^2 - (p-2)^2 \sigma^4 \mathbb{E}(\|X\|^{-2}).$$

On the other hand, it follows from (2.6) that

$$\mathbb{E}(\|X\|^{-2}) = \frac{1}{(2\pi)^{p/2}} \int_{\mathbb{R}^p} \|\mathbf{x}\|^{-2} e^{-\frac{1}{2}\|\mathbf{x}-\boldsymbol{\mu}\|^2} d\mathbf{x},$$

an integral which becomes finite if  $\int r^{p-3} dr$  converges near  $r = 0$ . Thus, we conclude that  $\text{mse}_{\boldsymbol{\mu}}(\hat{\boldsymbol{\mu}}_{JS}) < p\sigma^2$  if  $p \geq 3$  for *any*  $\boldsymbol{\mu}$ , which confirms that in those cases the unbiased MLE estimator  $\hat{\boldsymbol{\mu}}$  fails to be the most efficient one (if the “performance” is measured by the mean squared error); cf. Remark 7.28. Regarding this remarkable estimator, we add the following comments.

<sup>28</sup>Remarkably enough, the validity of (8.28) for all  $t$  varying in a suitable class of test functions completely characterizes  $\sigma^{-1}(X - \boldsymbol{\mu})$  as a standard normal random vector, which turns out to be the starting point of Stein's approach to the Berry-Esseen theorem discussed in Remark 6.7 [Che21].

- Since  $\hat{\mu} = X$  is simply the sample mean as we have just a single observation at our disposal, this is in sharp contrast with the result in Example 8.20, which says that the sample mean is the best estimator for  $\mu$  is  $p = 1$ . To reinforce this analogy, let us take a random sample  $X_j \sim \mathcal{N}(\mu, \sigma^2 \text{Id}_p)$ ,  $j = 1, \dots, n$ , so that  $n$  observations of the underlying multivariate normal population are available. Now, the log-likelihood function is

$$l(\mathbf{x}; \mu) = -\frac{np}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_j \|\mathbf{x}_j - \mu\|^2, \quad \mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{np},$$

so the MLE for  $\mu$  is  $\hat{\mu}^{(n)} = \bar{X}$ , where

$$\bar{X} = \frac{1}{n} \sum_j X_j \in \mathbb{R}^p.$$

Since  $\sigma^{-2}\|X_j - \mu\|^2 \sim \chi_p^2$  for each  $j$  and  $\{X_j - \mu\}_{j=1}^n$  is independent, we have

$$\begin{aligned} \text{mse}(\hat{\mu}^{(n)}) &= \mathbb{E}(\|\bar{X} - \mu\|^2) \\ &= \frac{1}{n^2} \mathbb{E} \left( \left\| \sum_j (X_j - \mu) \right\|^2 \right) \\ &= \frac{1}{n^2} \sum_j \mathbb{E}(\|X_j - \mu\|^2) \\ &= \frac{p}{n} \sigma^2, \end{aligned}$$

where again only the covariance contributes (since  $\hat{\mu}^{(n)}$  is unbiased). It turns out that essentially the same argument as above confirms that  $\hat{\mu}^{(n)}$  fails to be admissible if  $p \geq 3$ , as the corresponding James-Stein estimator

$$\hat{\mu}_{JS}^{(n)} = \left( 1 - \frac{(p-2)\sigma^2/n}{\|\hat{\mu}^{(n)}\|^2} \right) \hat{\mu}^{(n)} = \bar{X} - (p-2) \frac{\sigma^2}{n} \frac{\bar{X}}{\|\bar{X}\|^2}$$

satisfies  $\text{mse}(\hat{\mu}_{JS}^{(n)}) < p\sigma^2/n$ .

- In case  $\sigma^2$  is unknown, and restricting ourselves to the case  $n = 1$  for simplicity, let us replace (8.27) by

$$\hat{\mu}_{JS_l} = \left( 1 - \frac{(p-2)c_l \mathfrak{s}}{\|\hat{\mu}\|^2} \right) \hat{\mu} = X - (p-2) c_l \mathfrak{s} \frac{X}{\|X\|^2},$$

where  $p \geq 3$ ,  $c_l$  is a positive constant (depending on a positive integer  $l$  given in advance) to be determined below and  $\mathfrak{s}$  is the appropriate estimator of  $\sigma^2$  in the sense that  $\sigma^{-2}\mathfrak{s} \sim \chi_l^2$  and  $\{\mathfrak{s}, X\}$  is independent. Setting  $\mu_\bullet = \sigma^{-1}\mu$ ,  $\hat{\mu}_\bullet = \sigma^{-1}\hat{\mu}$ ,  $X_\bullet = \sigma^{-1}X$  and  $\mathfrak{s}_\bullet = \sigma^{-2}\mathfrak{s}$ , we compute

$$\begin{aligned} \text{mse}(\hat{\mu}_{JS_l}) &= \mathbb{E}(\|\hat{\mu}_{JS_l} - \mu\|^2) \\ &= \sigma^2 \mathbb{E} \left( \left\| \hat{\mu}_\bullet - \mu_\bullet - (p-2)c_l \mathfrak{s}_\bullet \frac{\hat{\mu}_\bullet}{\|\hat{\mu}_\bullet\|^2} \right\|^2 \right) \\ &= \sigma^2 \left( \text{mse}(\hat{\mu}_\bullet) - 2(p-2)c_l l \mathbb{E} \left( \frac{\langle X_\bullet, X_\bullet - \mu_\bullet \rangle}{\|X_\bullet\|^2} \right) + (p-2)^2 c_l^2 l(l+2) \mathbb{E}(\|X_\bullet\|^{-2}) \right), \end{aligned}$$

where we used the independence and that  $\mathbb{E}(\mathfrak{s}_\bullet) = l$  and  $\mathbb{E}(\mathfrak{s}_\bullet^2) = l(l+2)$  in the last step. Combining this with the obvious counterpart of (8.29) we end up with

$$\text{mse}(\hat{\mu}_{JS_l}) = \sigma^2 (p - (p-2)^2 l [2c_l - c_l^2(l+2)]) \mathbb{E}(\|X_\bullet\|^{-2}),$$

from which we see that the best choice is  $c_l = 1/(l + 2)$ , in which case

$$\hat{\mu}_{JS_l} = \sigma^2 \left( k - \frac{p-2}{l+2} \frac{\mathfrak{s}}{\|\hat{\mu}\|^2} \right) \hat{\mu}$$

certainly satisfies  $\text{mse}(\hat{\mu}_{JS_l}) < p\sigma^2$ .

- It is clear from the expressions above (as in (8.27), for instance) that the James-Stein estimator promotes a shrinkage of the sample mean towards the origin, which in particular introduces a small amount of bias into the most obvious estimate (by pulling it away from its observed value). At least in case the number  $p$  of entries of the underlying normal population mean is large enough, we have seen that by doing so we achieve a greater reduction in variance, resulting in an estimator with a lower total error as measured by the mean squared error<sup>29</sup>. We insist that this procedure (shrinkage) is not just helpful: it represented a paradigm shift in Statistics which directly led to the development of regularization techniques like Ridge and Lasso, which are indispensable tools for modern Data Science and Machine learning, especially in high-dimensional settings; see Subsection 9.3 for more on this point in the context of linear regression.
- Instead of shrinking toward the origin, it is sometimes more convenient to choose some  $\nu \in \mathbb{R}^k$  and replace (8.27) by

$$\hat{\mu}_{JS_\nu} = \left( 1 - \frac{(p-2)\sigma^2}{\|\hat{\mu} - \nu\|^2} \right) (\hat{\mu} - \nu) + \nu,$$

thus performing a shrinkage toward  $\nu$ , with the resulting estimator always satisfying  $\text{mse}(\hat{\mu}_{JS_\nu}) < \sigma^2 k$ . Although it is not clear which choice of  $\nu$  provides the better result, this certainly adds a lot of flexibility in performing the shrinkage, with a natural choice for  $\nu$  being the data-driven *grand mean vector*  $\bar{X}_{\text{gm}} \mathbf{1}$ , where  $\bar{X}_{\text{gm}}$  is the arithmetic mean of the components of the observed sample mean  $X$  (this is the choice used in the famous analysis of the baseball batting averages data set in [EM77]).  $\square$

**Remark 8.22.** By rewriting (8.24) as

$$\text{cov}(\hat{\theta}) \mathcal{F}(\theta) \geq \text{Id}_n,$$

it is patent the resemblance of the Cramér-Rao lower bound to the uncertainty principle in Quantum Mechanics.  $\square$

**8.3. Optimal asymptotic normality of ML estimators.** We now check that under suitable regularity assumptions (which are too complicated to reproduce here) the ML estimator achieves the Cramér-Rao lower bound as the sample size  $n$  grows indefinitely, which follows from the fact that consistent ML estimators are asymptotically normal (in the sense of Definition 7.12), with their asymptotic covariance  $\sigma_\theta^2$  determined by the (inverse of the) Fisher information matrix. As usual we consider an infinite family  $X_j \sim \psi_\theta$  of i.i.d. random variables, so that for each  $n$  the log-likelihood of  $X^{[n]} = (X_1, \dots, X_n)$  is given by

$$(8.30) \quad l^{(n)}(\mathbf{x}; \theta) = \sum_{j=1}^n \ln \psi_\theta(x_j),$$

where  $\theta \in \Theta$  is the true (but unknown) parameter. For simplicity, let us assume that  $\Theta \subset \mathbb{R}$  (the uni-dimensional case) so that  $\mathcal{F}(\theta) > 0$  is the Fisher information. For each  $n$  let  $\hat{\theta}_n$  be the corresponding (and unique!) ML estimator so that

$$(8.31) \quad \frac{d}{d\theta} l^{(n)}(\mathbf{x}; \hat{\theta}_n) = 0.$$

<sup>29</sup>This phenomenon also appears in Corollary 7.27, since we may view  $\hat{\sigma}_{(n+1)-1}^2$  as a shrinkage of both the unbiased estimator  $\hat{\sigma}_{(n-1)-1}^2$  and the ML estimator  $\hat{\sigma}_{n-1}^2$ .

**Theorem 8.23.** (Optimal asymptotic normality) Under the conditions above, if  $\hat{\theta}_n$  is consistent (in the sense of Definition 7.9) then

$$(8.32) \quad \sqrt{n}(\hat{\theta}_n - \theta) \rightarrow \mathcal{N}(0, \mathcal{F}_{(1)}(\theta)^{-1})$$

in distribution (with respect to  $\theta$ ), where  $\mathcal{F}_{(1)}$  is the Fisher information of a single observation (say,  $X_1$ ). As a consequence,

$$(8.33) \quad \hat{\theta}_n \approx_{n \rightarrow +\infty} \mathcal{N}(\theta, \mathcal{F}_{(n)}(\theta)^{-1}),$$

where  $\mathcal{F}_{(n)} = n\mathcal{F}_{(1)}$  is the Fisher information of the whole sample  $X^{[n]}$ .

**Remark 8.24.** (Asymptotic efficiency) It follows from (8.32) that the asymptotic variance of  $\{\hat{\theta}_n\}$  equals the Cramér-Rao lower bound (8.24) for the variance restricted to  $\mathcal{E}_{\text{id}}$ , where  $\text{id} : \Theta \rightarrow \mathbb{R}$  is the identity map. In other words, asymptotically the variance of  $\sqrt{n}(\hat{\theta}_n - \theta)$  is at least as small as the variance of any (sufficiently regular but not necessarily asymptotically normal) unbiased estimator. For obvious reasons, this is called *asymptotic efficiency*. Since  $\mathcal{F}_{(n)}(\theta)^{-1} = \mathcal{F}_{(1)}(\theta)^{-1}/n$ , in view of (8.33) this also means that the “fluctuation” of  $\hat{\theta}_n$  around  $\theta$ , as measured by its variance, decays with the sample size according to a rate which is inversely proportional to the reciprocal of the Fisher information of a single observation; compare with Example 8.12. More generally, if  $g : \Theta \subset \mathbb{R} \rightarrow \mathbb{R}$  is  $C^1$  with a nowhere vanishing derivative then Theorem 8.23 combines with Proposition 7.15 to yield

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) \rightarrow \mathcal{N}(0, |g'(\theta)|^2 \mathcal{F}_{(1)}(\theta)^{-1})$$

which by (8.23) means that the asymptotic variance of  $\{g(\hat{\theta}_n)\}$  approaches the Cramér-Rao lower bound (8.24) for the variance restricted to  $\mathcal{E}_g$ . Hence, this transformed estimator is asymptotically efficient in its own way.  $\square$

**Remark 8.25.** From (8.33) we know that

$$P(a \leq \hat{\theta}_n \leq b) \approx_{n \rightarrow +\infty} \sqrt{\frac{\mathcal{F}_{(n)}(\theta)}{2\pi}} \int_a^b e^{-\frac{\mathcal{F}_{(n)}(\theta)(x-\theta)^2}{2}} dx,$$

but we can rely on Theorem 2.23 to replace  $\theta$  by  $\hat{\theta}_n$  in the right-hand side because  $\hat{\theta}_n \rightarrow \theta$  in probability (consistency), so as to obtain

$$P(a \leq \hat{\theta}_n \leq b) \approx_{n \rightarrow +\infty} \sqrt{\frac{\mathcal{F}_{(n)}(\hat{\theta}_n)}{2\pi}} \int_a^b e^{-\frac{\mathcal{F}_{(n)}(\hat{\theta}_n)(x-\hat{\theta}_n)^2}{2}} dx.$$

The key point here is that the right-hand side depends solely on sample data, and only through the estimator  $\hat{\theta}$ . In the language of confidence intervals of Subsection 7.3, this translates into

$$(8.34) \quad \theta \in \left[ \hat{\theta}_n \mp \frac{z_{\delta/2}}{\sqrt{\mathcal{F}_{(n)}(\hat{\theta}_n)}} \right] \text{ with prob. } \approx 1 - \delta,$$

the “large sample” estimate for  $\theta$ .  $\square$

**Remark 8.26.** The consistency requirement in Theorem 8.23 may be often justified under suitable regularity assumptions on the underlying pdf's, which in particular apply to the ML estimator  $\hat{\sigma}_{n-1}^2$  in Example 8.4 [NM94, Theorem 2.5]. We may also directly retrieve the consistency of  $\hat{\sigma}_{n-1}^2$  as follows. First note from (7.23) that

$$(8.35) \quad \hat{\sigma}_{n-1}^2 = \frac{1}{n} \sum_{j=1}^n \sigma^2 \left( \frac{X_j - \mu}{\sigma} \right)^2 - (\bar{X}_n - \mu)^2.$$

Also, recalling that  $X_j$  is drawn from a normal population,  $\sigma^{-1}(X_j - \mu) \sim \mathcal{N}(0, 1)$  implies that  $\sigma^{-2}(X_j - \mu)^2 \sim \chi_1^2$  by Corollary 4.25 and hence

$$\mathbb{E} \left( \sigma^2 \left( \frac{X_j - \mu}{\sigma} \right)^2 \right) = \sigma^2$$

by Corollary 4.21. Thus, LLN (Theorem 6.2) applies to ensure that

$$\frac{1}{n} \sum_{j=1}^n \sigma^2 \left( \frac{X_j - \bar{X}_n}{\sigma} \right)^2 \xrightarrow{P} \sigma^2.$$

On the other hand, it also follows from LLN that  $(\bar{X}_n - \mu)^2 \xrightarrow{P} 0$ . Thus,  $\hat{\sigma}_{n-1}^2 \xrightarrow{P} \sigma^2$  by (8.35). Another approach to this same conclusion follows by taking  $c = n^{-1}$  in (7.25) to check that  $\text{mse}(\hat{\sigma}_{n-1}^2) \rightarrow 0$  as  $n \rightarrow +\infty$ , so that consistency follows by Proposition 7.11. In fact, any of these methods may be adapted to check that the ML estimators in Examples 8.6 and 8.7 above are consistent as well. We also note that, in general, Proposition 7.13 applies to ensure that consistency is a necessary condition for asymptotic normality.  $\square$

*Proof.* (of Theorem 8.23) Set  $\tilde{l}^{(n)} = n^{-1}l^{(n)}$  and note that by (8.31) and the Mean Value Theorem,

$$(8.36) \quad 0 = \frac{d}{d\theta} \tilde{l}^{(n)}(\hat{\theta}_n) = \frac{d}{d\theta} \tilde{l}^{(n)}(\theta) + \frac{d^2}{d\theta^2} \tilde{l}^{(n)}(\theta_n^\bullet)(\hat{\theta}_n - \theta),$$

for some  $\theta_n^\bullet$  lying between  $\hat{\theta}_n$  and  $\theta$ . A computation shows that for any  $\theta$  we have

$$\frac{d^2}{d\theta^2} \tilde{l}^{(n)}(\theta) = \frac{1}{n} \sum_{j=1}^n \left( \frac{d^2}{d\theta^2} \ln \psi_\theta(X_j) \right) \rightarrow \mathbb{E} \left( \frac{d^2}{d\theta^2} \ln \psi_\theta(X_1) \right),$$

where the convergence is in probability by LLN. Since  $\hat{\theta}_n \rightarrow \theta$  in probability, we conclude that

$$(8.37) \quad \frac{d^2}{d\theta^2} \tilde{l}^{(n)}(\theta_n^\bullet) \rightarrow \mathbb{E} \left( \frac{d^2}{d\theta^2} \ln \psi_\theta(X_1) \right) = -\mathcal{F}_{(1)}(\theta),$$

where the convergence is in probability and we used (8.21) in the last step. On the other hand,

$$\sqrt{n} \frac{d}{d\theta} \tilde{l}^{(n)}(\theta) = \sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n \frac{d}{d\theta} \ln \psi_\theta(X_j) \right),$$

which may be rewritten as

$$\sqrt{n} \frac{d}{d\theta} \tilde{l}^{(n)}(\theta) = \sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n \frac{d}{d\theta} \ln \psi_\theta(X_j) - \mathbb{E} \left( \frac{d}{d\theta} \ln \psi_\theta(X_1) \right) \right),$$

as the term within the expectation is a score (Corollary 8.14). Thus we may apply CLT to see that

$$(8.38) \quad \sqrt{n} \frac{d}{d\theta} \tilde{l}^{(n)}(\theta) \rightarrow \mathcal{N} \left( 0, \text{cov} \left( \frac{d}{d\theta} \ln \psi_\theta(X_1) \right) \right) = \mathcal{N}(0, \mathcal{F}_{(1)}(\theta)),$$

where the convergence is in distribution and the last step follows from the definition of  $\mathcal{F}_{(1)}$ . Since (8.36) leads to

$$\sqrt{n}(\hat{\theta}_n - \theta) = - \frac{\sqrt{n} \frac{d}{d\theta} \tilde{l}^{(n)}(\theta)}{\frac{d^2}{d\theta^2} \tilde{l}^{(n)}(\theta_n^\bullet)},$$

we may use Theorem 2.23, (8.37) and (8.38) to complete the proof.  $\square$

We now discuss the implications of this theory for some of the statistical models discussed earlier.

**Example 8.27.** We start with

$$\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2) = (\bar{X}_n, \hat{\sigma}_{n-1}^2),$$

the ML estimator for the bi-dimensional parameter  $(\theta_1, \theta_2) = (\mu, \sigma^2)$  coming from a normal population as in Example 8.4 above. We first look only at  $\hat{\theta}_2$ , which amounts to declaring that  $\mu$  is known. As observed in Remark 8.26, this estimator is consistent and hence asymptotically normal by Theorem 8.23. In order to determine the associated limiting normal distribution by means of Theorem 8.23 we need to recall the corresponding Fisher information. From (8.19) with  $n = 1$ ,

$$(8.39) \quad \mathcal{F}_{(1)}(\theta_2) = \mathcal{F}_{(1)}(\theta)_{22} = \frac{1}{2\theta_2^2}, \implies \mathcal{F}_{(n)}(\theta_2) = \frac{n}{2\theta_2^2},$$

so (8.32) and (8.33) apply to give

$$(8.40) \quad \sqrt{n}(\hat{\sigma}_{n-1}^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma^4)$$

and

$$\hat{\sigma}_{n-1}^2 \approx_{n \rightarrow +\infty} \mathcal{N}(\sigma^2, 2\sigma^4/n).$$

Thus, in view of (8.39), (8.34) translates into

$$(8.41) \quad \sigma^2 \in \left[ \left(1 - \sqrt{\frac{2}{n}} z_{\delta/2}\right) \hat{\sigma}_{n-1}^2, \left(1 + \sqrt{\frac{2}{n}} z_{\delta/2}\right) \hat{\sigma}_{n-1}^2 \right] \text{ with prob. at least } 1 - \delta.$$

We next consider the bi-dimensional case  $\theta = (\theta_1, \theta_2)$ . Again by (8.19),

$$\mathcal{F}_{(1)}(\theta_1, \theta_2) = \begin{pmatrix} \sigma^2 & 0 \\ 0 & 2\sigma^4 \end{pmatrix},$$

which gives

$$(8.42) \quad \sqrt{n} \left( \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & 2\sigma^4 \end{pmatrix} \right),$$

or equivalently,

$$(8.43) \quad \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} \approx_{n \rightarrow +\infty} \mathcal{N} \left( \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}, \begin{pmatrix} \sigma^2/n & 0 \\ 0 & 2\sigma^4/n \end{pmatrix} \right).$$

A key point here is that the asymptotic covariance matrix in (8.42) only depends on  $\theta_2 = \sigma^2$ , which allows us to proceed as in Remark 8.25: consistency allows us to replace  $\sigma^2$  by  $\hat{\theta}_2$  in the covariance matrix of (8.43) to obtain

$$\begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} \approx_{n \rightarrow +\infty} \mathcal{N} \left( \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix}, \begin{pmatrix} \hat{\theta}_2/n & 0 \\ 0 & 2\hat{\theta}_2^2/n \end{pmatrix} \right),$$

an asymptotic estimate in which the “fluctuation” around the center  $(\hat{\theta}_1, \hat{\theta}_2)$  of the “confidence region” where the unknown parameter  $(\mu, \sigma^2)$  is supposed to be (with probability at least  $1 - \delta$ ) depends on sample data, and only through the estimator  $\hat{\theta}_2 = \hat{\sigma}_{n-1}^2$ .  $\square$

**Example 8.28.** Let us refer to the notation and terminology of Examples 7.39 and 8.5, with the (simplifying and justifiable) assumption that  $\mu_X = \mu_Y = 0$ , so that  $\theta = (\sigma_X^2, \sigma_Y^2, \rho) \in \mathbb{R}^+ \times \mathbb{R}^+ \times (-1, 1)$ . In order to determine the asymptotic behaviour of the ML estimator  $\hat{\theta}_m = (\hat{\sigma}_{m-1}^2(X), \hat{\sigma}_{m-1}^2(Y), \hat{\rho}_m)$  for the unknown population parameter  $\theta$  as  $m \rightarrow +\infty$  (for a given jointly normal sample  $\{X_j, Y_j\}$ ) we start with (8.11) and, after a somewhat tedious computation, we end up with a complicated expression for the  $3 \times 3$  matrix  $\nabla_{\theta\theta} l(X, Y; \theta)$  whose entries

depend *linearly* on the symbols in (8.10) evaluated on the sample, with the corresponding coefficients being algebraic on the components of  $\theta$ . Using that

$$\mathbb{E}(A(X_j)) = \sigma_X^2, \quad \mathbb{E}(B(X_j, Y_j)) = \sigma_{XY} = \rho\sigma_X\sigma_Y, \quad \mathbb{E}(C(Y_j)) = \sigma_Y^2,$$

and (8.21) we conclude that the corresponding Fisher information matrix is

$$\mathcal{F} = \frac{1}{1 - \rho^2} \begin{pmatrix} \frac{2 - \rho^2}{4\sigma_X^4} & -\frac{\rho^2}{4\sigma_X^2\sigma_Y^2} & -\frac{\rho}{2\sigma_X^2} \\ -\frac{\rho}{4\sigma_X^2\sigma_Y^2} & \frac{2 - \rho^2}{4\sigma_Y^4} & -\frac{\rho}{2\sigma_Y^2} \\ -\frac{\rho}{2\sigma_X^2} & -\frac{\rho}{2\sigma_Y^2} & \frac{1 + \rho^2}{1 - \rho^2} \end{pmatrix},$$

so that

$$\mathcal{F}^{-1} = \begin{pmatrix} 2\sigma_X^4 & 2\rho^2\sigma_X^2\sigma_Y^2 & \rho(1 - \rho^2)\sigma_X^2 \\ 2\rho^2\sigma_X^2\sigma_Y^2 & 2\sigma_Y^4 & \rho(1 - \rho^2)\sigma_Y^2 \\ \rho(1 - \rho^2)\sigma_X^2 & \rho(1 - \rho^2)\sigma_Y^2 & (1 - \rho^2)^2 \end{pmatrix}.$$

From this and Theorem 8.23 we thus derive not only that

$$\sqrt{m}(\hat{\sigma}_{m-1}^2(X) - \sigma_X^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma_X^2), \quad \sqrt{m}(\hat{\sigma}_{m-1}^2(Y) - \sigma_Y^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma_Y^2),$$

which are fully compatible with (8.40), but also that

$$(8.44) \quad \sqrt{m}(\hat{\rho}_m - \rho) \xrightarrow{d} \mathcal{N}(0, (1 - \rho^2)^2),$$

which identifies the asymptotic variance of the sample correlation coefficient  $\hat{\rho}_m$  as being  $(1 - \rho^2)^2$ ; cf. Definition 7.12. Moreover, Remark 8.24 guarantees that this asymptotic invariance equals the Crámer-Rao lower bound (8.24) for the variances of all *unbiased* estimators for  $\rho$ . We remark, however, that  $\hat{\rho}_m$  itself is *not* unbiased even though there holds

$$F(\hat{\rho}_m) - \hat{\rho}_m \xrightarrow{p} 0$$

for any *unbiased* estimator of the form  $F(\hat{\rho}_m)$ , where  $F$  is assumed to be odd [OP58]. In particular,

$$\sqrt{m}(F(\hat{\rho}_m) - \rho) \xrightarrow{d} \mathcal{N}(0, (1 - \rho^2)^2).$$

As in Example 8.27 above, we may combine (8.44) with the consistency of  $\hat{\rho}_m$  to construct a large sample confidence interval for the unknown correlation coefficient  $\rho$ , namely,

$$\rho \in \left[ \hat{\rho}_m \mp z_{\delta/2} \frac{1 - \hat{\rho}_m^2}{\sqrt{m}} \right] \text{ with prob. } \approx 1 - \delta$$

whose “fluctuation” around its center  $\hat{\rho}_m$  depends only on sample data, and through the estimator  $\hat{\rho}_m$ . An alternate route is to apply the delta method (Proposition 7.15) with the *Fisher z-transformation*

$$z = g(\rho) = \frac{1}{2} \ln \left( \frac{1 + \rho}{1 - \rho} \right)$$

to (8.44) so as to get

$$\sqrt{m}(\hat{z}_m - z) \xrightarrow{d} \mathcal{N}(0, 1),$$

which allows us to obtain large sample estimates for  $z$  in terms of the familiar normal quantiles  $z_{\delta/2}$  and then transform them back to corresponding estimates for  $\rho = \tanh z$ ; see [Ken46, Section 14.18] and [And03, Subsection 4.2.3]  $\square$

**Example 8.29.** (The coefficient of variation of a normal population) Let  $g : \Theta \rightarrow \mathbb{R}$  be any smooth function satisfying  $\nabla g \neq \vec{0}$  everywhere, where  $\Theta = \mathbb{R} \times \mathbb{R}^+$  is the parameter space of a normal population as above; cf.



Example 8.4. Recall that  $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$  in this case. It then follows from (8.42) and the multi-dimensional version of the delta method in Proposition 7.15 that

$$(8.45) \quad \sqrt{n} \left( g \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - g \begin{pmatrix} \mu \\ \sigma^2 \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( 0, \nabla g(\theta)^t \begin{pmatrix} \sigma^2 & 0 \\ 0 & 2\sigma^4 \end{pmatrix} \nabla g(\theta) \right).$$

Assuming that  $\mu \neq 0$ , we may apply this to  $g(\theta) = \sqrt{\theta_2}/\theta_1 = \sigma/\mu$ , the *coefficient of variation*; cf. (4.14). Since

$$\nabla g(\mu, \sigma^2)^t = \left( -\frac{\sigma}{\mu^2}, \frac{1}{2\mu\sigma} \right)$$

we obtain that the corresponding estimator,  $\hat{\sigma}_{n-1}/\bar{X}_n$ , is asymptotically normal,

$$\sqrt{n} \left( \frac{\hat{\sigma}_{n-1}}{\bar{X}_n} - \frac{\sigma}{\mu} \right) \xrightarrow{d} \mathcal{N} \left( 0, \frac{\sigma^2}{\mu^2} \left( \frac{1}{2} + \frac{\sigma^2}{\mu^2} \right) \right),$$

with the asymptotic variance depending on  $\sigma/\mu$  itself. As usual, we may use consistence to get

$$\frac{\sigma}{\mu} \in \left[ \frac{\hat{\sigma}_{n-1}}{\bar{X}_n} \mp z_{\delta/2} \frac{\hat{\sigma}_{n-1}}{\sqrt{n}\bar{X}_n} \left( \frac{1}{2} + \frac{\hat{\sigma}_{n-1}^2}{\bar{X}_n^2} \right)^{1/2} \right] \text{ with prob. } \approx 1 - \delta,$$

a large sample confidence interval estimate for the coefficient of variation.  $\square$

**Example 8.30.** (MLE for a Gamma population) It follows from (4.22) that the log-likelihood function of a Gamma distribution  $\Gamma_{\alpha, \lambda}$  is

$$(8.46) \quad l(\mathbf{x}; \theta) = n \left( \lambda \ln \alpha - \ln \Gamma(\lambda) + (\lambda - 1) \overline{\ln x} - \alpha \bar{x} \right), \quad \theta = (\alpha, \lambda),$$

where  $\overline{\ln x}$  is the mean of  $(\ln x_1, \dots, \ln x_n)$ ; recall that  $x_j > 0$  for each  $j$ . Hence, the score vector is

$$(8.47) \quad s(\mathbf{x}; \theta) = n \begin{pmatrix} \frac{\lambda}{\alpha} - \bar{x} \\ \ln \alpha - \psi(\lambda) + \overline{\ln x} \end{pmatrix}, \quad \psi(\lambda) = \frac{d}{d\lambda} \ln \Gamma(\lambda),$$

so the ML estimator  $\hat{\theta} = (\hat{\alpha}, \hat{\lambda})$  satisfies

$$(8.48) \quad \begin{cases} \frac{\hat{\lambda}}{\hat{\alpha}} &= \bar{x} \\ \psi(\hat{\lambda}) - \ln \hat{\alpha} &= \overline{\ln x} \end{cases}$$

Note that trying to find a solution for this system in closed form is out of question so possible strategies here are:

- to use our favorite optimization package to find

$$\hat{\theta} = \operatorname{argmax}_{\theta} l(\mathbf{x}; \theta)$$

starting from (8.46);

- to use the first equation in (8.48) to eliminate  $\hat{\alpha}$  in the second equation, *numerically* solve for  $\hat{\lambda}$  in the resulting equation, namely,

$$(8.49) \quad \psi(\hat{\lambda}) - \ln \hat{\lambda} = \overline{\ln x} - \ln \bar{x},$$

and then replacing the result back in the first equation in order to get  $\hat{\alpha}$ <sup>30</sup>.

<sup>30</sup>That a unique solution  $\hat{\lambda}$  to (8.49) exists for any given  $x$  is a consequence of the facts that i)  $\overline{\ln x} < \ln \bar{x}$  if each  $x_j > 0$ ; ii) the function  $\lambda \mapsto \Psi(\lambda) = \psi(\lambda) - \ln \lambda$  is monotone continuous and satisfies

$$\lim_{\lambda \rightarrow 0} \Psi(\lambda) = -\infty \text{ and } \lim_{\lambda \rightarrow \infty} \Psi(\lambda) = 0.$$

In any case, with the ML estimator so determined, we may proceed to compute the associated Fisher information matrix by means of (8.47) and (8.21):

$$\mathcal{F}_{(n)}(\theta) = n \begin{pmatrix} \lambda/\alpha^2 & -1/\alpha \\ -1/\alpha & \psi_1(\lambda) \end{pmatrix} = n\mathcal{F}_{(1)}(\theta), \quad \psi_1 = d\psi/d\lambda.$$

It is not hard to check that  $\det \mathcal{F}_{(1)}(\theta) = (\lambda\psi_1(\lambda) - 1)/\alpha^2 > 0$ , so Theorem 8.23 gives asymptotic normality for  $\hat{\theta}_n$ :

$$\sqrt{n} \left( \begin{pmatrix} \hat{\alpha}_n \\ \hat{\lambda}_n \end{pmatrix} - \begin{pmatrix} \alpha \\ \lambda \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \vec{0}, \mathcal{F}_{(1)}^{-1} \right) = \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \frac{1}{\lambda\psi_1(\lambda) - 1} \begin{pmatrix} \alpha^2\psi_1(\lambda) & \alpha \\ \alpha & \lambda \end{pmatrix} \right).$$

As usual, may combine this with consistency in order to obtain

$$(8.50) \quad \begin{pmatrix} \hat{\alpha}_n \\ \hat{\lambda}_n \end{pmatrix} \approx_{n \rightarrow +\infty} \mathcal{N} \left( \begin{pmatrix} \alpha \\ \lambda \end{pmatrix}, \mathcal{F}_{(n)}^{-1} \right) = \mathcal{N} \left( \begin{pmatrix} \alpha \\ \lambda \end{pmatrix}, \frac{1}{n(\hat{\lambda}_n\psi_1(\hat{\lambda}_n) - 1)} \begin{pmatrix} \hat{\alpha}_n^2\psi_1(\hat{\lambda}_n) & \hat{\alpha}_n \\ \hat{\alpha}_n & \hat{\lambda}_n \end{pmatrix} \right),$$

which may be used to construct not only large sample confidence intervals for  $\alpha$  and  $\lambda$  (separately) but also large sample confidence regions for the whole vector parameter  $\theta$ ; see Remark 8.31 below for this latter kind of construction.

**Remark 8.31.** (Confidence region for the unknown parameter  $\theta$  via asymptotic normality of the ML estimator) Starting with the (possibly multivariate) version of (8.33), where  $\theta \in \mathbb{R}^p$ ,  $p \geq 1$ , and  $\mathcal{F}_{(n)}(\theta)$  is a  $p \times p$  symmetric, positive definite matrix, consistency of the ML estimator  $\hat{\theta}_n$  leads to the asymptotic normality relation

$$\hat{\theta}_n \approx_{n \rightarrow +\infty} \mathcal{N}(\theta, \mathcal{F}_{(n)}(\hat{\theta}_n)^{-1}),$$

from which (8.50) is a rather special case (with  $p = 2$ ). In order to extract from this a confidence region for the unknown vector parameter  $\theta$ , let us write  $\mathcal{F}_{(n)}(\hat{\theta}_n) = A^t A$  so that Remark 4.15 gives

$$A(\hat{\theta}_n - \theta) \sim \mathcal{N}(0, \text{Id}_p)$$

and hence

$$(\hat{\theta}_n - \theta)^t \mathcal{F}_{(n)}(\hat{\theta}_n) (\hat{\theta}_n - \theta) = \|A(\hat{\theta}_n - \theta)\|^2 \sim \chi_p^2.$$

In other words, the quadratic form in the left-hand side is a pivotal quantity to which the standard method may be applied: if  $\chi_{p,\alpha}^2$  is the quantile of  $\chi_p^2$  associated to  $0 < \alpha < 1$  then

$$P \left( (\hat{\theta}_n - \theta)^t \mathcal{F}_{(n)}(\hat{\theta}_n) (\hat{\theta}_n - \theta) \leq \chi_{p,\alpha}^2 \right) \approx 1 - \alpha.$$

Since  $\mathcal{F}_{(n)}(\hat{\theta}_n)$  is positive definite, the random confidence region where  $\theta$  is supposed to lie (within the given confidence level) is ellipsoidal in nature, with its size, shape, and orientation being completely determined by the totality of the elements of  $\mathcal{F}_{(n)}(\hat{\theta}_n)$ . Moreover, since its construction takes into account the possible correlations among the various components of  $\hat{\theta}_n$ , as encoded in the off-diagonal elements of the asymptotic covariance matrix  $\mathcal{F}_{(n)}(\hat{\theta}_n)^{-1}$ , in such cases it certainly encloses a much tighter volume than the  $p$ -cube which is the product of the separate confidence intervals for the entries of  $\theta$ .

**Remark 8.32.** When comparing Theorem 8.23 with Theorem 6.5 we see that “sample universality”, a treasured feature of this latter classical accomplishment, has been irremediably lost. Indeed, Theorem 8.23 roughly says that for each choice of the log-likelihood function as in (8.30), which by its turn is completely determined by the underlying pdf  $\psi(\cdot, \theta)$  via (8.7) and (8.2), the MLE method selects an estimator as in (8.5) for which a corresponding “limit theorem” holds as in (8.32). Thus, differently from CLT, Theorem 8.23 is model-dependent.  $\square$

## 9. THE METHOD OF LEAST SQUARES

If  $\theta_2 = \sigma^2$  is known, maximizing  $l$  in (8.9) is equivalent to minimizing

$$\theta_1 \mapsto \frac{1}{2} \sum_j (x_j - \theta_1)^2,$$

which furnishes a variational characterization of the arithmetic mean  $n^{-1} \sum_j x_j$ . This is of course a manifestation of the Method of Least Squares (MLS), a celebrated procedure which provides a solution to the following kind of problem. Let us arrange the outcome of  $n$  measurements of  $p$  features (regressors, independent/explanatory variables, etc.) of a population by means of the  $n \times (p+1)$ -matrix

$$\mathbf{x} = \begin{pmatrix} 1 & \vdots & \mathbf{x}_1 \\ & \vdots & \vdots \\ 1 & \vdots & \mathbf{x}_n \end{pmatrix}$$

where each

$$\mathbf{x}_j = (x_{j1}, \dots, x_{jp}), \quad j = 1, \dots, n,$$

is a row  $p$ -vector (representing the outcome of the  $j^{\text{th}}$  measurement) and  $\mathbf{1}$  is the column  $n$ -vector whose entries all equal 1. If we suspect that these features relate to a response (regressand, dependent/explained variable, etc.) which has also been measured, thus yielding an  $n$ -vector  $\mathbf{y}$ , we may try to “predict” the response at some unknown feature by “best fitting” a (possibly non-linear) functional dependence, say  $\mathbf{y} = F(\mathbf{x})$ , to the available data  $(\mathbf{x}, \mathbf{y})$ . The simplest choice is to postulate that  $F$  is *linear*, so that  $\mathbf{y} = \mathbf{x}\hat{\beta}$ , where  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$  is determined by minimizing the corresponding quadratic objective function:

$$(9.1) \quad \hat{\beta} = \operatorname{argmin}_{\beta} f(\beta), \quad f(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}\beta\|^2,$$

hence the “least squares” terminology.

This is a purely geometric problem (best fitting a hyperplane to a cloud of  $n$  points in  $\mathbb{R}^p \times \mathbb{R} = \mathbb{R}^{p+1}$ , where we usually assume that  $p+1 \ll n$ ), which can be solved by the methods of Calculus. Indeed, since

$$(\nabla f)(\beta) = -\mathbf{x}^t (\mathbf{y} - \mathbf{x}\beta),$$

where  $t$  means transpose, we obtain

$$(9.2) \quad \hat{\beta} = (\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t \mathbf{y},$$

where we assume that  $\mathbf{x}$  has full column-rank (this not only implies that  $p+1 \leq n$  but also that the *Gram matrix*  $\mathbf{x}^t \mathbf{x}$  is symmetric and positive definite, hence invertible). Moreover, since

$$\nabla^2 f = \mathbf{x}^t \mathbf{x},$$

we conclude that  $\hat{\beta}$  is the unique global minimum. Under these conditions, we then say that

$$(9.3) \quad \hat{\mathbf{y}} = \mathbf{x}\hat{\beta}$$

is the *fitted vector* (that is, the vector of fitted values). Finally, we observe that  $\hat{\beta}$  is *linear* in  $\mathbf{y}$  with coefficients depending on  $\mathbf{x}$ .

**9.1. The statistical model behind MLS.** In the examples below, we discuss the statistical rationale behind the purely data-driven minimization problem in (9.1). In alignment with the general estimation setting in Subsection 7.1, this involves imposing convenient assumptions on how the data array  $(\mathbf{y}, \mathbf{r})$  has been drawn from an underlying population, so as to be able to set up a statistical model in which  $\hat{\beta}$  gets confirmed as an efficient estimator for the corresponding population parameter, say  $\beta$ , which is usually regarded as unknown<sup>31</sup>.

