

# ABSTRACT

Proposal for an Adjusted Numerical Approximation to the Integrated Likelihood Function

Timothy Ruel

We propose an adapted version of a numerical approximation to the integral of a likelihood function with respect to a weight function over the space of its nuisance parameter. In doing so, we illustrate the overall utility of integrated likelihood functions as a method for eliminating nuisance parameters from a likelihood function. Chapter 1 introduces the concept of a statistical model and the various assumptions associated with it. Chapter 2 discusses the properties of the likelihood function and the maximum likelihood estimator. Chapter 3 motivates the use of pseudolikelihood functions as a tool for nuisance parameter elimination and considers the advantages and disadvantages of several types. Chapter 4 focuses on a procedure for approximating integrated likelihood functions that has been adapted from an earlier approximation technique by Severini (2007) based on the construction of a nuisance parameter that is unrelated to the parameter of interest. Chapter 5 gives examples of applications for this procedure.

## Table of Contents

ABSTRACT	1
Table of Contents	2
Chapter 1. Statistical Models	3
Chapter 2. The Likelihood Function	6
2.1. Definition	6
2.2. Transformations	7
2.3. Maximum Likelihood Estimation	10
2.4. The Bartlett Identities	14
2.5. One-Index Asymptotics	17
Chapter 3. Nuisance Parameter Elimination	24
3.1. Model Parameter Decomposition	24
3.2. Pseudolikelihood Functions	26
Chapter 4. Integrated Likelihood Analysis	31
4.1. Two-Index Asymptotics	31
4.2. The Zero-Score Expectation Parameter	36
Chapter 5. Applications	41
5.1. Multinomial Distribution	41
5.2. Standardized Mean Difference	42
References	47

	3
Appendix A. Chapter 2	49
A.1. Definitions and Notation	49
A.2. Theorems	53
Appendix B. Chapter 4	55
B.1. Desirable Properties of the Integrated Likelihood	55
B.2. Laplace's Method	56
Appendix C. Chapter 5	57

## CHAPTER 1

**Statistical Models**

The acquisition of knowledge regarding a population of interest has long been the impetus for the field of statistical inference. In all but the most basic of circumstances, limiting constraints such as time, accessibility, and cost make perfect knowledge of a population essentially impossible to obtain. It therefore becomes necessary to infer characteristics of the population based on a random and representative sample of observations drawn from it. The procedures by which these samples might be procured are themselves far from trivial, and indeed an entire branch of statistics has been dedicated to their study. However, we are primarily concerned in this paper with what occurs after the sample has been taken, and so we will generally take it for granted that a suitably representative sample of the population already exists.

Suppose  $(x_1, \dots, x_n)$  is one such sample. What information can we then glean from this sample about the population from which it has been drawn? Where is its point of central tendency located? Are its values clustered tightly around this point, or are they more diffuse? Are they distributed symmetrically or skewed to one side or the other? As the questions increase in complexity, so do the techniques required to answer them. Unfortunately the natural chaos of the real world all but guarantees there will never be an instrument capable of completely capturing the intricacies of a population whose properties we wish to infer. Hence, some amount of idealization will always be required in order to proceed.

This idealization typically comes in the form of additional assumptions that we impose on the population with the goal of sacrificing what we hope is only a small amount of accuracy in exchange for a large reduction in complexity. These assumptions are essentially never “true” in the sense that they are not a flawless representation of reality, but they may nevertheless serve as convenient approximations that are capable of producing answers with degrees of accuracy high enough to be useful in their own right. Taken as a whole, they form the basis for a statistical model.

The traditional framework for a statistical model begins by assuming that there exists an unknown probability distribution  $P$  over the population of interest that generates the data we observe from it. We choose to model this observed data as being the realized outcomes (“realizations”) of some random variable  $X$  that is distributed according to  $P$ . Let  $\mathcal{P}$  denote the set of all such distributions that we are willing to consider as candidates for the true distribution  $P$ . Out of necessity, we will proceed as though our choice of  $\mathcal{P}$  always contains  $P$  though in reality there is nothing specifically requiring it.

The next assumption we make is that  $\mathcal{P}$  is *parameterized*. That is, there exists a *parameter*  $\theta$  which indexes  $\mathcal{P}$ , acting as a label that allows us to differentiate between the distributions it contains. For a particular value of the parameter  $\theta$ , say  $\theta_1$ , we can refer to its corresponding distribution in  $\mathcal{P}$  with the notation  $P_{\theta_1}$ , and therefore  $\mathcal{P}$  itself may be written as  $\mathcal{P} = \{P_{\theta} | \theta \in \Theta\}$ .  $\Theta$  is called the *parameter space* and represents the set of all possible values  $\theta$  can take on.

We will restrict our attention in this paper to distributions that are absolutely continuous with respect to some  $\sigma$ -finite measure  $\mu$ , so that they admit a probability density function by the Radon-Nikodym Theorem. Let  $p_{\theta}$  denote the density function associated with the distribution  $P_{\theta}$ . This one-to-one correspondence between distribution and density allows us to define  $\mathcal{P}$  as  $\mathcal{P} = \{p_{\theta} | \theta \in \Theta\}$ . Going forward, we will use the notation  $p_{\theta}(x)$  and  $p(x; \theta)$  interchangeably to refer to a density function with parameter  $\theta$ . We will also simply write  $\int p_{\theta}(x)dx$  instead of  $\int p_{\theta}d\mu$  or  $\int p_{\theta}(x)d\mu(x)$  with the understanding that  $p_{\theta}$  is always defined with respect to some dominating measure  $\mu$ .

In general, a model’s parameterization is not unique, and for a given parameter  $\theta$ , we are free to choose any invertible function of  $\theta$  as a new parameter. Once we have made our choice of parameterization, we will assume that  $\theta$  does contain a singular true parameter value, which we will denote by  $\theta_0$ . The conventional interpretation of  $\theta_0$  is as a fixed but unknown constant that represents the value of the parameter corresponding to the true density function  $p_{\theta_0}$  in  $\mathcal{P}$ . Conversely,  $\theta$  represents an arbitrary parameter value that is allowed to range over all possible elements of  $\Theta$ , including  $\theta_0$ . In other words,  $\theta$  acts like a tuning dial for the population - rotate the dial and certain behaviors of the population (e.g. its location, scale, or shape) will change. Making inferences regarding  $\theta_0$  is like trying to figure out the particular value ( $\theta_0$ ) to which a population’s dial ( $\theta$ ) has been set. It is possible for the value of  $\theta_0$  itself

to change over time as well, depending on the population. In such cases, any estimate of  $\theta_0$  based on a cross-sectional sample drawn from the population is best thought of as an estimate of the true parameter value during the particular time in which the sample was collected.

The theoretical ability to identify the parameter correctly on the basis of the data we observe is another essential property that a model must satisfy. A model is considered *identifiable* if having perfect knowledge of the population would enable us to determine  $\theta_0$  with absolute certainty. This is equivalent to requiring that for some observed data  $x$  and any two parameters  $\theta_1, \theta_2 \in \Theta$ , if  $p_{\theta_1}(x) = p_{\theta_2}(x)$ , then it must follow that  $\theta_1 = \theta_2$ . A model that is not identifiable could potentially have two or more distinct parameter values that give rise to the same probability distribution. For example, suppose  $Y$  is distributed uniformly on the interval  $(0, \alpha + \beta)$ , where  $\alpha, \beta > 0$ . If we use  $\theta = (\alpha, \beta)$  as a parameter for the distribution of  $Y$ , then  $\theta$  is unidentifiable since, for instance, the case where  $\theta_1 = (0, 1)$  and  $\theta_2 = (1, 0)$  implies that  $p_{\theta_1}(y) = p_{\theta_2}(y)$  despite the fact that  $\theta_1 \neq \theta_2$ . This is clearly an undesirable property for a model to possess, and so we will consider only identifiable models in this paper as a means of avoiding it.

Finally, we must make a choice regarding the dimension of the parameter space  $\Theta$  when formulating our models. *Parametric* models are defined as having finite-dimensional parameter spaces. Any model that is not parametric is either *semi-parametric* or *nonparametric*. In this paper, we will consider only parametric models whose parameter spaces are subsets of the  $d$ -dimensional real coordinate space, i.e.,  $\Theta \subseteq \mathbb{R}^d$ , where  $d \in \mathbb{Z}^+$ .

## CHAPTER 2

## The Likelihood Function

## 2.1. Definition

Upon choosing a statistical model that we think best characterizes our population of interest, the obvious next step is to identify the true distribution in  $\mathcal{P}$  or at the very least, the one that best approximates the truth. This is equivalent to making inferences about  $\theta_0$  in the case where the model is parametric and identifiable. That is, given the particular form(s) we have chosen for the distributions in  $\mathcal{P}$ , the only unknown remaining is the value of  $\theta_0$  itself. Since this value is ultimately what controls the mechanism generating any sample of data  $\mathbf{x}_n = (x_1, \dots, x_n)$  that we might observe from the population, it stands to reason that information regarding  $\theta_0$  can be inferred from the specific values of  $x_1, \dots, x_n$  that we obtain. To make this notion more rigorous, we require some method of analyzing the joint probability of our sample as a function of our parameter  $\theta$ .

Given some observed data  $\mathbf{x}_n$ , the *likelihood function* for  $\theta$  is defined as

$$(2.1.1) \quad L(\theta) = L(\theta; \mathbf{x}_n) = p(\mathbf{x}_n; \theta), \quad \theta \in \Theta.$$

In other words, the value of the likelihood function evaluated at a particular  $\theta \in \Theta$  is simply equal to the output of the model's density function evaluated at the same inputs. However, while  $p(\mathbf{x}_n; \theta)$  is viewed primarily as a function of  $\mathbf{x}_n$  for fixed  $\theta$ , the reverse is actually true for  $L(\theta; \mathbf{x}_n)$ . Indeed, we regard the likelihood as being a function of the parameter  $\theta$  for fixed  $\mathbf{x}_n$ . The reversal of the order of the arguments  $\theta$  and  $x$  is a reflection of this difference in perspectives.

When  $X$  is discrete, we may interpret  $L(\theta; x)$  as the probability that  $X = x$  given that  $\theta$  is the true parameter value.<sup>1</sup> Crucially, this is *not* equivalent to the inverse probability that  $\theta$  is the true parameter

---

<sup>1</sup>Whichever value of  $\theta$  we choose to plug into  $L(\theta; x)$  is the value that we are currently “pretending” is the true one, regardless of whether or not it actually equals  $\theta_0$  in reality.

value given  $X = x$ . The likelihood does not directly tell us anything about the probability that  $\theta$  assumes any particular value at all. Though intuitively appealing, this interpretation constitutes a fundamental misunderstanding of what a likelihood function is, and great care must be taken to avoid it.

When  $X$  is continuous, the likelihood for  $\theta$  may still be defined as it is in Equation 2.1.1. However, we must forfeit our previous interpretation of  $L(\theta)$  as a probability since the probability that  $X$  takes on any particular value is now 0. We may however still think of the likelihood as being proportional to the probability that  $X$  takes on a value “close” to  $x$ , meaning that that  $X$  is within a tiny ball centered at  $x$ . Specifically, for two different observations  $x_1$  and  $x_2$ , if  $L(\theta; x_1) = c \cdot L(\theta; x_2)$ , where  $c > 1$ , then under this model we may conclude  $X$  is  $c$  times more likely to assume a value closer to  $x_1$  than  $x_2$  given that  $\theta$  is the true value of the parameter.

As in the discrete case, we must also be careful when  $X$  is continuous to avoid using  $L(\theta; \mathbf{x}_n)$  to make probabilistic assertions regarding  $\theta$ . Despite our use of probability in its definition, the likelihood itself is *not* a probability density function for the parameter  $\theta$  and is subject to neither the same rules nor interpretations as one.

## 2.2. Transformations

There are a few useful transformations of the likelihood function that we will define here for use in future sections. The first is the *log-likelihood function*, which is defined as the natural logarithm of the likelihood function:

$$(2.2.1) \quad \ell(\theta) = \ell(\theta; \mathbf{x}_n) = \log L(\theta; \mathbf{x}_n).$$

In practice, we will typically eschew direct analysis of the likelihood in favor of the log-likelihood due to the nice mathematical properties logarithms possess. Chief among these properties is the ability to turn products into sums (i.e.  $\log(ab) = \log(a) + \log(b)$  for  $a, b > 0$ ). Sums tend to be easier to differentiate than products, making this a particularly useful feature for likelihood functions, which are often expressed as the product of marginal density functions when the observations are independent.



The other key property of logarithms that makes the log-likelihood so useful is that they are strictly increasing functions of their arguments (i.e.  $\log x > \log y$  for  $x > y > 0$ ). This monotonicity ensures that the locations of a function's extrema are preserved when the function is passed to the argument of a logarithm. For example, for a positive function  $f$  with a global maximum,  $\operatorname{argmax}_x f(x) = \operatorname{argmax}_x \log f(x)$ .

In the general case in which  $\boldsymbol{\theta}$  is a  $d$ -dimensional vector, where  $d$  is an integer greater than 1, it follows that the first derivative of the log-likelihood with respect to  $\boldsymbol{\theta}$  will also be a  $d$ -dimensional vector, the second derivative will be a  $d \times d$  matrix, the third derivative will be a  $d \times d \times d$  array, and so forth. To emphasize the multidimensional nature of these results, we will use notation typically associated with partial derivatives involving functions of more than one variable (e.g.  $\nabla$ ,  $\mathbf{J}$ ,  $\mathbf{H}$ , etc.) along with subscripts that indicate the variable with respect to which the partial derivatives are being taken. See Appendix A for a review of this notation.

The gradient of  $\ell$  with respect to  $\boldsymbol{\theta}$  appears frequently enough in the analysis of likelihood functions that it has earned its own name - the *score function*, or just the *score*. Formally, it is defined as

$$(2.2.2) \quad \mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{x}_n) = \begin{pmatrix} \frac{\partial}{\partial \theta_1} \ell(\boldsymbol{\theta}; \mathbf{x}_n) \\ \vdots \\ \frac{\partial}{\partial \theta_d} \ell(\boldsymbol{\theta}; \mathbf{x}_n) \end{pmatrix} = \begin{pmatrix} \mathcal{S}_1(\boldsymbol{\theta}; \mathbf{x}_n) \\ \vdots \\ \mathcal{S}_d(\boldsymbol{\theta}; \mathbf{x}_n) \end{pmatrix},$$

where we think of each component as being a function  $S_j : \Theta \rightarrow \mathbb{R}$ .