**Example 9.1.** (The general regression model) We start by viewing  $(\mathbf{r}, \mathbf{y})$  as the realization of a  $\mathbb{R}^{n \times (p+1) + n}$ -valued random vector  $(\mathfrak{X}, \mathbf{Y})$  in  $L^2(\Omega)$  and satisfying

- $\mathfrak{X} = (\mathbf{1} \ \mathbf{X})$ , where  $\mathbf{X}$  is a random  $\mathbb{R}^{n \times p}$ -valued vector whose realization is  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^t$ . In other words,  $\mathfrak{X}$  is a random matrix whose first column is deterministic (non-random) and equals  $\mathbf{1}$ . Also, we assume that  $\mathfrak{X}$  has full column-rank a.s.
- $\{(\mathbf{X}_{j\bullet}, \mathbf{Y}_j)\}_{j=1}^n$ , where  $\bullet \in \{1, \dots, p\}$ , are i.i.d. copies of the same  $\mathbb{R}^p \times \mathbb{R}$ -valued random vector, say  $(\mathcal{X}, \mathcal{Y})$ .

Under these conditions, we may initially impose the (not necessarily linear) *regression model*

$$(9.4) \quad \mathbf{Y} = F(\mathfrak{X}) + \mathbf{e},$$

with the *regression function*  $F : \mathbb{R}^{n(p+1)} \rightarrow \mathbb{R}^n$  being defined by

$$(9.5) \quad F(\mathbf{r}) = \mathbb{E}(\mathbf{Y} | \mathfrak{X} = \mathbf{r}),$$

where we use here the notation of Subsection 3.1; see Remark 9.4 below for the justification of this choice of  $F$ , where it is shown that it minimizes the corresponding mean squared error:

$$\mathbb{E}(\|\mathbf{Y} - F(\mathfrak{X})\|^2) \leq \mathbb{E}(\|\mathbf{Y} - G(\mathfrak{X})\|^2),$$

for any  $G : \mathbb{R}^{n(p+1)} \rightarrow \mathbb{R}^n$ . Thus, we may view (9.4) as the definition of the random *error*  $\mathbf{e}$ , which by Proposition 3.14 may also be expressed as

$$\mathbf{e} = \mathbf{Y} - \mathbb{E}(\mathbf{Y} | \mathfrak{X}),$$

so that (3.14) easily implies *exogeneity*,

$$(9.6) \quad \mathbb{E}(\mathbf{e} | \mathfrak{X}) = 0,$$

and hence

$$(9.7) \quad \mathbb{E}(\mathbf{e}) = 0.$$

More generally, again by (3.14),

$$\mathbb{E}(\mathfrak{X}^t \mathbf{e}) = \mathbb{E}(\mathbb{E}(\mathfrak{X}^t \mathbf{e} | \mathfrak{X})) = \mathbb{E}(\mathfrak{X}^t \mathbb{E}(\mathbf{e} | \mathfrak{X})),$$

where we used Proposition 3.11 (7) in the last step, so that (9.6) applies to give

$$(9.8) \quad \mathbb{E}(\mathfrak{X}^t \mathbf{e}) = 0,$$

so (9.4) decomposes  $\mathbf{Y}$  as a sum of a term  $F(\mathfrak{X})$  which is “explained” by  $\mathfrak{X}$  and an error which has zero mean (conditioned to  $\mathfrak{X}$ ) and is uncorrelated to (any function of)  $\mathfrak{X}$ <sup>32</sup>. Also, if we define the *error covariance function* by

$$\sigma^2(\mathbf{r}) := \text{cov}(\mathbf{e} | \mathfrak{X} = \mathbf{r}) \stackrel{(9.6)}{=} \mathbb{E}(\mathbf{e}^t \mathbf{e} | \mathfrak{X} = \mathbf{r}), \quad \mathbf{r} \in \mathbb{R}^{n(p+1)},$$

<sup>31</sup>For a critical appraisal of model building and interpretation in Regression Analysis, with an emphasis on the distinction between the amount of information already present in the data and the inferential consequences of suitable assumptions on the sampling process by which these data have been drawn from a hypothetical population, and which in a sense applies to any other kind of statistical analysis, we refer to [Ber04]

<sup>32</sup>The implications of this remarkable decomposition to Regression Theory (and to Econometrics, in particular) are discussed at length in [AP09, Chapter 3].

then

$$\text{cov}(\mathbf{e}) \stackrel{(9.7)}{=} \mathbb{E}(\mathbf{e}^t \mathbf{e}) \stackrel{(3.14)}{=} \mathbb{E}(\mathbb{E}(\mathbf{e}^t \mathbf{e} | \mathfrak{X})),$$

and using Proposition 3.14,

$$\text{cov}(\mathbf{e}) = \mathbb{E}(\boldsymbol{\sigma}^2(\mathfrak{X})).$$

In words, the unconditioned error covariance equals the average of the conditioned error covariance.  $\square$

We now specialize the general setup above to the cases which appear more frequently in applications.

**Example 9.2.** (The linear regression model) The simplest of all choices for the regression function above is  $F(\mathfrak{x}) = \mathfrak{x}\beta$ , which gives rise to the *linear regression model*

$$(9.9) \quad \mathbf{Y} = \mathfrak{X}\beta + \mathbf{e},$$

where

$$(9.10) \quad \begin{aligned} \beta &= \operatorname{argmin}_{\beta' \in \mathbb{R}^{p+1}} \frac{1}{n} \mathbb{E}(\|\mathbf{Y} - \mathfrak{X}\beta'\|^2) \\ &= \operatorname{argmin}_{\beta' \in \mathbb{R}^{p+1}} \mathbb{E}(\left(\mathcal{Y} - \widetilde{\mathcal{X}}^t \beta'\right)^2), \quad \widetilde{\mathcal{X}} = (1, \mathcal{X}), \end{aligned}$$

provides the best linear fitting for  $\mathbf{Y}$  (or  $\mathcal{Y}$ ) in the  $L^2$  sense. In this setting, (9.8) should be interpreted as the “projection condition” that  $\mathbf{e} = \mathbf{Y} - \mathfrak{X}\beta$  should be “orthogonal” (again in the  $L^2$  sense) to any potential linear fitting; see Remark 9.22 for an elaboration of this viewpoint. Also, (9.6) is automatically satisfied due to Remark 3.9. Now, if we apply the usual first order test from Calculus to (9.10) we find that

$$\beta = \mathbb{E}(|\widetilde{\mathcal{X}}|^2)^{-1} \mathbb{E}(\mathcal{Y} \widetilde{\mathcal{X}}^t),$$

but this does not say much about the true nature of  $\beta$  because the joint distribution of  $(\mathcal{X}, \mathcal{Y})$  remains unknown, which makes the expectations intractable. Thus, this population parameter should somehow be estimated from data (a realization  $(\mathbf{x}, \mathbf{y})$  of  $(\mathbf{X}, \mathbf{Y})$ ) with  $\widehat{\beta}$  in (9.2) being the most obvious candidate for an estimator. As it is always the case with any estimator, its efficiency is only effectively certified by the establishment of good inferential properties (say, by requiring that its mse is minimized within a given class of estimators and/or that it is consistent and asymptotically normal, etc.; see the general discussion in Subsection 7.1), so with this purpose in mind it is convenient to add to (9.6) the assumption of *spherical error*, which means that there exists  $\sigma > 0$  such that

$$(9.11) \quad \text{cov}(\mathbf{e} | \mathfrak{x} = \mathfrak{x}) = \sigma^2 \text{Id}_n, \quad \text{independently of } \mathfrak{x}.$$

Thus,

$$\mathbb{E}(\mathbf{Y} | \mathfrak{x} = \mathfrak{x}) = \mathfrak{x}\beta \quad \text{and} \quad \text{cov}(\mathbf{Y} | \mathfrak{x} = \mathfrak{x}) = \sigma^2 \text{Id}_n$$

summarize the assumptions of the linear regression model<sup>33</sup>. In the language of Example 9.1, (9.11) means that the random matrix  $\boldsymbol{\sigma}^2(\mathfrak{X})$  is actually *constant* and equals  $\sigma^2 \text{Id}_n$ , an artifact also known as *homoscedasticity*.  $\square$

**Example 9.3.** (The linear regression model with a normal error) In the setting of the linear regression model (9.9), the “empirical” quadratic minimization in (9.1) may be justified via MLE under a normality assumption on the error<sup>34</sup>. Precisely, and in alignment with (9.6) and (9.11), let us further assume that the error  $\mathbf{e}$  is such that  $\{\mathbf{e}_j | \mathfrak{x}_j = \mathfrak{x}_j\}_{j=1}^n$  is independent and distributed according to

$$(9.12) \quad \mathbf{e}_j | \mathfrak{x}_j = \mathfrak{x}_j \sim \mathcal{N}(0, \sigma^2),$$

<sup>33</sup>Although the stronger assumption of the independence of  $\{\mathbf{e} | \mathfrak{x}_j = \mathfrak{x}_j\}_{j=1}^n$  may eventually be useful (as in Examples 9.23 and 9.24, for instance), we stress that only uncorrelatedness, as expressed by (9.11), is imposed at this point, as this already allows us to derive some nice inferential properties for  $\widehat{\beta}$ ; cf. Propositions 9.6 and 9.7 and Remark 9.19. In any case, if  $\mathbf{e} | \mathfrak{x} = \mathfrak{x}$  is normally distributed, as in Example 9.3, then these assumptions (uncorrelatedness and independence) are equivalent indeed (by Corollary 4.11).

<sup>34</sup>As it is well-known, this connection between MLS and the normal distribution has been first observed by Gauss and Laplace [Sti90, Chapter 4].

or equivalently,

$$(9.13) \quad \mathbf{e}|\mathbf{x}=\mathbf{r} \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_n),$$

by Proposition 4.10. It then follows from (9.9) that  $\{\mathbf{Y}_j|\mathbf{x}_j=\mathbf{r}_j\}_{j=1}^n$  is independent with

$$(9.14) \quad \mathbf{Y}_j|\mathbf{x}_j=\mathbf{r}_j \sim \mathcal{N}\left(\sum_{k=0}^p \mathbf{r}_{jk}\beta_k, \sigma^2\right).$$

Therefore, we have been able to construct an identifiable statistical model (albeit one “conditioned” on the observed value  $\mathbf{r}$  of  $\mathbf{x}$ ) in which  $\beta$  appears as the unknown parameter; compare with the extended notion of a statistical model in Remark 7.4 and note that for the moment we regard  $\sigma^2$  as known. In particular, we may apply MLE to (9.14), as in Definition 8.3, to find the corresponding estimator. Now, again by Proposition 4.10,

$$\mathbf{Y}|\mathbf{x}=\mathbf{r} \sim \mathcal{N}(\mathbf{r}\beta, \sigma^2 \text{Id}_n),$$

so the likelihood function of (9.14) is

$$(9.15) \quad L(\mathbf{y}; \beta) = (2\pi\sigma^2)^{-n/2} e^{-\frac{\|\mathbf{y}-\mathbf{r}\beta\|^2}{2\sigma^2}}.$$

Since the corresponding log-likelihood function to be maximized is

$$(9.16) \quad l(\mathbf{y}; \beta) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{\|\mathbf{y}-\mathbf{r}\beta\|^2}{2\sigma^2},$$

we see that, up to irrelevant constants (depending on  $\sigma^2$ ), solving this maximization problem is equivalent to finding  $\hat{\beta}$  as in (9.1), thus confirming that MLE implies MLS under the stated assumptions.  $\square$

**Remark 9.4.** The appearance of the regression function  $F$  in (9.5) may be justified by the fact that

$$\mathbb{E}(\|\mathbf{Y} - F(\mathbf{x})\|^2) = \inf_G \mathbb{E}(\|\mathbf{Y} - G(\mathbf{x})\|^2),$$

for any  $G : \mathbb{R}^{n(p+1)} \rightarrow \mathbb{R}^n$  measurable. To check this, first note that

$$\begin{aligned} \mathbb{E}(\|\mathbf{Y} - G(\mathbf{x})\|^2) &= \mathbb{E}(\|\mathbf{Y} - F(\mathbf{x}) + F(\mathbf{x}) - G(\mathbf{x})\|^2) \\ &= \mathbb{E}(\|\mathbf{Y} - F(\mathbf{x})\|^2) + \mathbb{E}(\|F(\mathbf{x}) - G(\mathbf{x})\|^2) \\ &\quad + 2\mathbb{E}(\langle \mathbf{Y} - F(\mathbf{x}), F(\mathbf{x}) - G(\mathbf{x}) \rangle). \end{aligned}$$

Also, by Proposition 3.11 (2) and Proposition 3.14,

$$\begin{aligned} \mathbb{E}(\langle \mathbf{Y} - F(\mathbf{x}), F(\mathbf{x}) - G(\mathbf{x}) \rangle) &= \mathbb{E}(\mathbb{E}(\langle \mathbf{Y} - F(\mathbf{x}), F(\mathbf{x}) - G(\mathbf{x}) \rangle | \mathbf{x})) \\ &= \mathbb{E}(\langle \mathbf{Y} - F(\mathbf{x}), F(\mathbf{x}) - G(\mathbf{x}) \rangle | \mathbf{x}=\mathbf{r}) \\ &= \langle F(\mathbf{r}) - G(\mathbf{r}), \mathbb{E}(\mathbf{Y} | \mathbf{x}=\mathbf{r}) - F(\mathbf{r}) \rangle \\ &= 0. \end{aligned}$$

Thus,

$$\mathbb{E}(\|\mathbf{Y} - G(\mathbf{x})\|^2) \geq \mathbb{E}(\|\mathbf{Y} - F(\mathbf{x})\|^2),$$

with the equality holding if and only if  $G(\mathbf{x}) = F(\mathbf{x})$  a.s.  $\square$

**Remark 9.5.** It is already clear from the discussion above that we may either view  $\mathbf{X}$  as random, thus definitely contributing to the randomness of  $\mathbf{Y}$ , or condition on any of its observations, so it becomes fixed (i.e. non-random) with the only contribution to the randomness of  $\mathbf{Y}$  coming from the error. Depending on the applications we have in mind, either viewpoint may be adopted. In an experimental field (such as agriculture), where the experimenter usually controls the variable  $\mathbf{x}$  and subsequently observes  $\mathbf{y}$ , the latter option is preferable. On the other hand, in a more socially oriented science (Econometrics, for instance) no such control is expected (or even acceptable) and  $\mathbf{X}$  should be regarded as random. Fortunately, the algebraic expressions

of the relevant inferential statistics seem to be insensitive to which perspective we take, notably under a suitable normality assumption [RS08, Chapter 10]. Since normality is mostly assumed in the following, we find it convenient to eliminate any reference to conditioning on  $\mathfrak{X} = \mathfrak{r}$ , although  $\mathbf{X}$  should be considered as random whenever issues of a more foundational nature are discussed (as in Remark 9.16). Thus, as a rule we use the same symbol to denote a random variable and its observed value, as this turns the notation much cleaner and (hopefully!) will cause no confusion.  $\square$

**9.2. Inference and goodness of fit for MLS.** With a statistical model for MLS at hand, we now proceed to the pertinent inferential analysis. We start by observing that although the normality assumption for the error in (9.13) is essential for interpreting MLS via MLE, good statistical properties of the associated estimator  $\hat{\beta}$  may be derived under the much less stringent assumptions of the linear regression model in Example 9.2.

**Proposition 9.6.** *Under the conditions of Example 9.2 there hold*

$$(9.17) \quad \mathbb{E}(\hat{\beta}) = \beta, \quad \text{cov}(\hat{\beta}) = \sigma^2(\mathbf{r}^t \mathbf{r})^{-1}.$$

*As a consequence,  $\hat{\beta}$  is unbiased and  $\text{mse}(\hat{\beta}) = \sigma^2 \text{tr}(\mathbf{r}^t \mathbf{r})^{-1}$ .*

*Proof.* With the simplifying notation suggested by Remark 9.5, we are assuming that

$$(9.18) \quad \mathbb{E}(\mathbf{e}_j) = 0 \quad \text{and} \quad \text{cov}(\mathbf{e}_j, \mathbf{e}_k) = \sigma^2 \delta_{jk}, \quad j, k = 1, \dots, n.$$

Hence, from (9.2),

$$\hat{\beta} = (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t (\mathbf{r} \beta + \mathbf{e}),$$

so that

$$(9.19) \quad \hat{\beta} = \beta + (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \mathbf{e},$$

which gives

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(\beta) + (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \mathbb{E}(\mathbf{e}) = \beta + (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t (\vec{0}) = \beta.$$

Also,

$$\begin{aligned} \text{cov}(\hat{\beta}) &= \mathbb{E}((\hat{\beta} - \beta)(\hat{\beta} - \beta)^t) \\ &= \mathbb{E}(((\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \mathbf{e})(\mathbf{e}^t \mathbf{r} (\mathbf{r}^t \mathbf{r})^{-1})) \\ &= (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \mathbb{E}(\mathbf{e} \mathbf{e}^t) \mathbf{r} (\mathbf{r}^t \mathbf{r})^{-1} \\ &= (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \text{cov}(\mathbf{e}) \mathbf{r} (\mathbf{r}^t \mathbf{r})^{-1} \\ &= \sigma^2 (\mathbf{r}^t \mathbf{r})^{-1}, \end{aligned}$$

as desired.  $\square$

As a first check on the efficiency of MLS estimator  $\hat{\beta}$  in (9.2), let us see how it competes with a general *linear* estimator

$$\bar{\beta} = C \mathbf{Y},$$

where  $C$  is a  $(p+1) \times n$  matrix which is allowed to depend on  $\mathfrak{r}$  but not on  $\mathbf{Y}$ . This leads to a remarkable result confirming that  $\hat{\beta}$  attains the best performance (as measured by the mse) within a natural class of estimators.

**Theorem 9.7.** *(Gauss-Markov) Let  $\bar{\beta}$  as above be unbiased with  $\mathbf{e}$  satisfying (9.18). Then  $\text{cov}(\bar{\beta}) \geq \text{cov}(\hat{\beta})$ .*

*Proof.* We write

$$\bar{\beta} = ((\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t + D) \mathbf{Y},$$

where  $D$  has the same properties as  $C$ . It follows that

$$\bar{\beta} = ((\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t + D)(\mathbf{x}\beta + \mathbf{e}) = \hat{\beta} + D\mathbf{x}\beta + D\mathbf{e}$$

with  $\mathbb{E}(D\mathbf{e}) = D\mathbb{E}(\mathbf{e}) = 0$  (by Proposition 4.9 (3)) so that

$$\mathbb{E}(\bar{\beta}) = \beta + D\mathbf{x}\beta,$$

and letting  $\beta$  vary we see that a vanishing bias for  $\bar{\beta}$  implies  $D\mathbf{x} = 0$ . Hence,

$$\text{cov}(\bar{\beta}) = \text{cov}(\hat{\beta}) + \text{cov}(D\mathbf{e}) + 2\text{cov}(\hat{\beta}, D\mathbf{e}).$$

Now note that  $\text{cov}(D\mathbf{e}) = D\text{cov}(\mathbf{e})D^t = \sigma^2 DD^t$ . Moreover, using (9.19) and the fact that  $\beta$  is non-random,

$$\begin{aligned} \text{cov}(\hat{\beta}, D\mathbf{e}) &= \text{cov}(\beta, D\mathbf{e}) + \text{cov}((\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t \mathbf{e}, D\mathbf{e}) \\ &= \text{cov}((\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t \mathbf{e}, D\mathbf{e}) \\ &= (\mathbf{x}^t \mathbf{x})^{-1} \mathbf{x}^t \text{cov}(\mathbf{e}, \mathbf{e}) D^t \\ &= \sigma^2 (\mathbf{x}^t \mathbf{x})^{-1} (D\mathbf{x})^t \\ &= 0. \end{aligned}$$

Thus,

$$\text{cov}(\bar{\beta}) = \text{cov}(\hat{\beta}) + \sigma^2 DD^t,$$

and the result follows because  $DD^t \geq 0$ .  $\square$

**Example 9.8.** If we take it for granted that  $\mathbf{X}$  has no influence whatsoever on  $\mathbf{Y}$  then we are actually dealing with the “intercept-only” case  $\beta = (\beta_0, 0, \dots, 0)$ , so we must impose  $\hat{\beta} = (\hat{\beta}_0, 0, \dots, 0)$  and, as expected, (9.2) gives

$$(9.20) \quad \hat{\beta}_0 = \bar{\mathbf{Y}} = \frac{1}{n} \sum_j \mathbf{Y}_j,$$

the sample mean of  $\mathbf{Y}$ ; in the simple linear regression case of Example 9.9 below, this is immediate from (9.21). If  $w = (w_1, \dots, w_n) \in \mathbb{R}^n$  we may consider the more general linear combination in the entries of  $\mathbf{Y}$  given by

$$\hat{\beta}_0^w = \sum_j w_j \mathbf{Y}_j = \hat{\beta}_0 + \sum_j \left( w_j - \frac{1}{n} \right) \mathbf{Y}_j,$$

so that

$$\begin{aligned} \mathbb{E}(\hat{\beta}_0^w) &= \beta_0 + \sum_j \left( w_j - \frac{1}{n} \right) \mathbb{E}(\mathbf{Y}_j) \\ &= \beta_0 + \sum_j \left( w_j - \frac{1}{n} \right) (\mathbf{x}\beta)_j \\ &= \beta_0 + \beta_0 \sum_j \left( w_j - \frac{1}{n} \right), \end{aligned}$$

and  $\hat{\beta}_0^w$  is unbiased if and only if  $w$  is a weight vector,  $\sum_j w_j = 1$ . Thus, Gauss-Markov applies to ensure that  $\hat{\beta}_0$  attains the best performance among all these *weighted* estimators of  $\beta_0$ . In particular, Gauss-Markov may be regarded as a generalization of Example 7.17.  $\square$



**Example 9.9.** (Simple linear regression) If  $p = 1$  in Example 9.2 then  $\mathfrak{X} = (\mathbf{1}, \mathbf{X})$ , where  $\mathbf{X} = (\mathbf{X}_{11}, \dots, \mathbf{X}_{n1})^t$ , so that

$$\mathfrak{X}^t \mathfrak{X} = \begin{pmatrix} n & n\bar{\mathbf{X}} \\ n\bar{\mathbf{X}} & \|\mathbf{X}\|^2 \end{pmatrix}, \quad \bar{\mathbf{X}} = \frac{1}{n} \sum_j \mathbf{X}_{j1}.$$

Since  $\mathbf{X}$  is supposed not to be a multiple of  $\mathbf{1}$  a.s., Cauchy-Schwartz implies that

$$\det \mathfrak{X}^t \mathfrak{X} = n\|\mathbf{X}\|^2 - n^2\bar{\mathbf{X}}^2 > 0,$$

so that  $\mathfrak{X}^t \mathfrak{X}$  is invertible and

$$(\mathfrak{X}^t \mathfrak{X})^{-1} = \frac{1}{n\|\mathbf{X}\|^2 - n^2\bar{\mathbf{X}}^2} \begin{pmatrix} \|\mathbf{X}\|^2 & -n\bar{\mathbf{X}} \\ -n\bar{\mathbf{X}} & n \end{pmatrix}.$$

A little computation using (9.2) then gives

$$(9.21) \quad \hat{\beta} := \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{Y}} - \hat{\beta}_1 \bar{\mathbf{X}} \\ S_{\mathbf{X}\mathbf{Y}}/S_{\mathbf{X}\mathbf{X}} \end{pmatrix},$$

where

$$\bar{\mathbf{Y}} = \frac{1}{n} \sum_j \mathbf{Y}_j$$

is the sample mean of  $\mathbf{Y}$  and

$$(9.22) \quad S_{\mathbf{X}\mathbf{Y}} = \sum_j (\mathbf{X}_{j1} - \bar{\mathbf{X}})(\mathbf{Y}_j - \bar{\mathbf{Y}}), \quad S_{\mathbf{X}\mathbf{X}} = \sum_j (\mathbf{X}_{j1} - \bar{\mathbf{X}})^2.$$

Thus, the second line in (9.21) is used to compute the *slope*  $\hat{\beta}_1$  from sample data whereas the first line determines the *intercept*  $\hat{\beta}_0$ . In this case, the fitted value is realized as

$$(9.23) \quad \hat{\mathbf{y}} = \hat{\beta}_0 + \hat{\beta}_1 \mathbf{x}.$$

Also,  $\hat{\beta}$  is unbiased with

$$\text{cov}(\hat{\beta}) = \frac{\sigma^2}{n\|\mathbf{x}\|^2 - n^2\bar{\mathbf{x}}^2} \begin{pmatrix} \|\mathbf{x}\|^2 & -n\bar{\mathbf{x}} \\ -n\bar{\mathbf{x}} & n \end{pmatrix},$$

under the assumptions of Proposition 9.6. □

We now briefly discuss the construction of confidence intervals for the entries of the unknown parameter  $\beta$ . Here we remain in the setting of Example 9.3, so we assume that, conditionally on  $\mathfrak{X} = \mathfrak{x}$ ,  $\{\mathbf{e}_j\}$  is independent with  $\mathbf{e}_j \sim \mathcal{N}(0, \sigma^2)$  as in (9.13). It follows from (9.17) and (9.19) that we may use the pivotal quantity

$$(9.24) \quad \frac{\hat{\beta}_j - \beta_j}{\mathbf{s}_j} \sim \mathcal{N}(0, 1), \quad \mathbf{s}_j := \sigma \sqrt{\mathfrak{s}_{jj}}, \quad \mathfrak{s} := (\mathfrak{x}^t \mathfrak{x})^{-1},$$

to exhibit confidence intervals for the unknown parameter  $\beta_j$  in case  $\sigma$  is known. Precisely, in the notation of Subsection 7.3,

$$(9.25) \quad \beta_j \in \left[ \hat{\beta}_j \mp z_{\delta/2} \mathbf{s}_j \right] \text{ with prob. } = 1 - \delta.$$

Otherwise, we proceed as follows. We define the *residual*

$$(9.26) \quad \hat{\mathbf{e}} := \mathbf{Y} - \hat{\mathbf{Y}},$$

where  $\hat{\mathbf{y}} = \mathfrak{x} \hat{\beta}$  is the fitted vector as in (9.3). As we shall see,  $\|\hat{\mathbf{e}}\|^2/(n - p - 1)$  qualifies as an appropriate estimator for the error variance  $\sigma^2$ .

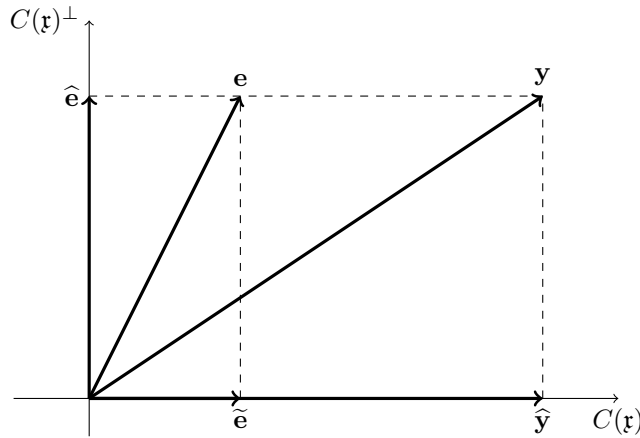


FIGURE 1. The geometry of the linear model

**Proposition 9.10.**  $\hat{\mathbf{e}} \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_{n-p-1})$  and  $\|\hat{\mathbf{e}}\|^2 / \sigma^2 \sim \chi_{n-p-1}^2$ ,  $n \geq p + 2$ .

*Proof.* We compute

$$\hat{\mathbf{e}} = (\text{Id} - \mathbf{r}(\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t) \mathbf{y} = Q(\mathbf{r} \beta + \mathbf{e}),$$

where  $Q = \text{Id}_n - \mathbf{r}(\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t$  is an idempotent, symmetric matrix satisfying

$$(9.27) \quad Q\mathbf{r} = 0,$$

which means that  $\text{rank } Q = n - p - 1 \geq 1$ . Therefore,  $\hat{\mathbf{e}} = Q\mathbf{e}$  is normally distributed as in the statement (either by Proposition 4.9 (3) or by rotational invariance (Corollary 4.12)). Moreover, since

$$(9.28) \quad \frac{\|\hat{\mathbf{e}}\|^2}{\sigma^2} = \left\langle \frac{\mathbf{e}}{\sigma}, Q \left( \frac{\mathbf{e}}{\sigma} \right) \right\rangle,$$

the last assertion follows from Proposition 4.27.  $\square$

**Remark 9.11.** (The geometry of the linear model and regression diagnostics) The projection matrix<sup>35</sup>  $H = \mathbf{r}(\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t$  appearing in the argument above is usually called the “hat matrix”, as it projects  $\mathbf{y}$  onto  $\hat{\mathbf{y}} = H\mathbf{y} \in C(\mathbf{r}) \equiv \mathbb{R}^{p+1}$ , the *design space*, which is the  $(p + 1)$ -subspace of  $\mathbb{R}^n$  spanned by the columns of the design matrix  $\mathbf{r}$ . On the other hand, its complementary projection matrix  $Q = \text{Id}_n - H$  projects  $\mathbf{y}$  (and also  $\mathbf{e}$ , because  $\mathbf{y} - \mathbf{e} = \mathbf{r}\beta \in C(\mathbf{r})$ ) onto the residual  $\hat{\mathbf{e}} \in C(\mathbf{r})^\perp \equiv \mathbb{R}^{n-p-1}$  lying in the orthogonal complement of  $C(\mathbf{r})$ . This nice “orthogonal” geometry, which hinges on the general setting of Example 9.2, is depicted in Figure 1, where  $\tilde{\mathbf{e}} = H\mathbf{e}$ . In particular, the orthogonal decomposition  $\mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}}$  clearly implies that the *sample correlation* between  $\hat{\mathbf{y}}$  and  $\hat{\mathbf{e}}$ , defined by

$$(9.29) \quad \text{corr}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) = \frac{\hat{\mathbf{e}}^t \hat{\mathbf{y}}}{\|\hat{\mathbf{e}}\| \|\hat{\mathbf{y}}\|},$$

vanishes, which justifies the common practice of using a scatterplot of the residuals against the fitted values in order to identify patterns of goodness of fit (or lack thereof) of a given linear model, as far as linearity and

<sup>35</sup>Recall that this means that  $H$  is symmetric and idempotent, hence defining an orthogonal projection onto its range.

homoscedasticity go [Far06, Section 6.1]. Now, if we further specialize to the setting of Example 9.3, which assumes  $\mathbf{e} \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_n)$ , then, by the projection property in (4.7),

$$(9.30) \quad \hat{\mathbf{e}} \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_{n-p-1}),$$

so in particular  $\hat{\mathbf{e}}/|\hat{\mathbf{e}}|$  is uniformly distributed in  $\mathbb{S}^{n-p-2} \subset C(\mathfrak{r})^\perp$  by Remark 4.26. Since  $\hat{\mathbf{e}}$  is accessible from data and adjusted values, graphical methods (say, a Q-Q plot) may be used to confirm the empirical validity of (9.30), which somehow works as an indirect checking of the theoretical assumption on the normality of errors underlying Example 9.3; again, see [Far06, Section 6.1]. A further gauging of the goodness of fit of the model may be implemented after properly combining the residual and the fitted vector. The simplest way of doing this, which leads to a sharpening of (9.29), is to look at the joint distribution of  $(\hat{\mathbf{y}}, \hat{\mathbf{e}})$  under error normality. Since

$$\begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} = \begin{pmatrix} H\mathbf{y} \\ Q\mathbf{y} \end{pmatrix} = \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix} \begin{pmatrix} \mathbf{y} \\ \mathbf{y} \end{pmatrix},$$

we see that  $(\hat{\mathbf{y}}, \hat{\mathbf{e}})$  is jointly normally distributed. To find the specific normal distribution we note that

$$\mathbb{E} \left( \begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} \right) = \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix} \mathbb{E} \begin{pmatrix} \mathbf{y} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix} \begin{pmatrix} \mathfrak{r}\beta \\ \mathfrak{r}\beta \end{pmatrix} = \begin{pmatrix} \mathfrak{r}\beta \\ 0 \end{pmatrix}$$

and

$$\text{cov} \left( \begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} \right) = \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix} \text{cov} \begin{pmatrix} \mathbf{y} \\ \mathbf{y} \end{pmatrix} \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix}^t = \sigma^2 \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix},$$

so that

$$\text{cov} \left( \Lambda \begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} \right) = \sigma^2 \Lambda \begin{pmatrix} H & 0 \\ 0 & Q \end{pmatrix} \Lambda^t,$$

with  $\Lambda$  being any  $2n \times 2n$  matrix. By choosing  $\Lambda$  orthogonal with the corresponding conjugation performing the appropriate diagonalization and viewing  $(\hat{\mathbf{y}}, \hat{\mathbf{e}})$  as an element of  $\mathbb{R}^{p+1} \times \mathbb{R}^{n-p-1} = \mathbb{R}^n$ , we find that there exists an orthogonal  $n \times n$  matrix  $\Lambda'$  such that

$$\text{cov} \left( \Lambda' \begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} \right) = \sigma^2 \text{Id}_n.$$

By Remark 4.15, and viewing  $\mathfrak{r}\beta$  as an element of  $\mathbb{R}^{p+1}$ ,

$$\Lambda' \begin{pmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{pmatrix} \sim \mathcal{N} \left( \Lambda' \begin{pmatrix} \mathfrak{r}\beta \\ 0 \end{pmatrix}, \sigma^2 \text{Id}_n \right),$$

and therefore  $\{\hat{\mathbf{y}}, \hat{\mathbf{e}}\}$  is independent by Corollary 4.12. Further uses of the residual in the art of quantifying the goodness of fit of the linear model may be found in Remark 9.16 below.  $\square$

We now come back to the business of finding confidence intervals for the entries of  $\beta$ , this time with  $\sigma^2$  regarded as unknown. In this case, the next result justifies the replacement of  $\mathbf{s}_j$  in (9.24) by

$$(9.31) \quad \hat{\mathbf{s}}_j := \hat{\sigma} \sqrt{\mathbf{s}_{jj}}, \quad \hat{\sigma} := \frac{\|\hat{\mathbf{e}}\|}{\sqrt{n-p-1}}.$$

Note that with this notation, Proposition 9.10 says that

$$(9.32) \quad (n-p-1) \frac{\hat{\sigma}^2}{\sigma^2} = \frac{\|\hat{\mathbf{e}}\|^2}{\sigma^2} \sim \chi_{n-p-1}^2.$$

**Proposition 9.12.**  $\hat{\sigma}^2$  is an unbiased and consistent estimator for  $\sigma^2$  (as  $n \rightarrow +\infty$  and  $p$  is held fixed). In particular,  $\hat{\sigma}$  is consistent for the error standard deviation  $\sigma$ . Moreover,  $\{\hat{\beta}, \hat{\sigma}^2\}$  is independent with

$$\frac{\hat{\beta}_j - \beta_j}{\hat{s}_j} \sim t_{n-p-1}.$$

*Proof.* From (9.32) and Corollary 4.21 we have  $\mathbb{E}(\hat{\sigma}^2) = \sigma^2$ , that is,  $\text{bias}(\hat{\sigma}^2) = 0$ . Also,

$$(9.33) \quad \text{var}(\hat{\sigma}^2) = \frac{2\sigma^4}{n-p-1},$$

so that  $\text{mse}(\hat{\sigma}^2) \rightarrow 0$  as  $n \rightarrow +\infty$  and consistency follows from Proposition 7.11. Moreover, since

$$\frac{\hat{\beta}_j - \beta_j}{\hat{s}_j} = \frac{\frac{\hat{\beta}_j - \beta_j}{s_j}}{\sqrt{\|\hat{e}\|^2/\sigma^2}},$$

the last assertion follows from (9.24), (9.32) and Proposition 4.30 once one verifies that  $\{\hat{\beta}, \|\hat{e}\|^2\}$  is independent. To check this, note that (9.19) gives

$$\mathbf{r}^t \mathbf{r} \left( \frac{\hat{\beta} - \beta}{\sigma} \right) = \mathbf{r}^t \left( \frac{\mathbf{e}}{\sigma} \right),$$

which together with (9.28), (9.27) and Proposition 4.28 implies that  $\{\mathbf{r}^t \mathbf{r} \hat{\beta}, \|\hat{e}\|^2\}$  is independent, from which the result follows (because  $\mathbf{r}^t \mathbf{r}$  is invertible).  $\square$

Thus, again using the notation of Subsection 7.3,

$$(9.34) \quad \beta_j \in \left[ \hat{\beta}_j \mp t_{n-p-1, \delta/2} \hat{s}_j \right] \text{ with prob. } \approx 1 - \delta,$$

a confidence interval estimate for  $\beta_j$  in case  $\sigma$  is unknown. Notice that if  $n - p - 1 \gg 0$  then we can replace  $t_{n-p-1, \delta/2}$  by  $z_{\delta/2}$  with a negligible error; this uses Remark 6.4.

**Example 9.13.** (Confidence region for the whole vector parameter  $\beta$ , with  $\sigma^2$  unknown) Pick  $\mathbf{p}$  so that  $\mathbf{p}^t \mathbf{p} = \mathbf{s}$  as in (9.24) and set  $\mathbf{n} = \sigma^{-1}(\mathbf{p}^t)^{-1}(\hat{\beta} - \beta)$ . It is immediate that  $\mathbf{n} \sim \mathcal{N}(\vec{0}, \text{Id}_{p+1})$  so that

$$\frac{(\hat{\beta} - \beta)^t \mathbf{s} (\hat{\beta} - \beta)}{\sigma^2} = \mathbf{n}^t \mathbf{n} \sim \chi_{p+1}^2$$

and is independent of  $\hat{\sigma}^2$ . Therefore, by Propositions 9.12 and 4.33,

$$\frac{(\hat{\beta} - \beta)^t \mathbf{s} (\hat{\beta} - \beta)}{(p+1)\hat{\sigma}^2} = \frac{\mathbf{n}^t \mathbf{n}/(p+1)}{\hat{\sigma}^2/\sigma^2} \sim F_{p+1, n-p-1},$$

which gives the “confidence region” estimate

$$(9.35) \quad \beta \in \mathcal{U}_{n,p,\delta}(\hat{\beta}; \hat{\sigma}^2) \text{ with prob. } 1 - \delta,$$

where

$$(9.36) \quad \mathcal{U}_{n,p,\delta}(\hat{\beta}; \hat{\sigma}^2) = \left\{ \beta' \in \mathbb{R}^{p+1}; (\hat{\beta} - \beta')^t \mathbf{s} (\hat{\beta} - \beta') \leq (p+1)\hat{\sigma}^2 t_{p+1, n-p-1, \delta} \right\}.$$

In other words, the random ellipsoidal region  $\mathcal{U}_{n,p,\delta}(\hat{\beta}; \hat{\sigma}^2)$ , which is fully specified by  $n, p, \delta$  and the sample data  $\hat{\beta}$  and  $\hat{\sigma}^2$ , covers the true vector parameter  $\beta$  with probability  $1 - \delta$ .  $\square$

**Example 9.14.** (Simultaneous confidence band for the mean response) According to [CW09, page 27], “the primary goal in a regression analysis is to understand, as far as possible with the available data, how the conditional distribution of the response varies across sub-populations determined by the possible values of the predictors”. Precisely, and using the notation of Example 9.1, if

$$\mathbf{R}^p = \{\mathbf{x}' = (1, x'_1, \dots, x'_p); x'_j \in \mathbb{R}, j = 1, \dots, p\}$$

and  $\mathbf{x} \in \mathbf{R}^p$  is given, this amounts to checking how much information on the conditioned random variable  $\mathcal{Y}|\widetilde{\mathcal{X}}=\mathbf{x}$  may be extracted from the data array  $(\mathbf{x}, \mathbf{y})$ . We insist that  $\mathbf{x}$  should be viewed as a future observation which has been consolidated from the same population *after* the data set has been drawn. From this perspective, the method of least squares in Example 9.2, according to which

$$(9.37) \quad \mathcal{Y}_{\mathbf{x}} := \mathcal{Y}|\widetilde{\mathcal{X}}=\mathbf{x} = \mathbf{x}^t \beta + \mathbf{e} \text{ with } \mathbb{E}(\mathbf{e}) = 0,$$

represents a first step toward this goal as it allows us to make use of the estimate  $\widehat{\beta}$ , which has been computed from  $(\mathbf{x}, \mathbf{y})$ , in order to retrieve information on the *mean response*

$$(9.38) \quad \mathbf{x}^t \beta = \mathbb{E}(\mathcal{Y}_{\mathbf{x}}),$$

a collection of summaries which, as  $\mathbf{x}$  varies, exhausts the realizations of the conditional expectation  $\mathbb{E}(\mathcal{Y}|\widetilde{\mathcal{X}})$  (by Proposition 3.14). Indeed, if we further specialize to the setting of Example 9.3, where

$$(9.39) \quad \mathbf{e} \sim \mathcal{N}(0, \sigma^2),$$

then  $\mathbf{x}^t \widehat{\beta}$  may be used to properly estimate the population parameter in (9.38): from (9.24) we have

$$(9.40) \quad \mathbf{x}^t (\widehat{\beta} - \beta) \sim \mathcal{N}(0, \sigma^2 \mathbf{x}^t \mathbf{s} \mathbf{x}), \quad \mathbf{s} = (\mathbf{r}^t \mathbf{r})^{-1},$$

that is,

$$\frac{\mathbf{x}^t (\widehat{\beta} - \beta)}{\sigma \sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}}} \sim \mathcal{N}(0, 1),$$

so that, by Propositions 9.12 and 4.30,

$$\frac{\mathbf{x}^t (\widehat{\beta} - \beta)}{\widehat{\sigma} \sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}}} \sim \mathbf{t}_{n-p-1},$$

and hence,

$$(9.41) \quad \mathbf{x}^t \beta \in \left[ \mathbf{x}^t \widehat{\beta} \mp \mathbf{t}_{n-p-1, \delta/2} \widehat{\sigma} \sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}} \right] \text{ with prob. } 1 - \delta.$$

If we allow for a bit more of dispersion, a “simultaneous” version of this pointwise bound is also available, with the corresponding estimate holding for *any*  $\mathbf{x} \in \mathbf{R}^n$ , as follows. Starting with (9.35) and (9.36) and using the notation of Example 9.13,

$$\begin{aligned} 1 - \delta &= P \left( \frac{(\widehat{\beta} - \beta)^t \mathbf{s} (\widehat{\beta} - \beta)}{(p+1) \widehat{\sigma}^2} \leq \mathbf{f}_{p+1, n-p-1, \delta} \right) \\ &= P \left( \frac{\|\mathbf{n}\|}{\widehat{\sigma}/\sigma} \leq \sqrt{(p+1) \mathbf{f}_{p+1, n-p-1, \delta}} \right), \end{aligned}$$

and recalling that, by Cauchy-Schwarz,

$$\sup_{\mathbf{x} \in \mathbf{R}^p} \frac{|(\mathbf{p}\mathbf{x})^t \mathbf{n}|}{\|\mathbf{p}\mathbf{x}\| \|\mathbf{n}\|} = 1,$$

we get

$$\begin{aligned} 1 - \delta &= P \left( \sup_{\mathbf{x} \in \mathbb{R}^p} \frac{|(\mathbf{p}\mathbf{x})^t \mathbf{n}|}{(\widehat{\sigma}/\sigma) \sqrt{(\mathbf{p}\mathbf{x})^t \mathbf{p}\mathbf{x}}} \leq \sqrt{(p+1)\mathbf{f}_{p+1, n-p-1, \delta}} \right) \\ &= P \left( \sup_{\mathbf{x} \in \mathbb{R}^p} \frac{|\mathbf{x}^t (\widehat{\beta} - \beta)|}{\widehat{\sigma} \sqrt{\mathbf{x}^t \mathbf{s}\mathbf{x}}} \leq \sqrt{(p+1)\mathbf{f}_{p+1, n-p-1, \delta}} \right), \end{aligned}$$

so that

$$(9.42) \quad \mathbf{x}^t \beta \in \left[ \mathbf{x}^t \widehat{\beta} \mp \sqrt{(p+1)\mathbf{f}_{p+1, n-p-1, \delta}} \widehat{\sigma} \sqrt{\mathbf{x}^t \mathbf{s}\mathbf{x}} \right] \forall \mathbf{x} \in \mathbb{R}^p \text{ with prob. } 1 - \delta,$$

As fully explained in [Liu10], this Scheffé-type simultaneous confidence band plays a fundamental role in the inference theory of MLS models; see also Example 11.14 for a generalization thereof.  $\square$

**Example 9.15.** (Simultaneous prediction band for the response) With the estimates for the mean response provided by (9.41) and (9.42) at hand, we may now “predict” where the response itself,

$$(9.43) \quad \mathcal{Y}_{\mathbf{x}} = \mathbf{x}^t \beta + \mathbf{e} \sim \mathcal{N}(\mathbf{x}^t \beta, \sigma^2),$$

is likely to fall, where, as in Example 9.14, we should think of  $\mathbf{x}$  as a new observation for the regressor which has taken place after the data  $(\mathbf{x}, \mathbf{y})$  has been gathered, hence the “prediction” terminology. In particular,  $\mathbf{e}$  is independent of  $\widehat{\beta}$ , which has been constructed out of  $(\mathbf{x}, \mathbf{y})$ , so that  $\mathcal{Y}_{\mathbf{x}}$  is independent of  $\widehat{\beta}$  as well. Combining this with (9.43), (9.40) and Proposition 4.7 (3) we see that

$$\mathcal{Y}_{\mathbf{x}} - \mathbf{x}^t \widehat{\beta} \sim \mathcal{N}(0, \sigma^2 (1 + \mathbf{x}^t \mathbf{s}\mathbf{x})),$$

which by the standard argument gives the sought-after pointwise “prediction interval” for the response,

$$(9.44) \quad \mathcal{Y}_{\mathbf{x}} \in \left[ \mathbf{x}^t \widehat{\beta} \mp t_{n-p-1, \delta/2} \widehat{\sigma} \sqrt{1 + \mathbf{x}^t \mathbf{s}\mathbf{x}} \right] \text{ with prob. } 1 - \delta,$$

with this new terminology being adopted because the random variable  $\mathcal{Y}_{\mathbf{x}}$  is *not* a parameter, so this fails to be a confidence interval in the ordinary sense. In any case, upon comparison with (9.41) we see that when passing from the mean response to the response itself, the point estimate  $\mathbf{x}^t \widehat{\beta}$  remains the same but the dispersion gets expanded by a factor that makes it at least as large as  $t_{n-p-1, \delta/2} \widehat{\sigma}$ , a lower bound which depends on the already observed data  $(\mathbf{x}, \mathbf{y})$  but *not* on the future observation  $\mathbf{x}$  for the regressor. Following [Car86, Theorem 1] and [SA90, Theorem 1], we may also contemplate a “simultaneous” version of (9.44), which is obtained by means of an easy generalization of Scheffé’s argument leading to (9.42). Indeed,

$$\mathbf{b} = \begin{pmatrix} \widehat{\beta} - \beta \\ \mathbf{e} \end{pmatrix} \sim \mathcal{N} \left( \vec{0}, \begin{pmatrix} \sigma^2 \mathbf{s} & 0 \\ 0 & \sigma^2 \end{pmatrix} \right)$$

and

$$\bar{\mathbf{s}} = \begin{pmatrix} \mathbf{s} & 0 \\ 0 & 1 \end{pmatrix}$$

are such that

$$\sigma^{-2} \mathbf{b}^t \bar{\mathbf{s}} \mathbf{b} = \sigma^{-2} (\widehat{\beta} - \beta)^t \mathbf{s} (\widehat{\beta} - \beta) + \sigma^{-2} \|\mathbf{e}\|^2 \sim \chi_{p+2}^2,$$

so that

$$\frac{\mathbf{b}^t \bar{\mathbf{s}} \mathbf{b}}{(p+2)\widehat{\sigma}^2} \sim F_{p+2, n-p-1},$$

and hence

$$\begin{aligned} 1 - \delta &= P \left( \frac{b^t \bar{s} b}{(p+2)\hat{\sigma}^2} \leq \mathbf{f}_{p+2, n-p-1, \delta} \right) \\ &= P \left( \frac{\|\bar{\mathbf{n}}\|}{\hat{\sigma}/\sigma} \leq \sqrt{(p+2)\mathbf{f}_{p+2, n-p-1, \delta}} \right), \end{aligned}$$

where  $\bar{\mathbf{n}} = \sigma^{-1}(\bar{\mathbf{p}}^t)^{-1}b$  with  $\bar{\mathbf{p}}$  satisfying  $\bar{\mathbf{p}}^t \bar{\mathbf{p}} = \bar{s}$ ; cf. the corresponding unbarred objects in Example 9.13. Thus, again using Cauchy-Schwarz,

$$\begin{aligned} 1 - \delta &= P \left( \sup_{\mathbf{z} \in \mathbb{R}^{p+1}} \frac{|(\bar{\mathbf{p}}\mathbf{z})^t \bar{\mathbf{n}}|}{(\hat{\sigma}/\sigma) \sqrt{(\bar{\mathbf{p}}\mathbf{z})^t \bar{\mathbf{p}}\mathbf{z}}} \leq \sqrt{(p+2)\mathbf{f}_{p+2, n-p-1, \delta}} \right) \\ &= P \left( \sup_{\mathbf{z} \in \mathbb{R}^{p+1}} \frac{|\mathbf{z}^t b|}{\hat{\sigma} \sqrt{\mathbf{z}^t \bar{s} \mathbf{z}}} \leq \sqrt{(p+2)\mathbf{f}_{p+2, n-p-1, \delta}} \right), \end{aligned}$$

where  $\mathbb{R}^{p+1} = \mathbb{R}^p \times \mathbb{R}$ , so if we choose  $\mathbf{z} = (\mathbf{x}, -1)$ , where  $\mathbf{x} \in \mathbb{R}^p$  is arbitrary, we get  $\mathbf{z}^t b = \mathbf{x}^t \hat{\beta} - \mathcal{Y}_{\mathbf{x}}$  and  $\mathbf{z}^t \bar{s} \mathbf{z} = 1 + \mathbf{x}^t \mathbf{s} \mathbf{x}$ , which finally gives

$$\mathcal{Y}_{\mathbf{x}} \in \left[ \mathbf{x}^t \hat{\beta} \mp \sqrt{(p+2)\mathbf{f}_{p+2, n-p-1, \delta}} \hat{\sigma} \sqrt{1 + \mathbf{x}^t \mathbf{s} \mathbf{x}} \right] \forall \mathbf{x} \in \mathbb{R}^p \text{ with prob. } 1 - \delta$$

as a simultaneous prediction band for the response.  $\square$

**Remark 9.16.** (Coefficient of determination and goodness of fit) As it is manifest from Figure 1, which depicts the residual  $\hat{\mathbf{e}}$  as sitting orthogonally to  $C(\mathbf{r})$ , the residual sum of squares

$$SS_{\text{Res}} := \|\hat{\mathbf{e}}\|^2 = \sum_j (\mathbf{Y}_j - \hat{\mathbf{Y}}_j)^2$$

considered above may be interpreted as a measure of how much variation in  $\mathbf{Y}$  is left unexplained by the linear regression model (which provides the fitted vector  $\hat{\mathbf{y}} \in C(\mathbf{r})$ ). Moreover, it appears in the decomposition

$$(9.45) \quad S_{\mathbf{Y}\mathbf{Y}} = SS_{\text{Reg}} + SS_{\text{Res}},$$

where

$$SS_{\text{Reg}} = \sum_j (\hat{\mathbf{Y}}_j - \bar{\mathbf{Y}})^2$$

corresponds to the amount of variation effectively explained by the model, and

$$S_{\mathbf{Y}\mathbf{Y}} = \sum_j (\mathbf{Y}_j - \bar{\mathbf{Y}})^2$$

is the total variation of  $\mathbf{Y}$  about its sample mean  $\bar{\mathbf{Y}}$  (the notation here is compatible with (9.22)). Since  $SS_{\text{Res}} \leq S_{\mathbf{Y}\mathbf{Y}}$ , it is natural to consider the *coefficient of determination*,

$$(9.46) \quad R^2 = \frac{SS_{\text{Reg}}}{S_{\mathbf{Y}\mathbf{Y}}} = 1 - \frac{SS_{\text{Res}}}{S_{\mathbf{Y}\mathbf{Y}}},$$

as a measure of *goodness of fit*: the closer its observed value is to its maximal value 1 the better we should regard the model. A simple argument justifying this assertion may be readily produced in the setting of the simple

linear regression in Example 9.9. Indeed,

$$\begin{aligned}
 \hat{\mathbf{e}} &= \mathbf{Y} - (\hat{\beta}_0 \mathbf{1} + \hat{\beta}_1 \mathbf{X}) \\
 &= \mathbf{Y} - \left( (\bar{\mathbf{Y}} - \hat{\beta}_1 \bar{\mathbf{X}}) \mathbf{1} + \hat{\beta}_1 \mathbf{X} \right) \\
 &= \mathbf{Y} - \bar{\mathbf{Y}} \mathbf{1} - \hat{\beta}_1 (\mathbf{X} - \bar{\mathbf{X}} \mathbf{1}) \\
 &= \mathbf{Y} - \bar{\mathbf{Y}} \mathbf{1} - \frac{S_{\mathbf{XY}}}{S_{\mathbf{XX}}} (\mathbf{X} - \bar{\mathbf{X}} \mathbf{1}),
 \end{aligned}$$

and from this we easily derive that

$$SS_{\text{Res}} = \frac{S_{\mathbf{XX}} S_{\mathbf{YY}} - S_{\mathbf{XY}}^2}{S_{\mathbf{XX}}}.$$

Eliminating  $SS_{\text{Res}}$  from this and (9.46) we see that

$$R^2 = \hat{\rho}^2,$$

where

$$\hat{\rho} = \frac{S_{\mathbf{XY}}}{\sqrt{S_{\mathbf{XX}} S_{\mathbf{YY}}}}$$

is the *sample correlation coefficient* of the random pair  $(\mathbf{X}, \mathbf{Y})$  underlying the model; cf. (9.29) and Example 9.2. If the degrees of freedom naturally associated to the quadratic statistics in (9.46) should be taken into account, in the line of what is formally done in Example 11.10 below under error normality, then we may consider instead the *adjusted* coefficient of determination

$$R_{\text{adj}}^2 = 1 - \frac{SS_{\text{Res}}/n - p - 1}{S_{\mathbf{YY}}/n - 1} = 1 - \frac{n - 1}{n - p - 1} (1 - R^2),$$

which may be more appropriate to compare the efficiency of models with a distinct number of regressors. One should be aware, however, that the mere observations of  $R^2$  and  $R_{\text{adj}}^2$  as measures of fit are rather disputable in general (essentially because no confidence level is attached to them), so whenever possible these observed values should be used in conjunction with the confidence interval estimates for the parameters in (9.34) and the machinery of hypothesis testing presented in Section 11. As an illustration of this latter method, which provides much more compelling (and flexible) arguments regarding the statistical significance of a given normal linear model, we refer to Example 11.10 below, which features a test based on whether the observed value of the statistics

$$\frac{SS_{\text{Reg}}/p}{SS_{\text{Res}}/n - p - 1} = \frac{n - p - 1}{p} \frac{SS_{\text{Reg}}/S_{\mathbf{YY}}}{SS_{\text{Res}}/S_{\mathbf{YY}}}$$

lies beyond a certain numerical threshold (the critical value of the corresponding F-test as in (11.12)).  $\square$

**Remark 9.17.** (Regression to the mean, again) Using the notation of Examples 9.9 and 9.16, we see from (9.21) that the slope of the regression line is

$$(9.47) \quad \hat{\beta}_1 = \frac{\sqrt{S_{\mathbf{YY}}}}{\sqrt{S_{\mathbf{XX}}}} \hat{\rho},$$

so that

$$\hat{\beta}_0 = \bar{\mathbf{y}} - \hat{\rho} \frac{\sqrt{S_{\mathbf{YY}}}}{\sqrt{S_{\mathbf{XX}}}} \bar{\mathbf{x}},$$

and from (9.23) we deduce that the fitted value is

$$\hat{\mathbf{y}} = \bar{\mathbf{y}} - \hat{\rho} \frac{\sqrt{S_{\mathbf{YY}}}}{\sqrt{S_{\mathbf{XX}}}} (\mathbf{x} - \bar{\mathbf{x}}).$$



Equivalently,

$$(9.48) \quad \frac{\hat{\mathbf{y}} - \bar{\mathbf{y}}}{\sqrt{S_{yy}}} = \hat{\rho} \frac{\mathbf{x} - \bar{\mathbf{x}}}{\sqrt{S_{xx}}},$$

and we conclude that, unless the sampling informs us that  $\mathbf{x}$  and  $\mathbf{y}$  are perfectly correlated ( $|\hat{\rho}| = 1$ ), we should regard the appropriated standardization of  $\hat{\mathbf{y}}$  as being strictly smaller (in absolute value) than the standardization of  $\mathbf{x}$ , a circumstance which certainly indicates a “regression to the mean”; compare with Remark 4.18.  $\square$

**Remark 9.18.** (MLS as a best unbiased estimator) It follows from (9.16) that the MLE for the parameter  $\theta = (\beta, \sigma^2)$  of the regression linear model under error normality is  $\hat{\theta} = (\hat{\beta}, \hat{\sigma}_{\bullet}^2)$ , where  $\hat{\beta}$  is the usual MLS estimator for  $\beta$  in (9.2), and  $\hat{\sigma}_{\bullet}^2 = \|\hat{\mathbf{e}}\|^2/n$  with  $\hat{\mathbf{e}} = \mathbf{Y} - \hat{\mathbf{Y}}$  being the residual as in (9.26), so that (9.17), (9.33) and the independence of  $\{\hat{\beta}, \hat{\sigma}_{\bullet}^2\}$  lead to

$$\text{cov}(\hat{\theta}) = \begin{pmatrix} \text{cov}(\hat{\beta}) & 0 \\ 0 & \text{var}(\hat{\sigma}_{\bullet}^2) \end{pmatrix} = \begin{pmatrix} \sigma^2(\mathbf{r}^t \mathbf{r})^{-1} & 0 \\ 0 & 2(n-p-1)\sigma^2/n^2 \end{pmatrix},$$

where we used that  $\hat{\sigma}_{\bullet}^2 = (n-p-1)\hat{\sigma}^2/n$  with  $\hat{\sigma}$  as in (9.31). On the other hand, applying (8.21) to (9.16) we easily compute that

$$(9.49) \quad \mathcal{F}(\theta) = \begin{pmatrix} \mathcal{F}(\beta) & 0 \\ 0 & \mathcal{F}(\sigma^2) \end{pmatrix} = \begin{pmatrix} \sigma^{-2} \mathbf{r}^t \mathbf{r} & 0 \\ 0 & n/2\sigma^4 \end{pmatrix}.$$

Since the Cramér-Rao lower bound is attained for  $\hat{\beta}$ , Corollary 8.18 implies that it is the best *unbiased* estimator for  $\beta$  (under error normality), which should be compared with Theorem 9.7 (Gauss-Markov), where normality is relaxed to (9.18) but the competing unbiased estimators are required to be linear. On the other hand,  $\text{var}(\hat{\sigma}_{\bullet}^2) < 2\sigma^2/n = \mathcal{F}(\sigma^2)^{-1}$ , a clear violation of the Cramér-Rao lower bound, but of course this poses no contradiction to Theorem 8.17 because  $\hat{\sigma}_{\bullet}^2$  is *not* unbiased. Incidentally, the unbiased estimator  $\hat{\sigma}^2$  for  $\sigma^2$  satisfies  $\text{var}(\hat{\sigma}^2) = 2\sigma^2/(n-p-1) > \mathcal{F}(\sigma^2)^{-1}$ , a strict inequality which suggests the existence of another unbiased estimator for  $\sigma^2$  with a better performance than  $\hat{\sigma}^2$ .  $\square$

**Remark 9.19.** (Asymptotic normality of the MLS estimator) The “small sample” computations leading to the confidence interval estimates (9.25) and (9.34) rely heavily on the normality of the error and should be compared to the corresponding “small sample” estimates for the population mean  $\mu$  in (7.35) and (7.36), respectively. Similarly to what occurred there, under the more general assumptions in (9.18) we must resort to the fundamental limit theorems in Section 6 in order to establish the asymptotic normality of the LSM estimator  $\hat{\beta}$  from which “large sample” estimates should be retrieved. We use the assumptions underlying the linear regression model in Example 9.2 and conveniently decompose the Gram matrix as

$$\mathbf{x}^t \mathbf{x} = \sum_{j=1}^n \mathbf{x}_j^t \mathbf{x}_j,$$

with a similar expression holding for  $\mathbf{x}^t \mathbf{e}$ . Thus,

$$\sqrt{n}(\hat{\beta} - \beta) = \left( \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^t \mathbf{x}_j \right)^{-1} \sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^t \mathbf{e}_j \right).$$

From LLN (Theorem 6.2) we know that

$$\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^t \mathbf{x}_j \xrightarrow{p} \mathbf{c} := \mathbb{E}(\mathbf{x}_j^t \mathbf{x}_j),$$

a symmetric and positive definite  $(p+1) \times (p+1)$  random matrix. Also, since  $\mathbb{E}(\mathfrak{X}_j^t \mathbf{e}) = \vec{0}$  by (9.8) and

$$\begin{aligned} \text{cov}(\mathfrak{X}_j^t \mathbf{e}_j) &\stackrel{(3.15)}{=} \mathbb{E}(\text{cov}(\mathfrak{X}_j^t \mathbf{e}_j | \mathfrak{X})) + \text{cov}(\mathbb{E}(\mathfrak{X}_j^t \mathbf{e}_j | \mathfrak{X})) \\ &= \mathbb{E}(\mathfrak{X}_j^t \text{cov}(\mathbf{e}_j | \mathfrak{X}) \mathfrak{X}_j) \\ &\stackrel{(9.11)}{=} \sigma^2 \mathfrak{C}, \end{aligned}$$

we may use CLT (Theorem 6.5) to check that

$$\sqrt{n} \left( \frac{1}{n} \sum_{j=1}^n \mathfrak{X}_j^t \mathbf{e}_j \right) \xrightarrow{d} \mathcal{N}(\vec{0}, \sigma^2 \mathfrak{C}).$$

Combining these calculations with Theorem 2.23 we conclude that, as  $n \rightarrow +\infty$  and  $p$  is held fixed,

$$(9.50) \quad \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(\vec{0}, \sigma^2 \mathfrak{C}^{-1}),$$

so that  $\hat{\beta}$  is asymptotically normal (and hence consistent) with asymptotic covariance  $\sigma_\beta = \sigma^2 \mathfrak{C}^{-1}$ , which should be reliably estimated in order to obtain the desired confidence regions. We refer to [Ame85, Hay11] for full accounts of the estimation theory of the linear regression model.  $\square$

**Remark 9.20.** (Asymptotic normality for the linear regression model under error normality) Using the results of Remark 9.18, notably the computation of the corresponding Fisher information matrix in (9.49), we find that the MLE estimator  $(\hat{\beta}, \hat{\sigma}_\bullet^2)$  for the parameter  $(\beta, \sigma^2)$  in the linear regression model (under error normality) satisfies, as  $n \rightarrow +\infty$  and  $p$  is held fixed,

$$\sqrt{n} \left( \begin{pmatrix} \hat{\beta} \\ \hat{\sigma}_\bullet^2 \end{pmatrix} - \begin{pmatrix} \beta \\ \sigma^2 \end{pmatrix} \right) \xrightarrow{d} \mathcal{N} \left( \begin{pmatrix} \vec{0} \\ 0 \end{pmatrix}, \begin{pmatrix} n\sigma^2 \mathfrak{s} & 0 \\ 0 & 2\sigma^4 \end{pmatrix} \right), \quad \mathfrak{s} = (\mathfrak{X}^t \mathfrak{X})^{-1},$$

which establishes its asymptotic normality. In particular,

$$(9.51) \quad \sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} \mathcal{N}(\vec{0}, n\sigma^2 \mathfrak{s}) \text{ and } \sqrt{n}(\hat{\sigma}_\bullet^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma^4),$$

which yield the large sample estimates

$$\beta_j \in \left[ \hat{\beta}_j \mp z_{\delta/2} \hat{\mathbf{s}}_{\bullet,j} \right] \text{ with prob. } \approx 1 - \delta, \quad \hat{\mathbf{s}}_{\bullet,j} = \hat{\sigma}_\bullet \sqrt{\mathfrak{s}_{jj}}, \quad \hat{\sigma}_\bullet = \frac{\|\hat{\mathbf{e}}\|}{\sqrt{n}},$$

which is consistent with (9.25) because  $\hat{\mathbf{s}}_{\bullet,j} \xrightarrow{P} \mathbf{s}_j$  by Proposition 9.12, and

$$\sigma^2 \in \left[ \hat{\sigma}_\bullet^2 \mp z_{\delta/2} \sqrt{\frac{2}{n} \hat{\sigma}_\bullet^2} \right] \text{ with prob. } \approx 1 - \delta.$$

Note also that, upon comparison of the first convergence of (9.51) with (9.50), we see that  $\mathfrak{C} = \mathbb{E}(\mathfrak{X}_j^t \mathfrak{X}_j)$  seems to be the right replacement for  $\mathfrak{X}^t \mathfrak{X}/n$  in case the error is not normal.  $\square$

**9.3. Regularization in high dimension, sparsity and the LASSO.** The presence of the root-squared diagonal terms  $\sqrt{\mathfrak{s}_{jj}}$ ,  $\mathfrak{s} = (\mathfrak{X}^t \mathfrak{X})^{-1}$ , in the confidence interval estimates above contributes to enlarge the corresponding dispersion in case  $\mathfrak{X}^t \mathfrak{X}$  is ill-conditioned (e.g. the ratio between its extremal eigenvalues is exceedingly large). Also, the analysis above uses that  $p+1 \leq n$  in a crucial way, so it is useless in the “high-dimensional” regime, where  $p \gg n$  and  $\mathfrak{X}^t \mathfrak{X}$  is no longer invertible. A possible way of remedying this is to work instead with the “penalized” regression estimator

$$\hat{\beta}_\lambda = \argmin_{\beta} \widehat{\mathcal{L}}_\lambda(\beta),$$

where

$$\widehat{\mathcal{L}}_\lambda(\beta) = \frac{1}{2} \|\mathbf{y} - \mathfrak{X}\beta\|^2 + \lambda \|\beta\|^2, \quad \lambda > 0.$$

It follows that

$$\widehat{\beta}_\lambda = (\mathbf{r}^t \mathbf{r} + 2\lambda I)^{-1} \mathbf{r}^t \mathbf{Y},$$

a solution that makes sense even if  $\mathfrak{X}^t \mathfrak{X}$  is not required to meet the full column rank assumption, hence the wording *ridge regularization* for this proposal [HTW15, HTFF09, Wai19, Led22]. Moreover, under the conditions of Proposition 9.6 one computes that

$$\mathbb{E}(\widehat{\beta}_\lambda) = (\mathbf{r}^t \mathbf{r} + 2\lambda I)^{-1} \mathbf{r}^t \mathbf{r} \beta$$

and

$$\text{cov}(\widehat{\beta}_\lambda) = \sigma^2 (\mathbf{r}^t \mathbf{r} + 2\lambda I)^{-1} \mathbf{r}^t \mathbf{r} (\mathbf{r}^t \mathbf{r} + 2\lambda I)^{-1}.$$

Thus, even though  $\widehat{\beta}_\lambda$  fails to be unbiased, it follows from these expressions the existence of  $\lambda_0 > 0$  such that  $\text{mse}(\widehat{\beta}_\lambda) < \text{mse}(\widehat{\beta})$  for  $0 < \lambda < \lambda_0$  [The74]; compare with Remark 7.28, where a similar phenomenon is described for the variance estimators  $\widehat{\sigma}_c^2$ ,  $c > 0$ . Since the ridge regression and its many variants are widely used in real-world applications, we see that a bit of bias is not that bad, specially if it comes with a substantial decrease in the variance. A modern incarnation of this perspective is presented in Example 9.25, where the basic prediction properties of the celebrated LASSO procedure, introduced in [Tib96] and widely used in modern Data Science [Tib96, HTFF09, JWHT13], are briefly discussed. For the sake of motivation, this material is preceded by the corresponding prediction bounds in Examples 9.21, 9.23 and 9.24 below for the classical “low dimensional” MLE ( $p \ll n$ ) under various assumptions on the error distribution.

**Example 9.21.** (High probability bounds for the prediction error under normality) In addition to the parameter recovery methods already discussed (based on the construction of confidence intervals for the unknown parameter  $\beta$ ), a possible way to evaluate the performance of the MLS is to look at

$$(9.52) \quad \mathbf{r} \widehat{\beta} - \mathbf{r} \beta = \mathbf{r} (\mathbf{r}^t \mathbf{r})^{-1} \mathbf{r}^t \mathbf{e},$$

where we assume as always that  $p \leq n$  and  $\mathbf{r}$  has full column rank and hence  $\mathbf{r}^t \mathbf{r}$  is invertible<sup>36</sup>. In the notation of Remark 9.11,

$$(9.53) \quad \mathbf{r} \widehat{\beta} - \mathbf{r} \beta = H \mathbf{e} = \mathbf{e} - \widehat{\mathbf{e}} \in C(\mathbf{r}),$$

the difference between the true error and the residual. In other words, rather than paying attention to the projection of  $\mathbf{e}$  onto  $C(\mathbf{r})^\perp$  under  $Q = \text{Id}_n - H$ , which defines the residual  $\widehat{\mathbf{e}}$ , we now focus on its projection onto  $C(\mathbf{r})$  under  $H$ ; in Figure 1,  $\mathbf{r} \widehat{\beta} - \mathbf{r} \beta$  is represented by  $\widetilde{\mathbf{e}}$ , so that

$$(9.54) \quad \|\widetilde{\mathbf{e}}\|^2 = \|\mathbf{r} \widehat{\beta} - \mathbf{r} \beta\|^2$$

is usually termed the *prediction error*<sup>37</sup>. Under the normality assumption  $\mathbf{e} \sim \mathcal{N}(\vec{0}, \sigma^2 I_{n \times n})$ , it follows from (9.53) and rotational invariance that

$$\sigma^{-2} \|\mathbf{r} \widehat{\beta} - \mathbf{r} \beta\|^2 \sim \chi_p^2,$$

so if

$$(9.55) \quad \widetilde{\text{mse}}(\mathbf{r} \widehat{\beta}) = \frac{\text{mse}(\mathbf{r} \widehat{\beta})}{n}$$

is the *average prediction risk* then

$$(9.56) \quad \widetilde{\text{mse}}(\mathbf{r} \widehat{\beta}) = \frac{\sigma^2 p}{n},$$

<sup>36</sup>Here and in the rest of this subsection we will assume, without loss of generality, that the intercept vanishes, so that  $\beta_0 = 0$ ,  $\mathbf{r} = (\mathbf{r}_1, \dots, \mathbf{r}_p)$ , where  $\mathbf{r}_j$  is the  $j^{\text{th}}$  column of  $\mathbf{r}$ ,  $j = 1, \dots, p$ , etc.

<sup>37</sup>The term “prediction” is used here in a somewhat different sense than in Example 9.15.

where we used Corollary 4.21<sup>38</sup>. Of course, Markov's inequality (2.13) allows us to pass from this "expectation bound" to the corresponding "high probability bound",

$$(9.57) \quad P\left(\frac{\|\mathfrak{r}\hat{\beta} - \mathfrak{r}\beta\|^2}{n} \leq \frac{\sigma^2 p}{\delta n}\right) \geq 1 - \delta, \quad \delta > 0.$$

As expected, this analysis only provides satisfactory prediction results for the classical MLS if either  $\sigma^2$ , which is assumed known, is very small or  $p \ll n$ .  $\square$

**Remark 9.22.** We have seen in Example 9.2 that statistical reasoning demands that the MLS estimator should be specified by solving the minimization problem

$$(9.58) \quad \hat{\beta} = \operatorname{argmin}_{\beta} \widehat{\mathcal{L}}(\beta), \quad \widehat{\mathcal{L}}(\beta) = \frac{1}{n} \|\mathbf{y} - \mathfrak{r}\beta\|^2.$$

At least if  $n$  is large, we may argue<sup>39</sup> that this is the "empirical" version of the more fundamental minimization problem

$$\beta_m := \operatorname{argmin}_{\beta} \mathcal{R}(\beta)$$

with

$$\mathcal{R}(\beta) = \mathbb{E}(\|\mathbf{Y} - \mathfrak{X}\beta\|^2)$$

being the associated *risk function*. Since  $\mathfrak{X}\beta_m$  geometrically corresponds to the orthogonal projection of  $\mathbf{Y}$  onto the subspace generated by the columns of  $\mathfrak{X}$ , we easily see that  $\mathbf{e}_m := \mathbf{Y} - \mathfrak{X}\beta_m$  satisfies

$$(9.59) \quad \mathbb{E}(\mathfrak{X}^t \mathbf{e}_m) = 0,$$

so that

$$\begin{aligned} \mathcal{R}(\hat{\beta}) &= \mathcal{R}(\beta_m + \hat{\beta} - \beta_m) \\ &= \mathbb{E}(\|\mathbf{Y} - \mathfrak{X}(\beta_m + \hat{\beta} - \beta_m)\|^2) \\ &= \mathbb{E}(\|\mathbf{e}_m - \mathfrak{X}(\hat{\beta} - \beta_m)\|^2) \\ &\stackrel{(9.59)}{=} \mathbb{E}(\|\mathbf{e}_m\|^2) + \mathbb{E}(\|\mathfrak{X}(\hat{\beta} - \beta_m)\|^2), \end{aligned}$$

which gives

$$\mathbb{E}(\|\mathfrak{X}(\hat{\beta} - \beta_m)\|^2) = \mathcal{R}(\hat{\beta}) - \mathcal{R}(\beta_m).$$

If we replace  $\beta_m$  by  $\beta$  (to comply with the notation of Example 9.21) and condition on  $\mathfrak{X} = \mathfrak{x}$  we see that

$$\widetilde{\text{mse}}(\mathfrak{x}\hat{\beta}) = \frac{\mathcal{R}(\hat{\beta}) - \mathcal{R}(\beta)}{n},$$

which justifies the terminology employed in (9.55).  $\square$

**Example 9.23.** (High probability bounds for the prediction error without normality) The calculations leading to the expectation and high probability bounds in (9.56) and (9.57) rely heavily on the usual normality assumption on the error. It turns out that we may still obtain a quite effective high probability bound for the prediction error  $\|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2$  in (9.54) by merely assuming that, besides (9.11), the errors  $\{\mathbf{e}_j\}_{j=1}^n$  are assumed to be independent and *sub-Gaussian* in the sense that

$$(9.60) \quad \mathbb{E}(e^{\mathbf{e}_j u}) \leq e^{\sigma^2 u^2 / 2}, \quad u \in \mathbb{R},$$

so that  $\mathbf{e}_j \in \text{SubG}(\sigma)$  as in Definition 5.1. The key point is that (9.52) leads to

$$(9.61) \quad \|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2 = \|D(D^t D)^{-1} D^t U^t \mathbf{e}\|^2,$$

<sup>38</sup>A justification for adopting (9.56) as a measure of accuracy for the prediction error in (9.54) appears in Remark 9.22 below.

<sup>39</sup>Say, by "freezing"  $\beta$  and applying the LLN.

where  $UDV^t$  is a singular value decomposition for  $\mathfrak{x}$  (in particular,  $U$  and  $V$  are both orthogonal). Using that  $D$  is diagonal, it is not hard to check that, under these conditions,

$$D(D^t D)^{-1} D^t = \begin{pmatrix} I_{p \times p} & \\ & 0_{(n-p) \times (n-p)} \end{pmatrix},$$

which gives

$$(9.62) \quad \|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2 = \sum_{j=1}^p |(U^t \mathbf{e})_j|^2.$$

Using that

$$(U^t \mathbf{e})_j = \sum_k U_{kj} \mathbf{e}_k, \quad \sum_k U_{kj}^2 = 1,$$

(9.60) and the independence one easily verifies that

$$\mathbb{E} \left( e^{(U^t \mathbf{e})_j u} \right) \leq e^{\sigma^2 u^2 / 2},$$

that is, each  $\sigma^{-1}(U^t \mathbf{e})_j \in \text{SubG}(1)$ , and from (9.62) we find that  $\sigma^{-2} \|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2 \in \text{SubE}(\nu, 1)$  is sub-exponential as in Definition 5.4; see Remarks 5.12 and 5.13. Using the concentration inequalities in Proposition 5.7 we thus conclude that

$$(9.63) \quad P \left( \frac{\|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2}{n} \leq t \sigma^2 \frac{p}{n} \right) = P \left( \sigma^{-2} \|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2 \leq pt \right) \geq 1 - 2e^{-t/2},$$

for  $t \geq \nu^2$ , which morally corresponds to (9.57) under the replacement  $t \rightarrow \delta^{-1}$ .  $\square$

**Example 9.24.** (High probability bounds for the prediction error without normality, again) An estimate similar to (9.63) may be obtained under the more general assumptions of Example 9.2), where no further knowledge of the error distribution is available besides (9.11). As we shall see, this ignorance will be counterbalanced by a precise control on the spectrum of the modified Gram matrix  $\hat{\Sigma} := \mathfrak{x}^t \mathfrak{x} / n$ , which is known to be positive definite. As in Remark 9.22, we identify  $\beta_m$  to  $\beta$  and explore the variational characterization of  $\hat{\beta}$  in (9.58) to get  $\widehat{\mathcal{L}}(\hat{\beta}) \leq \widehat{\mathcal{L}}(\beta)$ , which means that

$$\frac{\|\mathbf{y} - \mathfrak{x}\hat{\beta}\|^2}{n} \leq \frac{\|\mathbf{e}\|^2}{n}.$$

If we set  $\mathbf{y} = \mathfrak{x}\beta + \mathbf{e}$  in the right-hand side, expand the square and cancel out the terms which are quadratic in the errors we get

$$(9.64) \quad \frac{\|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2}{n} \leq 2 \frac{(\mathfrak{x}^t \mathbf{e})^t (\hat{\beta} - \beta)}{n} \leq 2 \frac{\|\mathfrak{x}^t \mathbf{e}\|}{n} \|\hat{\beta} - \beta\|,$$

where Cauchy-Schwarz has been used in the last step. If  $\lambda_{\min}(\hat{\Sigma}) \leq \lambda_{\max}(\hat{\Sigma})$  stand for the (positive) extremal eigenvalues of  $\hat{\Sigma}$  then we have

$$(9.65) \quad \frac{\|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2}{n} = \langle \hat{\Sigma}(\hat{\beta} - \beta), \hat{\beta} - \beta \rangle \geq \lambda_{\min}(\hat{\Sigma}) \|\hat{\beta} - \beta\|^2,$$

which may be viewed as a control on the sample correlation between the columns of  $\mathfrak{x}$  (because  $n\hat{\Sigma}_{jk} = \mathfrak{x}_j^t \mathfrak{x}_k = \|\mathfrak{x}_j\| \|\mathfrak{x}_k\| \text{corr}(\mathfrak{x}_j, \mathfrak{x}_k)$ ; cf (9.29)), so if we combine these estimates we get

$$\frac{\|\mathfrak{x}\hat{\beta} - \mathfrak{x}\beta\|^2}{n} \leq 4 \frac{\|\mathfrak{x}^t \mathbf{e}\|^2}{n^2 \lambda_{\min}(\hat{\Sigma})}.$$