Similarly, the Hessian matrix of the log-likelihood function with respect to  $\boldsymbol{\theta}$  (i.e. the transpose of the Jacobian matrix of the score) multiplied by  $-1$  is called the *observed information*, or just the *information*, and is denoted by

$$(2.2.3) \quad \mathcal{I}_{\mathbf{x}_n}(\boldsymbol{\theta}) = -\mathbf{H}_{\boldsymbol{\theta}}(\ell(\boldsymbol{\theta}; \mathbf{x}_n)) = -\mathbf{J}_{\boldsymbol{\theta}}(\mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n))^{\top} = - \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_d} \\ \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ell}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_d^2} \end{pmatrix}.$$

The use of the term “information” here derives from the fact that the second partial derivatives of  $\ell$  with respect to the components of  $\boldsymbol{\theta}$  are all related to the curvature of  $\ell$  near its maximum - the sharper the curve, the less uncertainty and therefore more information we have about  $\boldsymbol{\theta}$ .

Recall that  $L(\boldsymbol{\theta}; \mathbf{x}_n)$  is defined as a function of  $\boldsymbol{\theta}$  for a fixed sample of observations  $\mathbf{x}_n = (x_1, \dots, x_n)$ , where we think of each  $x_i$  as being a realization of a random variable  $X_i$ . We may therefore interpret  $L(\boldsymbol{\theta}; \mathbf{x}_n)$  as a random variable in the following sense: for a given  $\boldsymbol{\theta}$ , the value of  $L(\boldsymbol{\theta}; \mathbf{x}_n)$  depends entirely on the values of  $X_1, \dots, X_n$  that we happened to observe, and so  $L(\boldsymbol{\theta}; \mathbf{X}_n)$  is itself a random variable with respect to the joint probability distribution of  $\mathbf{X}_n = (X_1, \dots, X_n)$ . The same is also true for any function or estimate based on the likelihood, as they ultimately will all depend on the data through it as well. Going forward, we will use capital letters inside these functions when we want to emphasize this interpretation. For example,  $\mathcal{S}(\boldsymbol{\theta}; \mathbf{X}_n)$  is a random variable for which we have observed the value  $\mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n)$ .

The random nature of these likelihood-based quantities further implies that finding their expectations and variances with respect to  $p_{\boldsymbol{\theta}}(\mathbf{x}_n)$  is a well-defined, nontrivial task. The variance of the score function will be of particular importance, as it also relates to the amount of information pertaining to  $\boldsymbol{\theta}_0$  that is contained within the log-likelihood function of our model. Properly known as the *Fisher information* or the *expected information*, it is defined as

$$(2.2.4) \quad \mathcal{J}_{\mathbf{X}_n}(\boldsymbol{\theta}) = \text{Var}_{\boldsymbol{\theta}}[\mathcal{S}(\boldsymbol{\theta}; \mathbf{X}_n)].$$

Since we are working in the more general framework in which  $\mathcal{S}(\boldsymbol{\theta})$  is a  $d \times 1$  random vector, it would be more accurate to speak of the *Fisher information matrix*, which is equal to the variance-covariance

matrix of  $\mathcal{S}(\boldsymbol{\theta})$ . Hence, we have

$$\begin{aligned}
 \mathcal{J}_{\mathbf{X}_n}(\boldsymbol{\theta}) &= \text{Var}_{\boldsymbol{\theta}}[\mathcal{S}(\boldsymbol{\theta}; \mathbf{X}_n)] && \text{(by Eq. 2.2.4)} \\
 &= \text{Cov}_{\boldsymbol{\theta}} \left[ \left( \frac{\partial \ell}{\partial \theta_1}, \dots, \frac{\partial \ell}{\partial \theta_d} \right)^T \right] && \text{(by Eq. 2.2.2)} \\
 (2.2.5) \quad &= \begin{pmatrix} \text{Var}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_1} \right) & \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_2} \right) & \cdots & \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_d} \right) \\ \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_2}, \frac{\partial \ell}{\partial \theta_1} \right) & \text{Var}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_2} \right) & \cdots & \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_2}, \frac{\partial \ell}{\partial \theta_d} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_d}, \frac{\partial \ell}{\partial \theta_1} \right) & \text{Cov}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_d}, \frac{\partial \ell}{\partial \theta_2} \right) & \cdots & \text{Var}_{\boldsymbol{\theta}} \left( \frac{\partial \ell}{\partial \theta_d} \right) \end{pmatrix}.
 \end{aligned}$$

If the observations are independent, the Fisher information of the whole sample is equal to the sum of the Fisher information values for each of the observations individually. That is,

$$(2.2.6) \quad \mathcal{J}_{\mathbf{X}_n}(\boldsymbol{\theta}) = \sum_{i=1}^n \mathcal{J}_{X_i}(\boldsymbol{\theta}).$$

If the observations are also identically distributed according to the distribution of some random variable  $X$ , then  $\mathcal{J}_{X_i}(\boldsymbol{\theta}) = \mathcal{J}_X(\boldsymbol{\theta})$  for all  $i$ , and so the Fisher information for the entire sample is simply equal to the Fisher information for a single observation of  $X$  multiplied by a factor of  $n$ :

$$(2.2.7) \quad \mathcal{J}_{\mathbf{X}_n}(\boldsymbol{\theta}) = n \mathcal{J}_X(\boldsymbol{\theta}).$$

The above two equations hold true for the observed information under the same conditions as well.

## 2.3. Maximum Likelihood Estimation

### 2.3.1. Motivation

Maximum likelihood estimation is one of the most powerful and widespread techniques for obtaining point estimates of model parameters. The original intuition behind the method derives from the observation that when faced with a choice between two possible values of a parameter, the sensible choice is the one which makes the data we actually did observe more probable to have been observed. We have already defined the likelihood function as a means of capturing this probability, which makes expressing

this decision rule in terms of it very easy - we simply choose for our estimate the option that produces the higher value of the likelihood function. That is, if  $L(\theta_1; \mathbf{x}_n) > L(\theta_2; \mathbf{x}_n)$ , then under the preceding logic,  $\theta_1$  is the better estimate of the true parameter value.

This can be extended to include as many candidate parameter values as we would like. For  $n$  potential estimates of  $\theta_0$ , the best is the one that corresponds to the highest value of the likelihood function. Following this line of reasoning to its natural conclusion, a sensible choice for an estimate of  $\theta_0$  is any value that maximizes the likelihood function based on an observed dataset  $\mathbf{x}_n$ . Hence, let  $\hat{\theta} \in \Theta$  be any parameter value that renders the likelihood at the observed  $\mathbf{X}_n = \mathbf{x}_n$  as large as possible, i.e.,

$$(2.3.1) \quad L(\hat{\theta}; \mathbf{x}_n) = \sup_{\theta \in \Theta} L(\theta; \mathbf{x}_n).$$

We call such a value a *maximum likelihood estimate* of  $\theta_0$ . This definition of  $\hat{\theta}$  as a maximizer of  $L(\hat{\theta}; \mathbf{x}_n)$  necessarily makes it a function of the observed data. When this function is measurable, then we can further define the *maximum likelihood estimator* (MLE) of  $\theta_0$  as the statistic  $\hat{\theta}(\mathbf{X}_n)$  for which we observe the value  $\hat{\theta}(\mathbf{x}_n)$ .

### 2.3.2. Regularity Conditions

As a consequence of the random variable interpretation of likelihood-based quantities, a natural line of inquiry to investigate is how the behavior of these random variables changes as the sample size  $n$  increases. Of particular interest is the distribution to which the MLE converges, if any, as  $n$  tends toward infinity. To that end, it will be useful to establish some *regularity conditions* for our models. We can think of these conditions as being assumptions similar to those we discussed in the introduction to this paper that, when satisfied, endow our models with certain properties that enable us, among other things, to determine the aforementioned distribution.

For our purposes, we will call a model *regular* if it satisfies the following conditions:

- RC1)** Any observations  $x_1, \dots, x_n$  belonging to a sample that has been drawn from the model's sample space are independent and identically distributed (i.i.d.) realizations of a random variable  $X$  with density function  $p_{\theta}(x)$ .

- RC2)**  $P_{\theta_1} = P_{\theta_2} \implies \theta_1 = \theta_2$  for all  $\theta_1, \theta_2 \in \Theta$ .
- RC3)** The distributions in  $\mathcal{P}$  have a common support  $\mathcal{X} = \{x : p_{\theta}(x) > 0\} \subseteq \mathbb{R}$  not depending on  $\theta$ .
- RC4)** There exists an open set  $\Theta^* \subseteq \Theta$  of which  $\theta_0$  is an interior point.
- RC5)**  $p(x; \theta)$  is twice continuously differentiable with respect to  $\theta$  for all  $\theta$  in a neighborhood of  $\theta_0$ .
- RC6)** There exists a random function  $M(x)$  (that does not depend on  $\theta$ ) satisfying  $E[M(X)] < \infty$  such that each third partial derivative of  $\ell(\theta; x)$  is bounded in absolute value by  $M(x)$  uniformly in some neighborhood of  $\theta_0$ .
- RC7)** The integral  $\int_{\mathcal{X}} p(x; \theta) dx$  can be differentiated twice under the integral sign with respect to the components of  $\theta \in \Theta^*$ .
- RC8)**  $\mathcal{J}_X(\theta)$  is positive definite for all  $\theta \in \Theta$ .
- RC9)**  $\Theta$  is a compact and convex subset of  $\mathbb{R}^d$ .

While not strictly necessary, **RC1** is often assumed as a matter of convenience since it tends to simplify calculations greatly. We include it here for that purpose and to frame our discussion in the context of a standard case in which likelihood theory holds. Of course, it is possible to construct models lacking i.i.d. observations yet still possessing real world applications for which the results discussed in this paper hold.

The implication in **RC2** is simply the identifiability property we mentioned in Chapter 2. We repeat it here as it is necessary to guarantee the consistency of the MLE, i.e., that it converges in probability to  $\theta_0$  as  $n \rightarrow \infty$ .

**RC3** requires that the distributions in  $\mathcal{P}$  be supported on a common subset of the real line, and the definition of this subset cannot depend on  $\theta$ . This is to prevent situations in which, for example, the event  $\{X_i \leq x_i\}$  occurs with positive probability when  $\theta = \theta_1$  but not  $\theta = \theta_2$ .

**RC4** guarantees the existence of an open subset  $\Theta^*$  of  $\Theta$  containing  $\theta_0$  as an interior point. The fact that  $\theta_0$  is an interior point of  $\Theta^*$  further implies that it is possible to find a neighborhood of  $\theta_0$  that is contained in  $\Theta^*$ . **RC5** then goes on to assert the existence and continuity of the first two partial derivatives with respect to the components of  $\theta$  of  $p(x; \theta)$  in this neighborhood. This is a necessary requirement for defining a second-order Taylor series expansion of the score function around  $\theta_0$ .

Another way of stating **RC6** is that for all  $\boldsymbol{\theta}$  in a neighborhood  $N_{\boldsymbol{\theta}_0}$ , there exists a random function  $M(x)$  with finite expectation such that

$$\sup_{\boldsymbol{\theta} \in N_{\boldsymbol{\theta}_0}} \left| \frac{\partial^3 \ell(\boldsymbol{\theta}; x)}{\partial \theta_i \partial \theta_j \partial \theta_k} \right| \leq M(x)$$

for all integers  $1 \leq i, j, k \leq d$ . Equivalently, we could say that the entries of the Hessian matrix for each component of the score function are all bounded by  $M(x)$  as well. This ensures the remainder terms in a second-order Taylor series expansion of the score function around  $\boldsymbol{\theta}_0$  become negligible as the sample size increases to infinity.

**RC7** grants us the ability to freely interchange integration and second-order partial differentiation with respect to the components of  $\boldsymbol{\theta}$ , i.e.,

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} \int_{\mathcal{X}} p(x; \boldsymbol{\theta}) dx = \int_{\mathcal{X}} \frac{\partial^2}{\partial \theta_i \partial \theta_j} p(x; \boldsymbol{\theta}) dx$$

for all  $\boldsymbol{\theta} \in \Theta^*$  and  $i, j = 1, \dots, d$ . This implies first-order partial derivatives can be passed under the integral sign as well. This will prove useful in our discussion of the Bartlett identities in Section 3.5.

Finally, **RC8-9** play an important role in ensuring the existence and uniqueness of the maximum likelihood estimator of  $\boldsymbol{\theta}_0$ .

### 2.3.3. Properties

In general, there is no guarantee that an MLE for a model's parameter will exist, and even if it does, it will not necessarily be unique. However, since maximum likelihood estimation is critical to our discussion in Chapters 5 and 6, it would come as a great convenience if we possessed the ability to speak freely of *the* MLE of a model's parameter without having to clarify *which* MLE we mean or whether one even exists at all. Hence, some discussion of the conditions under which the MLE of a model's parameter exists and is unique is warranted.

The extreme value theorem (see Appendix A) implies that sufficient conditions for the existence of a model's MLE are that  $\Theta$  is compact, and  $\ell(\boldsymbol{\theta}; x)$  is continuous on  $\Theta$ . The former is satisfied directly

by **RC9** and the latter is implied through our assumption in **RC5** of the differentiability of  $p(x; \boldsymbol{\theta})$  in  $\boldsymbol{\theta}$ . Therefore, at least one MLE will always exist for the true parameter value of a regular model. These are only sufficient conditions, however, and not necessary; MLEs may exist for parameters of non-regular models as well.

Similarly, when the MLE does exist, a sufficient condition for its uniqueness is that  $\Theta$  is convex, and  $\ell(\boldsymbol{\theta}; x)$  is strictly concave on  $\Theta$ , as this ensures that it has exactly one global maximum, which it attains at the  $\hat{\boldsymbol{\theta}}$ . **RC9** directly satisfies the compactness criterion. Our assumption in **RC8** that the Fisher information matrix is positive definite forces the Hessian matrix of  $\ell(\boldsymbol{\theta}; x)$  to be negative definite. This in turn implies that  $\ell(\boldsymbol{\theta}; x)$  is strictly concave, so the requirement is met. Hence, a regular model will always have a unique MLE for its parameter.

One last useful property of the maximum likelihood estimator is its functional invariance. If  $\hat{\boldsymbol{\theta}}$  is an MLE of  $\boldsymbol{\theta}_0$ , then any function  $h(\boldsymbol{\theta}_0)$  will have  $h(\hat{\boldsymbol{\theta}})$  as its MLE. Hence, it is straightforward to find the MLE of a model that has undergone a reparameterization given that we know the MLE of the original parameter.

## 2.4. The Bartlett Identities

The Bartlett identities are a set of equations relating to the expectations of the derivatives of a log-likelihood function to one another. In general, there is no guarantee that an arbitrary function of a random variable  $X$  and its parameter  $\boldsymbol{\theta}$  will satisfy the Bartlett identities. It is guaranteed, however, that the log-likelihood function associated with  $X$  and  $\boldsymbol{\theta}$  will satisfy them, provided that the model is regular. This means we can think of any function that does satisfy the Bartlett identities (or at least some of them) as resembling that of a genuine log-likelihood.

Consider the case where a random variable  $X$  has density function  $p_\theta(x)$ , where  $\theta$  is a scalar. For a single observation  $X = x$ , the expectation of  $\frac{\partial}{\partial \theta} \ell(\theta; X)$  gives

$$\begin{aligned}
(2.4.1) \quad E_{\theta} \left[ \frac{\partial}{\partial \theta} \ell(\theta; X) \right] &= \int_{\mathbb{R}} \left[ \frac{\partial}{\partial \theta} \log p(x; \theta) \right] p(x; \theta) dx \\
&= \int_{\mathbb{R}} \frac{\frac{\partial}{\partial \theta} p(x; \theta)}{p(x; \theta)} p(x; \theta) dx \\
&= \int_{\mathbb{R}} \frac{\partial}{\partial \theta} p(x; \theta) dx \\
&= \frac{d}{d\theta} \int_{\mathbb{R}} p(x; \theta) dx \\
&= \frac{d}{d\theta} 1 \\
&= 0.
\end{aligned}$$

Equation 2.4.1 is called the first Bartlett identity. In words, it states that the expectation of the first partial derivative of the log-likelihood function of a statistical model with respect to the parameter will always be 0. Since the score is defined as  $\frac{\partial}{\partial \theta} \ell(\theta; x)$ , any function that satisfies the first Bartlett identity is said to be *score-unbiased*.

For any model with a log-likelihood satisfying the first Bartlett identity, the expected information for its parameter  $\theta$  may be rewritten as

$$\begin{aligned}
(2.4.2) \quad \mathcal{I}_X(\theta) &= \text{Var}_{\theta}[\mathcal{S}(\theta; X)] \\
&= \text{Var}_{\theta} \left[ \frac{\partial}{\partial \theta} \ell(\theta; X) \right] \\
&= \text{Var}_{\theta} \left[ \frac{\partial}{\partial \theta} \ell(\theta; X) \right] + \left( E_{\theta} \left[ \frac{\partial}{\partial \theta} \ell(\theta; X) \right] \right)^2 \quad (\text{by the first Bartlett identity}) \\
&= E_{\theta} \left[ \left( \frac{\partial}{\partial \theta} \ell(\theta; X) \right)^2 \right].
\end{aligned}$$



If we now consider the second partial derivative of  $\ell(\theta; x)$  with respect to  $\theta$ , we have

$$\begin{aligned}
\frac{\partial^2}{\partial \theta^2} \ell(\theta; x) &= \frac{\partial}{\partial \theta} \left[ \frac{\partial}{\partial \theta} \ell(\theta; x) \right] \\
&= \frac{\partial}{\partial \theta} \left[ \frac{\partial}{\partial \theta} \log p(x; \theta) \right] \\
&= \frac{\partial}{\partial \theta} \left[ \frac{\frac{\partial}{\partial \theta} p(x; \theta)}{p(x; \theta)} \right] \\
&= \frac{\left[ \frac{\partial^2}{\partial \theta^2} p(x; \theta) \right] p(x; \theta) - \left[ \frac{\partial}{\partial \theta} p(x; \theta) \right] \left[ \frac{\partial}{\partial \theta} p(x; \theta) \right]}{[p(x; \theta)]^2} \\
&= \frac{\frac{\partial^2}{\partial \theta^2} p(x; \theta)}{p(x; \theta)} - \left[ \frac{\frac{\partial}{\partial \theta} p(x; \theta)}{p(x; \theta)} \right]^2 \\
&= \frac{\frac{\partial^2}{\partial \theta^2} p(x; \theta)}{p(x; \theta)} - \left[ \frac{\partial}{\partial \theta} \log p(x; \theta) \right]^2 \\
&= \frac{\frac{\partial^2}{\partial \theta^2} p(x; \theta)}{p(x; \theta)} - \left[ \frac{\partial}{\partial \theta} \ell(\theta; x) \right]^2.
\end{aligned}$$

Rearranging terms and taking expectations yields

$$\begin{aligned}
\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta^2} \ell(\theta; X) \right] + \mathbb{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \ell(\theta; X) \right)^2 \right] &= \mathbb{E}_\theta \left[ \frac{\frac{\partial^2}{\partial \theta^2} p(X; \theta)}{p(X; \theta)} \right] \\
&= \int_{\mathbb{R}} \left[ \frac{\frac{\partial^2}{\partial \theta^2} p(x; \theta)}{p(x; \theta)} \right] p(x; \theta) dx \\
&= \int_{\mathbb{R}} \frac{\partial^2}{\partial \theta^2} p(x; \theta) dx \\
&= \frac{d^2}{d\theta^2} \int_{\mathbb{R}} p(x; \theta) dx \\
&= \frac{d^2}{d\theta^2} 1 \\
&= 0.
\end{aligned}$$

Therefore,

$$(2.4.3) \quad \mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta^2} \ell(\theta; X) \right] + \mathbb{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \ell(\theta; X) \right)^2 \right] = 0.$$

Equation 2.4.2 is called the second Bartlett identity. Any function that satisfies it is said to be *information-unbiased*. Regular models as we have defined them will automatically satisfy both the first and second Bartlett identities. Hence, for any regular model, the statements in Equation 2.2.4, Equation 2.4.2, and Equation 2.4.3 imply that the following definitions for its expected information regarding its parameter  $\theta$  are all equivalent:

$$(2.4.4) \quad \mathcal{I}_X(\theta) = \text{Var}_\theta \left[ \frac{\partial}{\partial \theta} \ell(\theta; X) \right] = \text{E}_\theta \left[ \left( \frac{\partial}{\partial \theta} \ell(\theta; X) \right)^2 \right] = \text{E}_\theta \left[ - \frac{\partial^2}{\partial \theta^2} \ell(\theta; X) \right] = \text{E}_\theta [\mathcal{J}_X(\theta)].$$

It is possible to derive further Bartlett identities by continuing in this manner for an arbitrary number of partial  $\theta$ -derivatives of the log-likelihood function, provided that they exist. However, the first two are sufficient for our purposes of evaluating the validity of approximations to genuine likelihoods so we will not go further here. While the above derivations were performed under the assumption that  $\theta$  is a scalar, the Bartlett identities also hold in the case where  $\theta$  is a  $d \times 1$  vector.

## 2.5. One-Index Asymptotics

The one-index asymptotics framework describes the behavior of likelihood-based statistics as the sample size ( $n$ ) grows to infinity. The aim of this section is to present an overview of some of the theory's classic results. This provides us with a readily available baseline against which to compare the results of the following section discussing the two-index asymptotics framework, in which the full model parameter is partitioned into a parameter of interest and a nuisance parameter, and the dimension of the nuisance parameter is allowed to increase with the sample size.

Let  $\theta$  be the parameter of a regular model with true value  $\theta_0$ . As a partial justification for our use of the MLE in estimating  $\theta_0$ , we will start by showing that with probability tending to 1 as  $n$  tends toward infinity, the likelihood for  $\theta$  is strictly larger at  $\theta_0$  than for any other  $\theta \in \Theta$ . **RC1** implies

$$(2.5.1) \quad L(\theta; \mathbf{x}_n) = p(\mathbf{x}_n; \theta) = \prod_{i=1}^n p(x_i; \theta)$$

and

$$(2.5.2) \quad \ell(\boldsymbol{\theta}; \mathbf{x}_n) = \log p(\mathbf{x}_n; \boldsymbol{\theta}) = \sum_{i=1}^n \log p(x_i; \boldsymbol{\theta}).$$

It follows that

$$(2.5.3) \quad \begin{aligned} L(\boldsymbol{\theta}; \mathbf{x}_n) < L(\boldsymbol{\theta}_0; \mathbf{x}_n) &\iff \ell(\boldsymbol{\theta}; \mathbf{x}_n) < \ell(\boldsymbol{\theta}_0; \mathbf{x}_n) \\ &\iff \sum_{i=1}^n \log p(x_i; \boldsymbol{\theta}) - \sum_{i=1}^n \log p(x_i; \boldsymbol{\theta}_0) < 0 \\ &\iff \sum_{i=1}^n [\log p(x_i; \boldsymbol{\theta}) - \log p(x_i; \boldsymbol{\theta}_0)] < 0 \\ &\iff \sum_{i=1}^n \log \frac{p(x_i; \boldsymbol{\theta})}{p(x_i; \boldsymbol{\theta}_0)} < 0 \\ &\iff \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i; \boldsymbol{\theta})}{p(x_i; \boldsymbol{\theta}_0)} < 0. \end{aligned}$$

**RC3** guarantees that the ratio  $p(x; \boldsymbol{\theta})/p(x; \boldsymbol{\theta}_0)$  is well-defined and finite for all  $x \in \mathcal{X}$ , the region of common support. Then by the Weak Law of Large Numbers,

$$(2.5.4) \quad \frac{1}{n} \sum_{i=1}^n \log \frac{p(X_i; \boldsymbol{\theta})}{p(X_i; \boldsymbol{\theta}_0)} \rightarrow \mathbb{E}_{\boldsymbol{\theta}_0} \left[ \log \frac{p(X; \boldsymbol{\theta})}{p(X; \boldsymbol{\theta}_0)} \right]$$

in probability as  $n \rightarrow \infty$ . Furthermore,

$$(2.5.5) \quad \mathbb{E}_{\boldsymbol{\theta}_0} \left[ \frac{p(X; \boldsymbol{\theta})}{p(X; \boldsymbol{\theta}_0)} \right] = \int_{\mathcal{X}} \left[ \frac{p(x; \boldsymbol{\theta})}{p(x; \boldsymbol{\theta}_0)} \right] p(x; \boldsymbol{\theta}_0) dx = \int_{\mathcal{X}} p(x; \boldsymbol{\theta}) dx = 1.$$

Since  $\log(x)$  is a strictly concave function, it follows from Jensen's inequality (see Appendix A) and Equation 2.5.5

$$(2.5.6) \quad \mathbb{E}_{\boldsymbol{\theta}_0} \left[ \log \frac{p(X; \boldsymbol{\theta})}{p(X; \boldsymbol{\theta}_0)} \right] < \log \mathbb{E}_{\boldsymbol{\theta}_0} \left[ \frac{p(X; \boldsymbol{\theta})}{p(X; \boldsymbol{\theta}_0)} \right] = \log 1 = 0.$$

Hence, the quantity on the left-hand side of Equation 2.5.4 is converging in probability to a constant that is less than 0 as  $n$  tends to infinity. From this and the equivalence we established in Equation 2.5.3, it

follows that

$$(2.5.7) \quad \lim_{n \rightarrow \infty} P_{\theta_0} [L(\boldsymbol{\theta}; \mathbf{x}_n) < L(\boldsymbol{\theta}_0; \mathbf{x}_n)] = \lim_{n \rightarrow \infty} P_{\theta_0} \left[ \frac{1}{n} \sum_{i=1}^n \log \frac{p(X_i; \boldsymbol{\theta})}{p(X_i; \boldsymbol{\theta}_0)} < 0 \right] = 1,$$

which proves the claim.

As a global maximizer of the log-likelihood function, the MLE  $\hat{\boldsymbol{\theta}}$  of a regular model must be a root of the log-likelihood function, i.e., it must satisfy the *likelihood equation*,

$$(2.5.8) \quad \nabla_{\boldsymbol{\theta}} \ell(\hat{\boldsymbol{\theta}}) = \mathbf{0}.$$

whenever it exists. For an arbitrary model, there may be other roots as well, even when the MLE doesn't exist. Assuming **RC2** and **RC5-8**, it can be shown that there will always be at least one sequence of roots  $\hat{\boldsymbol{\theta}}_n$  of its log-likelihood such that  $\hat{\boldsymbol{\theta}}_n$  tends to  $\boldsymbol{\theta}_0$  in probability as  $n \rightarrow \infty$  (Cf. Cramér 1945). The MLE will not necessarily be a part of this sequence though, even if it exists. Adding **RC9** is enough to ensure the MLE must be the unique solution to the likelihood equation, however, and therefore this sequence of roots will also be unique and for a given sample  $\mathbf{x}_n$ , the corresponding root  $\hat{\boldsymbol{\theta}}_n = \hat{\boldsymbol{\theta}}(\mathbf{x}_n)$  will be the unique MLE of  $\boldsymbol{\theta}_0$ . It follows that the MLE is a consistent estimator of  $\boldsymbol{\theta}_0$  for regular models.

We now turn our attention to the asymptotic distribution of the MLE. Similarly to the log-likelihood function, **RC1** implies the score function is equal to

$$(2.5.9) \quad \begin{aligned} \mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n) &= \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{x}_n) \\ &= \nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \ell(\boldsymbol{\theta}; x_i) \\ &= \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; x_i) \\ &= \sum_{i=1}^n \mathcal{S}(\boldsymbol{\theta}; x_i). \end{aligned}$$

where the last equality is true by . In other words, the score function for the parameter  $\boldsymbol{\theta}$  based on data  $x_1, \dots, x_n$  can be written as the sum of independent contributions  $\mathcal{S}(\boldsymbol{\theta}; x_i)$  ( $i = 1, \dots, n$ ) where each  $\mathcal{S}(\boldsymbol{\theta}; x_i)$  can be thought of as the score function for  $\boldsymbol{\theta}$  based only on observation  $x_i$ . This implies that

a Taylor series expansion of  $\mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n)$  will be equal to the sum of the Taylor series expansions of its individual contributions, plus a remainder term that depends on  $n$ . Since the observations are identically distributed, it suffices to consider the expansion for an arbitrary contribution,  $\mathcal{S}(\boldsymbol{\theta}; x_i)$ .