On the other hand, again using our standing assumptions (including (9.11)) we compute

$$\begin{aligned}\mathbb{E}(\|\mathbf{r}^t \mathbf{e}\|^2) &= \mathbb{E}(\text{tr}((\mathbf{r}^t \mathbf{e})(\mathbf{r}^t \mathbf{e})^t)) \\ &= \text{tr cov}(\mathbf{r}^t \mathbf{e}) \\ &= \sigma^2 \text{tr}(\mathbf{r}^t \mathbf{r}) \\ &\leq \sigma^2 p n \lambda_{\max}(\widehat{\mathbf{\Sigma}}),\end{aligned}$$

which gives the expectation bound

$$(9.66) \quad \widetilde{\text{mse}}(\mathbf{r}\widehat{\beta}) \leq 4\sigma^2 \lambda(\widehat{\mathbf{\Sigma}}) \frac{p}{n}, \quad \lambda(\widehat{\mathbf{\Sigma}}) := \frac{\lambda_{\max}(\widehat{\mathbf{\Sigma}})}{\lambda_{\min}(\widehat{\mathbf{\Sigma}})},$$

from which we obtain the high probability bound

$$(9.67) \quad P\left(\frac{\|\mathbf{r}\widehat{\beta} - \mathbf{r}\beta\|^2}{n} \leq \frac{4\sigma^2}{\delta} \lambda(\widehat{\mathbf{\Sigma}}) \frac{p}{n}\right) \geq 1 - \delta, \quad \delta > 0,$$

again via Markov. □

**Example 9.25.** (High dimensionality, sparsity and the LASSO) Although its derivation required only very mild assumptions on the error distribution, the high probability bound in (9.67) remains in essence quite similar to (9.63) and (9.57). In particular, its manifest dependence on the “dimensional” ratio  $p/n$  confirms that, with no further control on the error variance  $\sigma^2$  and on the condition number  $\lambda(\widehat{\mathbf{\Sigma}})$  of  $\widehat{\mathbf{\Sigma}}$ , the linear regression model only predicts reliably if  $p \ll n$ . Outside this regime, say if  $p \approx n$ , MLS gets plagued with at least two deficiencies, namely, *high variability* (although  $\mathbf{r}\widehat{\beta}$  is unbiased, estimates like (9.66) fail to control its variance) and *low interpretability* (the huge amount of variables makes it hard to select those which are indeed relevant for explaining the response). To remedy this, a possible course of action involves inserting a penalizing term in the classical model as we did in Example 9.3 for the ridge regression; see [Led22, Introduction] for a nice discussion of this approach, which has many applications in Supervised Learning, where being able to identify in which side of the threshold  $p \approx n$  the problem in hand lies is mandatory [Don00, HTFF09, BC11, BVDG11, HTW15, FBG<sup>+</sup>16, Ver18, Wai19, Led22]. Here we provide a modest introduction to this circle of ideas by showing how new insights are needed in the “high-dimensional” regime  $p \gg n$ , where in particular the key correlation assumption in (9.65) is no longer available (because  $\mathbf{r}^t \mathbf{r}$  fails to be invertible). As already suggested, we penalize a suitable multiple of the least squares objective function in order to obtain the *LASSO estimator*

$$\widehat{\beta}_L = \arg\min_{\beta'} f_L(\beta'), \quad \widetilde{\mathcal{L}}_L(\beta') = \frac{1}{2n} \|\mathbf{y} - \mathbf{r}\beta'\|^2 + \lambda \|\beta'\|_1,$$

where  $\lambda > 0$  is a tuning parameter to be chosen later and

$$\|\beta'\|_1 = \sum_{j=1}^p |\beta'_j|.$$

Since  $\widetilde{\mathcal{L}}_L(\widehat{\beta}_L) \leq \widetilde{\mathcal{L}}_L(\beta)$ , where  $\beta$  is the true parameter appearing in the model equation  $\mathbf{y} = \mathbf{r}\beta + \mathbf{e}$ , we thus get with a help from Hölder inequality,

$$\begin{aligned}\frac{1}{n} \|\mathbf{r}\widehat{\beta}_L - \mathbf{r}\beta\|^2 &\leq \frac{2}{n} (\mathbf{r}^t \mathbf{e})^t (\widehat{\beta}_L - \beta) + 2\lambda (\|\beta\|_1 - \|\widehat{\beta}_L\|_1) \\ &\leq \frac{2}{n} \|\mathbf{r}^t \mathbf{e}\|_\infty \|\widehat{\beta}_L - \beta\|_1 + 2\lambda (\|\beta\|_1 - \|\widehat{\beta}_L\|_1),\end{aligned}$$

an estimate which should be compared to (9.64), with its right-hand side effectively disentangling the contributions coming from the “effective error”  $2\|\mathbf{r}^t \mathbf{e}\|_\infty/n$  and the penalization. Now, sparsity enters the game precisely to handle this latter term, as it contemplates the belief, substantiated by an “omniscient oracle”, that

a considerable portion of regressors may be dispensed with, so the corresponding parameter entries may be set to vanish. Precisely, there exists  $S \subsetneq \{1, \dots, p\}$  with  $s := \#S \ll n$  such that  $\beta_j = 0$  exactly when  $j \notin S$ . Thus, if  $\beta_S$  is the “restriction” of  $\beta$  to  $S$ , so that  $\beta = \beta_S + \beta_{S^c}$ , and setting  $\hat{\delta} = \hat{\beta}_L - \beta$ , we have

$$\begin{aligned} \|\beta\|_1 - \|\hat{\beta}_L\|_1 &= \|\beta_S\|_1 - \|\beta + \hat{\delta}\|_1 \\ &= \|\beta_S\|_1 - \|\beta_S + \hat{\delta}_S + \hat{\delta}_{S^c}\|_1 \\ &= \|\beta_S\|_1 - \|\beta_S + \hat{\delta}_S\|_1 - \|\hat{\delta}_{S^c}\|_1 \\ &\leq \|\hat{\delta}_S\|_1 - \|\hat{\delta}_{S^c}\|_1, \end{aligned}$$

which gives

$$\frac{1}{n} \|\mathbf{r}\hat{\delta}\|^2 \leq \frac{2}{n} \|\mathbf{r}^t \mathbf{e}\|_\infty \|\hat{\delta}\|_1 + 2\lambda \left( \|\hat{\delta}_S\|_1 - \|\hat{\delta}_{S^c}\|_1 \right),$$

so if we further assume that the tuning parameter dominates the “effective error” according to

$$(9.68) \quad \frac{2}{n} \|\mathbf{r}^t \mathbf{e}\|_\infty \leq \lambda$$

we end up with

$$(9.69) \quad \frac{1}{n} \|\mathbf{r}\hat{\delta}\|^2 \leq \lambda \left( 3\|\hat{\delta}_S\|_1 - \|\hat{\delta}_{S^c}\|_1 \right).$$

As a direct consequence of this basic inequality we see that

$$\hat{\delta} \in \mathcal{C}(S) := \left\{ \beta' \in \mathbb{R}^{p+1}; \|\beta'_{S^c}\|_1 \leq 3\|\beta'_S\|_1 \right\},$$

which suggests that the appropriate replacement for (9.65) is to assume, for some  $\kappa > 0$ , that

$$(9.70) \quad \frac{1}{n} \|\mathbf{r}\beta'\|^2 \geq \kappa \|\beta'\|^2, \quad \beta' \in \mathcal{C}(S).$$

Under this *restricted eigenvalue* (RE) condition,

$$\begin{aligned} \frac{1}{n} \|\mathbf{r}\hat{\delta}\|^2 &\stackrel{(9.69)}{\leq} 3\lambda \|\hat{\delta}_S\|_1 \\ &\leq 3\lambda \sqrt{s} \|\hat{\delta}\| \\ &\stackrel{(9.70)}{\leq} \frac{3\lambda \sqrt{s}}{\sqrt{\kappa n}} \|\mathbf{r}\hat{\delta}\|, \end{aligned}$$

which finally gives the bound

$$(9.71) \quad \frac{\|\mathbf{r}\hat{\beta}_L - \mathbf{r}\beta\|^2}{n} \leq \frac{9\lambda^2 s}{\kappa}.$$

In order to estimate in terms of  $\lambda$  the probability of the event in (9.68), to which the validity of (9.71) is conditioned, let us assume for simplicity that  $\mathbf{e} \sim \mathcal{N}(\vec{0}, \sigma^2 \mathbf{I}_{n \times n})$ . By the projection property in (4.7),

$$2 \frac{\mathbf{r}_k^t \mathbf{e}}{n} \sim \mathcal{N} \left( 0, 4 \frac{\sigma^2}{n} \left\| \frac{\mathbf{r}_k}{\sqrt{n}} \right\|^2 \right), \quad k = 1, \dots, p,$$

so if the columns of the design matrix are normalized so that

$$\left\| \frac{\mathbf{r}_k}{\sqrt{n}} \right\| \leq C,$$

the standard Gaussian concentration inequality in (5.3) leads to

$$\begin{aligned} P\left(2\left\|\frac{\mathbf{r}^t \mathbf{e}}{n}\right\|_{\infty} \leq \lambda\right) &\geq 1 - 2pe^{-\frac{n\lambda^2}{8C^2\sigma^2}} \\ &= 1 - 2e^{-\frac{n\lambda^2}{8C^2\sigma^2} + \ln p}. \end{aligned}$$

This gives

$$P\left(2\left\|\frac{\mathbf{r}^t \mathbf{e}}{n}\right\|_{\infty} \leq \lambda\right) \geq 1 - 2e^{-\frac{t^2}{2}}$$

if

$$\lambda^2 = 8C^2\sigma^2 \left(\frac{\ln p}{n} + \frac{t^2}{2n}\right),$$

so with this choice of  $\lambda$ , (9.71) immediately yields the bound

$$(9.72) \quad \frac{\|\widehat{\beta}_L - \mathbf{r}\beta\|^2}{n} \leq \frac{72C^2\sigma^2}{\kappa} \frac{s}{n} \left(\ln p + \frac{t^2}{2}\right)$$

with at least the same probability. Upon comparison with (9.67) and not taking into account certain structural constants, we have been able to replace the dimensional ratio  $p/n$  by  $s \ln p/n$ , which is linear in the “sparsity index”  $s = \|\beta\|_0$  and scales logarithmically with  $p$ , added to another term which is driven by the “oracle rate”  $sn^{-1} = o(1)$ . Thus, it suffices to take  $n \gg s \ln p > s$  in order to have LASSO’s prediction nearly as accurate as if  $S = \text{supp } \beta$ , whose elements classify the relevant regressors, was known a priori. We mention that similar estimates hold true under much weaker assumptions on the error<sup>40</sup> and even for other kinds of penalizations; we refer to [BVDG11, Chapter 6], [HTW15, Chapter 11], [Wai19, Chapter 7] and [Led22, Chapter 6] for such generalizations and, more importantly, for the heuristics behind the crucial RE condition in (9.70) above. Finally, the practical question remains of fine-tuning the parameter  $\lambda$  so as to obtain the right balance between variability and interpretability. In this regard, the feasibility of the most adopted procedure, cross-validation, is theoretically confirmed in [CLC21], where it is shown that, under suitable conditions, its use only adds to the right-hand side of (9.72) a multiplicative factor which is  $O(\sqrt{\ln pn})$ , hence negligible for most realistic purposes.  $\square$

## 10. SUFFICIENCY

When passing, in a given statistical model, from the random sample

$$X = (X_1, \dots, X_n), \quad X_j \sim \psi_{\theta},$$

to an estimator  $\widehat{\theta}$  by means of a statistic  $h = h(X)$ , the important question arises of determining how much information has been gleaned in the process. A fully satisfactory answer to this question involves pondering on how to precisely measure the amount of information carried by the data, which is beyond the scope of these notes. A much less ambitious task would be to ensure that the given statistic indeed carries all the information provided by the sample data in the sense that no further knowledge is required to pinpoint the unknown parameter  $\theta$ . Thus, we seek to understand when the “extra randomness” in the data  $X$  not apprehended by  $h(X)$  fails to involve  $\theta$  and hence becomes irrelevant for its estimation. This admits a neat probabilistic formulation in terms of the conditioning notions introduced in Section 3 above.

**Definition 10.1.** A statistic  $h = h(X)$  is *sufficient* if for any realization  $\mathbf{x}$  of  $X$  the conditional probability distribution  $\psi_{\theta; X|h(X)=h(\mathbf{x})}$  evaluated at  $\mathbf{x}$  does not depend on  $\theta$ .

<sup>40</sup>For instance, if the error is sub-Gaussian then the corresponding concentration inequalities in Section 5 might be useful.



Using (3.3) we thus see that sufficiency of  $h$  leads to the existence of a function  $\xi = \xi(\mathbf{x})$  such that

$$(10.1) \quad \frac{\psi_{\theta; (h(X), X)}(h(\mathbf{x}), \mathbf{x})}{\psi_{\theta; h(X)}(h(\mathbf{x}))} = \xi(\mathbf{x}).$$

A key observation now is that the obvious inclusion of events  $\{X = \mathbf{x}\} \subset \{h(X) = h(\mathbf{x})\}$  implies that

$$(10.2) \quad \psi_{\theta; (h(X), X)}(h(\mathbf{x}), \mathbf{x}) = \psi_{\theta; X}(\mathbf{x}) = L(\mathbf{x}; \theta),$$

the likelihood function. Replacing this back in (10.1) leads to the following useful characterization of sufficiency, saying that it occurs precisely when the dependence of  $L(\mathbf{x}; \theta)$  on  $\theta$  gets confined in a term which only depends on  $\mathbf{x}$  through the given statistic. This clearly captures the essential content of the concept: all the information needed to estimate  $\theta$  is already contained in the estimator defined by the sufficient statistic  $h$ , with no need to further confer the data  $X$ .

**Theorem 10.2.** (Fisher-Neyman factorization)  $h$  is sufficient if and only if the likelihood function factorizes as

$$(10.3) \quad L(\mathbf{x}; \theta) = \eta(h(\mathbf{x}), \theta) \xi(\mathbf{x}),$$

for positive functions  $\eta$  and  $\xi$ .

*Proof.* We have already seen that sufficiency implies (10.3). For the converse we first note that (10.2) leads to

$$\begin{aligned} \psi_{\theta; h(X)}(h(\mathbf{x})) &= \int_{\{\mathbf{x}': h(\mathbf{x}') = h(\mathbf{x})\}} \psi_{\theta; (h(X), X)}(h(\mathbf{x}'), \mathbf{x}') d\mathbf{x}' \\ &= \int_{\{\mathbf{x}': h(\mathbf{x}') = h(\mathbf{x})\}} \psi_{\theta; X}(\mathbf{x}') d\mathbf{x}' \end{aligned}$$

so we may again use (3.3) to compute:

$$\begin{aligned} \psi_{\theta; X|h(X)=h(\mathbf{x})} &= \frac{\psi_{\theta; (h(X), X)}(h(\mathbf{x}), \mathbf{x})}{\psi_{\theta; h(X)}(h(\mathbf{x}))} \\ &= \frac{\eta(h(\mathbf{x}), \theta) \xi(\mathbf{x})}{\int_{\{\mathbf{x}': h(\mathbf{x}') = h(\mathbf{x})\}} \eta(h(\mathbf{x}'), \theta) \xi(\mathbf{x}') d\mathbf{x}'} \\ &= \frac{\eta(h(\mathbf{x}), \theta) \xi(\mathbf{x})}{\eta(h(\mathbf{x}), \theta) \int_{\{\mathbf{x}': h(\mathbf{x}') = h(\mathbf{x})\}} \xi(\mathbf{x}') d\mathbf{x}'} \end{aligned}$$

Thus,

$$\psi_{\theta; X|h(X)=h(\mathbf{x})} = \frac{\xi(\mathbf{x})}{\int_{\{\mathbf{x}': h(\mathbf{x}') = h(\mathbf{x})\}} \xi(\mathbf{x}') d\mathbf{x}'}$$

only depends on  $\mathbf{x}$ . □

**Corollary 10.3.** A unique ML estimator is a function of a sufficient statistic. More generally, if a ML estimator exists then an ML estimator may be chosen so as to be a function of a sufficient statistic.

*Proof.* Given that the ML estimator  $\hat{\theta}$  is obtained by maximizing the likelihood function  $L(\mathbf{x}; \theta)$  in  $\theta$  (for each  $\mathbf{x}$ ), this is an obvious consequence of (10.3). □

**Example 10.4.** (Sufficiency in a normal population) If  $X_j \sim \mathcal{N}(\mu, \sigma^2)$  we know from Example 8.4 that

$$L(\mathbf{x}; \theta) = (2\pi\theta_2)^{-n/2} e^{-\frac{1}{2\theta_2} \sum_{j=1}^n (x_j - \theta_1)^2}, \quad \mathbf{x} = (x_1, \dots, x_n),$$

where  $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+$ . We distinguish three cases:

- ( $\theta_2$  is known and  $\theta_1$  is the unknown parameter) Set

$$h_1(\mathbf{x}) = \frac{1}{n} \sum_j x_j$$

so that  $\sum_j (x_j - h_1(\mathbf{x})) = 0$  implies

$$\begin{aligned} \sum_j (x_j - \theta_1)^2 &= \sum_j (x_j - h_1(\mathbf{x}) + h_1(\mathbf{x}) - \theta_1)^2 \\ &= \sum_j (x_j - h_1(\mathbf{x}))^2 + n(h_1(\mathbf{x}) - \theta_1)^2, \end{aligned}$$

which leads to the factorization

$$(10.4) \quad L(\mathbf{x}; \theta_1) = \underbrace{(2\pi\theta_2)^{-n/2} e^{-\frac{\sum_j (x_j - h_1(\mathbf{x}))^2}{2\theta_2}}}_{\xi(\mathbf{x})} \underbrace{e^{-\frac{n(h_1(\mathbf{x}) - \theta_1)^2}{2\theta_2}}}_{\eta(h_1(\mathbf{x}), \theta_1)}.$$

This shows that  $h_1$  is a sufficient statistic for  $\theta_1$  (given  $\theta_2$ ).

- ( $\theta_1$  is known and  $\theta_2$  is the unknown parameter) Here,

$$h_2(\mathbf{x}) = \sum_j (x_j - \theta_1)^2$$

qualifies as a statistic and

$$L(\mathbf{x}; \theta_2) = \underbrace{(2\pi\theta_2)^{-n/2} e^{-\frac{h_2(\mathbf{x})}{2\theta_2}}}_{\eta(h_2(\mathbf{x}), \theta_2)} \times \underbrace{1}_{\xi(\mathbf{x})}$$

shows that  $h_2$  is a sufficient statistic for  $\theta_2$  (given  $\theta_1$ ).

- ( $\theta = (\theta_1, \theta_2)$  is the unknown bi-dimensional parameter). Here we set

$$\tilde{h}_2(\mathbf{x}) = \sum_j (x_j - h_1(\mathbf{x}))^2$$

so (10.4) gives

$$L(\mathbf{x}; \theta) = \underbrace{(2\pi\theta_2)^{-n/2} e^{-\frac{\tilde{h}_2(\mathbf{x}) + n(h_1(\mathbf{x}) - \theta_1)^2}{2\theta_2}}}_{\eta((h_1(\mathbf{x}), \tilde{h}_2(\mathbf{x}), \theta))} \times \underbrace{1}_{\xi(\mathbf{x})},$$

which shows that  $H(\mathbf{x}) = (h_1(\mathbf{x}), \tilde{h}_2(\mathbf{x}))$  is a sufficient statistic for  $\theta$ . We thus see that the common practice, which has been extensively used in Subsection 7.3, of regarding  $H$  as a sufficient statistic when sampling from a normal population, is fully justified.  $\square$

**Example 10.5.** (Sufficiency in an exponential population) If  $X_j \sim \text{Exp}(\lambda)$  then from Example 8.6 we get

$$L(\mathbf{x}; \lambda) = \underbrace{\lambda^n e^{-\lambda h(\mathbf{x})}}_{\eta(h(\mathbf{x}), \lambda)} \times \underbrace{1}_{\xi(\mathbf{x})},$$

where  $h(\mathbf{x}) = \sum_j x_j$  is a sufficient statistic for  $\lambda$ .  $\square$

**Example 10.6.** (Sufficiency in a Bernoulli or Poisson population) If  $X_j \sim \text{Ber}(p)$  then Example 8.7 gives

$$L(\mathbf{x}; p) = \underbrace{p^{h(\mathbf{x})} (1-p)^{n-h(\mathbf{x})}}_{\eta(h(\mathbf{x}), p)} \times \underbrace{1}_{\xi(\mathbf{x})},$$

which shows that  $h(\mathbf{x}) = \sum_j x_j$  is a sufficient statistic for estimating  $p$ . On the other hand, if  $X_j \sim \mathcal{P}(\rho)$  then, again by Example 8.7,

$$L(\mathbf{x}; \rho) = \underbrace{\rho^{k(\mathbf{x})} e^{-n\rho}}_{\eta(k(\mathbf{x}), \rho)} \times \underbrace{(\prod_j x_j!)^{-1}}_{\xi(\mathbf{x})},$$

which confirms that  $k(\mathbf{x}) = \sum_j x_j$  is a sufficient statistic for estimating  $\rho$ .  $\square$

As illustrated by the computations in Remark 8.25, at least in the regime of large samples it follows from (8.34) that the dependence on sample data of confidence intervals for ML estimators occurs only through the estimator itself. This general observation clearly aligns with Corollary 10.3 and is definitely confirmed by all the examples examined above, where a simple relationship of the given sufficient statistic with the corresponding ML estimator is manifest.

## 11. HYPOTHESIS TESTING

Our aim here is to discuss a bit more on the heuristics behind the choices of the rejection regions appearing in the F-tests implemented in Remark 7.37 and Example 7.40 above.

**11.1. A glimpse at the Neyman-Pearson setup.** As usual, we are given a parametric statistical model

$$X_1, \dots, X_n \sim \psi_\theta, \quad \theta \in \Theta \subset \mathbb{R}^p$$

as in Definition 7.2 and Remark 7.3, so that  $(\Omega, \mathcal{F}, \{\mathcal{P}_\theta\}_{\theta \in \Theta})$  is the underlying family of probability spaces and  $P_\theta = \psi_\theta dx$  is the common distribution of the components of the associated random vector  $X = (X_1, \dots, X_n) : \Omega \rightarrow \mathbb{R}^n$ . Given *disjoint* subsets  $\Theta_0, \Theta_a \subset \Theta$  with  $\Theta = \Theta_0 \cup \Theta_a$ , *hypothesis testing* concerns the prospect of using the available data in an observed value  $\mathbf{x}$  of  $X$  to provide statistical evidence for deciding between the *null hypothesis*

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

and the *alternative hypothesis*

$$H_a : \quad \theta \in \Theta_a.$$

One adheres to the usual asymmetry in regarding  $H_0$  as the *status quo* and then chooses a statistics  $h = h(X) : \Omega \rightarrow \mathbb{R}$  and a *rejection region*  $R \subset \mathbb{R}$  so that  $H_0$  gets *rejected* if the realization  $h(\mathbf{x})$  of  $h(X)$  takes value in  $R$ . A pair  $T = (h, R)$  as above is called a *test* for the given statistical model and we denote by  $\mathcal{T}$  the collection of all such tests<sup>41</sup>. As we will see in Remark 11.3 below, the eventual implementation of a test  $T \in \mathcal{T}$  necessarily involves the knowledge of the distribution of (a perhaps complicated function of)  $h(X)$  under the null hypothesis.

In order to quantify the possible types of errors in making such a decision, we consider the *power function*  $\pi : \Theta \rightarrow [0, 1]$  of  $(h, R)$ ,

$$\pi(\theta) = \mathcal{P}_\theta(h(X) \in R).$$

We then see that restriction to  $\Theta_0$ , namely,

$$\gamma(\theta) := \pi|_{\Theta_0}(\theta), \quad \theta \in \Theta_0,$$

quantifies the *type I error* of rejecting  $H_0$  when it is true, whereas restriction to  $\Theta_a$ ,

$$\delta(\theta) := \pi|_{\Theta_a}(\theta), \quad \theta \in \Theta_a,$$

is such that

$$1 - \delta(\theta) = \mathcal{P}_\theta(h(X) \notin R)$$

measures the *type II error* of *not* rejecting  $H_0$  when it is false. Ideally, one would seek for an strategy minimizing *both* errors at a time, but simple examples show that this is doomed to fail in general. The standard way to

<sup>41</sup>For instance, if  $R = [r, +\infty)$  then we say that  $r$  is a *critical value* for the test.

overcome this is to search for a test which minimizes type II error under the constraint that type I error remains uniformly bounded from above by a fixed amount given in advance.

**Definition 11.1.** Given  $\alpha \in (0, 1)$  we say that a test  $T \in \mathcal{T}$  has *confidence level*  $\alpha$  if

$$(11.1) \quad \sup_{\theta \in \Theta_0} \gamma(\theta) = \alpha,$$

and we denote by  $\mathcal{T}_\alpha$  the collections of all such tests.

**Definition 11.2.** A test  $T \in \mathcal{T}_\alpha$  is *uniformly most powerful* (UMP) if it satisfies (with self-explanatory notation)

$$\delta(\theta) \geq \delta^*(\theta), \quad \theta \in \Theta_a,$$

for any  $T^* \in \mathcal{T}_\alpha$ .

**Remark 11.3.** Note that (11.1), which may be rewritten as

$$(11.2) \quad \sup_{\theta \in \Theta_0} \mathcal{P}_\theta(h(X) \in R) = \alpha,$$

allows us to explicitly determine the rejection region  $R$  from the confidence level  $\alpha$  only in case a (perhaps approximate) knowledge of the distribution of  $h(X)$  under  $H_0$  is at hand, a procedure illustrated in the examples considered below. In other words, the ubiquitous “Problem of Distribution” in Parametric Statistics resurfaces in this setting as well, although here the relevant statistics  $h(X)$  gets restricted to the parametric region where the null hypothesis holds true.  $\square$

The celebrated Neyman-Pearson lemma [CB21, Theorem 8.3.12] exhibits a UMP test in the simple hypotheses case, where both  $\Theta_0$  and  $\Theta_a$  contain a single element. Unfortunately, such a test may not exist even for one of the simplest *composite* hypotheses cases, namely, a “two-sided” test of the form  $\Theta \subset \mathbb{R}$  some open interval,  $\Theta_0 = \{\theta_0\}$  for some  $\theta_0 \in \Theta$  and  $\Theta_1 = \Theta \setminus \Theta_0$  (as in Remark 7.37, for instance); see [CB21, Example 8.3.19] and the surrounding discussion for more on this rather delicate point. Of course, we may always restrict further the class of contenders where the ideal test should be sought (consistent, unbiased, etc.) but it seems that none of these strategies produces a test with optimal performance in *all* cases.

**11.2. Testing via likelihood ratios.** The state of affairs indicated in the previous paragraph suggests that, instead of *systematically* trying to find the best test in a given context, one should proceed *heuristically* so as to single out a family of tests which are relatively easy to implement, have nice asymptotic properties and reproduce most known composite tests for samples of any size.

Recall from Subsection 8.3 that the ML estimators, computed in terms of the likelihood function as in Definition 8.1, have many remarkable properties, including asymptotic normality. Moreover, given the available information contained in the realization  $\mathbf{x}$  of the random sample  $X$ , the discussion surrounding (8.4) justifies regarding  $\sup_{\theta \in \Theta_0} L(\mathbf{x}; \theta)$  as the best evidence in favor of  $H_0$  and  $\sup_{\theta \in \Theta_a} L(\mathbf{x}; \theta)$  as the best evidence in favor of  $H_a$ , which suggests formulating a hypothesis test based on the *likelihood ratio*

$$(11.3) \quad \mathbf{x} \in \mathbb{R}^n \mapsto \frac{\sup_{\theta \in \Theta_0} L(\mathbf{x}; \theta)}{\sup_{\theta \in \Theta_a} L(\mathbf{x}; \theta)} \in [0, +\infty].$$

At the risk of (over)simplifying the exposition, but at the same time remaining in a generality that will suffice for the applications we have in mind, we assume from now on that  $\Theta_0 \subset \Theta$  has a negligible size (as a subset of  $\Theta$ ) and  $\Theta_a = \Theta \setminus \Theta_0$ . Typically,  $\Theta \subset \mathbb{R}^p$  will be an open subset,  $\Theta_0$  the portion of an affine  $k$ -plane lying

in  $\Theta$  with  $k < p$  and  $\Theta_a$  its complement in  $\Theta$ . In this setting it seems reasonable to replace  $\Theta_a$  by  $\Theta$  in the denominator of (11.3), so that under suitable regularity assumptions the likelihood ratio becomes

$$(11.4) \quad \mathbf{x} \in \mathbb{R}^n \mapsto \Lambda(\mathbf{x}) := \frac{L(\mathbf{x}; \hat{\theta}_0)}{L(\mathbf{x}; \hat{\theta})} \in [0, 1],$$

where

$$\hat{\theta} = \sup_{\theta \in \Theta} L(\mathbf{x}; \theta)$$

is the MLE for  $\theta$  (as in Definition 8.3). and

$$\hat{\theta}_0 = \sup_{\theta \in \Theta_0} L(\mathbf{x}; \theta)$$

is the *null* MLE for  $\theta$ . This leads to a remarkable class of statistical tests.

**Definition 11.4.** Under the conditions above, the *likelihood ratio test*  $T = (h, R) \in \mathcal{T}_\alpha$  is performed by choosing

$$h(\mathbf{x}) = -2 \ln \Lambda(\mathbf{x})$$

and  $R = [r, +\infty)$ , where  $r > 0$  is determined by

$$(11.5) \quad \sup_{\theta \in \Theta_0} \mathcal{P}_\theta(h(X) \geq r) = \alpha.$$

**Remark 11.5.** Clearly, (11.5) is a special case of (11.2), where the rejection region  $R$  now takes the form  $[r, +\infty)$  for some critical value  $r > 0$ , and we are supposed to solve it for  $r$  given the confidence level  $\alpha$ . But notice that, as already observed in Remark 11.3, this requires knowing the distribution of  $h$  under  $H_0$ . In this regard, if  $h(X)$  is found to be  $\Theta_0$ -ancillary in the sense that its distribution does *not* depend on  $\theta \in \Theta_0$  then

$$(11.6) \quad \alpha = \mathcal{P}_\theta(h(X) \geq r) \text{ for any } \theta \in \Theta_0$$

determines  $r$  as a function of  $\alpha$ . □

Although we will restrict ourselves to normal populations, the examples below will suffice to illustrate the remarkable flexibility of this construction. Moreover, in all these examples the likelihood ratio statistics is  $\Theta_0$ -ancillary in the sense of Remark 11.5, so that (11.6) applies in order to solve for  $r$  in terms of  $\alpha$ .

**Example 11.6.** (*z-test for the mean of a normal population with known variance*) As in Example 8.4 we assume that  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ , where  $\theta_2 = \sigma^2$  is known. Thus,  $\Theta = \{\theta_1 = \mu\} = \mathbb{R}$ ,  $\Theta_0 = \{\mu_0\}$  for some  $\mu_0 \in \mathbb{R}$ ,  $\Theta_a = \mathbb{R} \setminus \{\mu_0\}$  and we want to test

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_a : \mu \neq \mu_0.$$

We recall that  $\hat{\theta}_1 = \bar{X}_n$  is the MLE for  $\mu$ . Using (8.8) and (7.23) we compute

$$\begin{aligned} h(\mathbf{x}) &= -2 \ln \frac{L(\mathbf{x}; \mu_0)}{L(\mathbf{x}; \hat{\theta}_1)} \\ &= -2 \ln \frac{(2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum_j (x_j - \mu_0)^2}}{(2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} \sum_j (x_j - \hat{\theta}_1)^2}} \\ &= -2 \ln e^{-\frac{n}{2\sigma^2} (\hat{\theta}_1 - \mu_0)^2} \\ &= \left( \frac{\bar{x}_n - \mu_0}{\sigma/\sqrt{n}} \right)^2. \end{aligned}$$

Thus, under  $H_0$  we see that

$$h(X) = Z(X)^2 \sim \chi_1^2,$$

where

$$Z(X) = \frac{\bar{X}_n - \mu_0}{\sigma/\sqrt{n}} \sim \mathcal{N}(0, 1).$$

Now, by (11.5) the rejection interval  $R = [r, +\infty)$  is determined by

$$\mathcal{P}_{\theta_0}(Z(X)^2 \geq r) = \alpha,$$

so we may take  $r = \chi_{1,\alpha}^2$ , the quantil of  $\chi_1^2$  whose tail probability is  $\alpha$  (that is,  $F_{\chi_k^2}(\chi_{k,\alpha}^2) = 1 - \alpha$ ). Since

$$\alpha = \mathcal{P}_{\theta_0}(Z(X)^2 \geq r) = \mathcal{P}_{\theta_0}(-\sqrt{r} \leq Z(X) \leq \sqrt{r}),$$

we may also use the standard normal quantiles to make sure that if

$$Z(\mathbf{x}) \in (-\infty, -z_{\alpha/2}] \cup [z_{\alpha/2}, +\infty)$$

then  $H_0$  gets rejected. □

**Example 11.7.** (*t*-test for the mean of a normal population with unknown variance) As in Example 8.4 we assume that  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ , where  $(\theta_1, \theta_2) = (\mu, \sigma^2)$  is unknown, so that  $\Theta = \mathbb{R} \times \mathbb{R}_+$ . Also, we fix  $\mu_0 \in \mathbb{R}$  and set  $\Theta_0 = \{\mu_0\} \times \mathbb{R}_+$ , a half-line contained in  $\Theta$ . As before, we want to test

$$H_0 : \mu = \mu_0 \quad \text{vs} \quad H_a : \mu \neq \mu_0.$$

We recall that  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)$ , where  $\hat{\theta}_1 = \bar{X}_n$  and

$$\hat{\theta}_2 = \frac{1}{n} \sum_j (X_j - \bar{\theta}_1)^2$$

is the MLE for  $\theta_2$ , so that

$$\sup_{\theta \in \Theta} L(\mathbf{x}; \theta) = L(\mathbf{x}; \hat{\theta}_1, \hat{\theta}_2).$$

On the other hand, one has

$$\sup_{\theta \in \Theta_0} L(\mathbf{x}; \theta) = L(\mathbf{x}; \mu_0, \hat{\theta}_{20}),$$

where the null MLE for  $\theta_2$  is

$$\hat{\theta}_{20} = \frac{1}{n} \sum_j (x_j - \mu_0)^2.$$

Thus, the likelihood ratio is

$$\begin{aligned} \Lambda(\mathbf{x}) &= \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_2} \right)^{-n/2} \frac{e^{-\frac{1}{2\hat{\theta}_{20}} \sum_j (x_j - \mu_0)^2}}{e^{-\frac{1}{2\hat{\theta}_2} \sum_j (x_j - \hat{\theta}_1)^2}} \\ &= \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_2} \right)^{-n/2} \frac{e^{-n/2}}{e^{-n/2}} \\ &= \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_2} \right)^{-n/2}, \end{aligned}$$

so that, again using (7.23),

$$h(\mathbf{x}) = n \ln \left[ 1 + \frac{1}{n-1} \left( \frac{\bar{x}_n - \mu_0}{s_n(x)/\sqrt{n}} \right)^2 \right].$$

Thus, under  $H_0$  we see that

$$(11.7) \quad h(X) = n \ln \left[ 1 + \frac{1}{n-1} T_{n-1}(X)^2 \right],$$

where

$$T_{n-1}(X) = \frac{\bar{X}_n - \mu_0}{S_n(X)/\sqrt{n}} \sim t_{n-1},$$

or also

$$T_{n-1}(X)^2 \sim F_{1,n-1}$$

by Corollary 4.34. Quite informally, we may expand (11.7) as  $n \rightarrow +\infty$  to obtain

$$\begin{aligned} h(X) &= \ln \left[ 1 + \frac{1}{n-1} T_{n-1}(X)^2 \right]^n \\ &= \ln \left[ 1 + \frac{n}{n-1} T_{n-1}(X)^2 + \dots \right] \\ &= \frac{n}{n-1} T_{n-1}(X)^2 + \dots \\ &\xrightarrow{P} \chi_1^2, \end{aligned}$$

where the dots represent lower order terms (which vanish as  $n \rightarrow +\infty$ ) and we used Remark 6.4 in the last step. Thus, for large samples we may take  $R = [\chi_{1,\alpha}^2, +\infty)$  as the “approximate” rejection interval. Otherwise, we use that

$$\alpha = \mathcal{P}_{\theta_0}(T_{n-1}(X)^2 \geq c_{n,r}) = \mathcal{P}_{\theta_0}(-\sqrt{c_{n,r}} \leq T_{n-1}(X) \leq \sqrt{c_{n,r}}),$$

where

$$c_{n,r} = (n-1)(e^{r/n} - 1),$$

to reject  $H_0$  if either

$$T_{n-1}(\mathbf{x})^2 \in [\mathbf{f}_{1,n-1,\alpha}, +\infty)$$

or equivalently

$$T_{n-1}(\mathbf{x}) \in (-\infty, -t_{n-1,\alpha/2}] \cup [t_{n-1,\alpha/2}, +\infty),$$

which is a more familiar presentation of the test.  $\square$

**Example 11.8.** (F-test for the equality of variances of independent normal populations) We will use the notation of Example 7.36 and Remark 7.37 with the aim of testing

$$H_0 : \sigma_X^2 = \sigma_Y^2 \quad \text{vs} \quad H_a : \sigma_X^2 \neq \sigma_Y^2.$$

Since we assume independence of the samples, Example 8.4 implies that the corresponding likelihood function is

$$L(\mathbf{x}, \mathbf{y}; \theta) = \frac{1}{(2\pi)^{(m+n)/2} \theta_{2X}^{m/2} \theta_{2Y}^{n/2}} e^{-\frac{1}{2} \left( \sum_{j=1}^m \frac{(x_j - \theta_{1X})^2}{\theta_{2X}} + \sum_{k=1}^n \frac{(y_k - \theta_{1Y})^2}{\theta_{2Y}} \right)},$$

where

$$\theta = (\theta_{1X}, \theta_{1Y}, \theta_{2X}, \theta_{2Y}) = (\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2) \in \mathbb{R}^2 \times \mathbb{R}_+^2,$$

so the MLE estimators are

$$\hat{\theta}_{1X} = \bar{X}_m, \quad \hat{\theta}_{1Y} = \bar{Y}_n$$

and

$$\hat{\theta}_{2X} = \frac{m-1}{m} S_X^2 = \frac{1}{m} \sum_j (X_j - \hat{\theta}_{1X})^2, \quad \hat{\theta}_{2Y} = \frac{n-1}{n} S_Y^2 = \frac{1}{n} \sum_k (Y_k - \hat{\theta}_{1Y})^2.$$

Hence,

$$\begin{aligned}
\sup_{\theta \in \Theta} L(\mathbf{x}, \mathbf{y}; \theta) &= L(\mathbf{x}, \mathbf{y}; \hat{\theta}_{1X}, \hat{\theta}_{1Y}, \hat{\theta}_{2X}, \hat{\theta}_{2Y}) \\
&= \frac{1}{(2\pi)^{(m+n)/2} \hat{\theta}_{2X}^{m/2} \hat{\theta}_{2Y}^{n/2}} e^{-\frac{1}{2} \left( \sum_{j=1}^m \frac{(x_j - \hat{\theta}_{1X})^2}{\hat{\theta}_{2X}} + \sum_{k=1}^n \frac{(y_k - \hat{\theta}_{1Y})^2}{\hat{\theta}_{2Y}} \right)} \\
&= \frac{e^{-(m+n)/2}}{(2\pi)^{(m+n)/2} \hat{\theta}_{2X}^{m/2} \hat{\theta}_{2Y}^{n/2}}.
\end{aligned}$$

On the other hand, restriction to  $\Theta_0$  gives

$$L(\mathbf{x}, \mathbf{y}; \theta) = \frac{1}{(2\pi)^{(m+n)/2} \theta_2^{(m+n)/2}} e^{-\frac{1}{2\theta_2} (\sum_j (x_j - \theta_{1X})^2 + \sum_k (y_k - \theta_{1Y})^2)},$$

where  $\theta_2 = \sigma_X^2 = \sigma_Y^2$  is the common variance, so that maximization over all possible values of  $(\theta_{1X}, \theta_{1Y}, \theta_2)$  is achieved at the null MLE  $(\hat{\theta}_{1X}, \hat{\theta}_{1Y}, \hat{\theta}_{20})$ , where

$$\hat{\theta}_{20} = \frac{1}{m+n} \left( \sum_j (x_j - \hat{\theta}_{1X})^2 + \sum_k (y_k - \hat{\theta}_{1Y})^2 \right).$$

It follows that

$$\begin{aligned}
\sup_{\theta \in \Theta_0} L(\mathbf{x}, \mathbf{y}; \theta) &= L(\mathbf{x}, \mathbf{y}; \hat{\theta}_{1X}, \hat{\theta}_{1Y}, \hat{\theta}_{20}) \\
&= \frac{1}{(2\pi)^{(m+n)/2} \hat{\theta}_{20}^{(m+n)/2}} e^{-\frac{1}{2\hat{\theta}_{20}} (\sum_j ((x_j - \hat{\theta}_{1X})^2 + \sum_k (y_k - \hat{\theta}_{1Y})^2))} \\
&= \frac{e^{-(m+n)/2}}{(2\pi)^{(m+n)/2} \hat{\theta}_{20}^{(m+n)/2}},
\end{aligned}$$

so the likelihood ratio is

$$\Lambda(\mathbf{x}, \mathbf{y}) = \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_{2X}} \right)^{-m/2} \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_{2Y}} \right)^{-n/2},$$

and from this we easily deduce that

$$h(\mathbf{x}, \mathbf{y}) = \ln \left[ c \left( 1 + a \left( \frac{S_X^2}{S_Y^2} \right)^{-1} \right)^m \left( 1 + b \frac{S_X^2}{S_Y^2} \right)^n \right],$$

where  $a$ ,  $b$  and  $c$  are certain constants depending only on  $m$  and  $n$  with  $ab = 1$ . Thus, as in Remark 7.37 we see that under  $H_0$ ,

$$h(X, Y) = \ln \left( c \left( 1 + aU(X, Y)^{-1} \right)^m \left( 1 + bU(X, Y) \right)^n \right),$$

where

$$U(X, Y) = \frac{S_X^2}{S_Y^2} \sim \mathbf{F}_{m-1, m-1}.$$

Now note that  $h(\mathbf{x}, \mathbf{y}) \geq e^r$  if and only if  $u := U(\mathbf{x}, \mathbf{y})$  satisfies  $f(u) \geq e^r/c$ , where

$$f(u) := (1 + au^{-1})^m (1 + bu)^n, \quad u > 0.$$

Also, a little Calculus shows that  $f$  is strictly convex with its unique minimal value achieved at  $u_0 = ma/n = m/nb$  (this analysis uses that  $ab = 1$  in a crucial way). Hence, there exist a *maximal*  $0 < \underline{u} < u_0$  and a *minimal*