**RC4-5** guarantee the existence of a neighborhood of  $\boldsymbol{\theta}_0$  on which the first two partial derivatives of  $\mathcal{S}(\boldsymbol{\theta}; x_i)$  with respect to  $\boldsymbol{\theta}$  exist and are continuous. Without loss of generality, we may assume this neighborhood, call it  $N_{\boldsymbol{\theta}_0}$ , is convex so that it contains all of the line segments connecting any two of its points. In particular, for any  $\boldsymbol{\theta} \in N_{\boldsymbol{\theta}_0}$ ,  $\text{LS}(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \subset N_{\boldsymbol{\theta}_0}$ . Then by Taylor's theorem with the Lagrange form of the remainder, there exists  $\tilde{\boldsymbol{\theta}}_j \in \text{LS}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)$  such that the  $j$ -th component of  $\mathcal{S}(\boldsymbol{\theta}; x_i)$  may be expanded as

$$\begin{aligned} \mathcal{S}_j(\boldsymbol{\theta}; x_i) &= \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H}_{\boldsymbol{\theta}}(\mathcal{S}_j(\tilde{\boldsymbol{\theta}}_j; x_i))(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \left[ \nabla_{\boldsymbol{\theta}} \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + \frac{1}{2} \mathbf{H}_{\boldsymbol{\theta}}(\mathcal{S}_j(\tilde{\boldsymbol{\theta}}_j; x_i)) \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \left[ \nabla_{\boldsymbol{\theta}} \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \right] \quad (\text{by RC6}) \\ &= \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + \left[ \nabla_{\boldsymbol{\theta}}^\top \mathcal{S}_j(\boldsymbol{\theta}_0; x_i) + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0). \end{aligned}$$

When we stack each of these individual equations into a system of equations, we get

$$(2.5.10) \quad \begin{pmatrix} \mathcal{S}_1(\boldsymbol{\theta}; x_i) \\ \vdots \\ \mathcal{S}_d(\boldsymbol{\theta}; x_i) \end{pmatrix} = \begin{pmatrix} \mathcal{S}_1(\boldsymbol{\theta}_0; x_i) \\ \vdots \\ \mathcal{S}_d(\boldsymbol{\theta}_0; x_i) \end{pmatrix} + \left[ \begin{pmatrix} \nabla^\top \mathcal{S}_1(\boldsymbol{\theta}_0; x_i) \\ \vdots \\ \nabla^\top \mathcal{S}_d(\boldsymbol{\theta}_0; x_i) \end{pmatrix} + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top.$$

The matrix of gradient vectors in the second term on the right-hand side of the above equation is simply the Jacobian of the score function evaluated at  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ , i.e.,  $\mathbf{J}_{\boldsymbol{\theta}}(\mathcal{S}(\boldsymbol{\theta}_0; x_i))$ . However, by Equation 2.2.3, this is just the negative transpose of the observed information matrix also evaluated at  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ ,  $\mathcal{J}(\boldsymbol{\theta}_0)$ . Furthermore, since we have assumed  $\ell$  is continuous in  $\boldsymbol{\theta}$ , we can freely swap the order of differentiation in all of its second partial derivatives with respect to  $\boldsymbol{\theta}$ . This implies the Jacobian of the score will be a symmetric matrix, and so we simply have  $\mathbf{J}_{\boldsymbol{\theta}}(\mathcal{S}(\boldsymbol{\theta}_0; x_i)) = -\mathcal{J}_i(\boldsymbol{\theta}_0)$ . Hence, using more

compact notation, Equation 2.5.10 becomes

$$\begin{aligned}
 \mathcal{S}(\boldsymbol{\theta}; x_i) &= \mathcal{S}(\boldsymbol{\theta}_0; x_i) + \left[ \mathbf{J}_{\boldsymbol{\theta}} \left( \mathcal{S}(\boldsymbol{\theta}_0; x_i) \right) + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \\
 (2.5.11) \quad &= \mathcal{S}(\boldsymbol{\theta}_0; x_i) - \left[ \mathcal{J}_i(\boldsymbol{\theta}_0) + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top
 \end{aligned}$$

The above represents the second-order Taylor expansion around  $\boldsymbol{\theta}_0$  for an individual observation  $x_i$ 's contribution to the score function. Summing over all of these contributions yields

$$\begin{aligned}
 \mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n) &= \sum_{i=1}^n \mathcal{S}(\boldsymbol{\theta}; x_i) \\
 (2.5.12) \quad &= \sum_{i=1}^n \left[ \mathcal{S}(\boldsymbol{\theta}_0; x_i) - \left[ \mathcal{J}_{X_i}(\boldsymbol{\theta}_0) + M(x_i) O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \right] \\
 &= \mathcal{S}(\boldsymbol{\theta}_0; \mathbf{x}_n) - \left[ \mathcal{J}_{\mathbf{x}_n}(\boldsymbol{\theta}_0) + \left\{ \sum_{i=1}^n M(x_i) \right\} O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \\
 &= \mathcal{S}(\boldsymbol{\theta}_0; \mathbf{x}_n) - \left[ \frac{1}{n} \mathcal{J}_{\mathbf{x}_n}(\boldsymbol{\theta}_0) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] n(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top.
 \end{aligned}$$

If we divide through by  $\sqrt{n}$ , we arrive at

$$(2.5.13) \quad \frac{1}{\sqrt{n}} \mathcal{S}(\boldsymbol{\theta}; \mathbf{x}_n) = \frac{1}{\sqrt{n}} \mathcal{S}(\boldsymbol{\theta}_0; \mathbf{x}_n) - \left[ \frac{1}{n} \mathcal{J}_{\mathbf{x}_n}(\boldsymbol{\theta}_0) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] \sqrt{n}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top.$$

We previously established that regular models will always have a sequence of MLEs  $\hat{\boldsymbol{\theta}}_n$  that converge in probability to  $\boldsymbol{\theta}_0$  as  $n \rightarrow \infty$ , and that each  $\hat{\boldsymbol{\theta}}_n$  in this sequence will satisfy  $\mathcal{S}(\hat{\boldsymbol{\theta}}_n; \mathbf{x}_n) = \mathbf{0}$ . Plugging  $\hat{\boldsymbol{\theta}}_n$  in for  $\boldsymbol{\theta}$  in Equation 2.5.13 gives

$$\frac{1}{\sqrt{n}} \mathcal{S}(\hat{\boldsymbol{\theta}}_n; \mathbf{x}_n) = \frac{1}{\sqrt{n}} \mathcal{S}(\boldsymbol{\theta}_0; \mathbf{x}_n) - \left[ \frac{1}{n} \mathcal{J}_{\mathbf{x}_n}(\boldsymbol{\theta}_0) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O(\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^\top$$

and therefore,

$$(2.5.14) \quad \frac{1}{\sqrt{n}} \mathcal{S}(\boldsymbol{\theta}_0; \mathbf{x}_n) = \left[ \frac{1}{n} \mathcal{J}_{\mathbf{x}_n}(\boldsymbol{\theta}_0) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O_p(\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|) \mathbf{1}_{d \times d} \right] \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^\top.$$

For the terms in the square brackets in the line above, we have the following observations:

- 1) By Equation 3.5.4,  $E_{\theta_0}[\mathcal{J}_X(\theta_0)] = \mathcal{J}_X(\theta_0)$ , and thus  $\frac{1}{n}\mathcal{J}_{\mathbf{X}_n}(\theta_0) = \frac{1}{n}\sum_{i=1}^n \mathcal{J}_{X_i}(\theta_0)$  is converging in probability to  $\mathcal{J}_X(\theta_0)$  as  $n \rightarrow \infty$  by the Weak Law of Large Numbers, i.e.,  $\mathcal{J}_X(\theta_0) = \frac{1}{n}\mathcal{J}_{\mathbf{X}_n}(\theta_0) + o_p(1)$ .
- 2)  $M(x_i)$  has finite expectation and does not depend on  $\theta$  by **RC6**, which implies through Markov's inequality that it is bounded in probability as  $n \rightarrow \infty$ , i.e.,  $M(x_i) = O_p(1)$ . It follows that  $\frac{1}{n}\sum_{i=1}^n M(x_i) = O_p(1)$  as well.
- 3) The fact that  $\hat{\theta}_n$  is converging in probability to  $\theta_0$  as  $n \rightarrow \infty$  implies that  $\|\hat{\theta}_n - \theta_0\|$  is  $o_p(1)$ .

From these three facts we can conclude that the entire term inside the square brackets is converging in probability to  $\mathcal{J}_X(\theta_0)$  as  $n \rightarrow \infty$ . This allows us to rewrite Equation 2.5.14 as

$$(2.5.15) \quad \frac{1}{\sqrt{n}}\mathcal{S}(\theta_0; \mathbf{x}_n) = \left[ \mathcal{J}_X(\theta_0) + o_p(1) \right] \sqrt{n}(\hat{\theta}_n - \theta_0)^\top.$$

This is useful because if we know the distribution to which the term on the left-hand side is converging as  $n \rightarrow \infty$ , we can deduce the asymptotic distribution of  $\hat{\theta}_n$  using Slutsky's theorem.

By definition,  $\text{Var}_{\theta_0}[\mathcal{S}(\theta_0; x_i)] = \mathcal{J}_X(\theta_0)$ . From Equation 2.5.9, we have that  $\frac{1}{n}\mathcal{S}(\theta_0; \mathbf{x}_n) = \frac{1}{n}\sum_{i=1}^n \mathcal{S}(\theta_0; x_i)$ . It follows from the Central Limit Theorem that

$$(2.5.16) \quad \sqrt{n}\left(\frac{1}{n}\mathcal{S}(\theta_0; \mathbf{x}_n) - E_{\theta_0}[\mathcal{S}(\theta_0; x_i)]\right) \xrightarrow{d} N(\mathbf{0}, \mathcal{J}_X(\theta_0)) \text{ as } n \rightarrow \infty.$$

But by the first Bartlett identity,  $E_{\theta_0}[\mathcal{S}(\theta_0; x_i)] = \mathbf{0}$ , and therefore

$$(2.5.17) \quad \frac{1}{\sqrt{n}}\mathcal{S}(\theta_0; \mathbf{x}_n) \xrightarrow{d} N(\mathbf{0}, \mathcal{J}_X(\theta_0)) \text{ as } n \rightarrow \infty.$$

Combining the results of Equation 2.5.15 and Equation 2.5.17, we see that

$$(2.5.18) \quad \left[ \mathcal{J}_X(\theta_0) + o_p(1) \right] \sqrt{n}(\hat{\theta}_n - \theta_0)^\top \xrightarrow{d} N(\mathbf{0}, \mathcal{J}_X(\theta_0)) \text{ as } n \rightarrow \infty.$$

Since the term in square brackets is converging in probability to  $\mathcal{J}(\theta_0)$ , which by **RC8** is a positive definite matrix, a minor extension of Slutsky's Theorem (see Appendix A) allows us to deduce that the

asymptotic distribution of the MLE for the true parameter value of a regular model is given by

$$\begin{aligned}
 (2.5.19) \quad \sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)^\top &\stackrel{d}{\rightarrow} \mathcal{I}_X(\boldsymbol{\theta}_0)^{-1} \mathbf{N}\left(\mathbf{0}, \mathcal{I}_X(\boldsymbol{\theta}_0)\right) \\
 &\stackrel{d}{=} \mathbf{N}\left(\mathbf{0}, \mathcal{I}_X(\boldsymbol{\theta}_0)^{-1}\right) \text{ as } n \rightarrow \infty.
 \end{aligned}$$



## CHAPTER 3

## Nuisance Parameter Elimination

## 3.1. Model Parameter Decomposition

It is often the case that we are not interested in estimating the full parameter  $\theta \in \Theta \subseteq \mathbb{R}^d$ , but rather a different parameter  $\psi$  taking values in a set  $\Psi \subseteq \mathbb{R}^p$ , where  $p < d$ . In such an event, we refer to  $\psi$  as the *parameter of interest*. Crucially, as we will see,  $\psi$  can always be expressed as a function of  $\theta$ .

Since  $\psi$  is of lower dimension than  $\theta$ , it necessarily follows that there is another parameter  $\lambda$ , taking values in a set  $\Lambda \subseteq \mathbb{R}^q$ , where  $p + q = d$ , that is made up of whatever is “left over” from the full parameter  $\theta$ . We refer to  $\lambda$  as the *nuisance parameter*, so named for its ability to complicate inference regarding the parameter of interest. Despite not being the object of study themselves, nuisance parameters are nevertheless capable of modifying the distributions of our observations and therefore must be accounted for when conducting inference or estimation regarding the parameter of interest.<sup>1</sup> The process by which this is accomplished is nontrivial and often represents a serious obstacle that must be overcome.

While not required, we will assume the parameter of interest  $\psi$  is always one-dimensional. That is,  $\Psi \subseteq \mathbb{R}$  and consequently  $\Lambda \subseteq \mathbb{R}^{d-1}$ . This restriction reflects the common habit of researchers to focus on scalar-valued summaries of vector quantities. For example, suppose we observe data  $Y = (y_1, \dots, y_n)$ , where each  $y_i$  is the outcome of some random variable  $Y_i \sim N(\mu_i, \sigma_i^2)$ , and we are interested in estimating the average of the population means,  $\frac{1}{n} \sum_{i=1}^n \mu_i$ . Rather than defining  $\psi = (\mu_1, \dots, \mu_n)$ , we can instead define  $\psi = \frac{1}{n} \sum_{i=1}^n \mu_i$  directly, bypassing the need to estimate each  $\mu_i$  individually before taking their average. This does carry the trade-off of increasing the dimension of the nuisance parameter, which must be dealt with before conducting inference or estimation on  $\psi_0$ , the true value of  $\psi$ . We will examine some of the issues posed by high-dimensional nuisance parameters in greater detail in the next chapter.

---

<sup>1</sup>Nuisance parameters are not always uniquely defined. In fact, depending on the choice of parameter of interest, there may be multiple or even infinite ways to define one.

### 3.1.1. Explicit Parameters

Parameters of interest and nuisance parameters can be broadly classified into two categories, explicit or implicit. For a given statistical model, both types of parameter must occupy the same category - it is not possible for  $\psi$  to be explicit and  $\lambda$  to be implicit, or vice versa.

Let us first consider the case in which  $\psi$  and  $\lambda$  are *explicit* parameters. This means that  $\psi$  is a sub-vector of  $\theta$ , so that all the components of  $\psi$  are also components of  $\theta$ . Then there exists a set  $I = \{I_1, \dots, I_p\} \subsetneq \{1, \dots, d\}$  such that

$$(3.1.1) \quad \psi = (\theta_{I_1}, \dots, \theta_{I_p}).$$

It immediately follows that  $\lambda$  is the sub-vector of all components of  $\theta$  that are not part of  $\psi$ . More precisely, if we let  $J = \{J_1, \dots, J_q\} \subsetneq \{1, \dots, d\}$  such that  $I \cup J = \{1, \dots, d\}$  and  $I \cap J = \emptyset$ , then

$$(3.1.2) \quad \lambda = (\theta_{J_1}, \dots, \theta_{J_q}).$$

$\theta$  can therefore be decomposed as  $\theta = (\psi, \lambda)$  when  $\psi$  and  $\lambda$  are explicit, provided we shuffle the indices appropriately.

### 3.1.2. Implicit Parameters

Now let us consider the case in which  $\psi$  and  $\lambda$  are *implicit* parameters. This means there exists some function  $\varphi : \theta \rightarrow \Psi$  for which the parameter of interest can be written as

$$(3.1.3) \quad \psi = \varphi(\theta).$$

As before,  $\Psi$  is still assumed to be a subset of  $\mathbb{R}^p$  where  $p$  is less than  $d$ . This reduction in dimension again implies the existence of a nuisance parameter  $\lambda \in \Lambda \subseteq \mathbb{R}^{k-m}$ . However, unlike in the explicit case, a closed form expression for  $\lambda$  in terms of the original components of  $\theta$  need not exist. For this reason, implicit nuisance parameters are in general more difficult to eliminate compared to their explicit counterparts.

When the parameter of interest and nuisance parameter are explicit, it is always possible to define a function  $\varphi$  such that

$$(3.1.4) \quad \varphi(\boldsymbol{\theta}) = (\theta_{I_1}, \dots, \theta_{I_p}) = \psi,$$

where  $\{I_1, \dots, I_p\}$  is defined as above. Hence, the first case is really just a special example of this more general one in which  $\psi = \varphi(\boldsymbol{\theta})$ . With this understanding in mind, we will use the notation  $\psi = \varphi(\boldsymbol{\theta})$  to refer to the parameter of interest in general, only making the distinction between implicitness and explicitness when the difference is relevant to the situation.

## 3.2. Pseudolikelihood Functions

### 3.2.1. Definition

The natural solution to the hindrance nuisance parameters pose to making inferences on the parameter of interest is to find a method for eliminating them from the likelihood function altogether. The result of this elimination is what is known as a pseudolikelihood function.

In general, a *pseudolikelihood function* for  $\psi$  is a function of the data and  $\psi$  only, having properties resembling that of a genuine likelihood function. Suppose  $\psi = \varphi(\boldsymbol{\theta})$  for some function  $\varphi$  and parameter  $\boldsymbol{\theta} \in \Theta$ . If we let  $\Theta(\psi) = \{\boldsymbol{\theta} \in \Theta : \varphi(\boldsymbol{\theta}) = \psi\}$ , then associated with each  $\psi \in \Psi$  is the set of likelihoods  $\mathcal{L}_\psi = \{L(\boldsymbol{\theta}) : \boldsymbol{\theta} \in \Theta(\psi)\}$ .

Any summary of the values in  $\mathcal{L}_\psi$  that does not depend on  $\lambda$  theoretically constitutes a pseudolikelihood function for  $\psi$ . There exist a variety of methods to obtain this summary but among the most popular are profiling (maximization), conditioning, and integration, each with respect to the nuisance parameter. None of these summaries come without a cost though, meaning some information about  $\psi_0$  is almost certainly sacrificed whenever a nuisance parameter is eliminated from a likelihood. One measure of a good pseudolikelihood, therefore, is how well it is able to retain information about  $\psi_0$  without becoming overly complex in its computation.

In the previous chapter, we introduced the Bartlett identities as being a set of equations relating the derivatives of the log-likelihood to one another. They can also be used to understand the difference

between likelihood and pseudolikelihood functions. A genuine likelihood function of  $\theta$  can be characterized as any nonnegative random function of  $\theta$  for which all of the Bartlett identities hold. Similarly, we can think of a pseudolikelihood function of  $\theta$  as being any nonnegative random function of  $\theta$  for which at least one of the Bartlett identities does not hold. Hence, the identities act as a litmus test of sorts for determining the validity of a pseudolikelihood as an approximation to the genuine likelihood from which it originated - the more identities it does satisfy, the better the approximation. A pseudolikelihood that satisfies the first Bartlett identity is called *score-unbiased*; one that satisfies the second is called *information-unbiased*. Historically, more attention has been given to constructing pseudo-likelihoods that are score-unbiased, at least asymptotically (Cf. Kalbfleisch and Sprott, 1970; De Bin et al., 2015; Schumann et al., 2021, 2023).

### 3.2.2. Types of Pseudolikelihoods

**3.2.2.1. Profile Likelihoods.** Profile likelihoods are one of the most straightforward methods for eliminating a nuisance parameter  $\lambda$  from a likelihood function. The idea is to summarize the set  $\mathcal{L}_\psi$  by its maximum value. Formally, the profile likelihood is defined as

$$L_p(\psi) = \sup_{\theta \in \Theta(\psi)} L(\theta) = L(\hat{\theta}_\psi),$$

where  $\hat{\theta}_\psi$  is the MLE of  $\theta_0$  for fixed  $\psi$ . If we define  $\hat{\lambda}_\psi$  to be the conditional MLE for  $\lambda$  given  $\psi$ , meaning the value of  $\lambda$  that maximizes the likelihood for a particular value of  $\psi$ , then we must have  $\hat{\theta}_\psi = (\psi, \hat{\lambda}_\psi)$ .

Much of the allure of a profile likelihood can be traced to its ease of computation. In the event of a regular model, finding the value of  $\hat{\lambda}_\psi$  corresponding to a particular  $\psi$  is equivalent to solving a convex optimization problem. Either an analytical solution to this problem will exist or a numerical one can be obtained using modern computational tools. In both cases,  $\hat{\lambda}_\psi$  can be found without much trouble, and from there it's just a matter of setting  $\lambda = \hat{\lambda}_\psi$  inside  $L(\psi, \lambda)$  to obtain  $L_p(\psi)$ .

The method is not without its drawbacks however. As noted by Kalbfleisch and Sprott (1970), a major disadvantage to profile likelihoods are their inability to take into account the dimensionality of the nuisance parameter. By effectively always assuming that  $\lambda = \hat{\lambda}_\psi$ , they fail to incorporate any uncertainty

we might have in its value into the resulting pseudolikelihood, leading to overly confident estimates of  $\psi$ . This effect is especially pronounced when the dimension of  $\Lambda$  is large. Berger, Liseo, and Wolpert (1999) also bring up the scenario in which the likelihood has a sharp ridge running in one direction as being one in which the profile likelihood will perform poorly. In such a situation, the profile of the likelihood along the ridge would not be representative of the shape of the likelihood elsewhere, and yet it is exactly that profile which will be obtained through maximization.

Nevertheless, a profile likelihood can still be a useful tool for conducting inference regarding  $\psi_0$ . At worst, it offers a baseline of sorts to assess the quality of other pseudolikelihoods. There is no point in using another pseudolikelihood that is harder to compute if it offers the same amount of information about  $\psi_0$  - it must justify this increase in complexity with a corresponding increase in features that lend themselves to a greater degree of accuracy in our inference, such as higher level of peakedness around the value of  $\psi_0$ .

**3.2.2.2. Marginal and Conditional Likelihoods.** Suppose we observe some data  $\mathbf{X} = \mathbf{x}$  such that the pair of statistics  $(\mathbf{T}(\mathbf{X}), \mathbf{S}(\mathbf{X}))$  is sufficient for  $\boldsymbol{\theta} = (\psi, \lambda)$ . Then, abusing notation slightly, we can write

$$(3.2.1) \quad p_{\boldsymbol{\theta}}(\mathbf{X}) = p_{\boldsymbol{\theta}}(\mathbf{T}, \mathbf{S}).$$

Since the right-hand side is a joint density for  $\mathbf{T}$  and  $\mathbf{S}$ , it can be factored into a product of a marginal and a conditional density as follows:

$$(3.2.2) \quad p_{\boldsymbol{\theta}}(\mathbf{T}, \mathbf{S}) = p_{\boldsymbol{\theta}}(\mathbf{T}|\mathbf{S})p_{\boldsymbol{\theta}}(\mathbf{S}).$$

If the right-hand term in this product doesn't depend on  $\lambda$ , i.e.,  $p_{\boldsymbol{\theta}}(\mathbf{S}) = p_{\psi}(\mathbf{S})$ , then we have

$$(3.2.3) \quad p_{\boldsymbol{\theta}}(\mathbf{T}, \mathbf{S}) = p_{\boldsymbol{\theta}}(\mathbf{T}|\mathbf{S})p_{\psi}(\mathbf{S}).$$

In this case, a pseudolikelihood for  $\psi$  is simply

$$(3.2.4) \quad L_m(\psi) = p_{\psi}(\mathbf{S}).$$

$L_m(\psi)$  is called a *marginal likelihood* for  $\psi$  as it is based on the marginal distribution of  $\mathbf{S} = \mathbf{S}(\mathbf{X})$ . The conditional part of the density,  $p_{\theta}(\mathbf{T}|\mathbf{S})$ , does still depend on  $\psi$ , however, and so by choosing to base our inferences regarding  $\psi_0$  solely on  $L_m(\psi)$ , we have ignored some relevant information for  $\psi_0$  that was contained in the data.

If instead it is the left-hand term in Equation 3.2.2 that doesn't depend on  $\lambda$ , i.e.,  $p_{\theta}(\mathbf{T}|\mathbf{S}) = p_{\psi}(\mathbf{T}|\mathbf{S})$ , then we have

$$(3.2.5) \quad p_{\theta}(\mathbf{T}, \mathbf{S}) = p_{\psi}(\mathbf{T}|\mathbf{S})p_{\theta}(\mathbf{S}).$$

Here, a pseudolikelihood for  $\psi$  may be given by

$$(3.2.6) \quad L_C(\psi) = p_{\psi}(\mathbf{T}|\mathbf{S}).$$

$L_C(\psi)$  is called a *conditional likelihood* for  $\psi$  as it is based on the conditional distribution of  $\mathbf{T}$  given  $\mathbf{S}$ . As with the marginal likelihood, some information has been disregarded through our omission of  $p_{\theta}(\mathbf{S})$  in our inference for  $\psi_0$ .

Thus, the use of either a marginal or a conditional likelihood as a pseudolikelihood for a parameter of interest is theoretically only justified when the benefits of eliminating the nuisance parameter through marginalization or conditioning outweigh the corresponding loss in information. In practice, however, the real limiting constraint of using a marginal or conditional likelihood tends not to be the lack of clarity they confer to our inferences, but rather their lack of existence in the first place. The ability to factor a density as in Equation 3.2.3 or Equation 3.2.5 is a rather strong condition, and there are plenty of models for which it is impossible to construct a marginal or conditional likelihood for its parameter of interest. When they do exist, it is typically worthwhile to use them.

**3.2.2.3. Integrated Likelihoods.** The *integrated likelihood* for  $\psi$  seeks to summarize  $\mathcal{L}_{\psi}$  by its average value with respect to some weight function  $\pi$  over  $\Theta(\psi)$ . From a theoretical standpoint, this is preferable to maximization as it incorporates (or at least is capable of incorporating) the uncertainty we have in the value of the nuisance parameter into the resulting pseudolikelihood in a very natural way. Formally,

the integrated likelihood function is defined as

$$(3.2.7) \quad \bar{L}(\psi) = \int_{\Lambda} L(\psi, \lambda) \pi(\lambda|\psi) d\lambda,$$

where  $\pi(\lambda|\psi)$  is a nonnegative function on  $\Lambda$ .  $\pi(\lambda|\psi)$  is sometimes called a conditional prior density for  $\lambda$  given  $\psi$ , though it need not satisfy the requirements of a genuine density function.

Note the similarity in form between the integral in Equation 3.2.7 and the expression for the normalizing constant of a posterior distribution:

$$\int_{\Theta} L(\theta; X) \pi(\theta) d\theta.$$

This similarity lends weight to the idea that Bayesian techniques used to obtain empirical approximations to posterior distributions, such as Markov Chain Monte Carlo, could also be used to approximate an integrated likelihood function, with the result being useful for Bayesian and frequentist inference alike.

In general, the selection of the weight function plays an important role in the properties of the resulting integrated likelihood. The obvious first choice to make for it, and therefore the one that could be considered the “default”, is simply  $\pi(\lambda|\psi) \propto 1$ , i.e., a uniform density. This yields what is known as the uniform-integrated likelihood:

$$(3.2.8) \quad \bar{L}^U(\psi) = \int_{\Lambda} L(\psi, \lambda) d\lambda.$$

In the next chapter, we will discuss a re-parameterization of the nuisance parameter developed by Severini (2007) that makes the integrated likelihood relatively insensitive to the exact weight function chosen. Using this new parameterization, we have great flexibility in choosing our weight function; as long as it does not depend on the parameter of interest, the integrated likelihood that is produced will enjoy many desirable frequency properties.

## CHAPTER 4

## Integrated Likelihood Analysis

## 4.1. Two-Index Asymptotics

Earlier we discussed the one-index asymptotics setting, in which the sample size ( $n$ ) of the model diverged to infinity while the dimension of the nuisance parameter ( $q$ ) remained fixed. Now we turn our attention to the two-index asymptotics setting which describes the behavior of likelihood and pseudolikelihood functions as  $n$  and  $q$  both tend to infinity, with  $q$  growing at least as fast as  $n$ . Under such a framework, De Bin, Sartori, and Severini (2015) showed that estimates for  $\psi$  based on a suitably constructed integrated likelihood function will outperform those coming from more traditional pseudolikelihoods, such as the profile likelihood. Such findings provide the motivation for our ensuing examination of two-index asymptotics theory, insofar as it relates to the performance of the integrated likelihood function as a method of inference regarding a parameter of interest.

To mirror the strategy used by Sartori (2003) and De Bin, Sartori, and Severini (2015), we will frame our discussion in terms of a stratified sample of data in which each stratum contributes one component to the overall nuisance parameter. Consider a model with parameter  $\theta = (\psi, \lambda)$  where  $\psi$  is the parameter of interest and  $\lambda = (\lambda_1, \dots, \lambda_q)$  is a  $q$ -dimensional nuisance parameter. For the sake of reducing complexity in our notation, we will only consider the case in which  $\psi$  and the individual components of  $\lambda$  are scalar parameters, though the results of this section should hold in the case where they are all vectors as well. Suppose that we have divided the model's population into  $q$  strata and collected a sample of size  $m_i$  from each stratum such that observation  $j$  from stratum  $i$  may be modeled as

$$(4.1.1) \quad X_{ij} \sim p_{ij}(x_{ij}; \psi, \lambda_i),$$



where  $i = 1, \dots, q$  and  $j = 1, \dots, m_i$ , making the total sample size  $n = \sum_{i=1}^q m_i$ .<sup>1</sup> Hence, there is a one-to-one correspondence between the strata and the components of  $\lambda$ ; assume that each  $\lambda_i \in \Lambda$ , where the space  $\Lambda$  is the same for all  $i$ , and they all have the same interpretation within their respective strata.

Assume that all the regularity conditions set forth in Section 3.4 apply, except possibly for **RC1** - it is not necessary to assume that the observations are i.i.d. here, and in fact it is perfectly acceptable for the  $p_{ij}$ 's in Equation 4.1.1 to differ from one another. We will also allow for the possibility of dependence among observations within a stratum, though not between them.

Let  $\mathbf{x}_i = (x_{i1}, \dots, x_{im})$  denote the sample of observations from stratum  $i$ , so that their joint density may be written as  $p_i(\mathbf{x}_i; \psi, \lambda_i)$ . Therefore, the likelihood and log-likelihood for the  $i$ th stratum are

$$(4.1.2) \quad L^{(i)}(\psi, \lambda_i) = p_i(\mathbf{x}_i; \psi, \lambda_i),$$

and

$$(4.1.3) \quad \ell^{(i)}(\psi, \lambda_i) = \log L^{(i)}(\psi, \lambda_i),$$

respectively. For a particular choice of weight function  $g(\lambda_i; \psi)$ , the integrated likelihood for  $\psi$  in stratum  $i$  is given by

$$(4.1.4) \quad \bar{L}^{(i)}(\psi) = \int_{\Lambda} L^{(i)}(\psi, \lambda_i) g(\lambda_i; \psi) d\lambda_i.$$

From here, we proceed by using Laplace's method as described by Tierney and Kadane (1986) (see Appendix B for a brief review) to obtain an analytic approximation to  $\bar{L}^{(i)}(\psi)$ . Setting  $h(\lambda_i) = -\frac{1}{m} \ell^{(i)}(\psi, \lambda_i)$  and  $f(\lambda_i) = g(\lambda_i; \psi)$ , we may rewrite the integral in Equation 4.1.4 as

$$\bar{L}^{(i)}(\psi) = \int_{\Lambda} f(\lambda_i) \exp[-mh(\lambda_i)] d\lambda_i.$$

---

<sup>1</sup>It will be convenient to work under the restriction that the stratum sample sizes are all identical, meaning there exists some positive integer  $m$  such that  $m_i = m$  for all  $i$ . However, as both Sartori (2003) and De Bin, Sartori, and Severini (2015) note, we could also assume a looser condition in which  $m_i = K_i m$  where  $0 < K_i < \infty$  without compromising our results.