$\bar{u} > u_0$  with the property that for any  $r$  such that

$$f\left(\frac{ma}{n}\right) = f\left(\frac{m}{nb}\right) = \left(1 + \frac{n}{m}\right)^m \left(1 + \frac{m}{n}\right)^n < \frac{e^r}{c}$$

there holds  $f(u) \geq e^r/a$  whenever either  $u \leq \underline{u}$  or  $u \geq \bar{u}$ . This means that we may reject  $H_0$  if  $u$  falls outside  $(\underline{u}, \bar{u})$ . More concretely, given a confidence level  $\alpha$  small enough, we may reject  $H_0$  if

$$u \in \left(0, f_{m-1, n-1, \alpha/2}^-\right] \cup \left[f_{m-1, n-1, \alpha/2}^+, +\infty\right),$$

which is consistent with (7.48).  $\square$

**Example 11.9.** (F-test for the equality of means of  $p \geq 2$  independent normal populations with a common but unknown variance) We will use the notation from Example 7.40, so we have *independent* random samples  $X_{jk} \sim \mathcal{N}(\mu_j, \sigma^2)$  for  $j = 1, \dots, p$ . Thus,  $\Theta = \mathbb{R}^p \times \mathbb{R}_+$  with  $\theta = (\theta_{11}, \dots, \theta_{1p}, \theta_2)$ , where  $\theta_{1j} = \mu_j$  and  $\theta_2 = \sigma^2$ , the common variance. Also,  $\Theta_0 = \{\theta \in \Theta; \theta_{11} = \dots = \theta_{1p}\}$  and we want to test

$$H_0 : \mu_1 = \dots = \mu_p \quad \text{vs} \quad H_a : \mu_j \neq \mu_{j'} \text{ for some } j \neq j'.$$

Notice that this is precisely the one way ANOVA test in Example 7.40. If  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$  with  $\mathbf{x}_j = (\mathbf{x}_{j1}, \dots, \mathbf{x}_{jn_j}) \in \mathbb{R}^{n_j}$  then the likelihood function is

$$\begin{aligned} L(\mathbf{x}; \theta) &= \prod_{j=1}^p \prod_{k=1}^{n_j} (2\pi\theta_2)^{-n_j/2} e^{-\frac{1}{2\theta_2} \sum_{k=1}^{n_j} (x_{jk} - \theta_{1j})^2} \\ &= (2\pi\theta_2)^{-n/2} e^{-\frac{1}{2\theta_2} \sum_{j=1}^p \sum_{k=1}^{n_j} (x_{jk} - \theta_{1j})^2}. \end{aligned}$$

By passing to the log likelihood function and maximizing over  $\Theta$  in the usual way we see that the MLE for  $\theta$  is  $\hat{\theta} = (\hat{\theta}_{11}, \dots, \hat{\theta}_{1p}, \hat{\theta}_2)$ , where

$$\hat{\theta}_{1j}(\mathbf{x}) = \bar{\mathbf{x}}_{j\bullet} = \frac{1}{n_j} \sum_{k=1}^{n_j} \mathbf{x}_{jk} \quad \text{and} \quad \hat{\theta}_2(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^p \sum_{k=1}^{n_j} (\mathbf{x}_{jk} - \bar{\mathbf{x}}_{j\bullet})^2$$

are realizations of  $\bar{X}_{j\bullet}$  and  $S_{\text{Within}}^2(X)/n$ , respectively. On the other hand, restricted to  $\Theta_0$  we have

$$L(\mathbf{x}; \theta) = (2\pi\theta_2)^{-n/2} e^{-\frac{1}{2\theta_2} \sum_{j=1}^p \sum_{k=1}^{n_j} (\mathbf{x}_{jk} - \theta_0)^2},$$

where  $\theta_0 = \mu_1 = \dots = \mu_p$  is the common mean, so that maximization over  $\Theta_0$  gives the corresponding null MLEs

$$\hat{\theta}_{00}(\mathbf{x}) = \mathbf{x}_{\bullet\bullet} = \frac{1}{n} \sum_{j=1}^p n_j \bar{\mathbf{x}}_{j\bullet} = \frac{1}{n} \sum_{j=1}^p \sum_{k=1}^{n_j} \mathbf{x}_{jk} \quad \text{and} \quad \hat{\theta}_{20}(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^p \sum_{k=1}^{n_j} (\mathbf{x}_{jk} - \bar{\mathbf{x}}_{\bullet\bullet})^2,$$

which are realizations of  $\bar{X}_{\bullet\bullet}$  and  $S_{\text{Total}}^2(X)/n$ , respectively. It follows that

$$\begin{aligned} \Lambda(\mathbf{x}) &= \frac{(2\pi\hat{\theta}_{20}(\mathbf{x}))^{-n/2} e^{-\frac{1}{2\hat{\theta}_{20}(\mathbf{x})} \sum_{j=1}^p \sum_{k=1}^{n_j} (\mathbf{x}_{jk} - \mathbf{x}_{\bullet\bullet})^2}}{(2\pi\hat{\theta}_2(\mathbf{x}))^{-n/2} e^{-\frac{1}{2\hat{\theta}_2(\mathbf{x})} \sum_{j=1}^p \sum_{k=1}^{n_j} (\mathbf{x}_{jk} - \bar{\mathbf{x}}_{j\bullet})^2}} \\ &= \left( \frac{\hat{\theta}_{20}(\mathbf{x})}{\hat{\theta}_2(\mathbf{x})} \right)^{-n/2} \frac{e^{-n/2}}{e^{-n/2}}, \end{aligned}$$

so that using (7.52),

$$h(\mathbf{x}) = n \ln \left( 1 + \frac{p-1}{n-p} \frac{s_{\text{Between}}^2(\mathbf{x})/(p-1)}{s_{\text{Within}}^2(\mathbf{x})/(n-p)} \right).$$

Hence, as in Example 7.40 we see that under  $H_0$ ,

$$h(X) = n \ln \left( 1 + \frac{p-1}{n-p} V \right),$$

where

$$V = \frac{S_{\text{Between}}^2/(p-1)}{S_{\text{Within}}^2/(n-p)} \sim F_{p-1, n-p}.$$

Again, we may expand this as  $n \rightarrow +\infty$  to find that

$$\begin{aligned} h(X) &= \ln \left( 1 + \frac{p-1}{n-p} V \right)^n \\ &= \ln \left( 1 + \frac{n(p-1)}{n-p} V + \dots \right) \\ &= (p-1)V + \dots \\ &\xrightarrow{d} \chi_{p-1}^2, \end{aligned}$$

where we used Remark 6.4 in the last step. Thus, for large samples we may take  $R = [\chi_{p-1, \alpha}^2, +\infty)$  as the rejection interval. Otherwise, we use that

$$\alpha = \sup_{\theta_0 \in \Theta_0} \mathcal{P}_{\theta_0} \left( V(\mathbf{x}) \geq \frac{n-p}{p-1} (e^{r/n} - 1) \right)$$

to reject  $H_0$  if

$$V(\mathbf{x}) \in [\mathbf{f}_{p-1, n-p, \alpha}, +\infty)$$

as in (7.60). □

**Example 11.10.** (F-test for statistical significance of the linear regression model) We consider here the linear regression model in (9.14), whose likelihood function  $L(\mathbf{y}; \beta, \sigma^2)$  is given by (9.15), in order to test the full “intercept-only” hypothesis appearing in Example 9.8:

$$(11.8) \quad H_0 : \beta_1 = \dots = \beta_p = 0 \quad \text{vs} \quad H_a : \beta_j \neq 0 \text{ for some } j.$$

In other words, the null hypothesis here says that  $\mathbf{X}$  has no influence whatsoever on  $\mathbf{Y}$  so its rejection provides statistical evidence for employing the model as it is posed in Example 9.3 (that is, with the full “slope”  $(\beta_1, \dots, \beta_n)$  included). We have  $\theta = (\beta, \theta_2)$ , where  $\theta_2 = \sigma^2$ , so the usual calculation implies that the corresponding MLE is  $(\hat{\beta}, \hat{\theta}_2)$ , where

$$\hat{\theta}_2 = \frac{1}{n} \|\mathbf{Y} - \mathbf{r}\hat{\beta}\|^2 = \frac{1}{n} \|\mathbf{r}\beta + \mathbf{e} - \mathbf{r}\hat{\beta}\|^2,$$

so that (9.53) gives

$$\hat{\theta}_2 = \frac{|\hat{\mathbf{e}}|^2}{n} = \frac{SS_{\text{Res}}}{n},$$

and we verify that

$$\sup_{\theta \in \Theta} L(\mathbf{y}; \beta, \theta_2) = L(\mathbf{y}; \hat{\beta}, \hat{\theta}_2) = (2\pi SS_{\text{Res}}/n)^{-n/2} e^{-n/2}.$$

On the other hand, under the null hypothesis,

$$L(\mathbf{y}; \beta_0, \theta_2) = (2\pi\theta_2)^{-n/2} e^{-\frac{\|\mathbf{y} - \beta_0 \mathbf{1}\|^2}{2\theta_2}},$$

so that the null MLE estimator for  $\theta_2$  is

$$(11.9) \quad \hat{\theta}_{20} = \frac{1}{n} \|\mathbf{Y} - \hat{\beta}_0 \mathbf{1}\|^2 \stackrel{(9.20)}{=} \frac{1}{n} \|\mathbf{Y} - \bar{\mathbf{Y}} \mathbf{1}\|^2 = \frac{SS_{\mathbf{Y}\mathbf{Y}}}{n},$$

which gives

$$\sup_{\theta \in \Theta_0} L(\mathbf{y}; \beta, \theta_2) = L(\mathbf{y}; \hat{\beta}_0, \hat{\theta}_{20}) = (2\pi SS_{\mathbf{Y}\mathbf{Y}}/n)^{-n/2} e^{-n/2}.$$

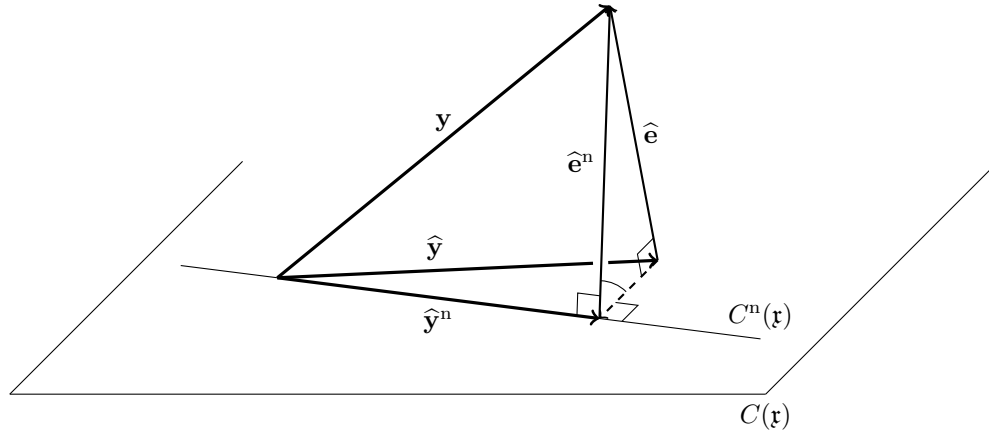


FIGURE 2. The geometry of nested models

It then follows that the likelihood ratio statistics is

$$(11.10) \quad h(\mathbf{y}) = \ln \left( \frac{\hat{\theta}_{20}}{\hat{\theta}_2} \right)^n = \ln \left( \frac{SS_{\mathbf{Y}\mathbf{Y}}}{SS_{\text{Res}}} \right)^n = \ln \left( 1 + \frac{SS_{\text{Reg}}}{SS_{\text{Res}}} \right)^n,$$

where we used that, as in (9.45),

$$(11.11) \quad SS_{\mathbf{Y}\mathbf{Y}} = SS_{\text{Reg}} + SS_{\text{Res}}.$$

We now proceed to the appropriate counting of degrees of freedom as we did in Example 11.9. We see from Propositions 9.10 and 9.12 that  $\sigma^{-2}SS_{\text{Res}} \sim \chi_{n-p-1}^2$  is independent of  $SS_{\text{Reg}}$  and hence of  $\hat{\mathbf{Y}} = \mathbf{x}\hat{\beta}$ , so  $SS_{\text{Res}}$  is independent of  $SS_{\text{Reg}}$  (recall that we are conditioning on  $\mathbf{x} = \mathbf{x}$ ). On the other hand, under  $H_0$  we have  $\mathbf{Y} \sim \mathcal{N}(\beta_0 \mathbf{1}, \sigma^2 I_{n \times n})$  and  $\hat{\beta}_0 = \bar{\mathbf{Y}}$ , which gives  $\sigma^{-2}SS_{\mathbf{Y}\mathbf{Y}} \sim \chi_{n-1}^2$  by Proposition 7.24. Thus, again under  $H_0$ , we get from (11.11) that  $\sigma^{-2}SS_{\text{Reg}} \sim \chi_p^2$  and we conclude that

$$h(\mathbf{Y}) = \ln \left( 1 + \frac{p}{n-p-1} W(\mathbf{Y}) \right)^n,$$

where

$$W(\mathbf{Y}) = \frac{SS_{\text{Reg}}/p}{SS_{\text{Res}}/(n-p-1)} \sim F_{p, n-p-1}.$$

As in Example 11.9,  $W(\mathbf{Y}) \xrightarrow{d} \chi_p^2$  as  $n \rightarrow +\infty$ , so for large samples we may take  $R = [\chi_{p, \alpha}^2, +\infty)$  as the rejection interval. Otherwise, we must reject  $H_0$  if

$$(11.12) \quad W(\mathbf{y}) \in [\mathbf{f}_{p, n-p-1, \alpha}, +\infty).$$

We mention that the theoretical procedure leading to the F-test above, based on a likelihood ratio test, is flexible enough to handle a general linear hypothesis test on the parameters, in which the null hypothesis may be expressed as  $B\beta = c$ , where  $B$  is a suitable  $q \times (p+1)$  matrix and  $c$  is a  $q$ -vector; see [Ame85, Subsection 1.5] and [SL03, Chapter 4]. For instance, if  $c = \vec{0}$  and  $B$  is suitably chosen then we can form the test

$$(11.13) \quad H_0^n : \beta_{q+1} = \cdots = \beta_p = 0 \quad \text{vs} \quad H_a^n : \beta_j \neq 0 \text{ for some } j \in \{q+1, \dots, p\},$$

which compares the full model and a new null model in which only the first  $q$  independent variables possibly appear as significant predictors; as indicated in Figure 2, the corresponding design spaces satisfy  $C^n(\mathbf{x}) \subset C(\mathbf{x})$  with  $\dim C(\mathbf{x}) \setminus C^n(\mathbf{x}) = p - q$ , this being the reason why the models are *nested*; see Remark 9.11. If we view

$SS_{\mathbf{Y}\mathbf{Y}}$  in (11.10) and (11.11) as the residual sum of squares (i.e. the norm squared residual) of the null model in (11.8) and proceed by analogy, it is not hard to check that the likelihood ratio statistics now is

$$h^n(\mathbf{Y}) = \ln \left( \frac{SS_{\text{Res}}^n}{SS_{\text{Res}}} \right)^n = \ln \left( 1 + \frac{p-q}{n-p-1} W^n(\mathbf{Y}) \right)^n,$$

where

$$W^n(\mathbf{Y}) = \frac{(SS_{\text{Res}}^n - SS_{\text{Res}})/(p-q)}{SS_{\text{Res}}/(n-p-1)}$$

and  $SS_{\text{Res}}^n$  is the residual sum of squares of the null model in (11.13). Since the usual counting of degrees of freedom shows that  $W^n(\mathbf{Y}) \sim F_{p-q, n-p-1}$  under  $H_0^n$ , we find that the null hypothesis in (11.13) gets rejected if

$$W^n(\mathbf{y}) \in [f_{p-q, n-p-1, \alpha}, +\infty),$$

the obvious extension of (11.12). Put in another way, if  $SS_{\text{Res}}^n$  and  $SS_{\text{Res}}$  are close to each other, which intuitively means that the null model fits as well as the full model, then  $W^n(\mathbf{y})$  is small and hence  $H_0^n$  should *not* be rejected. In any case, the geometry backing not only this latter assertion but also the whole argument above is fully discernible from Figure 2, where  $SS_{\text{Res}}^n = \|\hat{\mathbf{e}}^n\|^2$ ,  $SS_{\text{Res}}^n - SS_{\text{Res}} = \|\hat{\mathbf{e}}^n - \hat{\mathbf{e}}\|^2$ , the squared norm of the dashed vector, and so on.  $\square$

**Remark 11.11.** (*p*-value) As already observed, in all examples above the likelihood ratio statistics  $h(X)$  is  $\Theta_0$ -ancillary in the sense that its distribution does not depend on  $\theta \in \Theta_0$ . In those cases, an equivalent way of reporting the result of a likelihood ratio test is to look at the corresponding *p*-value

$$\mathfrak{p} = \mathcal{P}_\theta(h(X) \geq h(x)), \quad \theta \in \Theta_0,$$

where  $h(x)$  is the observed value of  $h(X)$ . Thus,  $\mathfrak{p}$  is the probability of finding, under  $H_0$ , an observed value at least as extreme as the one actually observed. With this terminology,  $H_0$  gets rejected if  $\mathfrak{p} \leq \alpha$ , which is just a rephrasing of the rejection condition  $h(x) \geq r$ . Although this seems to be the preferred way of summarizing the outcome of a test in Applied Statistics, it is argued that the common misinterpretation of regarding  $\mathfrak{p}$  as the probability that  $H_0$  is true, thus erroneously accepting the validity of the alternative hypothesis (with high probability) if  $\mathfrak{p}$  is found to be sufficiently small, may be a source of confusion leading to “*P*-hacking”, “the replication crisis”, etc.; see [HB03, WL16, FP15, Gib21] for more on this quite controversial issue.  $\square$

In all examples above where we have been able to directly carry out the corresponding computation, the *asymptotic* likelihood ratio statistics turned out to be  $\chi_l^2$ -distributed, where  $l = \dim \Theta - \dim \Theta_0$ . In fact, this is a general phenomenon which substantially simplifies the implementation of the test for large samples.

**Theorem 11.12.** [Wil38] *Under the conditions above, and requiring suitable regularity assumptions on the underlying statistical model as usual, there holds  $h(X) \xrightarrow{d} \chi_l^2$  as  $n \rightarrow +\infty$  and under  $H_0$ .*

*Proof.* We only sketch the argument, which relies on the (multi-dimensional version) of the proof of Theorem 8.23 on the asymptotic normality of ML estimators (and the simplifying assumption that we may choose rectangular coordinates  $(\theta_1, \dots, \theta_p)$  on  $\Theta$  so that  $\Theta_0$  is singled out by  $\theta_{k+1} = \dots = \theta_p = 0$ ). Now, under  $H_0$  the true parameter value, say  $\theta$ , lies in  $\Theta_0$ . Moreover, if  $n$  is large enough then both

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} L(\mathbf{x}; \theta) \quad \text{and} \quad \hat{\theta}_0 = \operatorname{argmax}_{\theta \in \Theta_0} L(\mathbf{x}; \theta)$$

are close to  $\theta$  and hence close to each other, so we may expand about  $\hat{\theta}$ ,

$$\begin{aligned} h(X) &= -2 \left( \ln L(X; \hat{\theta}) - \ln L(X; \hat{\theta}_0) \right) \\ &\approx \langle \hat{\theta} - \hat{\theta}_0, (-\nabla_{\theta\theta}^2)(\ln L(X; \hat{\theta}))(\hat{\theta} - \hat{\theta}_0) \rangle, \end{aligned}$$

where we have discarded terms of order at least three and used that  $\nabla_\theta \ln L(\mathbf{x}; \hat{\theta}) = 0$  by the definition of  $\hat{\theta}$ ; compare with (8.31). From (8.37) we also know that

$$(-\nabla_{\theta\theta}^2)(\ln L(X; \hat{\theta})) \xrightarrow{p} \mathcal{F}(\theta),$$

where  $\mathcal{F}(\theta)$  is the Fisher information matrix of a single observation. Moreover,

$$W := \sqrt{n} \nabla_\theta \ln L(X; \theta) \xrightarrow{d} \mathcal{N}(\vec{0}, \mathcal{F}(\theta))$$

by (8.38), so that (8.32) may be rewritten as

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{F}(\theta)^{-1}W.$$

Also, by examining the maximization problem restricted to  $\Theta_0$  which yields  $\hat{\theta}_0$ , it is not hard to check that

$$\sqrt{n}(\hat{\theta}_0 - \theta) \xrightarrow{d} \mathcal{G}(\theta)W,$$

where  $\mathcal{G}(\theta)$  is a certain symmetric matrix with rank  $k = \dim \Theta_0$  and satisfying

$$\mathcal{G}(\theta)\mathcal{F}(\theta)\mathcal{G}(\theta) = \mathcal{G}(\theta).$$

Thus, eliminating  $\theta$  in the convergences above and using the expansion we get

$$h(X) \approx \langle W, (\mathcal{F}(\theta)^{-1} - \mathcal{G}(\theta))W \rangle$$

Now, if  $\mathcal{F}(\theta) = B^2$  then using Remark 4.15 we have that  $Z = B^{-1}W \approx \mathcal{N}(\vec{0}, \text{Id}_p)$ , a standard normal vector, so that

$$\begin{aligned} h(X) &\approx \langle BZ, (\mathcal{F}(\theta)^{-1} - \mathcal{G}(\theta))BZ \rangle \\ &= \langle Z, B(\mathcal{F}(\theta)^{-1} - \mathcal{G}(\theta))BZ \rangle \\ &= \langle Z, (\text{I}_p - B\mathcal{G}(\theta)B)Z \rangle, \end{aligned}$$

where  $B\mathcal{G}(\theta)B$  is easily seen to be idempotent with the same rank as  $\mathcal{G}(\theta)$ . Hence,  $\text{I}_p - B\mathcal{G}(\theta)B$  is idempotent as well with rank  $l = p - k = \dim \Theta - \dim \Theta_0$  and the result follows from Proposition 4.27.  $\square$

We refer to [Ame85, Subsection 4.5.1] and [Sha08, Theorem 6.5] for those interested in filling out the omitted details in the argument above. Also, it is worthwhile mentioning that Wilks' original proof in [Wil38] is equally elegant as it involves checking that  $\phi_{h(X)}$ , the characteristic function of  $h(X)$ , asymptotically approaches  $\phi_{\chi^2}$ . In any case, we stress that it is not required in Theorem 11.12 that  $h(X)$  is  $\Theta_0$ -ancillary so in a sense the result guarantees that this property gets restored in the asymptotic regime. Finally, we note that the material above by no means exhausts the rich literature on hypothesis testing and extensive treatments may be found elsewhere [Ame85, DM88, Wel96, LR05, Cas08, Sha08, Hay11, DS14].

We now discuss a non-conventional hypothesis testing and its connection with a major result in Differential Geometry, namely, Weyl's formula for the volume of tubes [Wey39, Gra03].

**Example 11.13.** (Testing for an additional term in the linear model and Weyl's formula for the volume of tubes) Let us consider a (possibly non-linear) perturbation of the linear model with a normal error from Example 9.3,

$$\mathbf{Y}_j = \mathbf{x}_j\beta + cf_j(\mathbf{X}_j, \gamma) + \mathbf{e},$$

a setting first considered in a seminal paper by H. Hotelling [Hot39]. Here,  $f$  is known but  $c \in \mathbb{R}$  and  $\gamma \in \mathbb{R}^k$  are unknown parameters and our aim is to test

$$H_0 : c = 0 \quad \text{vs} \quad H_a : c \neq 0,$$

so that not being able to reject  $H_0$  indicates statistical evidence for ignoring  $f$  in the model design. In the geometric language of Remark 9.11, the null hypothesis says that  $\mathbb{E}(\mathbf{Y}|\mathbf{x}=\mathbf{x}) \in C(\mathbf{x}) \equiv \mathbb{R}^{p+1}$ , whereas the alternative adds a multiple of  $f_\gamma = f(\cdot, \gamma)$  to this vector. Without loss of generality, we may assume that

$f_\gamma \in C(\mathfrak{r})^\perp \equiv \mathbb{R}^{n-p-1}$  and, in order to make the model identifiable, that  $f_\gamma$  is not a multiple of  $f_{\gamma'}$  if  $\gamma \neq \gamma'$ , but notice that the model still gets *non*-identifiable under  $H_0$ , so the usual dimensional counting that would allow us to determine the asymptotic distribution of the likelihood ratio statistics  $h$  via Theorem 11.12 does not apply. In particular, there is no point here in working with  $h$  so we turn our attention to the corresponding likelihood ratio  $\Lambda = e^{-h/2}$ , thus rejecting  $H_0$  if this ratio, when observed, is conveniently small. Now, the full likelihood function is

$$L(\mathbf{y}; \beta, \theta_2, c, \gamma) = (2\pi\theta_2)^{-n/2} e^{-\frac{\|\mathbf{y} - \mathfrak{r}\beta - cf_\gamma\|^2}{2\theta_2}}, \quad \theta_2 = \sigma^2,$$

which under  $H_0$  reduces to the usual likelihood function of the linear model treated in Example 11.10:

$$L(\mathbf{y}; \beta, \theta_2) = (2\pi\theta_2)^{-n/2} e^{-\frac{\|\mathbf{y} - \mathfrak{r}\beta\|^2}{2\theta_2}}.$$

Hence,

$$\sup_{\beta, \theta_2} L(\mathbf{y}; \beta, \theta_2) = (2\pi\hat{\theta}_{20})^{-n/2} e^{-n/2},$$

where

$$\hat{\theta}_{20} = \frac{1}{n} \|\mathbf{Y} - \mathfrak{r}\hat{\beta}\|^2 = \frac{1}{n} \|\hat{\mathbf{e}}\|^2$$

and  $\hat{\mathbf{e}}$  is the residual. It follows that the likelihood ratio is

$$\Lambda = \left( \frac{\hat{\theta}_{20}}{\theta_2} \right)^{-n/2},$$

where

$$\hat{\theta}_2 = \frac{1}{n} \|\mathbf{Y} - \mathfrak{r}\hat{\beta} - \hat{c}f_{\hat{\gamma}}\|^2$$

and

$$(\hat{\beta}, \hat{c}, \hat{\gamma}) = \operatorname{argmax}_{\beta, c, \gamma} L(\mathbf{y}; \beta, c, \gamma) = \operatorname{argmin}_{\beta, c, \gamma} \|\mathbf{y} - \mathfrak{r}\beta - cf_\gamma\|^2.$$

Now, using that  $f_\gamma^t \mathfrak{r}\beta = 0$  (in particular,  $f_\gamma^t \mathfrak{r}\hat{\beta} = 0$ ) we compute

$$\begin{aligned} (\hat{\beta}, \hat{c}, \hat{\gamma}) &= \operatorname{argmin}_{\beta, c, \gamma} \|\mathbf{y} - \mathfrak{r}\beta\|^2 - 2cf_\gamma^t \mathbf{y} + c^2 \|f_\gamma\|^2 \\ &= \operatorname{argmin}_{c, \gamma} \|\mathbf{y} - \mathfrak{r}\hat{\beta}\|^2 - 2cf_\gamma^t \mathbf{y} + c^2 \|f_\gamma\|^2 \\ &= \operatorname{argmin}_{c, \gamma} \|\mathbf{y} - \mathfrak{r}\hat{\beta}\|^2 - 2cf_\gamma^t (\mathbf{y} - \mathfrak{r}\hat{\beta}) + c^2 \|f_\gamma\|^2 \\ &= \operatorname{argmin}_{c, \gamma} \|\hat{\mathbf{e}} - cf_\gamma\|^2, \end{aligned}$$

and since

$$\hat{c} = \operatorname{argmin}_c \|\hat{\mathbf{e}} - cf_\gamma\|^2 = \frac{f_\gamma^t \hat{\mathbf{e}}}{\|f_\gamma\|^2},$$

we see that

$$\begin{aligned} \Lambda^{2/n} &= \inf_{\gamma} \frac{\|\hat{\mathbf{e}} - \hat{c}f_\gamma\|^2}{\|\hat{\mathbf{e}}\|^2} \\ &= \inf_{\gamma} \left( 1 - \left( \frac{f_\gamma^t \hat{\mathbf{e}}}{\|f_\gamma\| \|\hat{\mathbf{e}}\|} \right)^2 \right) \\ &= 1 - \sup_{\gamma} (\tilde{f}_\gamma^t \mathbf{U})^2. \end{aligned}$$

where, as  $\gamma$  varies,  $\tilde{f}_\gamma = f_\gamma / \|f_\gamma\|$  traces a subset  $M_\gamma \subset \mathbb{S}^{n-p-2}$ , the unit sphere in  $C(\mathfrak{r})^\perp$ , and  $\mathbf{U} = \hat{\mathbf{e}} / \|\hat{\mathbf{e}}\|$  is a random vector also taking values in  $\mathbb{S}^{n-p-2}$ . Since

$$\tilde{f}_\gamma^t \mathbf{u} = \cos \operatorname{dist}(f_\gamma, \mathbf{u}),$$

where  $\text{dist}$  is the intrinsic distance in  $\mathbb{S}^{n-p-2}$ , we may choose as rejection region the “tubular neighborhood”

$$B_\rho(M_\gamma) = \{\vartheta \in \mathbb{S}^{n-p-2}; \text{dist}(\vartheta, M_\gamma) \leq \rho\}$$

of radius  $\rho > 0$  around  $M_\gamma$ . Now, under  $H_0$  we know from Remark 9.11 that  $\mathbf{U}$  is uniformly distributed in  $\mathbb{S}^{n-p-2}$  and we conclude that the significance level  $\alpha$  of the test satisfies

$$\alpha = P(\mathbf{U} \in B_\rho(M_\gamma)) = \text{vol}_{P_U}(B_\rho(M_\gamma)),$$

where  $\text{vol}_{P_U}$  is the (normalized) intrinsic volume (so that  $\text{vol}_{P_U}(\mathbb{S}^{n-p-2}) = 1$ ). Thus, in order to determine the rejection “tube” associated to a given confidence level, an explicit formula for the volume of the tube is required, at least for  $\rho$  small enough. This turns out to be a rather formidable geometric problem which has been completely solved by H. Weyl [Wey39] in case  $M$  is a closed submanifold of a space form (a Riemannian manifold with constant sectional curvature)<sup>42</sup>. Unfortunately, this Weyl’s formula does not directly apply to this problem (as  $M_\gamma$  may be only piecewise smooth or carry a boundary, etc.) so adjustments, mainly based on suitable approximations, are required [Nai90].  $\square$

**Example 11.14.** (Scheffé-type simultaneous band for the mean response in a normal linear model and the volume of tubes, again) One is often interested in obtaining simultaneous confidence bands for the mean response  $\mathbf{x}^t \beta$  in a normal regression model (as in Example 9.3) with  $\mathbf{x}$  varying in some subset  $\mathcal{S} \subset \mathbb{R}^p$ , say  $\mathcal{S}$  diffeomorphic to an interval or a rectangle and so on; here we retain the notation of Examples 9.13 and 9.14. Although the Scheffé-type band in (9.42) may be applied to this end, it certainly provides a wider band than required for a given confidence level, so a sensible strategy here is to seek for  $c > 0$  satisfying

$$(11.14) \quad \mathbf{x}^t \beta \in \left[ \mathbf{x}^t \hat{\beta} \mp c \hat{\sigma} \sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}} \right] \forall \mathbf{x} \in \mathcal{S} \text{ with prob. } 1 - \delta.$$

Now, in terms of

$$(11.15) \quad \mathbf{m} = \sigma \mathbf{n} = (\mathbf{p}^t)^{-1}(\hat{\beta} - \beta) \sim \mathcal{N}(\vec{0}, \sigma^2 \text{Id}_{p+1}),$$

$\varepsilon(\mathbf{x}) = \mathbf{p}\mathbf{x}/|\mathbf{p}\mathbf{x}| \in \mathbb{S}^p \subset \mathbb{R}^{p+1}$ ,  $\mathbf{x} \in \mathcal{S}$ , and  $\mathbf{u} = \mathbf{m}/\|\mathbf{m}\|$ , a uniformly distributed random vector in  $\mathbb{S}^p$ , we have

$$\frac{\mathbf{x}^t(\hat{\beta} - \beta)}{\sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}}} = \frac{(\mathbf{p}^t)^{-1} \mathbf{x} \mathbf{m}}{\sqrt{(\mathbf{p}\mathbf{x})^t \mathbf{p}\mathbf{x}}} = (\varepsilon(\mathbf{x})^t \mathbf{u}) \|\mathbf{m}\|,$$

so if  $T = \hat{\sigma}/\|\mathbf{m}\|$  then (11.14) expresses the corresponding tail probability as

$$\begin{aligned} \delta &= P \left( \sup_{\mathbf{x} \in \mathcal{S}} \left| \frac{\mathbf{x}^t(\hat{\beta} - \beta)}{\sqrt{\mathbf{x}^t \mathbf{s} \mathbf{x}}} \right| \geq c \hat{\sigma} \right) \\ &= P \left( \sup_{\mathbf{x} \in \mathcal{S}} |\varepsilon(\mathbf{x})^t \mathbf{u}| \geq cT \right). \end{aligned}$$

Hence, using a notation similar to the one of the previous example, we see that the coverage probability in (11.14) is

$$1 - \delta = \text{vol}_{P_u}(B_{cT}(M_\varepsilon \cup M_{-\varepsilon})),$$

<sup>42</sup>In case  $M \subset \mathbb{R}^l$  has dimension  $m$ , Weyl’s formula says that

$$\text{vol}(B_\rho(M)) = c_{m,l} \rho^{l-m} \sum_{q=0}^{[m/2]} \left( d_{m,l,q} \int_M \kappa_{2q} dM \right) \rho^{2q},$$

where  $c_{m,l}$  and  $d_{m,l,q}$  are positive constants (only depending on the indicated natural parameters) and  $\kappa_{2q}$  is a certain (universal) polynomial expression which is homogeneous of degree  $q$  in the curvature tensor of  $M$ . In particular,  $\text{vol}(B_\rho(M))$  depends only on  $\rho$  and the *intrinsic* geometry of  $M$  and not on the specific way it is embedded in  $\mathbb{R}^l$ . For a masterly account of this remarkable result and its many applications (including a proof, in full generality, of the Chern-Gauss-Bonnet formula in Riemannian Geometry) we refer to [Gra03]. Also, for a brief overview of the ubiquitous role the Gauss-Bonnet curvatures  $\kappa_{2q}$  play in Riemannian Geometry and related areas, see [Lab07] and the references therein.

so that being able to compute the (normalized) volume of certain tubes around  $M_\varepsilon \cup M_{-\varepsilon} \subset \mathbb{S}^p$  intervenes in determining the critical value  $c$ . Since  $T$  is independent of  $\mathbf{u}$ , this may be rewritten as

$$1 - \delta = \int_0^{1/c} \text{vol}_{P_{\mathbf{u}}} (B_{ct}(M_\varepsilon \cup M_{-\varepsilon})) \psi_T(t) dt,$$

where  $\psi_T$  is the pdf of  $T$ . It follows from (9.32), (11.15) and the independence of  $\mathbf{u}$  and  $\hat{\sigma}$  that  $(p+1)T^2 \sim F_{n-p-1, p+1}$ , so (2.8) applies to give

$$\psi_T(t) = 2(p+1)t \psi_{F_{n-p-1, p+1}}((p+1)t^2), \quad t \geq 0,$$

and the substitution  $t = \cos \theta / c$  in the previous integral leads to

$$1 - \delta = \int_0^{\pi/2} \text{vol}_{P_{\mathbf{u}}} (B_{\cos \theta}(M_\varepsilon \cup M_{-\varepsilon})) \psi(\theta) d\theta,$$

where

$$\psi(\theta) = \frac{2(p+1) \sin \theta \cos \theta}{c^2} \psi_{F_{n-p-1, p+1}}\left(\frac{(p+1) \cos^2 \theta}{c^2}\right).$$

To see how this implies Scheffé's original contribution, note that if  $\mathcal{S} = \mathcal{R}^p$  then the volume function within this integral clearly equals 1 identically, so the substitution  $\tau = (p+1) \cos^2 \theta / c^2$  gives

$$1 - \delta = \int_0^{(p+1)/c^2} \psi_{F_{n-p-1, p+1}}(\tau) d\tau,$$

and using Corollary 4.34 with  $\tau' = 1/\tau$ ,

$$1 - \delta = \int_0^{c^2/(p+1)} \psi_{F_{p+1, n-p-1}}(\tau') d\tau'.$$

But this means that

$$\frac{c^2}{p+1} = \mathbf{f}_{p+1, n-p-1, \delta},$$

which recovers (9.42) as promised. In general, when  $\mathcal{S}$  is a proper subset of  $\mathcal{R}^p$ , “approximate” versions of Weyl's formula are needed in order to establish simultaneous confidence bands based on the computations above [SL94, Liu10].  $\square$

## 12. A BRIEF OVERVIEW OF “CLASSICAL” PARAMETRIC ESTIMATION

Looking in retrospect, we may now briefly describe some of the main mathematical underpinnings in the “classical” approach to Parametric Estimation Theory. Given a (say, uni-dimensional) statistical model with random sample

$$(12.1) \quad X_1, \dots, X_n \sim \psi_\theta, \quad \theta \in \mathbb{R},$$

of size  $n$ , we may form a statistic  $h(X_1, \dots, X_n)$  with the purpose of designing an estimator  $\hat{\theta}$  for  $\theta$  as in (7.2). We stress that  $h$  should *not* depend on the unknown parameter  $\theta$ . As a first step toward probing the efficiency of  $\hat{\theta}$ , its mean squared error  $\text{mse}(\hat{\theta})$  should be computed (or at least reliably estimated). As explained in detail above for the case  $\hat{\theta} = \hat{\sigma}_c^2$ , this involves computing the associated variance, an usually demanding task which has been substantially simplified here upon the requirement that  $\psi_\theta$  is already normal (this somewhat restrictive assumption is sometimes justified on heuristical grounds by combining CLT (Theorem 6.5) and the “hypothesis of elementary errors” [Fis11, Chapter 3]). If  $\text{mse}(\hat{\theta})$  is shown to be sufficiently small (so that a good performance for  $\hat{\theta}$  is guaranteed), further information about the distribution of  $\hat{\theta}$  may be needed depending on the goal we have in mind. As illustrated in Subsection 7.3 for the sample mean, if our (perhaps too ambitious) aim is to provide “small sample” confidence intervals estimates for  $\theta$ , we should be able to find, for any  $n$ , an



explicit expression for the sampling distribution  $\psi_{\theta}^{(n)}$  of  $h(X_1, \dots, X_n)$ , or of a pivotal quantity thereof. This turns out to be a rather delicate matter, even under sample normality, given that  $h$  might depend *non-linearly* on the sample (as in the case of the sample variance). In this regard, we mention that an earlier breakthrough in the theory is Student's determination of the sampling distribution of his pivotal quantity  $T_{n-1}$  in (7.26), which immediately yields "small sample" confidence intervals for the population mean  $\mu$ , but notice that this has been accomplished only under the assumption that the original sample already follows a normal. In any case, this remarkable contribution certainly inspired R. Fisher to pursue his "geometric method" which not only allowed him to put Student's result on solid mathematical grounds (Remark 7.31) but also led to his celebrated calculation of the sampling distribution of the correlation coefficient (Example 7.39)<sup>43</sup>.

A much less challenging task is to seek for *asymptotic* information, with a preliminary step in this direction being checking that  $\text{mse}(\hat{\theta}_n) \rightarrow 0$  as  $n \rightarrow +\infty$ , as this guarantees that  $\hat{\theta}_n$  is consistent (cf. Proposition 7.11, which may be viewed as a version of LLN in this broader setting). Since it is highly desirable to aggregate to this "asymptotic point estimate" the corresponding "dispersion analysis", a sensible strategy here would be to explicitly determining a distribution  $\psi$  such that an appropriate standardization of  $\hat{\theta}_n$  converges in distribution to  $\psi$ . Hopefully, this should provide estimates for  $\sqrt{n}(\hat{\theta}_n - \theta)$  to the effect that its sampling distribution gets tightly confined around the unknown parameter  $\theta$ , which makes it possible the construction of reliable confidence intervals in terms of "tail" probabilities of  $\psi$  in this asymptotic regime. This approach has been outlined in Subsection 7.3 for  $\hat{\theta}_n = \bar{X}_n$ , the sample mean, and more generally, it is precisely the content of the "asymptotic normality" in Theorem 8.23, which applies to a rather general class of consistent ML estimators, thus confirming that they are asymptotically efficient in the sense that, when the sample size grows indefinitely, they display the smallest dispersion (as measured by their standard deviation) allowed by the Cramér-Rao lower bound in Theorem 8.17; compare with the discussion in Remark 8.25<sup>44</sup>. Full accounts of the most basic aspects of this "large sample statistics" may be found in [NM94, Leh99, LC06, Fer17, VdV00].

It is remarkable that the modern strategy outlined above essentially reproduces the proposal first put forward by R. Fisher in his foundational paper [Fis22] written a century ago. There, after declaring that "the object of statistical methods is the reduction of data", Fisher elects as its primary task to solve the "Problems of Specification" which arise "in the choice of the right mathematical form of the population": in modern language this means finding the pertinent statistical model (as in (12.1)) to work with<sup>45</sup>. Next he moves to the "Problems of Estimation", which seek to determine the right statistic "designed to estimate the values of the parameters of the hypothetical population". At this point he advances three criteria to confirm the optimal character of an estimator, namely, "the Criterion of Consistency" (in the long run the estimator hits the unknown parameter<sup>46</sup>), "the Criterion of Efficiency" (in the regime of large samples, when the distributions of the estimators "tend to normality, that statistic is to be chosen which has the least probable error"<sup>47</sup>) and "the Criterion of Sufficiency", ensuring "that the chosen statistic should summarize the whole of the relevant information supplied by the sample". Finally, in order to make sure that "the theoretical aspect of the treatment of any particular body of data has been completely elucidated", he proposes the "Problems of Distributions", which amounts to explicitly computing the exact form of the law of the relevant estimators. As Fisher emphasizes, estimators meeting both the consistency and efficiency criteria may differ in their performance on finite samples, which

<sup>43</sup>Perhaps due to its indisputable appeal to spacial intuition, Fisher's geometric approach to the calculation of small sample distributions is hardly reproduced in modern textbooks, where it is usually replaced by its analytical counterpart involving the manipulation of the Jacobians associated to the underlying coordinate transformations.

<sup>44</sup>This ubiquitous dichotomy between the small and large sample regimes also transpires in Section 11, where hypothesis tests have been treated.

<sup>45</sup>According to Fisher, "these are entirely a matter for the practical statistician".

<sup>46</sup>But notice that Fisher consistency is essentially distinct from the modern one presented here.

<sup>47</sup>In modern terms, the dispersion, as measured by the variance, attains the Cramér-Rao lower bound in (8.24), where Fisher information, also discussed in [Fis22], plays a key role; compare with Remark 8.24.

justifies the consideration of the sufficiency criterion as a sort of tiebreaker. Also, whereas the asymptotic questions underlying the “Problems of Estimations” are quite manageable to a analytical treatment, the “Problems of Distribution” definitely involve “small samples” and hence, as Fisher himself acknowledges, are of “great mathematical difficulty”<sup>48</sup>. Finally, Fisher proposes, in this same paper, the MLE method as a systematic procedure of constructing statistics which solve the estimation problem in each given model, thus going a long way toward firmly establishing the conceptual framework upon which the “frequentist” approach to Statistical Estimation ultimately rests<sup>49</sup>.

### 13. THE BAYESIAN PATHWAY

As a way of comparison with the frequentist approach developed above, we now briefly comment on the *Bayesian approach* to estimation, where it is assumed that the parameter  $\theta$  in the i.i.d. measurements  $X_j \sim \psi(\cdot; \theta)$  is random with a pdf  $\psi_\vartheta$ , so that probabilities are assigned to parameters as well as to observations (here, we pretend that the random variable  $\vartheta$  on  $\Theta$  is the “capital” version of  $\theta$ ). It follows from Theorem 3.7 (Bayes rule) that

$$\psi_{\vartheta|X=\mathbf{x}}(\theta) = \frac{\psi_{X|\vartheta=\theta}(\mathbf{x})\psi_\vartheta(\theta)}{\psi_X(\mathbf{x})}, \quad \psi_X(\mathbf{x}) = \int_{\Theta} \psi_{X|\vartheta=\theta}(\mathbf{x})\psi_\vartheta(\theta)d\theta,$$

where  $X = (X_1, \dots, X_n)$  and  $\mathbf{x} = (x_1, \dots, x_n)$  is a realization of  $X$ . In the Bayesian jargon,  $\psi_\vartheta(\theta)$  is the *prior*, reflecting our knowledge of the underlying parameter  $\theta$  previous to any measurement (and hence viewed as a *hypothesis*) and  $\psi_{X|\vartheta=\theta}(\mathbf{x})$  is the *likelihood*, which indicates the compatibility of the *evidence*  $X$  with the given hypothesis. Note that

$$\psi_{X|\vartheta=\theta}(\mathbf{x}) = L(\mathbf{x}; \theta),$$

the likelihood function in (8.2), hence the terminology; here we are momentarily coming back to the “frequentist” setting of Section 12 and thus regarding  $\theta$  as deterministic (i.e. non-random). The prior and the likelihood combine to yield the *posterior*  $\psi_{\vartheta|X=\mathbf{x}}(\theta)$  through the proportionality

$$(13.1) \quad \psi_{\vartheta|X=\mathbf{x}}(\theta) \propto L(\mathbf{x}; \theta)\psi_\vartheta(\theta),$$

which provides an update of the probability distribution of the hypothesis as more observed evidence becomes available.

**Example 13.1.** For a normal sample  $X_j \sim \mathcal{N}(\mu, \sigma^2)$ , with  $\sigma$  known, we find that the likelihood is

$$L(\mathbf{x}; \mu) = \frac{1}{(2\pi)^{n/2}\sigma^n} e^{-\frac{1}{2\sigma^2} \sum_j (x_j - \mu)^2}.$$

Now assume that the prior, which expresses our initial degree of belief on the unknown parameter  $\mu$ , follows the normal  $\mathcal{N}(\mu_{\text{pr}}, \sigma_{\text{pr}}^2)$ , so that

$$\psi_\vartheta(\mu) = \frac{1}{\sqrt{2\pi}\sigma_{\text{pr}}} e^{-\frac{(\mu - \mu_{\text{pr}})^2}{2\sigma_{\text{pr}}^2}}.$$

A direct computation using (13.1) confirms that the posterior also follows a normal, namely,

$$\psi_{\vartheta|X=\mathbf{x}} \sim \mathcal{N}(\mu_{\text{pos}}, \sigma_{\text{pos}}^2),$$

where

$$(13.2) \quad \mu_{\text{pos}} = (1 - \lambda)\mu_{\text{pr}} + \lambda \frac{\sum_j x_j}{n}, \quad \lambda = \frac{\sigma_{\text{pr}}^2}{\sigma_{\text{pr}}^2 + \sigma^2/n},$$

<sup>48</sup>This partially explains his lifelong appreciation of Student’s fundamental contribution described in Remark 7.29. Also, as explained in Remark 11.3, this “Problem of Distribution” also emerges in the hypothesis testing context.

<sup>49</sup>As mentioned in [Sti05], “The paper is an astonishing work: It announces and sketches out a new science of statistics, with new definitions, a new conceptual framework and enough hard mathematical analysis to confirm the potential and richness of this new structure.”

and

$$\sigma_{\text{pos}}^2 = \frac{\sigma_{\text{pr}}^2 \sigma^2 / n}{\sigma_{\text{pr}}^2 + \sigma^2 / n}.$$

Hence, the Bayesian recipe confines the posterior mean  $\mu_{\text{pos}}$  somewhere between the prior mean  $\mu_{\text{pr}}$  and the realization  $\sum_j x_j / n$  of the sample mean, with a higher degree of belief than before (since  $\sigma_{\text{pos}} < \min\{\sigma_{\text{pr}}, \sigma\}$ ). We thus see that the data-gathering provided by the sample has the net effect of fine-tuning our initial subjective knowledge regarding  $\mu$ .  $\square$

**Example 13.2.** (Laplace's rule of succession) What chances are that the sun will rise tomorrow given that it has been so for the last  $n$  days? To ponder on this, consider a Bernoulli sample  $X_j \sim \text{Ber}(p)$  assigning probability  $p$  to a successful outcome corresponding to the event  $\{1\}$ . The question above is a special case (with  $s = n$ ) of the general problem of computing

$$P(X_{n+1} = 1 | X^{(n)} = s) = P_{X_{n+1} | X^{(n)} = s}(\{1\}), \quad X^{(n)} = X_1 + \cdots + X_n,$$

the probability that success occurs at the  $(n+1)^{\text{th}}$  outcome given that it has occurred  $s$  times previously; here we use the notation of (3.2). From Example 8.7 we know that the likelihood is

$$L(\mathbf{x}; p) = p^s (1-p)^{n-s},$$

where  $s = x_1 + \cdots + x_n$  is the realization of  $X^{(n)}$ . The simplest choice for the prior distribution of  $\mathbf{p}$ , the random variable associated to the Bayesian parameter  $p$ , appeals to the "Principle of Insufficient Reason": we declare that

$$\psi_{\mathbf{p}}(p) = \mathbf{1}_{[0,1]}(p),$$

the *uniform distribution* supported on the unit interval  $[0, 1]$ . Using (13.1) we see that the posterior is

$$\begin{aligned} \psi_{\mathbf{p} | X^{(n)} = s}(p) &= \frac{p^s (1-p)^{n-s} \mathbf{1}_{[0,1]}(p)}{\int_0^1 p^s (1-p)^s dp} \\ &= \frac{(n+1)!}{s!(n-s)!} p^s (1-p)^{n-s} \mathbf{1}_{[0,1]}(p), \end{aligned}$$

the Beta distribution  $\text{Beta}(s+1, n-s+1)$ ; cf. Definition 4.35. It follows that

$$\begin{aligned} P_{X_{n+1} | X^{(n)} = s}(\{1\}) &= \mathbb{E}(\mathbf{p} | X^{(n)} = s) \\ &= \int_0^1 p \psi_{\mathbf{p} | X^{(n)} = s}(p) dp, \end{aligned}$$

and using the previous expression for the posterior we get

$$P(X_{n+1} = 1 | X^{(n)} = s) = \frac{s+1}{n+2},$$

which in particular gives  $(n+1)/(n+2)$  as the solution for Laplace's sunrise problem<sup>50</sup>.  $\square$

**Example 13.3.** If  $\Theta \subset \mathbb{R}^q$  has a finite volume then the "Principle of Insufficient Reason" leads to

$$\psi_{\vartheta}(\theta) = \frac{1}{\text{vol}(\Theta)} \mathbf{1}_{\Theta}(\theta), \quad \theta \in \Theta,$$

as the choice for the prior, so the corresponding posterior is

$$\psi_{\vartheta | X = \mathbf{x}}(\theta) = \frac{L(\mathbf{x}; \theta) \mathbf{1}_{\Theta}(\theta)}{\int_{\Theta} L(\mathbf{x}; \theta) d\theta},$$

<sup>50</sup>In the end of [Lap98, Chapter III], Laplace takes  $n$  corresponding to five thousand years and finds that "it is a bet of 1820214 to one that it will rise again tomorrow". But as Laplace himself recognizes in the sequel, this should be taken with a salt of grain, specially in regard to the choice of prior, as possibilities other than the uniform are certainly available; see Example 13.7.

a suitable normalization of the likelihood. This confirms that, viewed as a function of  $\theta$ , the likelihood in general does not qualify as a pdf, which is consistent with the fact that the prescription for the MLE estimator in Definition 8.3 is insensitive to replacing  $L$  by  $cL$ ,  $c > 0$  a constant. In other words, any multiple of the likelihood function carries the same information as far as selecting the ML estimator is concerned.  $\square$

The examples above illustrate the Bayesian credo according to which probability is nothing but a measure of our degree of belief on the underlying parameter  $\theta$ , which thus should be random in nature. In any case, we may proceed with the corresponding estimation theory as follows. Given  $\tilde{\theta}$  define its *Bayes risk* as

$$\mathcal{R}(\tilde{\theta}) = \mathbb{E}_{\psi_{\tilde{\theta}}}(\mathcal{L}(\tilde{\theta}, \theta)),$$

where  $\mathcal{L}$  is a (previously chosen) *loss function* (for instance, the quadratic loss  $\mathcal{L}(\tilde{\theta}, \theta) = |\tilde{\theta} - \theta|^2$  gives rise to the Bayesian analogue of the mse in (7.3), but be aware of a crucial difference: here we average against the prior  $\psi_{\tilde{\theta}}(\theta)d\theta$  in alignment with the Bayesian philosophy according to which  $\theta$  is random, whereas there we integrate against  $dP_{\theta}$  since  $\theta$  is regarded as deterministic (i.e. non-random); see Remark 7.3.

**Definition 13.4.** A *Bayes estimator* is any  $\hat{\theta}$  that minimizes the Bayes risk.

Notice that this only depends on the prior distribution (and the given loss function) and hence involves no observation. In any case, given this setup we are now in a position to implement the Bayesian updating paradigm relying on the subsequent measurement  $X = x$  via (13.1). This leads to the following result, which provides a method for constructing Bayes estimators by solving a minimization problem formulated in terms of the posterior distribution  $\psi_{\tilde{\theta}|X=x}$ .

**Theorem 13.5.** Assume that for almost all  $\mathbf{x}$  there exists  $\hat{\theta}(\mathbf{x})$  minimizing

$$\tilde{\theta} \mapsto \mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(\mathcal{L}(\tilde{\theta}(\mathbf{x}), \theta)),$$

where  $\tilde{\theta}$  runs over the set of estimators with a finite risk. Then  $\hat{\theta} = \hat{\theta}(X)$  is a Bayes estimator.

*Proof.* By assumption we have, for almost all  $\mathbf{x}$  and any  $\tilde{\theta}$ ,

$$\mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(\mathcal{L}(\tilde{\theta}(\mathbf{x}), \theta)) \geq \mathbb{E}_{\psi_{\hat{\theta}|X=\mathbf{x}}}(\mathcal{L}(\hat{\theta}(\mathbf{x}), \theta)).$$

By Proposition 3.14, this may be expressed in terms of conditional expectations as

$$\mathbb{E}_{\psi_{\tilde{\theta}}}(\mathcal{L}(\tilde{\theta}(X), \vartheta)|\mathcal{F}_X) \geq \mathbb{E}_{\psi_{\hat{\theta}}}(\mathcal{L}(\hat{\theta}(X), \vartheta)|\mathcal{F}_X).$$

By applying  $\mathbb{E}_{\psi_{\tilde{\theta}}}$  to both sides and using Proposition 3.11 (2) we conclude that  $\mathcal{R}(\tilde{\theta}(X)) \geq \mathcal{R}(\hat{\theta}(X))$ .  $\square$

**Corollary 13.6.** The Bayes estimator associated to the weighted quadratic loss  $\mathcal{L}(\tilde{\theta}, \theta) := w(\theta)|\tilde{\theta} - g(\theta)|^2$ ,  $w > 0$ , is given by

$$\hat{\theta}(\mathbf{x}) = \frac{\mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(wg)}{\mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(w)}.$$

*Proof.* The Cauchy-Schwartz inequality

$$\mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(w(\theta)g(\theta))^2 < \mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(w(\theta))\mathbb{E}_{\psi_{\tilde{\theta}|X=\mathbf{x}}}(w(\theta)g(\theta)^2)$$

implies that

$$\begin{aligned}\mathbb{E}_{\psi_{\theta|X=\mathbf{x}}}(\mathcal{L}(\tilde{\theta}(\mathbf{x}), \theta)) &= \mathbb{E}_{\psi_{\theta|X=\mathbf{x}}}(w(\theta))\tilde{\theta}(\mathbf{x})^2 \\ &\quad - 2\mathbb{E}_{\psi_{\theta|X=\mathbf{x}}}(w(\theta)g(\theta))\tilde{\theta}(\mathbf{x}) + \mathbb{E}_{\psi_{\theta|X=\mathbf{x}}}(w(\theta)g(\theta)^2),\end{aligned}$$

viewed as a quadratic expression in  $\tilde{\theta}(\mathbf{x})$ , has a negative discriminant and hence is minimized at  $\tilde{\theta}(\mathbf{x}) = \hat{\theta}(\mathbf{x})$ .  $\square$

We refer to [Rob07, Chapter 2] for an extensive discussion on this “decision-theoretic” approach to Bayes estimation.

**Example 13.7.** A much more flexible choice for the prior in the sunrise problem of Example 13.2 is

$$\psi_{\mathbf{p}}(p) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1} \mathbf{1}_{[0,1]}, \quad \alpha, \beta > 0,$$

the Beta( $\alpha, \beta$ ) distribution (the case  $\alpha = \beta = 1$  corresponds to the uniform distribution). We thus calculate that the posterior is

$$\psi_{\mathbf{p}|X^{(n)=s}} \sim \text{Beta}(\alpha + s, \beta + n - s),$$

so in this case we obtain

$$(13.3) \quad P(X_{n+1} = 1 | X^{(n)=s}) = \frac{s + \alpha}{n + \alpha + \beta}.$$

This illustrates how sensitive the Bayesian machinery is to the choice of the prior. Moreover, taking into account that the right-hand side of (13.3) equals the expected value of the posterior, if we choose the loss function to be  $\mathcal{L}(\tilde{\theta}, \theta) = |\tilde{\theta} - \theta|^2$  in Corollary 13.6 we find that the corresponding Bayes estimator is

$$\hat{p}_{(n)}(X) = \frac{X^{(n)} + \alpha}{n + \alpha + \beta} = (1 - \gamma) \frac{\alpha}{\alpha + \beta} + \gamma \bar{X}_n, \quad \gamma = \frac{n}{n + \alpha + \beta},$$

which interpolates between  $\alpha/(\alpha + \beta)$ , the expected value of  $\psi_{\mathbf{p}}$  (the natural estimator prior to any observation) and the sample mean  $\bar{X}_n$  (the “frequentist” estimator that completely ignores the Bayesian paradigm incarnated in the prior). In particular, for large samples the prior mean plays a negligible role as  $\hat{p}_{(n)}(X)$  becomes indistinguishable from the ML estimator  $\bar{X}_n$ . Finally, if we apply this same recipe to the normal setting of Example 13.1, it follows from (13.2) that the corresponding Bayes estimator is

$$\hat{\mu}_{\text{pos},n}(X) = (1 - \lambda)\mu_{\text{pr}} + \lambda\bar{X}_n, \quad \lambda = \frac{\sigma_{\text{pr}}^2}{\sigma_{\text{pr}}^2 + \sigma^2/n}.$$

Again, this interpolates between the prior mean and the sample mean with

$$\hat{\mu}_{\text{pos},n}(X) \approx_{n \rightarrow +\infty} \bar{X}_n,$$

so the asymptotic behavior completely disregards the prior mean.  $\square$

**Remark 13.8.** (asymptotic efficiency of Bayes estimators) As illustrated in Example 13.7 above, Corollary 13.6 shows that the determination of a Bayes estimator for  $\theta$  under a quadratic loss boils down to computing the expectation of the posterior, a quite feasible task in some cases. Similarly to the course of action taken in the “frequentist” setting, with those estimators at hand we may then examine their asymptotic efficiency. In the cases treated above, this may be easily reduced to CLT. Indeed, in the Bernoulli case we compute that

$$\sqrt{n}(\hat{p}_{(n)}(X) - p) = \sqrt{n}(\bar{X}_n - p) + \frac{\sqrt{n}}{\alpha + \beta + n}(\alpha - (\alpha + \beta)\bar{X}_n),$$

where  $p$  is the true value of the unknown parameter, so that Theorem 2.23 and CLT apply to conclude that

$$\sqrt{n}(\hat{p}_{(n)}(X) - p) \xrightarrow{d} \mathcal{N}(0, p(1-p)).$$

Similarly, in the normal setting,

$$\sqrt{n}(\hat{\mu}_{\text{pos},n}(X) - \mu) = \sqrt{n}\lambda(\bar{X}_n - \mu) + \sqrt{n}(1 - \lambda)(\mu_{\text{pr}} - \mu),$$

and since  $\lambda \rightarrow 1$  and  $\sqrt{n}(1 - \lambda) = O(n^{-1/2}) \rightarrow 0$  we see that

$$\sqrt{n}(\hat{\mu}_{\text{pos},n}(X) - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$

Thus, in each case the limiting distribution of the appropriate standardization of the Bayes estimator is normal with a dispersion independent of the parameters of the prior distribution. This turns out to be a quite general phenomenon. Indeed, results in [LC06, Section 6.8] guarantee, under suitable regularity conditions and in the regime of large samples, that:

- the posterior distribution becomes asymptotically normal, and hence insensitive to the chosen prior, with a variance depending on the true value  $\theta_0$  of the unknown parameter only through its Fisher information:

$$(13.4) \quad \sqrt{n}(\psi_{\vartheta|X=x} - \theta_0) \xrightarrow{d} \mathcal{N}(0, \mathcal{F}(\theta_0)^{-1}).$$

- as a consequence, the limiting distribution associated to the Bayes estimator  $\hat{\theta}_n$  (under quadratic loss) is normal as well with the same asymptotic variance:

$$(13.5) \quad \sqrt{n}(\hat{\theta}_n(X) - \theta_0) \xrightarrow{d} \mathcal{N}(0, \mathcal{F}(\theta_0)^{-1}).$$

In particular,  $\hat{\theta}_n$  is asymptotically efficient.

We note that (13.5) follows from (13.4) and the fact that

$$\sqrt{n}(\hat{\theta}_n(X) - \psi_{\vartheta|X=x}) \xrightarrow{d} 0.$$

Although in the long run these results eventually succeed in altogether eliminating the effect of the subjective choice of the prior, they remain a bit extraneous in regard to the Bayesian tenet according to which by its very nature the posterior is conditional on the sample, whose size has been fixed once and for all.  $\square$

## APPENDIX A. BROWNIAN MOTION AND SOME OF ITS APPLICATIONS

In this rather long appendix, we turn our attention to Brownian motion, an important example of a stochastic process, and its most elementary properties. Although the main motivation here is to provide a proof of the Gaussian concentration inequality (5.18) with the sharp constant  $C = 1/2$ , which is presented in Section A.4, we also include a few other applications of the associated Itô's calculus, a cornerstone in the modern theory of stochastic processes.

**A.1. Brownian motion: its construction and basic regularity properties.** Since Brownian motion is the prototypical example of a stochastic process, we start by recalling the definition of this fundamental concept.

**Definition A.1.** A stochastic process on a probability space  $(\Omega, \mathcal{F}, P)$  is a one-parameter family of random variables  $X_t : \Omega \rightarrow \mathbb{R}^n, t \geq 0$ .

The map  $\omega \in \Omega \mapsto X_t(\omega) \in \mathbb{R}^n$  allows us to think of  $\Omega$  as a subset of  $(\mathbb{R}^n)^{[0,+\infty)}$ . Thus, to each  $\omega \in \Omega$  the process defines a path in  $\mathbb{R}^n$ . In general, the regularity of the process is expressed in terms of the regularity of these paths. For instance, we say that the process is continuous if  $X_t(\omega)$  is continuous for almost any  $\omega \in \Omega$ . Here we only deal with processes which are at least continuous. In any case, this pathwise description of stochastic processes motivates the following definition.

**Definition A.2.** Given a stochastic process  $X_t$ , its probability distributions in  $\mathbb{R}^{nk}$ ,  $k = 1, 2, \dots$ , are given by

$$\mu_{t_1, \dots, t_k}^X(F_1 \times \dots \times F_k) = P(X_{t_1} \in F_1, \dots, X_{t_k} \in F_k),$$

where  $t_i \geq 0$  and  $F_i \in \mathcal{B}^n$ ,  $i = 1, \dots, k$ .

In other words,  $\mu_{t_1, \dots, t_k}^X = P_{(X_{t_1}, \dots, X_{t_k})}$ , the joint distribution of  $(X_{t_1}, \dots, X_{t_k})$ ; cf. Definition 2.6. The next result shows that a stochastic process can be reconstructed from their probability distributions given that a couple of compatibility conditions are satisfied.

**Theorem A.3.** (Kolmogorov's extension) Assume that for any  $t_1, \dots, t_k \geq 0$  there exists a probability measure  $\nu_{t_1, \dots, t_k}$  in  $\mathbb{R}^{nk}$  such that:

- $(K_1)$   $\nu_{t_{\tau(1)}, \dots, t_{\tau(k)}}(F_1 \times \dots \times F_k) = \nu_{t_1, \dots, t_k}(F_{\tau^{-1}(1)}, \dots, F_{\tau^{-1}(k)})$ , for any permutation  $\tau$ .
- $(K_2)$   $\nu_{t_1, \dots, t_k}(F_1 \times \dots \times F_k) = \nu_{t_1, \dots, t_k, t_{k+1}, \dots, t_{k+m}}(F_1 \times \dots \times F_k \times \mathbb{R}^n \times \dots \times \mathbb{R}^n)$ , for any  $m \geq 1$ .

Then there exists a probability space  $(\Omega, \mathcal{F}, P)$  and a stochastic process  $X_t : \Omega \rightarrow \mathbb{R}^n$  such that

$$\nu_{t_1, \dots, t_k} = \mu_{t_1, \dots, t_k}^X,$$

for  $(t_1, \dots, t_k)$ .

*Proof.* See [Tao11, Section 2.4]. □

Now we can construct Brownian motion in  $\mathbb{R}^n$  following an approach due to Kolmogorov; for other possibilities see [SP14]. If  $0 \leq t_1 < \dots < t_k$  and  $y = (y_1, \dots, y_k) \in \mathbb{R}^{nk}$  define, if  $t_1 > 0$ ,

$$\nu_{t_1, \dots, t_k}(F_1 \times \dots \times F_k) = \frac{1}{(2\pi)^{nk/2} \sqrt{\det C}} \int_{F_1 \times \dots \times F_k} e^{-\frac{1}{2} \langle C^{-1} y, y \rangle} dy,$$

where each  $F_i \in \mathcal{B}^n$  and  $C$  is  $nk \times nk$ -matrix whose  $ij$ -block is  $C_{ij} = t_i \wedge t_j I_n$ <sup>51</sup>. If  $t_1 = 0$  we use instead  $\delta_{\vec{0}} \otimes \nu_{t_2, \dots, t_k}$ , where  $\delta_{\vec{0}}$  is the Dirac measure centered at the origin. This may be extended to all  $(t_1, \dots, t_k)$  so that  $(K_1)$  is satisfied. Moreover,  $(K_2)$  is satisfied as well because of Proposition 4.2. Thus, by means of Theorem A.3 we establish the following foundational existence result.

**Theorem A.4.** There exists a probability space  $(\Omega, \mathcal{F}, P)$  and a stochastic process  $b_t : \Omega \rightarrow \mathbb{R}^n$  so that

$$P(b_{t_1} \in F_1, \dots, b_{t_k} \in F_k) = \frac{1}{(2\pi)^{nk/2} \sqrt{\det C}} \int_{F_1 \times \dots \times F_k} e^{-\frac{1}{2} \langle C^{-1} x, x \rangle} dx.$$

This is called Brownian motion (BM) in  $\mathbb{R}^n$  (starting at  $\vec{0}$ ).

The next proposition lists the characterizing properties of BM.

**Proposition A.5.** BM in  $\mathbb{R}^n$  satisfies the following properties:

- (1)  $b_0 = 0$  a.s.;

<sup>51</sup>Here,  $a \wedge b = \min\{a, b\}$ . Also, note that the symmetric matrix  $C$  is positive definite.

(2) it has stationary normal increments, i.e. for any  $0 \leq s < t$ ,  $h \geq -s$ ,  $b_{t+h} - b_{s+h}$  and  $b_t - b_s$  are identically distributed with

$$(A.1) \quad b_t - b_s \sim \mathcal{N}(0, (t-s)\mathbf{I}_n);$$

(3) it has independent increments, that is, for any  $0 = t_0 < t_1 < \dots < t_k$ ,  $\{b_{t_1} - b_{t_0}, \dots, b_{t_k} - b_{t_{k-1}}\}$  are independent random vectors;

(4)  $t \mapsto b_t(\omega)$  is continuous for any  $\omega$ .

*Proof.* From the construction, (1) follows immediately. To approach (2) and (3) we take  $u = (u_1, \dots, u_k) \in \mathbb{R}^{nk}$ ,  $v_j = u_j + \dots + u_k$ ,  $j = 1, \dots, k$ , and  $b = (b_{t_1}, \dots, b_{t_k})$ , so if we use the language of characteristic functions in Definition 2.24 and the explicit computation of this object for normally distributed random vectors in Proposition 4.4 we have (recalling that  $b_{t_0} = 0$ )

$$\begin{aligned} \phi_{(b_{t_1}-b_{t_0}, \dots, b_{t_k}-b_{t_{k-1}})}(v_1, \dots, v_k) &= \mathbb{E}(e^{i \sum_{j=1}^k \langle b_{t_j} - b_{t_{j-1}}, v_j \rangle}) \\ &= \mathbb{E}(e^{i \langle b, u \rangle}) \\ &= \phi_b(u) \\ &= e^{-\frac{1}{2} \langle Cu, u \rangle}. \end{aligned}$$

But

$$\begin{aligned} \langle Cu, u \rangle &= \sum_{j=1}^k \sum_{l=1}^k (t_j \wedge t_l) \langle u_j, u_l \rangle \\ &= t_k \|u_k\|^2 + \sum_{j=1}^{k-1} t_j \langle u_j, u_j + 2u_{j+1} + \dots + 2u_k \rangle \\ &= t_k \|u_k\|^2 + \sum_{j=1}^{k-1} t_j (\|u_j + \dots + u_k\|^2 - \|u_{j+1} + \dots + u_k\|^2) \\ &= \sum_{j=1}^k t_j \|u_j + \dots + u_k\|^2 - \sum_{j=1}^k t_{j-1} \|u_j + \dots + u_k\|^2 \\ &= \sum_{j=1}^k (t_j - t_{j-1}) \|v_j\|^2, \end{aligned}$$

so that

$$(A.2) \quad \phi_{(b_{t_1}-b_{t_0}, \dots, b_{t_k}-b_{t_{k-1}})}(v_1, \dots, v_k) = \prod_{j=1}^k e^{-\frac{1}{2}(t_j - t_{j-1})\|v_j\|^2}.$$

Notice that for  $0 \leq s < t$  this specializes to

$$\phi_{b_t - b_s}(v) = e^{-\frac{1}{2}(t-s)\|v\|^2}, \quad v \in \mathbb{R}^n,$$

so that Corollary 4.6 applies to ensure that (A.1) holds, which proves (2). As for (3), note that (A.2) may be rewritten as

$$\phi_{(b_{t_1}-b_{t_0}, \dots, b_{t_k}-b_{t_{k-1}})}(v_1, \dots, v_k) = \prod_{j=1}^k \phi_{b_{t_j}-b_{t_{j-1}}}(v_j),$$

so we may proceed as in the last step of the proof of Proposition 4.10 and use the standard Fourier inversion formula to confirm that the joint distribution of the random vector of increments decomposes as

$$\psi_{(b_{t_1}-b_{t_0}, \dots, b_{t_k}-b_{t_{k-1}})}(x_1, \dots, x_k) = \prod_{j=1}^k \psi_{b_{t_j}-b_{t_{j-1}}}(x_j),$$



which proves (3) by Proposition 2.13. The proof of (4) is presented in the next section; see Proposition A.10.  $\square$

**Proposition A.6.** *If  $t \leq s$  then  $\mathbb{E}(\|b_s - b_t\|^2) = n(s - t)$ .*

*Proof.* We have seen that  $\mathbb{E}(b_t) = 0$  and  $\text{cov}(b_s, b_t) = s \wedge t I_n$ . Thus, if  $b_t = (b_t^{(1)}, \dots, b_t^{(n)})$  is the coordinate expression of  $b_t$  we have

$$\begin{aligned} \mathbb{E}(\langle b_s, b_t \rangle) &= \sum_i \mathbb{E}(b_s^{(i)} b_t^{(i)}) \\ &= \sum_i \text{cov}(b_s, b_t)_{ii} \\ &= ns \wedge t. \end{aligned}$$

If  $t \leq s$  we then have

$$\begin{aligned} \mathbb{E}(\|b_s - b_t\|^2) &= \mathbb{E}(\|b_s\|^2 - 2\langle b_s, b_t \rangle + \|b_t\|^2) \\ &= n(s - 2t + t), \end{aligned}$$

as desired.  $\square$

**Remark A.7.** It follows from (A.1) that  $\text{cov}(b_t)_{ij} = t\delta_{ij}$ , so that by Proposition 4.10 we see that the coordinate components  $\{b_t^{(i)}\}_{i=1}^n$  of  $b_t$  form an independent family of BMs in  $\mathbb{R}$ . Conversely, we may first construct BM  $b_t$  in  $\mathbb{R}$  by using the Kolmogorov's argument above (note that in this case the matrix  $C$  has a much simpler structure) and then take  $n$  independent copies of  $b_t$ , say  $\{b_t^{(1)}, \dots, b_t^{(n)}\}$ , in order to exhibit BM in  $\mathbb{R}^n$  as  $(b_t^{(1)}, \dots, b_t^{(n)})$ .  $\square$

We now turn to the basic regularity properties of Brownian motion and we start by proving Proposition A.5, (4). To simplify matters, we only consider the case  $n = 1$ ; cf. Remark A.7. The proof is based on the following general regularity result for stochastic processes. We recall that saying that  $X'_t$  is a *modification* of  $X_t$  means that  $P(X_t = X'_t) = 1$  for any  $t$ .

**Theorem A.8.** (Kolmogorov's continuity) *If  $X_t : \Omega \rightarrow \mathbb{R}$  is a stochastic process satisfying*

$$\mathbb{E}(|X_s - X_t|^\alpha) \leq C|s - t|^{\beta+1}, \quad s, t \geq 0,$$

*then there exists a modification  $X'_t$  of  $X_t$  whose paths are locally  $\gamma$ -Hölder continuous, where  $0 < \gamma < \beta/\alpha$ . In particular,  $X'_t \in C^0$ .*

*Proof.* See [LG13, Section 2.2].  $\square$

The regularity of BM now follows from the following fact.

**Proposition A.9.** *BM in  $\mathbb{R}$  satisfies*

$$\mathbb{E}(|b_s - b_t|^{2k}) = \frac{(2k)!}{2^k k!} |s - t|^k, \quad k \geq 1.$$

*Proof.* Since  $b_s - b_t \sim \mathcal{N}(0, s - t)$ , this follows from the discussion in Example 4.8.  $\square$

**Proposition A.10.** *Eventually passing to a modification, BM is locally  $(\frac{1}{2} - \epsilon)$ -Hölder continuous, for any  $\epsilon > 0$ .*

*Proof.* Apply the results above with  $\alpha = 2k$  and  $\beta = k - 1$  and send  $k \rightarrow +\infty$ . □

This is in a sense the best regularity we can have. To check this we need a definition.

**Definition A.11.** If  $X_t : \Omega \rightarrow \mathbb{R}$  and  $p > 0$ , we define its  $p^{\text{th}}$  variation by

$$\langle X \rangle_t^{(p)}(\omega) = \lim_{\Delta t_k \rightarrow 0} \sum_{t_k \leq t} |X_{t_{k+1}}(\omega) - X_{t_k}(\omega)|^p,$$

where  $\Delta t_k = t_{k+1} - t_k = t/k$  and the limit is taken in probability.

It turns out that the quadratic variation of BM can be explicitly computed.

**Proposition A.12.** *BM  $b_t$  in  $\mathbb{R}$  satisfies*

$$\langle b \rangle_t^{(2)} = t,$$

*with convergence in  $L^2$ -mean.*

*Proof.* Note that

$$\begin{aligned} \mathbb{E} \left( \left( \sum_{t_k \leq t} (b_{t_{k+1}} - b_{t_k})^2 - t \right)^2 \right) &= \mathbb{E} \left( \left( \sum_{t_k \leq t} (b_{t_{k+1}} - b_{t_k})^2 - (t_{k+1} - t_k) \right)^2 \right) \\ &= I + II, \end{aligned}$$

where

$$I = \mathbb{E} \left( \sum_{t_k \leq t} ((b_{t_{k+1}} - b_{t_k})^2 - (t_{k+1} - t_k))^2 \right)$$

and

$$II = 2 \sum_{t_j < t_k \leq t} \underbrace{\mathbb{E} \left( ((b_{t_{j+1}} - b_{t_j})^2 - (t_{j+1} - t_j)) ((b_{t_{k+1}} - b_{t_k})^2 - (t_{k+1} - t_k)) \right)}_{III}$$

By Proposition A.6 we know that

$$\mathbb{E} ((b_{t_{j+1}} - b_{t_j})^2) = t_{j+1} - t_j,$$

which implies that

$$III = \text{cov} ((b_{t_{j+1}} - b_{t_j})^2, (b_{t_{k+1}} - b_{t_k})^2).$$

But by Proposition A.5, (3),

$$b_{t_{j+1}} - b_{t_j} \perp b_{t_{k+1}} - b_{t_k} \Rightarrow (b_{t_{j+1}} - b_{t_j})^2 \perp (b_{t_{k+1}} - b_{t_k})^2,$$

and hence  $II = 0$  by Corollary 2.5. On the other hand,

$$\begin{aligned}
 I &= \sum_{t_k \leq t} (\mathbb{E}((b_{t_{k+1}} - b_{t_k})^4) - 2(t_{k+1} - t_k)\mathbb{E}((b_{t_{k+1}} - b_{t_k})^2) + (t_{k+1} - t_k)^2) \\
 &\stackrel{\text{Prop. A.9}}{=} \sum_{t_k \leq t} (3(t_{k+1} - t_k)^2 - 2(t_{k+1} - t_k)^2 + (t_{k+1} - t_k)^2) \\
 &= 2 \sum_{t_k \leq t} (t_{k+1} - t_k)^2 \\
 &\leq 2 \frac{t^2}{k} \rightarrow 0,
 \end{aligned}$$

as  $k \rightarrow +\infty$ . □

**Proposition A.13.** *One has  $\langle b \rangle_t^{(1)} = +\infty$  a.s. In other words, the total variation of  $b_t$  blows up in any interval.*

*Proof.* For  $\omega \in \Omega$  we have

$$\begin{aligned}
 \sum_{t_k \leq t} (b_{t_{k+1}}(\omega) - b_{t_k}(\omega))^2 &\leq \sum_{t_k \leq t} (b_{t_{k+1}}(\omega) - b_{t_k}(\omega)) \sup_{t_k} \sum_{t_k \leq t} |b_{t_{k+1}}(\omega) - b_{t_k}(\omega)| \\
 &\leq \langle b \rangle_t^{(1)}(\omega) \sup_{t_k} \sum_{t_k \leq t} |b_{t_{k+1}}(\omega) - b_{t_k}(\omega)|.
 \end{aligned}$$

From Proposition A.12, and possibly passing to a subsequence along the given partitions of  $[0, t]$ , we may assume that the left-hand side converges to  $t$  a.s. But the supremum goes to 0 as  $b_t(\omega)$  is uniformly continuous in  $[0, t]$ , which yields a contradiction if  $\langle b \rangle_t^{(1)}(\omega)$  is finite. □

A similar argument yields the following result.

**Proposition A.14.** *The paths of  $b_t$  are nowhere  $\gamma$ -Hölder continuous for  $\gamma > 1/2$ .*

*Proof.* Assume  $|b_{s'} - b_{t'}| \leq K|s' - t'|^\gamma$ ,  $0 \leq t' \leq s' \leq t$ . It follows that

$$\sum_{t_k \leq t} (b_{t_{k+1}}(\omega) - b_{t_k}(\omega))^2 \leq K^2 t \sup_k |t_{k+1} - t_k|^{2\gamma-1}.$$

As above, we may assume that the left-hand side converges to  $t$  a.s. But the supremum goes to 0 if  $\gamma > 1/2$ , which gives a contradiction. □

**Remark A.15.** (The Wiener space) From Proposition A.5 (4) we may conveniently identify the sample space  $\Omega$  underlying the construction of Brownian motion in Theorem A.4 to  $C_{\vec{0}}$ , the space of continuous functions  $\omega : [0, +\infty) \rightarrow \mathbb{R}^n$  with  $\omega(0) = \vec{0}$  by the rule  $\omega(t) = b_t(\omega)$ . It then follows from Proposition A.14 that the support of the underlying probability measure  $P$  fails to contain any  $\omega$  which is sufficiently regular (i.e.  $\gamma$ -Hölder continuous for  $\gamma > 1/2$ ). From this perspective, we call  $C_{\vec{0}}$  endowed with the induced probability measure, still denoted by  $P$ , as the *Wiener space* (starting at  $\vec{0}$ ) and any  $\omega$  lying in the support of  $P$  (the *Wiener measure*) as a *Brownian path* (again, starting at  $\vec{0}$ ).

**A.2. Martingales.** We now isolate another central notion in the theory.

**Definition A.16.** Let  $b_t : \Omega \rightarrow \mathbb{R}^n$  be BM in  $\mathbb{R}^n$  (starting at  $x$ ) and let  $\mathcal{F}_t$  be the  $\sigma$ -algebra generated by  $\{b_{t'}\}_{t' \leq t}$ . A *martingale* (rel. to  $b_t$ ) is a stochastic process  $M_t : \Omega \rightarrow \mathbb{R}^n$  such that

- $M_t$  is  $\mathcal{F}_t$ -measurable for any  $t > 0$ ;
- $\mathbb{E}(\|M_t\|) < +\infty$ ;
- $\mathbb{E}(M_s | \mathcal{F}_t) = M_t$  whenever  $t \leq s$ .