One consequence of the regularity conditions we have assumed is that  $L^{(i)}(\psi, \lambda_i)$  is that an MLE for  $\theta_0$ ,  $\hat{\theta}$ , exists and is unique. This further implies the existence and uniqueness of a conditional MLE for the true value of each stratum-specific nuisance parameter given  $\psi$  - denote this value for the  $i$ th stratum by  $\hat{\lambda}_{i\psi}$ . By definition this value maximizes  $\ell^{(i)}(\psi, \lambda_i)$  as a function of  $\lambda_i$ , and so it also maximizes  $-h(\lambda_i)$  since the two functions differ only by a multiplicative constant  $\frac{1}{m}$ .

The Laplace approximation to  $\bar{L}^{(i)}(\psi)$  is then given by

$$(4.1.5) \quad \hat{\bar{L}}^{(i)}(\psi) = f(\hat{\lambda}_{i\psi}) \sqrt{\frac{2\pi}{m}} \sigma \exp[-mh(\hat{\lambda}_{i\psi})],$$

where

$$\begin{aligned} \sigma &= \left[ \frac{\partial^2 h}{\partial \lambda_i^2} \Big|_{\lambda_i = \hat{\lambda}_{i\psi}} \right]^{-1/2} \\ &= \left[ -\frac{1}{m} \frac{\partial^2 \ell^{(i)}(\psi, \lambda_i)}{\partial \lambda_i^2} \Big|_{\lambda_i = \hat{\lambda}_{i\psi}} \right]^{-1/2} \\ &= \left[ \frac{1}{m} \mathcal{J}(\hat{\lambda}_{i\psi}) \right]^{-1/2}. \end{aligned}$$

Here,  $\mathcal{J}(\hat{\lambda}_{i\psi})$  denotes the observed information function for  $\lambda_i$  only (i.e. the negative second partial derivative of the log-likelihood with respect to  $\lambda_i$ ) evaluated at  $\hat{\lambda}_{i\psi}$ . Plugging in the appropriate quantities for  $f$ ,  $h$ , and  $\sigma$  into Equation 4.1.5, we arrive at

$$\begin{aligned} (4.1.6) \quad \hat{\bar{L}}^{(i)}(\psi) &= f(\hat{\lambda}_{i\psi}) \sqrt{\frac{2\pi}{m}} \sigma \exp[-mh(\hat{\lambda}_{i\psi})] \\ &= g(\hat{\lambda}_{i\psi}; \psi) \sqrt{\frac{2\pi}{m}} \left[ \frac{1}{m} \mathcal{J}(\hat{\lambda}_{i\psi}) \right]^{-1/2} \exp \left\{ -m \cdot -\frac{1}{m} \ell^{(i)}(\psi, \hat{\lambda}_{i\psi}) \right\} \\ &= \frac{\sqrt{2\pi}}{m} L^{(i)}(\psi, \hat{\lambda}_{i\psi}) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2} \\ &= \frac{\sqrt{2\pi}}{m} L_P^{(i)}(\psi) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2}, \end{aligned}$$

where  $L_P^{(i)}(\psi) = L^{(i)}(\psi, \hat{\lambda}_{i\psi})$  is the profile likelihood for  $\psi$ . The error in this approximation is

$$(4.1.7) \quad \bar{L}^{(i)}(\psi) = \hat{\bar{L}}^{(i)}(\psi) \left\{ 1 + O\left(\frac{1}{m}\right) \right\} \quad \text{as } m \rightarrow \infty.$$

Let

$$(4.1.8) \quad \bar{\ell}^{(i)}(\psi) = \log \bar{L}^{(i)}(\psi)$$

denote the integrated log-likelihood for  $\psi$ . Putting the results in Equation 4.1.6, Equation 4.1.7, and Equation 4.1.8 together, we have

$$\begin{aligned} \bar{\ell}^{(i)}(\psi) &= \log \bar{L}^{(i)}(\psi) && \text{(by Equation 4.2.8)} \\ &= \log \left( \hat{\bar{L}}^{(i)}(\psi) \left\{ 1 + O\left(\frac{1}{m}\right) \right\} \right) && \text{(by Equation 4.2.7)} \\ &= \log \hat{\bar{L}}^{(i)}(\psi) + \log \left\{ 1 + O\left(\frac{1}{m}\right) \right\} && \text{(by Equation 4.2.6)} \\ &= \log \left\{ \frac{\sqrt{2\pi}}{m} L_P^{(i)}(\psi) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2} \right\} + O\left(\frac{1}{m}\right) \\ &= \frac{1}{2} \log(2\pi) - \log(m) + \log L_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{1}{m}\right) \\ &= \ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}) + \frac{1}{2} \log(2\pi) - \log(m) + O\left(\frac{1}{m}\right) \quad \text{as } m \rightarrow \infty, \end{aligned}$$

where  $\ell_P^{(i)}(\psi) = \ell^{(i)}(\psi, \hat{\lambda}_{i\psi})$  is the profile log-likelihood for  $\psi$ . Since log-likelihoods are equivalent up to additive constants, we can discard the  $\frac{1}{2} \log(2\pi)$  and  $\log(m)$  terms in the final line above to arrive at our final approximation for the integrated log-likelihood in stratum  $i$ :

$$(4.1.9) \quad \hat{\bar{\ell}}^{(i)}(\psi) = \ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}).$$

The error in this approximation is given by

$$(4.1.10) \quad \bar{\ell}^{(i)}(\psi) = \hat{\bar{\ell}}^{(i)}(\psi) + O\left(\frac{1}{m}\right).$$

Since the observations between the strata are independent, we may write the likelihood and log-likelihood functions for the entire model as

$$(4.1.11) \quad L(\psi, \lambda) = \prod_{i=1}^q L^{(i)}(\psi, \lambda_i)$$

and

$$(4.1.12) \quad \ell(\psi, \lambda) = \sum_{i=1}^q \ell^{(i)}(\psi, \lambda_i),$$

respectively. If we define the weight function

$$(4.1.13) \quad G(\lambda; \psi) \equiv \prod_{i=1}^q g(\lambda_i; \psi)$$

then the integrated likelihood function for  $\psi$  becomes separable. That is,

$$\begin{aligned} \bar{L}(\psi) &= \int_{\Lambda^q} L(\psi, \lambda) G(\lambda; \psi) d\lambda \\ &= \int_{\Lambda} \cdots \int_{\Lambda} \left[ \prod_{i=1}^q L^{(i)}(\psi, \lambda_i) \right] \left[ \prod_{i=1}^q g(\lambda_i; \psi) \right] d\lambda_1 \cdots d\lambda_q \\ (4.1.14) \quad &= \prod_{i=1}^q \int_{\Lambda} L^{(i)}(\psi, \lambda_i) g(\lambda_i; \psi) d\lambda_i \\ &= \prod_{i=1}^q \bar{L}^{(i)}(\psi). \end{aligned}$$

Let  $\bar{\ell}(\psi) = \log \bar{L}(\psi)$  denote the integrated log-likelihood function for  $\psi$ . Taking the logarithm of both sides in Equation 4.1.14, we have

$$(4.1.15) \quad \bar{\ell}(\psi) = \sum_{i=1}^q \bar{\ell}^{(i)}(\psi).$$

Plugging in our approximation to  $\bar{\ell}^{(i)}(\psi)$  and its error term in Equation 4.1.9 and Equation 4.1.10, respectively, yields

$$\begin{aligned}
 \bar{\ell}(\psi) &= \sum_{i=1}^q \left[ \ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{1}{m}\right) \right] \\
 (4.1.16) \quad &= \ell_P(\psi) + \sum_{i=1}^q \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \sum_{i=1}^q \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{q}{m}\right) \quad \text{as } m \rightarrow \infty.
 \end{aligned}$$

## 4.2. The Zero-Score Expectation Parameter

Severini (2007) considered the problem of selecting a weight function  $\pi(\lambda|\psi)$  such that when the likelihood function is integrated with respect to this density over the nuisance parameter space, the result is useful for non-Bayesian inference. To do this, he outlined four properties (see Appendix B) that an integrated likelihood must satisfy if it is to be of any use and went on to show that such a function could be obtained by doing the following:

- 1) Find a reparameterization  $(\psi, \lambda) \mapsto (\psi, \phi)$  of the model such that the new nuisance parameter  $\phi$  is unrelated to  $\psi$  in the sense that  $\hat{\phi}_\psi = \hat{\phi}$ ; that is, the conditional maximum likelihood estimate of  $\phi$  given  $\psi$  is simply equal to the unrestricted maximum likelihood estimate of  $\phi$ .
- 2) Select a prior density  $\pi(\lambda|\psi)$  such that when the model undergoes the above reparameterization, the resulting prior density  $\pi(\phi)$  does not depend on  $\psi$ .

An integrated likelihood function for  $\psi$  that possesses the desired properties will then be given by

$$(4.2.1) \quad \bar{L}(\psi) = \int_{\Phi} \tilde{L}(\psi, \phi) \pi(\phi) d\phi,$$

where  $\tilde{L}(\psi, \phi)$  is the likelihood function for the model after it has been reparameterized in terms of  $\phi$ . The exact choice of prior density for  $\phi$  is not particularly important; the only restriction placed upon it is that it must not depend on  $\psi$ . Hence, the crux of the matter really lies in completing the first step. The approach taken by Severini (2007) is to define a new nuisance parameter  $\phi$  as the solution to the

equation

$$(4.2.2) \quad \mathbb{E}_{(\psi_0, \lambda_0)} \left[ \nabla_{\lambda} \ell(\psi, \lambda; \mathbf{X}_n) \right] \bigg|_{(\psi_0, \lambda_0) = (\hat{\psi}, \phi)} = \mathbf{0},$$

where  $\psi_0$  and  $\lambda_0$  denote the true values of  $\psi$  and  $\lambda$ , and  $\hat{\psi}$  is the unrestricted MLE for  $\psi_0$ . The expectation here is being taken with respect to the data  $\mathbf{X}_n = \mathbf{x}_n$  and not the parameters themselves. Equation 4.2.2 can thus be rewritten as

$$I(\psi, \lambda, \hat{\psi}, \phi) = \mathbf{0},$$

where

$$I(\psi, \lambda, \psi_0, \lambda_0) = \int_{\mathbb{R}^n} [\nabla_{\lambda} \ell(\psi, \lambda; \mathbf{x}_n)] p(\mathbf{x}_n; \psi_0, \lambda_0) d\mathbf{x}_n.$$

Assuming  $I$  is invertible, for a particular value of  $(\psi, \lambda, \hat{\psi})$ , there will be a unique value of  $\phi$  that solves Equation 4.2.2.  $\phi$  is called the *zero-score expectation* (ZSE) parameter because it is defined as the value that makes the expectation of the score function (in terms of  $\lambda$ , not the full parameter) with respect to  $p(\mathbf{x}_n; \psi_0, \lambda_0)$  evaluated at the point  $(\psi_0, \lambda_0) = (\hat{\psi}, \phi)$  equal to zero. This means that  $\phi$  is really a function of  $(\psi, \lambda, \hat{\psi})$ , i.e.,  $\phi = \phi(\psi, \lambda, \hat{\psi})$ . This in turn implies that  $\phi$  is a function of the data through  $\hat{\psi}$ . Normally we try to avoid creating such dependencies in our parameters as it renders them useless for the purpose of parameterizing a statistical model. However, from the perspective of the likelihood function, once the data have been collected they are considered fixed in place and there is no issue with using a quantity such as  $\phi$  that depends on the data to parameterize it.

For a given value of  $\phi$ , the corresponding value of  $\lambda$  can be found by

$$(4.2.3) \quad \lambda(\psi, \phi) = \operatorname{argmax}_{\lambda \in \Lambda} \mathbb{E}_{(\hat{\psi}, \phi)} [\ell(\psi, \lambda; \mathbf{X}_n)].$$

For a certain choice of prior density  $\pi(\phi)$ , this allows us to write Equation 4.2.1 in terms of  $L(\psi, \lambda)$ :

$$(4.2.4) \quad \bar{L}(\psi) = \int_{\Phi} L(\psi, \lambda(\psi, \phi)) \pi(\phi) d\phi.$$

#### 4.2.1. Approximating the Integrated Likelihood for an Implicit Parameter of Interest

The procedure described in the previous section is based on the assumption that  $\lambda$  is an explicit nuisance parameter so that taking partial derivatives of  $\ell$  with respect to its components is a well-defined operation. However, Severini (2018) proved that reparameterizing the model in terms of the ZSE parameter yields the same nice properties in the subsequent integrated likelihood when  $\psi$  and  $\lambda$  are implicit as well. In this section, we consider an approach to approximating the integrated likelihood function that has been adapted from this method.

Consider a model with parameter  $\theta \in \Theta$  and implicit parameter of interest  $\psi = \varphi(\theta)$  for some function  $\varphi : \Theta \rightarrow \mathbb{R}$ . Equation 4.2.4 tells us that we can calculate the integrated likelihood for  $\psi$  if we know the value of  $\theta$  corresponding to a given value of the ZSE parameter, and this will be true for models with both explicit and implicit parameters. Let  $\hat{\psi} = \varphi(\hat{\theta})$  denote the unrestricted MLE for  $\psi$  and define the set

$$(4.2.5) \quad \Omega_{\psi} = \left\{ \omega \in \Theta : \varphi(\omega) = \psi \right\}.$$

Then an element of  $\Omega_{\hat{\psi}}$  for a model without an explicit nuisance parameter is functionally equivalent to a value  $(\hat{\psi}, \phi)$  for a model with an explicit nuisance parameter.

Generalizing the result in Equation 4.2.3, for a given element  $\omega \in \Omega_{\hat{\psi}}$ , the corresponding value of  $\theta$  is that which maximizes  $E_{\omega}[\ell(\theta; \mathbf{X}_n)]$  subject to the restriction that  $\varphi(\theta) = \psi$ . This allows us to define a function  $T_{\psi} : \Omega_{\hat{\psi}} \rightarrow \Theta(\psi)$  such that  $T_{\psi}(\omega) = \underset{\theta \in \Theta(\psi)}{\operatorname{argmax}} E_{\omega}[\ell(\theta; \mathbf{X}_n)]$ . The integrated likelihood for  $\psi$  is then given by

$$(4.2.6) \quad \bar{L}(\psi) = \int_{\Omega_{\hat{\psi}}} L(T_{\psi}(\omega)) \pi(\omega) d\omega,$$

where  $\pi(\omega)$  is a density defined on  $\Omega_{\hat{\psi}}$ .