**Proposition A.17.**  $b_t$  is a martingale.

*Proof.* If  $t \leq s$  write

$$\mathbb{E}(b_s | \mathcal{F}_t) = \mathbb{E}(b_s - b_t | \mathcal{F}_t) + \mathbb{E}(b_t | \mathcal{F}_t).$$

Proposition A.5, (3), implies that  $b_s - b_t \perp \mathcal{F}_t$ . Hence, by Proposition 3.11, (4),  $\mathbb{E}(b_s - b_t | \mathcal{F}_t) = \mathbb{E}(b_s - b_t) = 0$ . On the other hand, since  $b_t$  is (obviously)  $\mathcal{F}_t$ -measurable, Proposition 3.11, (3), implies that  $\mathbb{E}(b_t | \mathcal{F}_t) = b_t$ .  $\square$

For our purposes, a basic property of a martingale is that its expectation is preserved in time. This confirms that martingales are “pure fluctuation” processes.

**Proposition A.18.** If  $M_t$  is a martingale then  $\mathbb{E}(M_t) = \mathbb{E}(M_s)$ , for any  $s, t$ .

*Proof.* By Proposition 3.11, (2), if  $t \leq s$  we have

$$\mathbb{E}(M_t) = \mathbb{E}(\mathbb{E}(M_s | \mathcal{F}_t)) = \mathbb{E}(M_s).$$

$\square$

**A.3. Itô’s integral and Itô’s formula.** Consider a partition  $0 = t_0 < t_1 \cdots < t_k = t$ . Let  $b_t$  be BM in  $\mathbb{R}$  with  $b_0 = 0$ . By Proposition A.5, (3),  $b_{t_{j+1}} - b_{t_j} \perp b_{t_j}$  and hence

$$(A.3) \quad \mathbb{E} \left( \sum_j b_{t_j} (b_{t_{j+1}} - b_{t_j}) \right) = 0.$$

On the other hand,

$$\begin{aligned} \mathbb{E} \left( \sum_j b_{t_{j+1}} (b_{t_{j+1}} - b_{t_j}) \right) &\stackrel{(A.3)}{=} \mathbb{E} \left( \sum_k (b_{t_{j+1}} - b_{t_j})^2 \right) \\ &= \sum_j (t_{j+1} - t_j) \\ &= t, \end{aligned}$$

where we used Proposition A.6 in the second step. This simple computation, which reflects the already known fact that  $db_t$  can not be interpreted as a classical Lebesgue-Stieltjes integrator since  $b_t$  has infinite total variation by Proposition A.13, illustrates the difficulty of making sense of stochastic integrals like  $\int_0^t b_s db_s$  by standard methods. Put in another way, each choice of  $\{\hat{t}_j\}$  such that  $t_j \leq \hat{t}_j \leq t_{j+1}$  yields its own output for the “approximate” stochastic integral  $\sum_j b_{\hat{t}_j} (b_{t_{j+1}} - b_{t_j})$ . Among the many possibilities available, Itô’s integration corresponds to choosing the first option (A.3) above. The basics of this kind of stochastic integration may be found in many sources [KS12, LG13, Bau14] and our presentation below follows [Oks13] closely, a text very

much oriented to applications (in particular, to Mathematical Finance) to which we refer for the detailed proofs of most of the results on Itô's calculus described in the sequel.

Recall that a filtration  $\mathcal{F}_t$  of a  $\sigma$ -algebra is a nested family of  $\sigma$ -subalgebras of  $\mathcal{F}$ . Here we consider the filtration  $\mathcal{F}_t = \mathcal{F}_{\{b_{t'}\}_{t' \leq t}}$ .

**Definition A.19.** We say that a process  $f : \mathbb{R}_+ \times \Omega \rightarrow \mathbb{R}$  is *adapted* if  $\omega \mapsto f(t, \omega)$  is  $\mathcal{F}_t$ -measurable for any  $t$ .

**Definition A.20.** For  $0 \leq S < T$  we denote by  $\mathcal{V}(S, T)$  the class of all processes  $f : \mathbb{R}_+ \times \Omega \rightarrow \mathbb{R}$  such that:

- (1)  $f$  is  $\mathcal{B} \times \mathcal{F}$ -measurable;
- (2)  $f$  is adapted to  $\mathcal{F}_t$ , the filtration defined by  $b_t$ ;
- (3)  $\mathbb{E}(\int_S^T f(t, \omega)^2 dt) < +\infty$ .

**Definition A.21.** We say that  $f \in \mathcal{V}(S, T)$  is *elementary* if

$$f(t, \omega) = \sum_j f_j(\omega) \mathbf{1}_{[t_j, t_{j+1})}(t).$$

Note that if  $f$  is elementary then  $f_j$  is  $\mathcal{F}_{t_j}$ -measurable.

**Definition A.22.** (Itô's integral for elementary processes) If  $f \in \mathcal{V}(S, T)$  is elementary we define

$$\int_S^T f(t, \omega) db_t(\omega) = \sum_j f_j(\omega) (b_{t_{j+1}}(\omega) - b_{t_j}(\omega)).$$

Notice that this depends measurably on  $\omega$  and hence defines a random variable.

**Proposition A.23.** (Itô's isometry) If  $f \in \mathcal{V}(S, T)$  is elementary then

$$\mathbb{E} \left( \left( \int_S^T f(t, \omega) db_t(\omega) \right)^2 \right) = \mathbb{E} \left( \int_S^T f(t, \omega)^2 dt \right).$$

*Proof.* Let  $f = \sum_j f_j \mathbf{1}_{[t_j, t_{j+1})}$ . Since  $f_j$  is  $\mathcal{F}_{t_j}$ -measurable, Proposition A.5, (3), implies that  $f_j \perp b_{t_{j+1}} - b_{t_j}$ . Hence,  $f_j^2 \perp (b_{t_{j+1}} - b_{t_j})^2$  and we have

$$\mathbb{E}(f_j^2 (b_{t_{j+1}} - b_{t_j})^2) = \mathbb{E}(f_j^2) \mathbb{E}((b_{t_{j+1}} - b_{t_j})^2) = \mathbb{E}(f_j^2) (t_{j+1} - t_j),$$

where we used Proposition A.6 in the last step. On the other hand, if  $j < k$  we have  $f_j f_k (b_{t_{j+1}} - b_{t_j}) \perp b_{t_{k+1}} - b_{t_k}$  and hence,

$$\mathbb{E}(f_j f_k (b_{t_{j+1}} - b_{t_j})(b_{t_{k+1}} - b_{t_k})) = \mathbb{E}(f_j f_k (b_{t_{j+1}} - b_{t_j})) \mathbb{E}(b_{t_{k+1}} - b_{t_k}) = 0$$

It follows that

$$\begin{aligned} \mathbb{E} \left( \left( \int_S^T f(t, \omega) db_t(\omega) \right)^2 \right) &= \sum_{jk} \mathbb{E} (f_j f_k (b_{t_{j+1}} - b_{t_j})(b_{t_{k+1}} - b_{t_k})) \\ &= \sum_j \mathbb{E} (f_j^2) (t_{j+1} - t_j) \\ &= \mathbb{E} \left( \int_S^T f(t, \omega)^2 dt \right), \end{aligned}$$

as desired.  $\square$

**Proposition A.24.** (approximation) For any  $f \in \mathcal{V}(S, T)$  there exists  $\{f_i\}_{i=1}^{+\infty} \subset \mathcal{V}(S, T)$ ,  $f_i$  elementary, so that

$$(A.4) \quad \lim_{i \rightarrow +\infty} \mathbb{E} \left( \int_S^T |f - f_i|^2 dt \right) = 0.$$

*Proof.* [Oks13, pg. 27-28].  $\square$

**Definition A.25.** (Itô's integral in  $\mathcal{V}(S, T)$ ) If  $f \in \mathcal{V}(S, T)$  we define

$$\int_S^T f(t, \omega) db_t(\omega) \stackrel{L^2}{=} \lim_{i \rightarrow +\infty} \int_S^T f_i(t, \omega) db_t(\omega),$$

for some  $\{f_i\}$  as in (A.4).

Notice that, by Proposition A.23, the limit exists and does not depend on the sequence  $\{f_i\}$  chosen to approximate  $f$ .

We now list the basic properties of Itô's integral.

**Proposition A.26.** The Itô's integral satisfies the following properties:

- (1)  $\int_S^T f db_t = \int_S^U f db_t + \int_U^T f db_t$ ;
- (2)  $\int_S^T (af + bg) db_t = a \int_S^T f db_t + b \int_S^T g db_t$ ,  $a, b \in \mathbb{R}$ ;
- (3)  $\mathbb{E}(\int_S^T f db_t) = 0$ ;
- (4)  $\int_S^T f db_t$  is  $\mathcal{F}_T$ -measurable;
- (5) (Itô's isometry) There holds

$$\mathbb{E} \left( \left( \int_S^T f(t, \omega) db_t(\omega) \right)^2 \right) = \mathbb{E} \left( \int_S^T f(t, \omega)^2 dt \right).$$

- (6) If

$$\lim_{n \rightarrow +\infty} \mathbb{E} \left( \int_S^T (f_n(t, \omega) - f(t, \omega))^2 dt \right) = 0$$

then

$$\int_S^T f_n db_t \xrightarrow{L^2} \int_S^T f db_t.$$

(7) Any Itô's integral has a continuous modification.

*Proof.* The proofs of (1)-(5) follow the same method, namely, we first check the property for elementary processes and then pass the limit. Also, (6) follows immediately from (5). Finally, the proof of (7) can be found in [Oks13, Theorem 3.2.5].  $\square$

**Example A.27.** Let

$$f_n(s, \omega) = \sum_j b_{t_j}(\omega) \mathbf{1}_{[t_j, t_{j+1})}(s),$$

where  $\Delta t_j = t_{j+1} - t_j = t/n$ . We have

$$\begin{aligned} \mathbb{E} \left( \int_0^t (f_n - b_s)^2 ds \right) &= \mathbb{E} \left( \sum_j \int_{t_j}^{t_{j+1}} (f_n - b_s)^2 ds \right) \\ &= \mathbb{E} \left( \sum_j \int_{t_j}^{t_{j+1}} (b_{t_j} - b_s)^2 ds \right) \\ &= \sum_j \int_{t_j}^{t_{j+1}} (s - t_j) ds \\ &= \sum_j \frac{1}{2} (t_{j+1} - t_j)^2 \xrightarrow{n \rightarrow +\infty} 0. \end{aligned}$$

Thus, by Proposition A.26, (6),

$$\int_0^t b_s db_s = \lim_{n \rightarrow +\infty} \int_0^t f_n db_s = \lim_{\Delta t_j \rightarrow 0} \sum_j b_{t_j} (b_{t_{j+1}} - b_{t_j}).$$

But, since  $b_0 = 0$ ,

$$\begin{aligned} b_t^2 &= \sum_j (b_{t_{j+1}}^2 - b_{t_j}^2) \\ &= \sum_j (b_{t_{j+1}} - b_{t_j})^2 + 2 \sum_j b_{t_j} (b_{t_{j+1}} - b_{t_j}), \end{aligned}$$

that is,

$$\sum_j b_{t_j} (b_{t_{j+1}} - b_{t_j}) = \frac{1}{2} b_t^2 - \frac{1}{2} \sum_j (b_{t_{j+1}} - b_{t_j})^2.$$

By passing the limit and using Proposition A.12 to handle the last term in the right-hand side we conclude that

$$\int_0^t b_s db_s = \frac{1}{2} b_t^2 - \frac{t}{2}.$$

At least formally, we can rewrite this as

$$db_t^2 = 2b_t db_t + dt.$$

Setting  $f(x) = x^2$  we have

$$df(b_t) = f'(b_t) db_t + \frac{1}{2} f''(b_t) dt.$$

This rather special case of the famous Itô's formula illustrates the appearance of an extra term in the chain rule in the stochastic chain rule. In fact, if we interpret Proposition A.12 as saying that  $db_t^2 = dt$ , we have

$$df(b_t) = f'(b_t)db_t + \frac{1}{2}f''(b_t)db_t^2,$$

which means that we must expand up to second order in  $db_t$  to obtain the correct version of the chain rule.  $\square$

We now prove that Itô's integrals are martingales.

**Proposition A.28.** *If  $f \in \mathcal{V}(0, t)$  consider the process*

$$M_t = \int_0^t f(\rho, \omega) db_\rho(\omega).$$

*Then  $M_t$  is martingale.*

*Proof.* If  $t \leq s$  we have

$$\mathbb{E}(M_s | \mathbb{F}_t) = \mathbb{E}(M_t | \mathbb{F}_t) + \mathbb{E}\left(\int_t^s f(\rho, \omega) db_\rho(\omega)\right).$$

Since  $M_t$  is  $\mathcal{F}_t$ -measurable (Proposition A.26, (4)), we have  $\mathbb{E}(M_t | \mathbb{F}_t) = M_t$ . Moreover,  $\int_t^s f(\rho, \omega) db_\rho(\omega)$  is 'independent' of  $\mathcal{F}_t$  in the sense that

$$(A.5) \quad \mathbb{E}\left(\int_t^s f(\rho, \omega) db_\rho(\omega)\right) = 0,$$

which completes the proof except for the checking of (A.5), which needs to be carried out only for  $f$  of the type  $f = \sum_j f_j \mathbf{1}_{[t_j, t_{j+1})}$ . In this case,

$$\begin{aligned} \mathbb{E}\left(\int_t^s f(\rho, \omega) db_\rho(\omega)\right) &= \mathbb{E}\left(\sum_j f_j(b_{t_{j+1}} - b_{t_j}) | \mathcal{F}_t\right) \\ &\stackrel{\mathcal{F}_t \subset \mathcal{F}_{t_j} + \text{Prop. 3.11, (5)}}{=} \sum_j \mathbb{E}\left(\mathbb{E}(f_j(b_{t_{j+1}} - b_{t_j}) | \mathcal{F}_{t_j}) | \mathcal{F}_t\right) \\ &\stackrel{\text{Prop. 3.11, (2)}}{=} \sum_j \mathbb{E}(f_j \mathbb{E}((b_{t_{j+1}} - b_{t_j}) | \mathcal{F}_{t_j}) | \mathcal{F}_t), \end{aligned}$$

and this vanishes because  $\mathbb{E}((b_{t_{j+1}} - b_{t_j}) | \mathcal{F}_{t_j}) = 0$ .  $\square$

We now discuss a multi-dimensional version of Itô's integral which will suffice for our applications. We first recall from Remark A.7 that if  $b_t = (b_t^{(1)}, \dots, b_t^{(n)})$  is Brownian motion in  $\mathbb{R}^n$  then  $\{b_t^{(i)}\}_{i=1}^n$  is an independent family of BMs on  $\mathbb{R}$  (and conversely). We will use this to define the integral

$$(A.6) \quad \int_S^T v db_t = \int_S^T \begin{pmatrix} v_{11} & \cdots & v_{1n} \\ \vdots & \ddots & \vdots \\ v_{m1} & \cdots & v_{mn} \end{pmatrix} \begin{pmatrix} db_t^{(1)} \\ \vdots \\ db_t^{(n)} \end{pmatrix}$$

as a process in  $\mathbb{R}^n$  for a suitable  $v_{ij} = v_{ij}(t, \omega)$ .



**Definition A.29.** Let  $\mathcal{V}_{\mathcal{H}}(S, T)$  be the collection of functions  $f : [0, +\infty) \times \Omega \rightarrow \mathbb{R}$  such that

- (1)  $f$  is  $\mathcal{B} \times \mathcal{F}$ -measurable;
- (2) there exists a filtration  $\mathcal{H}_t \subset \mathcal{F}$  such that:
  - $b_t$  is a martingale with respect to  $\mathcal{H}_t$  ( $b_t$  is BM in  $\mathbb{R}$ );
  - $f(t, \cdot)$  is adapted to  $\mathcal{H}_t$ ,  $t > 0$ .
- (3)  $\mathbb{E}(\int_S^T f(t, \omega)^2 dt) < +\infty$ .

Since  $\mathcal{F}_t \subset \mathcal{H}_t$  and  $\mathbb{E}(b_s - b_t | \mathcal{H}_t) = 0$  we can proceed as before and define

$$\int_S^T f db_t, \quad f \in \mathcal{V}_{\mathcal{H}}(S, T),$$

which is a martingale (see Proposition A.28). Coming back to  $b = (b^{(1)}, \dots, b^{(n)}) \in \mathbb{R}^n$ , let  $\mathcal{F}_t^{(n)}$  be the  $\sigma$ -algebra generated by  $b_{s_1}^{(1)}, \dots, b_{s_n}^{(n)}$ , where  $s_k \leq t$ ,  $k = 1, \dots, n$ . Using the componentwise independence mentioned above, we see that  $t < s$  implies that  $b_s^{(k)} - b_t^{(k)} \perp \mathcal{F}_t^{(n)}$ , so that by a previous argument each  $b_t^{(k)}$  is a martingale with respect to  $\mathcal{F}_t^{(n)}$ . This allows us to define integrals like

$$\int_S^T f(t, b_t^{(1)}, \dots, b_t^{(n)}) db_t^{(k)}, \quad f \in \mathcal{V}_{\mathcal{F}^{(n)}}(S, T), \quad k = 1, \dots, n.$$

Thus, if we set  $\mathcal{V}_{\mathcal{F}^{(n)}}^{m,n}(S, T)$  to be the collection of all  $\{v_{ij}\}_{i=1, \dots, m; j=1, \dots, n}$  such that  $v_{ij} \in \mathcal{V}_{\mathcal{F}^{(n)}}(S, T)$  then the multi-dimensional Itô's integral in (A.6) above is well-defined and has the expected properties (in particular, it is a martingale).

The Itô's integral considered above turns out to be the main ingredient in defining an important class of stochastic processes which, as we shall see, are quite amenable to formal manipulations resembling those available from the ordinary calculus.

**Definition A.30.** Let  $b_t$  be BM in  $\mathbb{R}$  (with probability space  $(\Omega, \mathcal{F}, P)$ ). Then an *Itô process* (or *diffusion*) is a stochastic process in  $(\Omega, \mathcal{F}, P)$  of the type

$$(A.7) \quad X_t = X_0 + \int_0^t u(s, \omega) ds + \int_0^t v(s, \omega) db_s, \quad v \in \mathcal{V}_{\mathcal{H}},$$

for some  $\mathcal{H}$  as in Definition A.29.

Formally, we can rewrite (A.7) as

$$(A.8) \quad dX_t = u dt + v db_t,$$

where  $u$  is the *drift coefficient* and  $v$  is the *diffusion coefficient*. Itô's formula in (A.10) below shows that reasonable functions of Itô's processes are Itô's processes as well. It provides the correct change of variables formula in the setting of Stochastic Calculus.

**Proposition A.31.** If  $X_t$  is an Itô's process as in (A.8) and  $g = g(t, \omega) \in C^{2,1}([0, +\infty) \times \mathbb{R})$  then  $Y_t(t, \omega) = g(t, X_t(\omega))$  is an Itô's process as well. More precisely,

$$(A.9) \quad dY_t(t, X_t) = \frac{\partial g}{\partial t}(t, X_t) dt + \frac{\partial g}{\partial x}(t, X_t) dX_t + \frac{1}{2} \frac{\partial^2 g}{\partial x^2}(t, X_t) dX_t^2,$$

where in handling the quadratic term  $dX_t^2$  we should use the multiplication table

	$dt$	$db_t$
$dt$	0	0
$db_t$	0	$dt$

As a consequence,

$$(A.10) \quad dY_t = \left( \frac{\partial g}{\partial t} + u \frac{\partial g}{\partial x} + \frac{1}{2} v^2 \frac{\partial^2 g}{\partial x^2} \right) dt + v \frac{\partial g}{\partial x} db_t.$$

*Proof.* [Oks13, Theorem 4.1.2]. □

We now discuss the multi-dimensional version of this result. Let  $b_t = (b_t^{(1)}, \dots, b_t^{(n)})$  be BM in  $\mathbb{R}^n$ ,  $b_0 = 0$ . We can consider a multi-dimensional Itô's process

$$dX_t = udt + vdb_t,$$

where  $X = (X_1, \dots, X_n)^t \in \mathbb{R}^n$ ,  $u = (u_1, \dots, u_n)^t \in \mathbb{R}^n$ ,  $u_i \in \mathcal{V}_{\mathcal{F}(n)}$ , and  $v \in \mathcal{V}_{\mathcal{F}(n)}^{n,m}$ . Thus,

$$dX_{ti} = u_i dt + \sum_{j=1}^m v_{ij} db_t^{(j)}, \quad i = 1, \dots, n.$$

**Proposition A.32.** *If  $X$  is as above and  $Y(t, \omega) = g(t, X_t(\omega))$ , where  $g : [0, +\infty) \times \mathbb{R}^n \rightarrow \mathbb{R}^p$ , then*

$$dY_k = \frac{\partial g_k}{\partial t} dt + \sum_{i=1}^n \frac{\partial g_k}{\partial x_i} dX_i + \frac{1}{2} \sum_{i,j=1}^n \frac{\partial^2 g_k}{\partial x_i \partial x_j} dX_i dX_j,$$

where in handling the quadratic terms  $dX_i dX_j$  we should use the multiplication table

	$dt$	$db_t^{(i)}$
$dt$	0	0
$db_t^{(j)}$	0	$\delta_{ij} dt$

As a consequence,

$$(A.11) \quad dY_k = \left( \frac{\partial g_k}{\partial t} + \sum_{i=1}^n \frac{\partial g_k}{\partial x_i} u_i + \frac{1}{2} \sum_{i,j=1}^n \underbrace{\sum_{l=1}^m v_{il} v_{jl}}_{=(v^* v)_{ij}} \frac{\partial^2 g_k}{\partial x_i \partial x_j} \right) dt + \sum_{i=1}^n \sum_{l=1}^m \frac{\partial g_k}{\partial x_i} v_{il} db_t^{(l)}.$$

*Proof.* [Oks13, Theorem 4.2.1]. □

**Example A.33.** (The geometric Brownian) We assume that  $Y_t(t, \omega) = g(t, b_t(\omega))$  in (A.9), so that (A.10) becomes

$$(A.12) \quad dY_t = \left( \frac{\partial g}{\partial t} + \frac{1}{2} \frac{\partial^2 g}{\partial x^2} \right) dt + \frac{\partial g}{\partial x} db_t.$$

If  $Y_t = e^{\sigma b_t}$ ,  $\sigma \neq 0$ , we then get

$$(A.13) \quad dY_t = \frac{\sigma^2}{2} Y_t dt + \sigma Y_t db_t,$$

which shows that, in the stochastic setting, the exponential fails to satisfy the self-reproducing property under differentiation. We may cancel out the annoying first term in the right-hand side above by considering the *geometric Brownian process*

$$(A.14) \quad Z_t = e^{(\mu - \frac{\sigma^2}{2})t + \sigma b_t}, \quad \mu \in \mathbb{R},$$

which satisfies

$$(A.15) \quad dZ_t = \mu Z_t dt + \sigma Z_t db_t,$$

or equivalently,

$$(A.16) \quad \frac{dZ_t}{Z_t} = \mu dt + \sigma db_t.$$

As we shall see in Subsection A.6, this kind of process plays a central role in the Black-Scholes strategy in Finance.  $\square$

**Remark A.34.** (Diffusion processes) We discuss here one of the most notable motivation behind Itô's construction of his integral. Roughly, this accomplishment allowed him to properly interpret solutions of a large class of stochastic differential equations, which in particular led to a pathwise approach to diffusion processes. We start by recalling that a *transition function* is a map  $\mathfrak{R} : [0, +\infty) \times \mathbb{R}^n \times \mathcal{B}_n \rightarrow [0, 1]$  such that:

- $x \mapsto \mathfrak{R}(t, x, B)$  is measurable;
- $B \mapsto \mathfrak{R}(t, x, B)$  is a probability measure on  $\mathbb{R}^n$ ;
- $\mathfrak{R}(0, x, \cdot) = \delta_x$ ;
- The *Chapman-Kolmogorov equation* holds:

$$(A.17) \quad \mathfrak{R}(t + s, x, B) = \int_{\mathbb{R}^n} \mathfrak{R}(t, x, dy) \mathfrak{R}(s, y, B).$$

An application of Theorem A.3 guarantees the existence of a measurable space  $(\Omega, \mathcal{F})$  and, for each  $x \in \mathbb{R}^n$ , a probability measure  $\mathbb{P}^x$  on  $(\Omega, \mathcal{F})$  and a stochastic process  $X_t^x : (\Omega, \mathcal{F}, \mathbb{P}^x) \rightarrow \mathbb{R}^n$  such that  $\mathfrak{R}(t, x, B) = \mathbb{P}^x(X_t^x \in B)$ . Equivalently,  $\mathfrak{R}(t, x, \cdot) = X_t^x \# \mathbb{P}^x$ . In particular,  $X_0^x \# \mathbb{P}^x = \delta_x$ . Moreover, the following *Markov property* holds:

$$(A.18) \quad \mathbb{E}^x(f(X_{t+s}^x) | \mathcal{F}_s^X) = \mathbb{E}^x(f(X_t^x)), \quad \mathbb{P}^x \text{ a.s.},$$

for any  $x \in \mathbb{R}^n$  and any  $f$  as above. Intuitively, this means  $X_t^x$  is memoryless. Attached to any  $\mathfrak{R}$  as above is the associated semigroup  $t \mapsto \mathfrak{P}_t$  given by

$$(\mathfrak{P}_t f)(x) = \int_{\mathbb{R}^n} f(y) \mathfrak{R}(t, x, dy),$$

for  $f$  a bounded, measurable function on  $\mathbb{R}^n$ . Notice that  $(\mathfrak{P}_t f)(x) = \mathbb{E}^x f(X_t^x)$ . Now, any *Markov process*  $X_t^x$  as above has an *infinitesimal generator*  $L$ , which is a linear operator defined on the space of all functions  $f$  such that

$$\lim_{t \rightarrow 0} \frac{\mathfrak{P}_t f - f}{t}$$

exists. We then define, for any such  $f$ ,

$$(A.19) \quad (Lf)(x) = \lim_{t \rightarrow 0} \frac{(\mathfrak{P}_t f)(x) - f(x)}{t}.$$

Under certain regularity assumptions, the generator is a *diffusion operator*, that is,

$$(A.20) \quad (Lf)(x) = \frac{1}{2} \sum_{ij} a_{ij}(x) \frac{\partial^2 f}{\partial x_i \partial x_j} + \sum_i u_i(x) \frac{\partial f}{\partial x_i}, \quad x \in \mathbb{R}^n,$$

where  $a$  is a symmetric, non-negative matrix. The problem now is how to recover  $X_t^x$  starting from  $L$  (or, more precisely, from the coefficients  $a$  and  $u$  defining it). It turns out that Itô's calculus may be used to solve this problem as follows. Write  $a = v^t v$  and form the stochastic differential equation

$$dX_t = u(X_t)dt + v(X_t)db^t.$$

Under mild conditions on the coefficients, it is shown that this equation has a unique solution  $X_t^x$  with  $X_0^x = x$ . This of course means that

$$X_t^x = x + \int_0^t u(X_s^x)ds + \int_0^t v(X_s^x)db^s,$$

so a notion of stochastic integral is required here in order to properly interpret the last term above. It is now immediate to check that  $X_t^x$  solves the problem in the sense that (A.17), (A.18) and (A.19) are satisfied if we set  $\mathfrak{R}(t, x, B) = P(X_t^x \in B)$ , where  $P$  is Wiener measure. Here we only check that (A.19) holds. From Itô's formula (A.11), for any  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  we have

$$f(X_t^x) = f(x) + \int_0^t (Lf)(X_s^x)ds + M_t^f,$$

where  $M_t^f$  is a martingale. By taking expectation we see that

$$(\mathfrak{P}_t f)(x) = \mathbb{E}^x(f(X_t^x)) = f(x) + \mathbb{E}^x \left( \int_0^t (Lf)(X_s^x)ds \right),$$

and (A.19) follows. In fact, this procedure of solving suitable stochastic differential equations yields a systematic way of associating a diffusion process  $X_t^x$  to any operator  $L$  as above. Indeed, the resulting process is completely characterized by the fact that for any  $f$  the martingale

$$M_t^{f,x} = f(X_t^x) - f(x) - \int_0^t (Lf)(X_s^x)ds$$

has quadratic variation given by

$$\langle M^{f,x} \rangle_t^{(2)} = \int_0^t \left( \sum_{ij} a_{ij} \frac{\partial f}{\partial x_i} \frac{\partial f}{\partial x_j} \right) (X_s^x) ds.$$

We then say that  $X_t^x$  is the *diffusion process* driven by  $L$ . To see that we are in the right track, let us take  $L = \frac{1}{2}\Delta$  and let us set  $X^i = M^{x^i, x}$  for simplicity. It follows that

$$\begin{aligned} \langle X^i, X^j \rangle_t^{(2)} &:= \frac{1}{2} \left( \langle X^i + X^j \rangle_t^{(2)} - \langle X^i \rangle_t^{(2)} - \langle X^j \rangle_t^{(2)} \right) \\ &= \delta_{ij}t, \end{aligned}$$

so a celebrated result due to P. Lévy [KS12, Theorem 3.16] implies that  $X_t = b_t$ . Thus, as expected, BM is the diffusion process driven by the Laplacian  $\frac{1}{2}\Delta$ . Finally, we indicate how the diffusion process can be directly defined in terms of the coefficients defining  $L$ . It suffices to take  $f(x) = x^i$  and observe that

$$(A.21) \quad N_t^i := X_t^i - x_i - \int_0^t u_i(X_s)ds$$

is a martingale with

$$\langle N^i, N^j \rangle_t^{(2)} = \int_0^t a_{ij}(X_s)ds.$$

In particular, if  $L = \sum_i u_i \partial_i$  is a vector field then  $\langle N^i, N^j \rangle_t^{(2)} = 0$  and hence  $N_t = 0$ . Thus, (A.21) says that  $L$  integrates to a deterministic flow (no fluctuation!). In this way we recover the classical result on integration of vector fields. In a nutshell, whereas standard calculus allows us to integrate vector fields (giving rise to deterministic dynamical systems), Itô's calculus allows us to integrate diffusion operators (giving rise to random dynamical systems). For a modern take on the theory of diffusion processes we refer to [Bau14]  $\square$

**A.4. The Gaussian concentration inequality (again).** With the basics of Itô's calculus at hand, we present in the rest of this Appendix some of its most glamorous applications. We start by providing here an elegant proof of the optimal version of the Gaussian concentration inequality (5.18) which is attributed to B. Maurey in [Pis06, Chapter 2] and relies on the full power of the Stochastic Calculus developed in the previous section.

For  $0 \leq t \leq 1$  we consider the “reversed” heat semigroup  $P_t = e^{\frac{1}{2}(1-t)\Delta}$ , so that for any (smooth and Lipschitz)  $F : \mathbb{R}^k \rightarrow \mathbb{R}$  there holds

$$\frac{\partial}{\partial t}(P_t F) + \frac{1}{2}\Delta(P_t F) = 0.$$

We now use Itô's formula (A.11) with  $Y(t, \cdot) = (P_t F)(b_t)$ , where  $b_t$  is a standard BM in  $\mathbb{R}^k$  (recall that  $b_t - b_{t'} \sim \mathcal{N}(\vec{0}, (t - t')\mathbf{I}_k)$ ,  $t' < t$ ). Since  $u = 0$  and  $v_{ij} = \delta_{ij}$  it simplifies to

$$dY_t = \langle (\nabla P_t F)(b_t), db_t \rangle,$$

and integrating this from  $t = 0$  to  $t = 1$ ,

$$\begin{aligned} F(b_1) &= (e^{\frac{1}{2}\Delta} F)(0) + \int_0^1 \langle (\nabla P_t F)(b_t), db_t \rangle \\ &= \mathbb{E}(F(b_1)) + \int_0^1 \langle (\nabla P_t F)(b_t), db_t \rangle, \end{aligned}$$

where we used Propositions A.18 and A.28 in the last step. We may now adapt the Cramér-Chernoff method in Section 5 to this setting: for  $\tau > 0$  and  $w \geq 0$ ,

$$(A.22) \quad P(|F(b_1) - \mathbb{E}(F(b_1))| > \tau) \leq 2e^{-w\tau} \mathbb{E} \left( e^{w \int_0^1 \langle (\nabla P_t F)(b_t), db_t \rangle} \right),$$

so it remains to estimate the expectation in the right-hand side.

The key observation at this point is that  $\text{Lip}(P_t F) = \text{Lip}(F)$  and hence  $|\nabla P_t F| \leq \text{Lip}(F)$  a.s. Now let  $\pi = \{t_0 = 0 < t_1 < \dots < t_n = 1\}$  be a partition of the interval  $[0, 1]$  with  $|\pi| = \max_l |t_l - t_{l-1}|$  its width. As  $|\pi| \rightarrow 0$  the Itô's integral inside the expectation, by its very definition, may be arbitrarily approximated (say, in probability) by  $S_n$ , where

$$S_j = \sum_{l=1}^j \langle V_l, b_{t_l} - b_{t_{l-1}} \rangle, \quad 1 \leq j \leq n,$$

$V_l = (\nabla P_{t_{l-1}} F)(b_{t_{l-1}})$  is  $\mathcal{F}_{t_{l-1}}$ -measurable (where  $\{\mathcal{F}_t\}$  is the filtration associated to  $b_t$ ) and satisfies  $|V_l| \leq \text{Lip}(F)$ . We have

$$S_j = S_{j-1} + \langle V_j, b_{t_j} - b_{t_{j-1}} \rangle,$$

a decomposition into independent factors, and since the inner product follows the normal  $\mathcal{N}(0, |V_j|^2(t_j - t_{j-1}))$  by Proposition (4.7) (3), assuming of course that  $V_j \neq \vec{0}$ , we obtain

$$\begin{aligned} \mathbb{E}(e^{wS_j}) &= \mathbb{E}(e^{wS_{j-1}}) \mathbb{E} \left( e^{w \langle V_j, b_{t_j} - b_{t_{j-1}} \rangle} \right) \\ &\stackrel{(4.5)}{=} \mathbb{E}(e^{wS_{j-1}}) e^{\frac{1}{2}w^2 |V_j|^2 (t_j - t_{j-1})} \\ &\leq \mathbb{E}(e^{wS_{j-1}}) e^{\frac{1}{2}w^2 \text{Lip}(f)^2 (t_j - t_{j-1})}. \end{aligned}$$

Note that this obviously remains true if  $V_j = \vec{0}$ . In any case, if we iterate this starting with  $j = n$  we get

$$\mathbb{E}(e^{wS_n}) \leq e^{\frac{1}{2}w^2 \text{Lip}(f)^2},$$

where the right-hand side, remarkably enough, does not depend on  $\pi$ . Passing the limit as  $|\pi| \rightarrow 0$  on the left-hand side we thus obtain

$$\mathbb{E} \left( e^{w \int_0^1 \langle (\nabla P_t F)(b_t), db_t \rangle} \right) \leq e^{\frac{1}{2}w^2 \text{Lip}(f)^2},$$

which is the same as saying that

$$(A.23) \quad \int_0^1 \langle (\nabla P_t F)(b_t), db_t \rangle \in \text{SubG}(\text{Lip}(f)).$$

Leading this to (A.22) we get

$$P(|F(b_1) - \mathbb{E}(F(b_1))| > \tau) \leq 2e^{\frac{1}{2}w^2 \text{Lip}(F)^2 - w\tau},$$

so if we minimize the right-hand side over  $w \geq 0$  we find that

$$P(|F(b_1) - \mathbb{E}(F(b_1))| > \tau) \leq 2e^{-\frac{\tau^2}{2\text{Lip}(F)^2}}, \quad \tau > 0,$$

which is equivalent to (5.18) with the optimal constant  $C = 1/2$  because the normal random vector appearing there has been chosen so that  $X \sim \mathcal{N}(\vec{0}, I_k) \sim b_1$ .

**A.5. The Feynman-Kac formula and the path integral representation of the heat kernel.** Let us consider

$$g(t, X_t) = e^{-\int_0^t V(X_t)dt} w(T-t, X_t), \quad 0 \leq t \leq T,$$

where  $V = V(x)$ ,  $x \in \mathbb{R}^n$ , is a (well-behaved) potential function and we assume that  $X_t$  is an Itô's diffusion as in (A.8):

$$dX_t = u(X_t)dt + v(X_t)db_t,$$

where  $b_t$  is BM in  $\mathbb{R}^n$ . We compute that

$$\frac{\partial g}{\partial t} = -e^{-\int_0^t V(X_t)dt} \left( V(X_t)w(T-t, X_t) + \frac{\partial w}{\partial t}(T-t, X_t) \right),$$

$$\frac{\partial g}{\partial x_i} = e^{-\int_0^t V(X_t)dt} \frac{\partial w}{\partial x_i}(T-t, X_t),$$

and

$$\frac{\partial^2 g}{\partial x_i \partial x_j} = e^{-\int_0^t V(X_t)dt} \frac{\partial^2 w}{\partial x_i \partial x_j}(T-t, X_t),$$

so that Itô's formula in (A.11) applies to give

$$\begin{aligned} dg &= e^{-\int_0^t V(X_t)dt} \left( -\frac{\partial w}{\partial t}(T-t, X_t) + \sum_i u_i \frac{\partial w}{\partial x_i}(T-t, X_t) - \right. \\ &\quad \left. -V(X_t)w(T-t, X_t) + \frac{1}{2}\mathcal{L}'w(T-t, X_t) \right) dt \\ &\quad + e^{-\int_0^t V(X_t)dt} \sum_i \frac{\partial g}{\partial x_i} db_t^{(i)}, \end{aligned}$$

where

$$\mathcal{L}'w = \sum_{ij} (v^* v)_{ij} \frac{\partial^2 w}{\partial x_i \partial x_j}.$$

Putting all the pieces of this computation together we obtain a remarkable stochastic (or path integral) representation of solutions of certain heat-type equations.

**Proposition A.35.** (Feynman-Kac formula I) If  $w = w(t, x)$  satisfies the heat-type equation

$$\begin{cases} \frac{\partial w}{\partial t} &= \frac{1}{2} \mathcal{L}' w + \langle u, \nabla w \rangle - V w \\ w(0, x) &= f(x) \end{cases}$$

then the following holds:

$$(A.24) \quad w(t, x_0) = \mathbb{E}_{x_0} \left( e^{-\int_0^t V(X_s) ds} f(X_t) \right),$$

where  $\mathbb{E}_{x_0}$  refers to the law  $P_{x_0}$  of BM in  $\mathbb{R}^n$  starting at  $x_0$ .

*Proof.* From the computation above,

$$dg = e^{-\int_0^t V(X_s) ds} \sum_i \frac{\partial g}{\partial x_i} db_t^{(i)},$$

which shows that the process  $g(t, X_t)$  is a martingale (with respect to  $\mathcal{F}^{(n)}$ ); see Proposition A.28. Since

$$\mathbb{E}_{x_0} (g(t, X_t)) |_{t=0} = \mathbb{E}_{x_0} (w(T, X_0)) = w(T, x_0),$$

and

$$\mathbb{E}_{x_0} (g(t, X_t)) |_{t=T} = \mathbb{E}_{x_0} \left( e^{-\int_0^t V(X_s) ds} w(0, X_T) \right) = \mathbb{E}_{x_0} \left( e^{-\int_0^t V(X_s) ds} f(X_T) \right),$$

the result follows in view of Proposition A.18 and the fact that  $T$  is arbitrary.  $\square$

**Corollary A.36.** (Exponential control) Under the conditions above, if  $|f| \leq M$  and  $V \geq c$ ,  $c \in \mathbb{R}$ , then

$$|u(t, x_0)| \leq M e^{-ct}.$$

An important special case of Proposition A.24 occurs when  $u = 0$  and  $v_{ij} = \delta_{ij}$ , so that

$$\mathcal{L}' = \frac{1}{2} \Delta,$$

where  $\Delta$  is the Laplacian. We then see that any solution of

$$\begin{cases} \frac{\partial w}{\partial t} &= \frac{1}{2} \Delta w - V w \\ w(0, x) &= f(x) \end{cases}$$

satisfies

$$(A.25) \quad w(t, x) = \mathbb{E}_x \left( e^{-\int_0^t V(b_\tau) d\tau} f(b_t) \right).$$

On the other hand, we know from Analysis [Paz12] that this can be rewritten as

$$(A.26) \quad w(t, x) = (e^{t\mathcal{L}} f)(x) = \int_{\mathbb{R}^n} K_{\mathcal{L}}(t; x, y) f(y) dy,$$

where  $e^{t\mathcal{L}}$  is the heat semigroup generated by  $\mathcal{L} = \frac{1}{2} \Delta - V$  and  $K_{\mathcal{L}}$  is the associated *heat kernel*, i.e.  $K_{\mathcal{L}}$  satisfies

$$\begin{cases} \frac{\partial K_{\mathcal{L}}}{\partial t} &= \frac{1}{2} \Delta K_{\mathcal{L}} - V K_{\mathcal{L}} \\ K_{\mathcal{L}}(0; x, y) &= \delta(x - y) \end{cases}$$

This suggests the existence of a stochastic representation for  $K_{\mathcal{L}}$ , thus pointing toward a version of a Feynman-Kac formula working at the more fundamental level of heat kernels.