Equation 4.2.6 can also be written as

$$(4.2.7) \quad \bar{L}(\psi) = E_{\omega} \left[ L(T_{\psi}(W)) \right],$$

where  $W$  represents a random variable with density  $\pi(\omega)$ . We can further define a function  $Q : \mathcal{U} \rightarrow \Omega_{\hat{\psi}}$  for some set  $\mathcal{U}$  such that if  $U$  is a random variable taking values in  $\mathcal{U}$ , then  $Q(U)$  will be a random variable in  $\Omega_{\hat{\psi}}$  with a distribution that is completely determined by our choice of  $U$ . Therefore,

$$(4.2.8) \quad \bar{L}(\psi) = E_{\omega} \left[ L(T_{\psi}(Q(U))) \right]$$

is an integrated likelihood for  $\psi$  with a weight function corresponding to the density of  $Q(U)$ .

Define

$$(4.2.9) \quad \tilde{L}(u; \psi) = L(T_{\psi}(Q(u))), \quad u \in \mathcal{U}.$$

Then we simply have  $\tilde{L}(u; \psi) = L(\theta)$ , with  $\theta$  taken to be  $T_{\psi}(Q(u))$ .  $\tilde{L}(u; \psi)$  will be a genuine likelihood for the parameter  $(u, \psi)$  provided that there always exists a value  $(u, \psi)$  such that  $T_{\psi}(Q(u)) = \theta_0$  for any possible value of  $\theta_0$ . A sufficient condition for this to occur is that for any  $\psi \in \Psi$ ,  $T_{\psi}(Q(\mathcal{U})) = \Omega_{\psi}$ .

We can then write the integrated likelihood for  $\psi$  in terms of  $\tilde{L}(u; \psi)$  as follows:

$$(4.2.10) \quad \bar{L}(\psi) = \int_{\mathcal{U}} \tilde{L}(u; \psi) \tilde{\pi}(u) du,$$

where  $\tilde{\pi}(u)$  is a density on  $\mathcal{U}$  of our choosing. If  $\tilde{\pi}(u)$  doesn't depend on  $\psi$ , the integrated likelihood given by Equation 4.2.10 will have the properties we desire.

We can further rewrite Equation 4.2.10 as

$$(4.2.11) \quad \bar{L}(\psi) = \int_{\Theta} \frac{\tilde{L}(u; \psi)}{\check{L}(u)} \check{L}(u) \tilde{\pi}(u) du,$$

where  $\check{L}(u)$  represents an arbitrary likelihood function for the “parameter”  $u$ . From a Bayesian point of view, the quantity  $\check{L}(u) \tilde{\pi}(u)$  can then be thought of as a posterior density for  $u$  up to some proportionality constant. If  $\tilde{\pi}(u)$  is chosen to be a conjugate prior for  $\check{L}(u)$  such that the posterior density  $\check{L}(u) \tilde{\pi}(u)$  has the same form as  $\check{L}$  itself, and  $\check{L}$  is chosen to be a known distribution, then random samples can be drawn directly from this posterior density using modern statistical computing packages. Alternatively, it may



also be possible to obtain random samples from the posterior density through Monte Carlo methods such as importance sampling or MCMC.

In either case, once random variates  $u_1, \dots, u_R$  have been sampled from  $\check{L}(u)\check{\pi}(u)$ , a simple empirical estimate to  $\bar{L}(\psi)$  at a particular value of  $\psi$  is given by

$$(4.2.12) \quad \hat{\bar{L}}(\psi) = \frac{1}{R} \sum_{i=1}^R \frac{\check{L}(u_i; \psi)}{\check{L}(u_i)}.$$

Repeating this procedure for every value of  $\psi \in \Psi$  will give an overall shape for  $\hat{\bar{L}}(\psi)$ , which can be used to conduct inference for  $\psi_0$  without interference from a nuisance parameter.

## CHAPTER 5

## Applications

## 5.1. Multinomial Distribution

For  $n$  independent trials each of which leads to a success for exactly one of  $d$  categories, with the  $j$ th category having a fixed probability of success  $\theta_j$ , the multinomial distribution allows us to calculate the probability of any particular combination of numbers of successes for the various categories. If  $N_j$  represents the number of successes in the  $j$ th category, then the random vector  $(N_1, \dots, N_d)$  will follow a multinomial distribution with parameter  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$  such that  $E_{\boldsymbol{\theta}}(N_j) = n\theta_j$  for  $j = 1, \dots, d$ , where  $n = \sum_{j=1}^d N_j$ . Since each  $\theta_j$  is meant to be interpreted as a probability, the parameter space  $\Theta$  must be the probability simplex in  $\mathbb{R}^d$  so that  $\theta_j \geq 0$  for all  $j$  and  $\sum_{j=1}^d \theta_j = 1$ .

Suppose we are interested in the entropy of this distribution, so that the parameter of interest is  $\psi = \varphi(\boldsymbol{\theta})$ , where

$$\varphi(\boldsymbol{\theta}) = -\sum_{j=1}^d \theta_j \log(\theta_j).$$

Let  $(n_1, \dots, n_d)$  denote the observed counts of  $(N_1, \dots, N_d)$  so that the likelihood function is given by

$$L(\boldsymbol{\theta}) = \prod_{j=1}^d \theta_j^{n_j}.$$

The unrestricted MLE of  $\theta_j$  is given by  $\hat{\theta}_j = \frac{n_j}{n}$ . Similarly, due to the invariance property of the MLE, the unrestricted MLE of  $\psi$  will simply be  $\hat{\psi} = \varphi(\hat{\boldsymbol{\theta}})$ .

To obtain an integrated likelihood for  $\psi$  alone, we can use the procedure described in the previous chapter to find an approximation to the integral in Equation 4.2.11 where  $\tilde{L}(u; \psi)$  is the likelihood function reparameterized in terms of the ZSE parameter, and  $\tilde{L}(u)$  and  $\tilde{\pi}(u)$  are chosen such that  $\tilde{\pi}$  is a conjugate

prior for  $\tilde{L}$ .<sup>1</sup> In this case, a natural choice exists due to the fact that the Dirichlet distribution is a conjugate prior for the multinomial distribution. Since our original likelihood function  $L$  is based on a multinomial distribution, we can simply set  $\tilde{L}(u) := L(u)$  and  $\tilde{\pi}(u) \sim \text{Dir}(\boldsymbol{\alpha})$ , where  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)$ . Then the posterior distribution for  $u$  based on data  $\mathbf{n} = (n_1, \dots, n_d)$  is given by  $L(u)\tilde{\pi}(u) \sim \text{Dir}(\mathbf{n} + \boldsymbol{\alpha})$ . Consequently, we will take  $\tilde{\pi}(u)$  to be the symmetric Dirichlet distribution on the probability simplex in  $\mathbb{R}^d$  with  $\alpha_j = 1$  for all  $j$  so that random variates of  $u$  can be sampled from a  $\text{Dir}(\mathbf{n} + 1)$  distribution.

From Equation 4.2.9, finding  $\tilde{L}(u; \psi)$  is a matter of defining two functions,  $Q$  and  $T_\psi$  such that  $\tilde{L}(u; \psi) = L(T_\psi(Q(u)))$ .  $Q$  maps a random variate  $u$  sampled from the above posterior to an element  $\omega$  in  $\Omega_{\hat{\psi}}$ , and  $T_\psi$  then maps  $\omega$  to an element  $\boldsymbol{\theta}$  in  $\Theta(\psi)$ . Since in this situation, these quantities  $u$ ,  $\omega$ , and  $\boldsymbol{\theta}$  are all members of the probability simplex, the maps  $Q$  and  $T_\psi$  can be defined as returning the elements in their respective output spaces that are closest to the input element they have each received. That is,  $Q$  returns the element  $\omega$  that minimizes the distance to an input  $u$ , subject to the constraints that the sum of the components of  $\omega$  equal 1 and  $\varphi(\omega) = \hat{\psi}$ . Similarly,  $T_\psi$  returns the element  $\boldsymbol{\theta}$  that minimizes the distance to an input  $\omega$ , subject to the constraints that the sum of the components of  $\boldsymbol{\theta}$  equal 1 and  $\varphi(\boldsymbol{\theta}) = \psi$ .

From here, we sample  $R$  random variates from the appropriate Dirichlet distribution, calculate  $\tilde{L}(u; \psi)$  and  $L(u)$  for each one, and use Equation 4.2.12 to approximate the integrated likelihood at a particular value of  $\psi$ . We then repeat this process over a finely spaced sequence of the possible values of  $\psi$  in order to get a shape of the overall integrated likelihood. See Appendix C for a graph comparing the plot of one such integrated likelihood to the profile likelihood, for observed data  $(n_1, \dots, n_6) = (1, 1, 4, 7, 10)$  and 250 samples of the appropriate Dirichlet distribution drawn for each value of  $\psi$ .<sup>2</sup>

## 5.2. Standardized Mean Difference

Broadly speaking, meta-analysis seeks to synthesize the results of multiple studies, typically with the goal of estimating with heightened precision the effect of a shared treatment among the studies in

<sup>1</sup>This procedure is an adapted version of another algorithm developed by Severini (2022) for approximating the integrated likelihood for the entropy of a multinomial distribution.

<sup>2</sup>The samples were obtained using the ‘LaplacesDemon’ R package. The distance minimizations needed for  $Q$  and  $T_\psi$  were computed numerically using the ‘nloptr’ R package.

question. An effect size frequently used in a meta-analysis when the outcome variable is continuous is the *standardized mean difference* (SMD), defined as the difference in means between the treatment and control group divided by their common standard deviation. The aim of this section is to provide a blueprint for obtaining the integrated likelihood function of an SMD in a meta-analysis.

Consider a sample of  $q$  independent studies, each comparing a treatment group to a control group. Suppose the  $j$ th observations of the treatment and control groups in the  $i$ th study are given by

$$Y_{ij}^T \sim N(\mu_i^T, \sigma_i^2)$$

and

$$Y_{ij}^C \sim N(\mu_i^C, \sigma_i^2),$$

respectively, where  $j = 1, \dots, m_i$  observations and  $i = 1, \dots, q$  studies. The number of observations in study  $i$  may be further decomposed as  $m_i = m_i^T + m_i^C$ . The true SMD effect size for study  $i$  is defined as

$$(5.2.1) \quad \delta_i = \frac{\mu_i^T - \mu_i^C}{\sigma_i}.$$

Standard practice is to estimate  $\delta_i$  by

$$(5.2.2) \quad d_i = \frac{\bar{y}_i^T - \bar{y}_i^C}{s_i},$$

where  $\bar{y}$  is the traditional sample mean and  $s$  is the pooled sample standard deviation. Hedges and Olkin (1985) showed that as the number of observations in the study tends to infinity,  $d_i$  converges to a  $N(\delta_i, \sigma_{d_i}^2)$  distribution, where

$$(5.2.3) \quad \sigma_{d_i}^2 = \frac{m_i^T + m_i^C}{m_i^T m_i^C} + \frac{\delta_i^2}{2(m_i^T + m_i^C - 2)}.$$

$\sigma_{d_i}^2$  represents the variance of all observations within the  $i$ th study.

Under a linear fixed-effect model,

$$(5.2.4) \quad d_i = \delta + \epsilon_i,$$

with  $\epsilon_i \sim N(0, \sigma_{d_i}^2)$  for  $i = 1, \dots, q$ , where  $\delta = \frac{1}{q} \sum_{i=1}^q \delta_i$  is the common mean. The overarching objective of a meta-analysis using this framework is to find an estimate of  $\delta$  and its associated standard error. Hence,  $\delta$  can be considered an implicit parameter of interest in this model, and so it makes sense to consider the efficacy of an integrated likelihood function as a tool for its estimation.

Let  $\theta_i = (\mu_i^T, \mu_i^C, \sigma_i)$  denote the vector of parameters associated with study  $i$  so that the overall model parameter is given by  $\theta = (\theta_1, \dots, \theta_q) = ((\mu_1^T, \mu_1^C, \sigma_1), \dots, (\mu_q^T, \mu_q^C, \sigma_q))$ . The likelihood function for  $\theta$  is

$$\begin{aligned} L(\theta) &= \left[ \prod_{i=1}^q \prod_{j=1}^{m_i^T} p(y_{ij}^T; \mu_i^T, \sigma_i) \right] \left[ \prod_{i=1}^q \prod_{j=1}^{m_i^C} p(y_{ij}^C; \mu_i^C, \sigma_i) \right] \\ &\propto \left[ \prod_{i=1}^q \prod_{j=1}^{m_i^T} \frac{1}{\sigma_i} \exp \left\{ -\frac{1}{2} \left( \frac{y_{ij}^T - \mu_i^T}{\sigma_i} \right)^2 \right\} \right] \left[ \prod_{i=1}^q \prod_{j=1}^{m_i^C} \frac{1}{\sigma_i} \exp \left\{ -\frac{1}{2} \left( \frac{y_{ij}^C - \mu_i^C}{\sigma_i} \right)^2 \right\} \right] \\ &= \left[ \prod_{i=1}^q \sigma_i^{-m_i} \right] \exp \left\{ -\frac{1}{2} \sum_{i=1}^q \frac{1}{\sigma_i^2} \left[ \sum_{j=1}^{m_i^T} (y_{ij}^T - \mu_i^T)^2 + \sum_{j=1}^{m_i^C} (y_{ij}^C - \mu_i^C)^2 \right] \right\}. \end{aligned}$$

Let

$$(5.2.5) \quad \varphi(\theta) = \frac{1}{q} \sum_{i=1}^q \frac{\mu_i^T - \mu_i^C}{\sigma_i}$$

so that  $\delta = \varphi(\theta)$ . It follows that the MLE of  $\delta$  is simply  $\hat{\delta} = \varphi(\hat{\theta})$ .