To find this representation we fix  $t > 0$  and consider the process  $\{B_s\}_{0 \leq s < t}$  satisfying

$$dB_s = db_s - \frac{B_s - y}{t - s} ds, \quad B_0 = x,$$

or equivalently,

$$B_s = x + b_s - \int_0^s \frac{B_\tau - y}{t - \tau} d\tau.$$

We will now show that the law of this process can be computed in terms of  $K_{\frac{1}{2}\Delta}$ , the heat kernel of the Laplacian  $\frac{1}{2}\Delta$ , and  $P_x$ , the law of BM starting at  $x$ .

We first note that the discussion above gives

$$(A.27) \quad (e^{\frac{1}{2}\tau\Delta}f)(x) = \int_{\mathbb{R}^n} K_{\frac{1}{2}\Delta}(\tau; x, y) f(y) dy = \mathbb{E}_x(f(b_\tau)), \quad \tau \geq 0.$$

In particular,

$$(A.28) \quad \int_{\mathbb{R}^n} K_{\frac{1}{2}\Delta}(t; x, y) dy = 1,$$

which also follows from Proposition 4.2 because, as is well-known,

$$K_{\frac{1}{2}\Delta}(t; x, y) = (2\pi t)^{-n/2} e^{-|x-y|^2/2t},$$

From this we see that

$$\nabla_x \ln K_{\frac{1}{2}\Delta}(t; x, y) = -\frac{x-y}{t},$$

and hence

$$(A.29) \quad dB_s = db_s + \nabla_x \ln K_{\frac{1}{2}\Delta}(t-s; B_s, y) ds, \quad s < t.$$

Thus, the Brownian bridge  $B_s$  is just the Brownian motion  $b_s$  with an added drift involving the logarithmic derivative of  $K_{\frac{1}{2}\Delta}$ . We note however that the drift is singular at  $s = t$ . Fortunately, careful first order estimates of  $K_{\frac{1}{2}\Delta}$  [Hsu02, Section 5.5] allow us to bypass this difficulty and confirm not only that this is well defined for  $s = t$  but also that  $B_s \rightarrow y$  as  $s \rightarrow t$ . Thus, we call  $\{B_s\}_{0 \leq s \leq t}$  the *Brownian bridge* connecting  $x$  to  $y$  with lifetime  $t$ .

We should think of  $B_s$  as a process on the *bridge space*  $C_{t;x,y} \subset C_x$  of all Brownian paths starting at  $x$  and conditioned to hit  $y$  at time  $t$ ; cf. Remark A.15. To find the law  $B_s$  we first note that

$$\frac{\partial}{\partial t} \ln K_{\frac{1}{2}\Delta} = \frac{1}{2} \Delta \ln K_{\frac{1}{2}\Delta} + \frac{1}{2} \|\nabla \ln K_{\frac{1}{2}\Delta}\|^2,$$

so if we apply Itô's formula to

$$E_s = \ln \frac{K_{\frac{1}{2}\Delta}(t-s; B_s, y)}{K_{\frac{1}{2}\Delta}(t; x, y)}$$

we find that

$$dE_s = \langle F_s, dB_s \rangle - \frac{1}{2} \|F_s\|^2 ds,$$

where  $F_s = \nabla_x K_{\frac{1}{2}\Delta}(t-s; B_s, y)$ . Now define a measure  $Q$  in  $C_{t;x,y}$  by

$$\frac{dQ}{dP_x}|_{\mathcal{F}_s} = \exp \left( \int_0^s \langle F_u, db_u \rangle - \frac{1}{2} \int_0^s \|F_u\|^2 du \right).$$

By Girsanov's theorem [Oks13, Theorem 8.6.4], under  $Q$  the process

$$B_s - \int_0^s F_u du$$

is a BM. Thus, we see that  $P_{t;x,y} := Q$  is the law of the Brownian bridge  $B_s$  and there holds

$$(A.30) \quad \frac{dP_{t;x,y}}{dP_x}|_{\mathcal{F}_s} = \frac{K_{\frac{1}{2}\Delta}(t-s; b_s, y)}{K_{\frac{1}{2}\Delta}(t; x, y)}.$$



With these preliminaries at hand, finally we will be able to provide a path integral representation for  $K_{\mathcal{L}}$ .

**Proposition A.37.** (Feynman-Kac formula II) One has

$$K_{\mathcal{L}}(t; x, y) = K_{\frac{1}{2}\Delta}(t; x, y) \mathbb{E}_{t;x,y} \left( e^{-\int_0^t V(B_\tau) d\tau} \right).$$

In other words, if we define the conditional Wiener measure on  $C_{t;x,y}$  by

$$(A.31) \quad \mu_{t;x,y} = K_{\frac{1}{2}\Delta}(t; x, y) P_{t;x,y}$$

then

$$(A.32) \quad K_{\mathcal{L}}(t; x, y) = \int_{C_{t;x,y}} e^{-\int_0^t V(B_\tau) d\tau} d\mu_{t;x,y}.$$

*Proof.* First we have from (A.27) with  $\tau = 0$  that

$$e^{-\int_0^t V(b_\tau) d\tau} f(b_t) = \int_{\mathbb{R}^n} K_{\frac{1}{2}\Delta}(0; b_t, y) e^{-\int_0^t V(b_\tau) d\tau} f(y) dy,$$

and taking expectation we get

$$\begin{aligned} \mathbb{E}_x \left( e^{-\int_0^t V(b_\tau) d\tau} f(b_t) \right) &= \int_{\mathbb{R}^n} K_{\frac{1}{2}\Delta}(t; x, y) \mathbb{E}_x \left( \frac{K_{\frac{1}{2}\Delta}(0; b_t, y)}{K_{\frac{1}{2}\Delta}(t; x, y)} e^{-\int_0^t V(b_\tau) d\tau} f(y) \right) dy \\ &\stackrel{(A.30) \text{ with } s=t}{=} \int_{\mathbb{R}^n} K_{\frac{1}{2}\Delta}(t; x, y) \mathbb{E}_{t;x,y} \left( e^{-\int_0^t V(B_\tau) d\tau} f(B_t) \right) dy. \end{aligned}$$

On the other hand, we know from (A.25) and (A.26) that

$$\mathbb{E}_x \left( e^{-\int_0^t V(b_\tau) d\tau} f(b_t) \right) = \int_{\mathbb{R}^n} K_{\mathcal{L}}(t; x, y) f(y) dy.$$

Since  $B_t = y$  and  $f$  is arbitrary, the result follows.  $\square$

**Remark A.38.** (The Laplacian heat kernel as a transition probability) It follows from the formalism above that the conditioned Wiener measure  $\mu_{t;x,y}$  in (A.31) may be characterized by

$$\int_{\mathbb{R}^n} \left( \int_{C_{t;x,y}} F(\omega) d\mu_{t;x,y}(\omega) \right) dy = \int_{C_x} F(\omega) dP_x(\omega),$$

where  $F$  varies over the set all bounded functions on  $(C_x, P_x)$ , the Wiener space starting at  $x$ ; cf. Remark A.15. By taking  $F \equiv 1$  we thus see that

$$\int_{\mathbb{R}^n} \left( \int_{C_{t;x,y}} d\mu_{t;x,y}(\omega) \right) dy = 1,$$

which is just a restatement of (A.28), as it follows either from (A.31) or from (A.32) with  $V \equiv 1$  that

$$K_{\frac{1}{2}\Delta}(t; x, y) = \int_{C_{t;x,y}} d\mu_{t;x,y}(\omega),$$

the total measure of the Brownian bridge  $C_{t;x,y}$  endowed with  $\mu_{t;x,y}$ . It then follows that:

- for each  $t \geq 0$  and  $x \in \mathbb{R}^n$  the function

$$y \mapsto K_{\frac{1}{2}\Delta}(t; x, y) = \mu_{t;x,y}(C_{t;x,y})$$

defines a probability density in  $\mathbb{R}^n$ ;

- For each  $U \in \mathcal{B}^n$  the quantity

$$P_{t;x}(U) := \int_U K_{\frac{1}{2}\Delta}(t; x, y) dy$$

may be interpreted as the probability that a Brownian path passes through  $U$  when  $s = t$  given that it has started at  $x$  when  $s = 0$ . By shrinking  $U$  to  $\{y\}$  we thus conclude that  $K_{\frac{1}{2}\Delta}(t; x, y)$  may be viewed as the *transition probability* that a Brownian path hits  $y$  at  $s = t$  given that it has started at  $x$  when  $s = 0$ .

**Remark A.39.** (Weyl's law and its modern incarnations in Index Theory) In case  $e^{t\mathcal{L}}$  is trace class, its trace can be computed by integrating (A.32) along the diagonal of  $\mathbb{R}^n \times \mathbb{R}^n$  defined by  $x = y$ :

$$\mathrm{Tr} e^{t\mathcal{L}} = \int_{\mathbb{R}^n} \left( \int_{C_{t;x,x}} e^{-\int_0^t V(X_\tau) d\tau} d\mu_{t;x,x} \right) dx.$$

Thus, if we define a “measure”  $d\mu_t$  on the space  $C_t = \cup_{x \in \mathbb{R}^n} C_{t;x,x}$  of all *Brownian loops* on  $\mathbb{R}^n$  with lifetime  $t$  by  $d\mu_t = d\mu_{t;x,x} dx$ , then

$$\mathrm{Tr} e^{t\mathcal{L}} = \int_{C_t} e^{-\int_0^t V(X_\tau) d\tau} d\mu_t.$$

This result holds *verbatim* if we replace  $\mathbb{R}^n$  by a compact Riemannian manifold  $(M, g)$  and  $\Delta_g$  is the corresponding Laplace-Beltrami operator. In particular,

$$\mathrm{Tr} e^{\frac{1}{2}t\Delta_g} = \int_{C_t} d\mu_t$$

where the measure  $\mu_t$  is now defined locally by using coordinate charts and then pasted together in the usual manner; see [Hsu02] for a careful discussion of Brownian motion on Riemannian manifolds and its basic properties. Using the well-known fact that, as  $t \rightarrow 0$ , the typical Brownian loop in  $C_t$  shrinks to its base point while still remaining is a geodesic ball whose radius vanishes as  $t \rightarrow 0$  [Hsu02, Lemma 7.7], we see that the path integral on the right-hand side “localizes” around  $M \subset C_t$ . Combining this with the “principle of not feeling the curvature”, according to which  $K_{\frac{1}{2}\Delta_g}(t; x, x) \sim (2\pi t)^{-n/2}$  as  $t \rightarrow 0$ , we conclude that

$$\mathrm{Tr} e^{\frac{1}{2}t\Delta_g} \sim (2\pi t)^{-n/2} \mathrm{vol}(M, g).$$

Since  $\mathrm{Tr} e^{\frac{1}{2}t\Delta_g} = \sum_i e^{-\frac{1}{2}\lambda_i t}$ , where the  $\{\lambda_i\}$  are the (positive) eigenvalues of  $\Delta_g$ , we thus obtain Weyl's celebrated result that we can “hear” the volume of  $(M, g)$  from the asymptotic behavior of its spectrum. A much more sophisticated incarnation of this argument, exploring the short time heat kernel asymptotics of a certain “supersymmetric” version of the Dirac operator acting on spinors, eventually leads to a probabilistic proof of the Atiyah-Singer index formula [Bis84, Hsu02]. Further applications of this approach, which involves a careful analysis of the asymptotics of an appropriate Feynman-Kac formula for the heat kernel of certain “Hodge Laplacians” acting on sections of geometric vector bundles over (not necessarily compact) Riemannian manifolds (possibly carrying a non-empty boundary), we refer to [dL17a, dL17b, dL20] and their references.  $\square$

**A.6. The Black-Scholes strategy in Finance.** Here we derive the celebrated Black-Scholes option pricing formula<sup>52</sup>. From the outset, this involves a risky asset  $S_t$ , a *stock*, evolving in time according to a geometric Brownian as in Example A.33:

$$(A.33) \quad \frac{dS_t}{S_t} = \mu dt + \sigma db_t.$$

Here,  $\mu > 0$  is the *mean rate of return* and  $\sigma > 0$  is the *volatility*. Recall that

$$S_t = S_0 e^{(\mu - \frac{\sigma^2}{2})t + \sigma b_t}$$

<sup>52</sup>As it is well-known, this has been worth a Nobel Prize in 1997.

provides the explicit solution of (A.33). In particular,

$$\ln S_t = \ln S_0 + \left( \mu - \frac{\sigma^2}{2} \right) t + \sigma b_t \sim \mathcal{N} \left( \underbrace{\ln S_0 + \left( \mu - \frac{\sigma^2}{2} \right) t}_{=:m}, \underbrace{\sigma^2 t}_{=: \nu^2} \right),$$

so that, by Example 4.9,

$$(A.34) \quad S_t = e^{\ln S_t} \sim \mathcal{LN} \left( S_0 e^{\mu t}, S_0^2 e^{2\mu t} (e^{\sigma^2 t} - 1) \right) = \Lambda(m, \nu^2).$$

On the other hand, we have an investor's *portfolio*  $(A_t, B_t)$  whose value is

$$(A.35) \quad V_t = A_t S_t + B_t \gamma_t,$$

with the risk-less *bond*  $\gamma_t$  satisfying  $d\gamma_t = r\gamma_t dt$ , where  $r > 0$  is the associated *interest rate*. The *option pricing problem* addressed by Black-Scholes consists in adjusting the trading strategy  $(A_t, B_t)$  to the underlying asset  $S_t$  by (deterministically!) finding a function  $u$  such that

$$(A.36) \quad V_t = u(t, S_t), \quad 0 \leq t \leq T,$$

where  $T > 0$  is the *expiration time* for the option. The key point here is that  $u = u(t, x)$  should satisfy a certain PDE. In order to find it, we start with (A.33) and apply Itô's formula to (A.36) to check that

$$(A.37) \quad dV_t = \left( \frac{\partial u}{\partial t} + \mu S_t \frac{\partial u}{\partial x} + \frac{\sigma^2 S_t^2}{2} \frac{\partial^2 u}{\partial x^2} \right) dt + \sigma S_t \frac{\partial u}{\partial x} db_t.$$

On the other hand, if we assume that our portfolio is *self-financing* in the sense that

$$dV_t = A_t dS_t + B_t d\gamma_t,$$

we get

$$(A.38) \quad dV_t = (\mu A_t S_t + r B_t \gamma_t) dt + \sigma A_t S_t db_t.$$

By comparing the diffusion and drift coefficients in the expressions for  $dV_t$  above we get

$$(A.39) \quad A_t = \frac{\partial u}{\partial x}$$

and hence

$$(A.40) \quad \frac{\partial u}{\partial t} + \frac{\sigma^2 S_t^2}{2} \frac{\partial^2 u}{\partial x^2} = r B_t \gamma_t.$$

Now note that from (A.35), (A.36) and (A.39),

$$(A.41) \quad B_t \gamma_t = V_t - A_t S_t = u - S_t \frac{\partial u}{\partial x},$$

so if we replace this in the right-hand side of (A.40) and make  $S_t = x$  we conclude that  $u$  must satisfy the *Black-Scholes equation*

$$(A.42) \quad \frac{\partial u}{\partial t} + \frac{\sigma^2 x^2}{2} \frac{\partial^2 u}{\partial x^2} + r x \frac{\partial u}{\partial x} - r u = 0.$$

Notice that the coefficients of this PDE depend on  $\sigma$  and  $r$  but not on  $\mu$ . Also, as written the PDE fails to be of heat type because the coefficients of  $\partial u / \partial t$  and  $\partial^2 u / \partial x^2$  have the same sign. This suggests that we should try to solve it by imposing the "terminal condition"

$$(A.43) \quad u(T, S_T) = V_T.$$

In fact, the choice

$$(A.44) \quad V_T = \max\{0, S_T - K\}, \quad K > 0,$$

corresponds to the investor holding at time  $t = 0$  the option (but not the obligation) of buying the stock by a fixed price  $K$  at the expiration time  $T$ . Hence, in this *European call*, if  $S_T > K$  then the owner of the option will obtain the payoff  $S_T - K$  whereas if  $S_T \leq K$  the owner will not exercise his option, thus obtaining a null payoff.

In terms of the cumulative normal distribution

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-y^2/2} dy, \quad x \in \mathbb{R},$$

the Black-Scholes equation (A.42) with  $u(T, S_T) = \max\{0, S_T - K\}$  may be explicitly solved as

$$(A.45) \quad u(t, x) = x\Phi(g(t, x)) - Ke^{-r(T-t)}\Phi(h(t, x)),$$

where

$$g(t, x) = \frac{\ln(x/K) + (r + \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{T-t}}$$

and

$$h(t, x) = g(t, x) - \sigma\sqrt{T-t} = \frac{\ln(x/K) + (r - \frac{1}{2}\sigma^2)(T-t)}{\sigma\sqrt{T-t}}.$$

Notice that, as expected, the solution depends on  $\sigma$  and  $r$  but not on  $\mu$ . We conclude that the *Black-Scholes option pricing formula*

$$V_0 = u(0, S_0) = S_0\Phi(g(0, S_0)) - Ke^{-rT}\Phi(h(0, S_0))$$

provides the rational price to hold at the initial time  $t = 0$  a European call option with price  $K$ . Also, from (A.39) and (A.41) we see that the corresponding self-financing strategy is

$$(A_t, B_t) = \left( \frac{\partial u}{\partial x}, \gamma_0^{-1} e^{-rt} \left( u - S_t \frac{\partial u}{\partial x} \right) \right).$$

The explicit solution (A.45) to (A.42)-(A.44) may be obtained in many ways. For instance, a “deterministic” approach may be pursued upon successive changes of variables so as to transform (A.42) into the standard heat equation, which can then be explicitly solved by the usual methods [KK12, Section 10.3]. Alternatively, we may appeal to the full power of the Stochastic Calculus as follows. Let us set  $\theta = (\mu - r)/\sigma$  and consider the process

$$M_t = e^{-\theta b_t - \frac{1}{2}\theta^2 t}, \quad 0 \leq t \leq T.$$

From Itô’s formula we obtain

$$dM_t = -\theta M_t db_t,$$

so  $M_t$  is a  $b_t$ -martingale. Hence, by Proposition A.18,

$$\mathbb{E}^P(M_T) = \mathbb{E}^P(M_0) = \mathbb{E}^P(1) = 1,$$

where  $P$  is Wiener measure. Thus, if we define a new measure  $P^\bullet$  on Wiener space by requiring that  $dP^\bullet = M_T dP$ , it is immediate that  $P^\bullet$  is a probability measure. A version of Girsanov’s theorem [Oks13, Theorem 8.6.4] applies here and we conclude that  $b_t^\bullet := b_t + \theta t$  is a standard Brownian with respect to  $P^\bullet$  (so that  $b_t^\bullet \sim_{P^\bullet} \mathcal{N}(0, t)$ ) and, moreover,

$$(A.46) \quad dS_t = rS_t dt + \sigma S_t db_t^\bullet.$$

Thus, we have been able to modify the drift of the dynamics of the stock (from  $\mu S_t$  to  $r S_t$ ) at the cost of changing the underlying measure (from  $P$  to  $P^\bullet$ ) and the driving Brownian (from  $b_t$  to  $b_t^\bullet$ ). The reason for doing this is now obvious: in terms of the infinitesimal generator of (A.46), namely,

$$L = \frac{\sigma^2 x^2}{2} \frac{\partial^2}{\partial x^2} + r x \frac{\partial}{\partial x},$$

(A.42) may be rewritten as

$$\frac{\partial u}{\partial t} + Lu - ru = 0,$$

whose solution may be obtained by the method leading to the Feynman-Kac formula discussed in Section A.5 (with  $r$  playing the role of a constant potential). Indeed, if we apply Itô's formula to  $v(t, S_t) := e^{r(T-t)} u(t, S_t)$ , we easily see that

$$dv(t, S_t) = \sigma S_t e^{r(T-t)} \frac{\partial u}{\partial x}(t, S_t) db_t^\bullet,$$

which means that  $v(t, S_t)$  is a  $b_t^\bullet$ -martingale. Thus, if we calculate the (identical!) expectations at the endpoints of the interval  $[t, T]$  and use (A.44) we end up with

$$u(t, S_t) = e^{-r(T-t)} \mathbb{E}^{P^\bullet}(\max\{S_T - K, 0\}).$$

We now observe that, due to (A.46) and similarly to (A.34), we now have

$$(A.47) \quad S_T \sim_{P^\bullet} \mathcal{LN}\left(S_t e^{r(T-t)}, S_t^2 e^{2r(T-t)} (e^{\sigma^2(T-t)} - 1)\right),$$

or equivalently,

$$(A.48) \quad \ln S_T \sim_{P^\bullet} \mathcal{N}\left(\underbrace{\ln S_t + \left(\mu - \frac{\sigma^2}{2}\right)(T-t)}_{=:m}, \underbrace{\sigma^2(T-t)}_{=:v^2}\right) = \Lambda(m, v^2).$$

Now,

$$\begin{aligned} \mathbb{E}^{P^\bullet}(\max\{S_T - K, 0\}) &= \int_K^{+\infty} (S_T - K) dF_{S_T} \\ &= \int_K^{+\infty} S_T dF_{S_T} - K \int_K^{+\infty} dF_{S_T}, \end{aligned}$$

where  $dF_{S_T} = \psi_{S_T} dx$  is the cdf of  $S_T$ , and these integrals may be computed in terms of  $\Phi$  by means of the recipe in Example 4.9. For instance,

$$\begin{aligned}
\int_K^{+\infty} dF_{S_T} &= 1 - \int_{-\infty}^K dF_{S_T} \\
&= 1 - F_{S_T}(K) \\
&= 1 - \Phi \left( \frac{\overbrace{\ln K - \ln S_t - \left(r - \frac{\sigma^2}{2}\right)(T-t)}^{=-m}}{\underbrace{\sigma\sqrt{T-t}}_{=\nu}} \right) \\
&= 1 - \Phi \left( -\frac{\ln S_t/K + \left(r - \frac{\sigma^2}{2}\right)(T-t)}{\sigma\sqrt{T-t}} \right) \\
&= \Phi \left( \frac{\ln S_t/K + \left(r - \frac{\sigma^2}{2}\right)(T-t)}{\sigma\sqrt{T-t}} \right) \\
&= \Phi(h(t, S_t)),
\end{aligned}$$

where we used (4.17), our choices for  $m$  and  $\nu$  as in (A.48) and the fact that  $\Phi(x) + \Phi(-x) = 1$ . Similarly,

$$\begin{aligned}
\int_K^{+\infty} S_T dF_{S_T} &= \mathbb{E}^{P^\bullet}(S_T) - \int_{-\infty}^K S_T dF_{S_T} \\
&= S_t e^{r(T-t)} - \int_{-\infty}^K x \psi_{S_T}(x) dx \\
&= S_t e^{r(T-t)} - \frac{1}{\sqrt{2\pi}\sigma\sqrt{T-t}} \int_{-\infty}^K e^{-\frac{1}{2} \left( \frac{\ln x - \ln S_t - \left(r - \frac{\sigma^2}{2}\right)(T-t)}{\sigma\sqrt{T-t}} \right)^2} dx,
\end{aligned}$$

where we used (A.47) and (4.16). After an appropriate change of variables and using again that  $\Phi(x) + \Phi(-x) = 1$ , we get

$$\begin{aligned}
\int_K^{+\infty} S_T dF_{S_T} &= e^{r(T-t)} S_t \Phi \left( \frac{-\ln K + \ln S_t + \left(r - \frac{\sigma^2}{2}\right)(T-t) + \sigma^2(T-t)}{\sigma\sqrt{T-t}} \right) \\
&= e^{r(T-t)} S_t \Phi(g(t, X_t)).
\end{aligned}$$

Putting all the pieces of this computation together we find that

$$u(t, S_t) = S_t \Phi(g(t, S_t)) - K e^{-r(T-t)} \Phi(h(t, S_t)),$$

which matches (A.45) if we make  $S_t = x$ .

## REFERENCES

- [AB69] John Aitchison and James Alan Calvert Brown. *The lognormal distribution, with special reference to its uses in Economics*. Cambridge University Press, 1969.
- [Ame85] Takeshi Amemiya. *Advanced econometrics*. Harvard university press, 1985.
- [And03] Theodore Wilbur Anderson. *An introduction to multivariate statistical analysis*. Wiley & Sons, 2003.

- [AP09] Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics: An empiricist's companion*. Princeton university press, 2009.
- [AS16] Noga Alon and Joel H Spencer. *The probabilistic method*. John Wiley & Sons, 2016.
- [Bau13] Johannes Bausch. On the efficient calculation of a linear combination of chi-square random variables with an application in counting string vacua. *Journal of Physics A: Mathematical and Theoretical*, 46(50):505202, 2013.
- [Bau14] Fabrice Baudoin. *Diffusion processes and stochastic calculus*. EMS, 2014.
- [BC11] Alexandre Belloni and Victor Chernozhukov. *High dimensional sparse econometric models: An introduction*. Springer, 2011.
- [Ber04] Richard A Berk. *Regression analysis: A constructive critique*, volume 11. Sage, 2004.
- [BHK20] Avrim Blum, John Hopcroft, and Ravindran Kannan. *Foundations of data science*. Cambridge University Press, 2020.
- [Bis84] Jean-Michel Bismut. The Atiyah—Singer theorems: a probabilistic approach. i. the index theorem. *Journal of functional analysis*, 57(1):56–99, 1984.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*, (2013). OUP: Oxford, 2013.
- [Bor75] Christer Borell. The Brunn-Minkowski inequality in Gauss space. *Inventiones mathematicae*, 30(2):207–216, 1975.
- [Bul03] Michael George Bulmer. *Francis Galton: pioneer of heredity and biometry*. JHU Press, 2003.
- [BVDG11] Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- [CA06] Kai Lai Chung and Farid AitSahlia. *Elementary probability theory: with stochastic processes and an introduction to mathematical finance*. Springer Science & Business Media, 2006.
- [Car86] E Carlstein. Simultaneous confidence regions for predictions. *The American Statistician*, 40(4):277–279, 1986.
- [Cas08] George Casella. *Statistical design*. Springer, 2008.
- [CB21] George Casella and Roger L Berger. *Statistical inference*. Cengage Learning, 2021.
- [Che21] Louis HY Chen. Stein's method of normal approximation: Some recollections and reflections. *The Annals of Statistics*, 49(4):1850–1863, 2021.
- [CLC21] Denis Chetverikov, Zhipeng Liao, and Victor Chernozhukov. On cross-validated lasso in high dimensions. *The Annals of Statistics*, 49(3):1300–1317, 2021.
- [Cow72] Ruth Schwartz Cowan. Francis galton's statistical ideas: the influence of eugenics. *Isis*, 63(4):509–528, 1972.
- [CW09] R Dennis Cook and Sanford Weisberg. *Applied regression including computing and graphics*. John Wiley & Sons, 2009.
- [DF87] Persi Diaconis and David Freedman. A dozen de Finetti-style results in search of a theory. *Annales de l'IHP Probabilités et statistiques*, 23(S2):397–423, 1987.
- [dL17a] Levi de Lima. A Feynman–Kac formula for differential forms on manifolds with boundary and geometric applications. *Pacific Journal of Mathematics*, 292(1):177–201, 2017.
- [dL17b] Levi Lopes de Lima. A probabilistic proof of the Gauss-Bonnet formula for manifolds with boundary. *arXiv:1709.03772*, 2017.
- [dL20] Levi L. de Lima. Heat conservation for generalized Dirac Laplacians on manifolds with boundary. *Annali di Matematica Pura ed Applicata (1923-)*, 199(3):997–1021, 2020.
- [dL25] Levi Lopes de Lima. Rmarkdown labs in Github. <https://github.com/levilopesdelima/stat-inference-labs>, 2025.
- [DM88] Edward J Dudewicz and Satya Mishra. *Modern mathematical statistics*. John Wiley & Sons, Inc., 1988.
- [Don00] David L Donoho. High-dimensional data analysis: The curses and blessings of dimensionality. *AMS math challenges lecture*, 1(2000):32, 2000.
- [DS14] Morris H DeGroot and Mark J Schervish. *Probability and statistics*, volume 563. Pearson Education London, UK., 2014.
- [EM77] Bradley Efron and Carl Morris. Stein's paradox in statistics. *Scientific American*, 236(5):119–127, 1977.
- [Far06] Julian J Faraway. *Linear models with R*. Chapman and Hall/CRC, 2006.
- [FBG<sup>+</sup>16] Arnoldo Frigessi, Peter Bühlmann, Ingrid K Glad, Sylvia Richardson, and Marina Vannucci. Some themes in high-dimensional statistics. In *Statistical Analysis for High-Dimensional Data: The Abel Symposium 2014*, pages 1–13. Springer, 2016.
- [Fer17] Thomas S Ferguson. *A course in large sample theory*. Routledge, 2017.
- [FG13] Bert E Fristedt and Lawrence F Gray. *A modern approach to probability theory*. Springer Science & Business Media, 2013.
- [Fis15] Ronald A Fisher. Frequency distribution of the values of the correlation coefficient in samples from an indefinitely large population. *Biometrika*, 10(4):507–521, 1915.
- [Fis21] R Fisher. On the “probable error” of a coefficient of correlation deduced from a small sample. *Metron*, pages 3–32, 1921.
- [Fis22] Ronald A Fisher. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604):309–368, 1922.
- [Fis25] Ronald A. Fisher. Applications of “student's” distribution. *Metron*, pages 90–104, 1925.
- [Fis11] Hans Fischer. *A history of the central limit theorem: from classical to modern probability theory*, volume 4. Springer, 2011.
- [FK16] Alan Frieze and Michal Karonski. *Introduction to random graphs*. Cambridge University Press, 2016.
- [FP15] Juliana Carvalho Ferreira and Cecilia Maria Patino. What does the  $p$  value really mean? *Jornal Brasileiro de Pneumologia*, 41(5):485, 2015.

- [Gib21] Eric W Gibson. The role of  $p$ -values in judging the strength of evidence and realistic replication expectations. *Statistics in Biopharmaceutical Research*, 13(1):6–18, 2021.
- [GKPS99] Mikhael Gromov, Misha Katz, Pierre Pansu, and Stephen Semmes. *Metric structures for Riemannian and non-Riemannian spaces*, volume 152. Springer, 1999.
- [Gne18] Boris V Gnedenko. *Theory of probability*. CRC Press, 2018.
- [Gor16] Prakash Gorroochurn. *Classic topics on the history of modern mathematical statistics: From Laplace to more recent times*. John Wiley & Sons, 2016.
- [Gra03] Alfred Gray. *Tubes*, volume 221. Springer Science & Business Media, 2003.
- [Gut06] Allan Gut. *Probability: a graduate course*. Springer, 2006.
- [Hay11] Fumio Hayashi. *Econometrics*. Princeton University Press, 2011.
- [HB03] Raymond Hubbard and María Jesús Bayarri. Confusion over measures of evidence ( $p$ 's) versus errors ( $\alpha$ 's) in classical statistical testing. *The American Statistician*, 57(3):171–178, 2003.
- [HCP22] Serim Hong, Carlos A Coelho, and Junyong Park. An exact and near-exact distribution approach to the Behrens–Fisher problem. *Mathematics*, 10(16):2953, 2022.
- [Hil73] Victor Hilt. Statistics and social science. *Foundations of scientific method: the nineteenth century*, pages 206–233, 1973.
- [HL51] JL Hodges and EL Lehmann. Some applications of the cramer-rao inequality. In *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, volume 2, pages 13–23. University of California Press, 1951.
- [HN64] Takeyuki Hida and Hisao Nomoto. Gaussian measure on the projective limit space of spheres. *Proceedings of the Japan Academy*, 40(5):301–304, 1964.
- [Hot39] Harold Hotelling. Tubes and spheres in  $n$ -spaces, and a class of statistical problems. *American Journal of Mathematics*, 61(2):440–460, 1939.
- [Hsu02] Elton P Hsu. *Stochastic analysis on manifolds*. American Mathematical Soc., 2002.
- [HTFF09] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- [HTW15] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity: the lasso and generalizations*. CRC press, 2015.
- [JLR11] Svante Janson, Tomasz Luczak, and Andrzej Rucinski. *Random graphs*. John Wiley & Sons, 2011.
- [JS61] William James and Charles Stein. Estimation with quadratic loss. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 361–379. University of California Press, 1961.
- [JWHT13] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.
- [KC98] Seock-Ho Kim and Allan S Cohen. On the Behrens-Fisher problem: A review. *Journal of Educational and Behavioral Statistics*, 23(4):356–377, 1998.
- [Ken46] Maurice George Kendall. *The advanced theory of statistics*. Charles Griffin and Co., 1946.
- [Kit04] Charles Kittel. *Elementary statistical physics*. Courier Corporation, 2004.
- [KK12] Gopinath Kallianpur and Rajeeva L Karandikar. *Introduction to option pricing theory*. Springer Science & Business Media, 2012.
- [Kle13] Achim Klenke. *Probability theory: a comprehensive course*. Springer Science & Business Media, 2013.
- [Kol18] Andrei Nikolaevich Kolmogorov. *Foundations of the theory of probability: Second English Edition*. Courier Dover Publications, 2018.
- [Kre11] Ulrich Krengel. *Ergodic theorems*, volume 6. Walter de Gruyter, 2011.
- [KS12] Ioannis Karatzas and Steven Shreve. *Brownian motion and stochastic calculus*, volume 113. Springer Science & Business Media, 2012.
- [Lab07] Mohammed Larbi Labbi. On gauss-bonnet curvatures. *SIGMA. Symmetry, Integrability and Geometry: Methods and Applications*, 3:118, 2007.
- [Lap98] Pierre-Simon Laplace. *Pierre-Simon Laplace philosophical essay on probabilities: translated from the fifth french edition of 1825 with notes by the translator*, volume 13. Springer Science & Business Media, 1998.
- [LC06] Erich L Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.
- [Led01] Michel Ledoux. *The concentration of measure phenomenon*, volume 89. American Mathematical Soc., 2001.
- [Led06] Michel Ledoux. Isoperimetry and gaussian analysis. *Lectures on Probability Theory and Statistics: Ecole d'Eté de Probabilités de Saint-Flour XXIV—1994*, pages 165–294, 2006.
- [Led22] Johannes Lederer. *Fundamentals of high-dimensional statistics*. Springer, 2022.
- [Leh99] Erich Leo Lehmann. *Elements of large-sample theory*. Springer, 1999.
- [LG13] Jean-Francois Le Gall. *Mouvement brownien, martingales et calcul stochastique*, volume 71. Springer, 2013.
- [Liu10] Wei Liu. *Simultaneous inference in regression*. CRC Press, 2010.
- [LR05] Erich Leo Lehmann and Joseph P Romano. *Testing statistical hypotheses*. Springer, 2005.
- [LT13] Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.
- [Luk42] Eugene Lukacs. A characterization of the normal distribution. *The Annals of Mathematical Statistics*, 13(1):91–93, 1942.



- [Luk70] Eugene Lukacs. *Characteristic functions*. Charles Griffin & Company, 1970.
- [Mac81] Donald A MacKenzie. Statistics in Britain, 1865-1930: The social construction of scientific knowledge. (*No Title*), 1981.
- [McK73] Henry P McKean. Geometry of differential space. *The Annals of Probability*, 1(2):197–206, 1973.
- [Mos85] Peter G Moschopoulos. The distribution of the sum of independent gamma random variables. *Annals of the Institute of Statistical Mathematics*, 37(1):541–544, 1985.
- [Nai90] Daniel Q Naiman. Volumes of tubular neighborhoods of spherical polyhedra and statistical inference. *The Annals of Statistics*, pages 685–716, 1990.
- [NM94] Whitney K Newey and Daniel McFadden. Large sample estimation and hypothesis testing. *Handbook of econometrics*, 4:2111–2245, 1994.
- [Oks13] Bernt Oksendal. *Stochastic differential equations: an introduction with applications*. Springer Science & Business Media, 2013.
- [ONe14] Ben O'Neill. Some useful moment results in sampling problems. *The American Statistician*, 68(4):282–296, 2014.
- [OP58] Ingram Olkin and John W Pratt. Unbiased estimation of certain correlation coefficients. *The Annals of Mathematical Statistics*, pages 201–211, 1958.
- [Paz12] Amnon Pazy. *Semigroups of linear operators and applications to partial differential equations*, volume 44. Springer Science & Business Media, 2012.
- [Pis06] Gilles Pisier. Probabilistic methods in the geometry of Banach spaces. In *Probability and Analysis: Lectures given at the 1st 1985 Session of the Centro Internazionale Matematico Estivo (CIME) held at Varenna (Como), Italy May 31–June 8, 1985*, pages 167–241. Springer, 2006.
- [RCC10] Christian P Robert, George Casella, and George Casella. *Introducing monte carlo methods with R*, volume 18. Springer, 2010.
- [RF10] Halsey Royden and Patrick Michael Fitzpatrick. *Real analysis*. China Machine Press, 2010.
- [Rob07] Christian P. Robert. *The Bayesian choice: from decision-theoretic foundations to computational implementation*, volume 2. Springer, 2007.
- [RP61] WD Ray and AENT Pitman. An exact distribution of the Fisher–Behrens–Welch statistic for testing the difference between the means of two normal populations with unknown variances. *Journal of the Royal Statistical Society: Series B (Methodological)*, 23(2):377–384, 1961.
- [RS08] Alvin C Rencher and G Bruce Schaalje. *Linear models in statistics*. John Wiley & Sons, 2008.
- [SA90] SM Sadooghi-Alvandi. Simultaneous prediction intervals for regression models with an intercept. *Communications in Statistics-Theory and Methods*, 19(4):1433–1441, 1990.
- [Saz81] Vjačeslav V Sazonov. *Normal approximation: some recent advances*, volume 879. Springer, 1981.
- [Sch47] Henry Scheffé. A useful convergence theorem for probability distributions. *The Annals of Mathematical Statistics*, 18(3):434–438, 1947.
- [Sha08] Jun Shao. *Mathematical statistics*. Springer Science & Business Media, 2008.
- [Shi16] Takashi Shioya. *Metric measure geometry - Gromov's theory of convergence and concentration of metrics and measures*. European Mathematical Society, 2016.
- [SL94] Jiayang Sun and Clive R Loader. Simultaneous confidence bands for linear regression and smoothing. *The Annals of Statistics*, pages 1328–1345, 1994.
- [SL03] George AF Seber and Alan J Lee. *Linear regression analysis*. John Wiley & Sons, 2003.
- [SP14] René L Schilling and Lothar Partzsch. *Brownian motion: an introduction to stochastic processes*. Walter de Gruyter, 2014.
- [ST78] Vladimir N Sudakov and Boris S Tsirel'son. Extremal properties of half-spaces for spherically invariant measures. *Journal of Soviet Mathematics*, 9(1):9–18, 1978.
- [Sti90] Stephen M Stigler. *The history of statistics: The measurement of uncertainty before 1900*. Harvard University Press, 1990.
- [Sti97] Stephen M Stigler. Regression towards the mean, historically considered. *Statistical methods in medical research*, 6(2):103–114, 1997.
- [Sti05] Stephen Stigler. Fisher in 1921. *Statistical Science*, pages 32–49, 2005.
- [Stu08a] Student. Probable error of a correlation coefficient. *Biometrika*, pages 302–310, 1908.
- [Stu08b] Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- [Tal96] Michel Talagrand. A new look at independence. *The Annals of probability*, pages 1–34, 1996.
- [Tao11] Terence Tao. *An introduction to measure theory*, volume 126. American Mathematical Soc., 2011.
- [Tao12] Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.
- [The74] Chris M Theobald. Generalizations of mean square error applied to ridge regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 36(1):103–106, 1974.
- [Tib96] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- [Ton90] Y L Tong. *The multivariate normal distribution*. Springer, 1990.
- [VDH24] Remco Van Der Hofstad. *Random graphs and complex networks*. Cambridge university press, 2024.
- [VdV00] A W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge University Press, 2000.
- [Ver18] Roman Vershynin. *High-dimensional probability: An introduction with applications in Data Science*, volume 47. Cambridge university press, 2018.

- [Wai19] Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- [Wel96] Alan H Welsh. *Aspects of statistical inference*, volume 246. John Wiley & Sons, 1996.
- [Wey39] Hermann Weyl. On the volume of tubes. *American Journal of Mathematics*, 61(2):461–472, 1939.
- [Wil38] Samuel S Wilks. The large-sample distribution of the likelihood ratio for testing composite hypotheses. *The annals of mathematical statistics*, 9(1):60–62, 1938.
- [Wil91] David Williams. *Probability with martingales*. Cambridge University Press, 1991.
- [WL16] Ronald L Wasserstein and Nicole A Lazar. The ASA statement on  $p$ -values: context, process, and purpose. *The American Statistician*, 70(2):129–133, 2016.

UNIVERSIDADE FEDERAL DO CEARÁ, DEPARTAMENTO DE MATEMÁTICA, CAMPUS DO PICI, R. HUMBERTO MONTE, S/N, 60455-760, FORTALEZA/CE, BRAZIL (levi@mat.ufc.br).