The MLEs of  $\mu_i^T$  and  $\mu_i^C$  are

$$(5.2.6) \quad \hat{\mu}_i^T = \bar{y}_i^T = \frac{1}{m_i^T} \sum_{j=1}^{m_i^T} y_{ij}^T$$

and

$$(5.2.7) \quad \hat{\mu}_i^C = \bar{y}_i^C = \frac{1}{m_i^C} \sum_{j=1}^{m_i^C} y_{ij}^C.$$

Since  $\sigma_i$  represents the common standard deviation of both the treatment and control groups in study  $i$ , the standard way to estimate it is by pooling the sample variances of the two groups and then taking the

square root. Let  $s_i^{T^2}$  and  $s_i^{C^2}$  denote the sample variances of the treatment and control groups, so that

$$(5.2.8) \quad s_i^{T^2} = \frac{1}{m_i^T - 1} \sum_{j=1}^{m_i^T} (y_{ij}^T - \bar{y}_i^T)^2$$

and

$$(5.2.9) \quad s_i^{C^2} = \frac{1}{m_i^C - 1} \sum_{j=1}^{m_i^C} (y_{ij}^C - \bar{y}_i^C)^2.$$

Then the pooled estimator of the common standard deviation  $\sigma_i$  is

$$(5.2.10) \quad s_i = \left[ \frac{(m_i^T - 1)s_i^{T^2} + (m_i^C - 1)s_i^{C^2}}{m_i^T + m_i^C - 2} \right]^{\frac{1}{2}}.$$

Note, however, that these sample standard deviation estimators include Bessel's correction to account for the bias in the estimates given by maximum likelihood estimation, and hence are technically not MLEs themselves (though they are equivalent asymptotically). Since the ZSE parameter depends on the MLE of the parameter of interest, we will undo this correction in our estimate of  $\sigma_i$  in order to obtain its MLE. Therefore, we have the following:

$$(5.2.11) \quad \begin{aligned} \hat{\sigma}_i^{T^2} &= \frac{m_i^T - 1}{m_i^T} s_i^{T^2}; \\ \hat{\sigma}_i^{C^2} &= \frac{m_i^C - 1}{m_i^C} s_i^{C^2}; \\ \hat{\sigma}_i^2 &= \frac{m_i^T + m_i^C - 2}{m_i^T + m_i^C} s_i^2, \end{aligned}$$

and so

$$(5.2.12) \quad \hat{\sigma}_i = \sqrt{\frac{m_i - 2}{m_i}} s_i.$$

This allows us to rewrite Equation 5.2.2 as

$$(5.2.13) \quad d_i = \sqrt{\frac{m_i - 2}{m_i}} \frac{\bar{y}_i^T - \bar{y}_i^C}{\hat{\sigma}_i}.$$

Going forward, one of our objectives will be to find appropriate choices for  $\check{L}(u)\check{\pi}(u)$ ,  $Q(u)$ , and  $T_\delta(\omega)$  that allow us to form an approximation to the integrated likelihood function

$$\bar{L}(\delta) = \int_{\mathcal{U}} \frac{L(T_\delta(Q(u)))}{\check{L}(u)} \check{L}(u)\check{\pi}(u) du$$

using the procedure described in the previous chapter. We will also explore a similar procedure under the framework of a random-effects model in which we form approximations to the integrated likelihood for an SMD, as well as one for the heterogeneity variance parameter  $\tau^2$ .

## References

- Barndorff-Nielsen, O. E., and David R. Cox. 1996. “Prediction and Asymptotics.” *Bernoulli* 2 (4): 319–40. <http://www.jstor.org/stable/3318417>.
- Basu, Debabrata. 1977. “On the Elimination of Nuisance Parameters.” *Journal of the American Statistical Association* 72 (358): 355–66. <http://www.jstor.org/stable/2286800>.
- Berger, James O., Brunero Liseo, and Robert L. Wolpert. 1999. “Integrated Likelihood Methods for Eliminating Nuisance Parameters.” *Statistical Science* 14 (1): 1–22. <http://www.jstor.org/stable/2676641>.
- Cramér, Harald. 1945. “Section 33.3. Asymptotic Properties of Maximum Likelihood Estimates.” In *Mathematical Methods of Statistics*, 500–506. Princeton University Press.
- De Bin, Riccardo, Nicola Sartori, and Thomas A. Severini. 2015. “Integrated likelihoods in models with stratum nuisance parameters.” *Electronic Journal of Statistics* 9 (1): 1474–91. <https://doi.org/10.1214/15-EJS1045>.
- Hedges, Larry V., and I. Olkin. 1985. *Statistical Methods for Meta-Analysis*. Edited by Larry V. Hedges. Academic Press.
- Johnson, Steven G. 2021. “The NLOpt Nonlinear-Optimization Package.” <https://github.com/stevengi/nlopt>.
- Kalbfleisch, J. D., and D. A. Sprott. 1970. “Application of Likelihood Methods to Models Involving Large Numbers of Parameters.” *Journal of the Royal Statistical Society. Series B (Methodological)* 32 (2): 175–208. <http://www.jstor.org/stable/2984524>.
- . 1973. “Marginal and Conditional Likelihoods.” *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)* 35 (3): 311–28. <http://www.jstor.org/stable/25049882>.



- Liseo, Brunero. 1993. “Elimination of Nuisance Parameters with Reference Priors.” *Biometrika* 80 (2): 295–304. <http://www.jstor.org/stable/2337200>.
- Sartori, N. 2003. “Modified Profile Likelihoods in Models with Stratum Nuisance Parameters.” *Biometrika* 90 (3): 533–49. <http://www.jstor.org/stable/30042064>.
- Schumann, Martin, Thomas A. Severini, and Gautam Tripathi. 2021. “Integrated Likelihood Based Inference for Nonlinear Panel Data Models with Unobserved Effects.” *Journal of Econometrics* 223 (1): 73–95. <https://doi.org/10.1016/j.jeconom.2020.10.001>.
- . 2023. “The Role of Score and Information Bias in Panel Data Likelihoods.” *Journal of Econometrics* 235 (2): 1215–38. <https://doi.org/10.1016/j.jeconom.2022.08.011>.
- Severini, Thomas A. 2000. *Likelihood Methods in Statistics*. Oxford University Press.
- . 2007. “Integrated Likelihood Functions for Non-Bayesian Inference.” *Biometrika* 94 (3): 529–42. <http://www.jstor.org/stable/20441394>.
- . 2018. “Integrated Likelihoods for Functions of a Parameter.” *Stat* 7 (1): e212. <https://doi.org/10.1002/sta4.212>.
- . 2022. “Integrated Likelihood Inference in Multinomial Distributions.” *Metron*. <https://doi.org/10.1007/s40300-022-00236-x>.
- Statisticat, and LLC. 2021. *LaplacesDemon: Complete Environment for Bayesian Inference*. Bayesian-Inference.com. <https://web.archive.org/web/20150206004624/http://www.bayesian-inference.com/software>.
- Tierney, Luke, and Joseph B. Kadane. 1986. “Accurate Approximations for Posterior Moments and Marginal Densities.” *Journal of the American Statistical Association* 81 (393): 82–86. <http://www.jstor.org/stable/2287970>.
- Zejnnullahi, Rrita. 2021. “Topics in Meta-Analysis with Few Studies.” PhD thesis.

## APPENDIX A

## Chapter 2

## A.1. Definitions and Notation

## A.1.1. Open and Closed Balls

**A.1.1.1. Open Ball.** The *open ball of radius  $r > 0$*  centered at a point  $p \in \mathbb{R}^d$ , is the set of all points  $x \in \mathbb{R}^d$  such that the distance between  $p$  and  $x$  is less than  $r$ . We denote this set using the notation

$$B_r(\mathbf{p}) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{p}\| < r\},$$

where  $\|\cdot\|$  indicates the Euclidean norm, i.e.,  $\|\mathbf{x}\| = \left(\sum_{i=1}^d x_i^2\right)^{1/2}$  for  $\mathbf{x} \in \mathbb{R}^d$ .

**A.1.1.2. Closed Ball.** The *closed ball of radius  $r > 0$*  centered at a point  $\mathbf{p} \in \mathbb{R}^d$ , is the set of all points  $\mathbf{x} \in \mathbb{R}^d$  such that the distance between  $\mathbf{p}$  and  $\mathbf{x}$  is less than or equal to  $r$ . We denote this set using the notation

$$B_r[\mathbf{p}] = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{p}\| \leq r\}.$$

## A.1.2. Open and Closed Sets

**A.1.2.1. Open Set.** A subset  $S$  of  $\mathbb{R}^d$  is called an *open set* of  $\mathbb{R}^d$  if every point in  $S$  is the center of an open ball entirely contained in  $S$ . That is,  $S$  is open if and only if for any  $\mathbf{p} \in S$ , there exists a radius  $r > 0$  such that  $B_r(\mathbf{p}) \subseteq S$ .

**A.1.2.2. Closed Set.** A subset  $S$  of  $\mathbb{R}^d$  is called a *closed set* of  $\mathbb{R}^d$  if its complement  $S^c = \mathbb{R}^d \setminus S$  is open.

## A.1.3. Neighborhoods

A set  $N_{\mathbf{p}} \subseteq \mathbb{R}^d$  is called a *neighborhood* of a point  $\mathbf{p} \in \mathbb{R}^d$  if it contains an open ball centered at  $\mathbf{p}$ , i.e., for some radius  $r > 0$  there exists an open ball  $B_r(\mathbf{p})$  such that  $B_r(\mathbf{p}) \subseteq N_{\mathbf{p}}$ .

#### A.1.4. Boundedness

**A.1.4.1. Bounded Set.** A set  $S \subset \mathbb{R}^d$  is called *bounded* if there exists some radius  $r > 0$  such that  $B_r(\mathbf{0}) \subset S$ .

**A.1.4.2. Bounded Function.** A function  $f : X \rightarrow \mathbb{R}$  is called *bounded* if there exists a real number  $M$  such that  $|f(x)| \leq M$  for all  $x \in X$ .

#### A.1.5. Compact Set

A subset of  $\mathbb{R}^d$  is called *compact* if it is closed and bounded.

#### A.1.6. Interiority

**A.1.6.1. Interior Point.** A point  $\mathbf{p} \in S \subseteq \mathbb{R}^d$  is called an *interior point* of  $S$  if there exists some radius  $r > 0$  such that  $B_r(\mathbf{p}) \subseteq S$ .

**A.1.6.2. Interior of a Set.** The *interior of a set*  $S \subset \mathbb{R}^d$ , denoted by  $\text{int } S$ , is the set of all interior points of  $S$ .

#### A.1.7. Line Segments

For two points  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$ , a third point  $\tilde{\mathbf{x}}$  is said to be on the *line segment* connecting  $\mathbf{x}_1$  and  $\mathbf{x}_2$  if there exists  $\omega \in [0, 1]$  such that  $\tilde{\mathbf{x}} = \omega\mathbf{x}_1 + (1 - \omega)\mathbf{x}_2$ . We use the following notation to refer to such line segments:

$$\text{LS}(\mathbf{x}_1, \mathbf{x}_2) = \{\omega\mathbf{x}_1 + (1 - \omega)\mathbf{x}_2 : \omega \in [0, 1]\}.$$

#### A.1.8. Convexity and Concavity

**A.1.8.1. Convex Set.** A set  $S \subseteq \mathbb{R}^d$  is called *convex* if for any two points  $\mathbf{x}_1, \mathbf{x}_2 \in S$ , the line segment connecting  $\mathbf{x}_1$  and  $\mathbf{x}_2$  is entirely contained within  $S$ , i.e.,

$$\text{LS}(\mathbf{x}_1, \mathbf{x}_2) \subseteq S \text{ for all } \mathbf{x}_1, \mathbf{x}_2 \in S.$$

**A.1.8.2. Convex Function.** Let  $f : X \rightarrow \mathbb{R}$ , where  $X$  is a convex set.  $f$  is called a *convex function* if for all  $t \in [0, 1]$  and all  $x_1, x_2 \in X$ ,

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2).$$

If it is possible to graph the function on the coordinate plane, this is equivalent to saying that the line segment between any two distinct points on the graph of the function lies above the graph.

$f$  is called *strictly convex* if the equality is tightened, i.e.,

$$f(tx_1 + (1-t)x_2) < tf(x_1) + (1-t)f(x_2).$$

**A.1.8.3. Concave Function.** Let  $f : X \rightarrow \mathbb{R}$ , where  $X$  is a convex set.  $f$  is called a *concave function* if for all  $t \in [0, 1]$  and all  $x_1, x_2 \in X$ ,

$$f(tx_1 + (1-t)x_2) \geq tf(x_1) + (1-t)f(x_2).$$

If it is possible to graph the function on the coordinate plane, this is equivalent to saying that the line segment between any two distinct points on the graph of the function lies below the graph.

$f$  is called *strictly concave* if the equality is tightened, i.e.,

$$f(tx_1 + (1-t)x_2) > tf(x_1) + (1-t)f(x_2).$$

#### A.1.9. Positive Definiteness

A  $d \times d$  symmetric real matrix  $M$  is called *positive definite* if  $\mathbf{x}^\top M \mathbf{x} > 0$  for all non-zero  $\mathbf{x} \in \mathbb{R}^d$ .

### A.1.10. Derivatives of Multivariable Functions

**A.1.10.1. Gradient.** The *gradient* of a multivariable scalar-valued function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is given by its  $d \times 1$  vector of partial derivatives:

$$\nabla f(x_1, x_2, \dots, x_d) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{pmatrix}.$$

**A.1.10.2. Jacobian Matrix.** The *Jacobian matrix* of a multivariable vector-valued function  $\mathbf{f} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  given by

$$\mathbf{f}(x_1, x_2, \dots, x_d) = \begin{pmatrix} f_1(x_1, x_2, \dots, x_d) \\ f_2(x_1, x_2, \dots, x_d) \\ \vdots \\ f_k(x_1, x_2, \dots, x_d) \end{pmatrix}.$$

is defined as its  $k \times d$  matrix of partial derivatives:

$$\mathbf{J}(\mathbf{f}) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_d} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_k}{\partial x_1} & \frac{\partial f_k}{\partial x_2} & \dots & \frac{\partial f_k}{\partial x_d} \end{pmatrix} = \begin{pmatrix} \nabla^\top f_1 \\ \vdots \\ \nabla^\top f_k \end{pmatrix}$$

In the case where  $k = 1$ , the Jacobian matrix simply reduces to  $\nabla^\top f$ , the transpose of the gradient of  $\mathbf{f}$ .

**A.1.10.3. Hessian Matrix.** The Hessian matrix of a multivariable scalar-valued function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is given by its  $d \times d$  matrix of second partial derivatives:

$$\mathbf{H}(f) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_d} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1} & \frac{\partial^2 f}{\partial x_d \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_d^2} \end{pmatrix}$$

The *Hessian matrix* of  $f$  is also equal to the transpose of the Jacobian matrix of the gradient of  $f$ , i.e.,  $\mathbf{H}(f) = \mathbf{J}(\nabla f)^\top$ . When the order of differentiation does not matter, which occurs if and only if all of  $f$ 's second partial derivatives are continuous, both matrices become symmetric and we simply have  $\mathbf{H}(f) = \mathbf{J}(\nabla f)$ .

A twice differentiable function of several variables is strictly convex (concave) on a convex set if and only if its Hessian matrix is positive (negative) definite on the interior of the convex set.

## A.2. Theorems

### A.2.1. Jensen's Inequality

For a real-valued random variable  $X$  with finite expectation and a strictly concave function  $\varphi$ ,

$$\mathbb{E}[\varphi(X)] < \varphi(\mathbb{E}[X]).$$

### A.2.2. The Extreme Value Theorem

If  $K$  is a compact set and  $f : K \rightarrow \mathbb{R}$  is a continuous function, then  $f$  is bounded and there exists  $p, q \in K$  such that

$$f(p) = \sup_{x \in K} f(x)$$

and

$$f(q) = \inf_{x \in K} f(x).$$

### A.2.3. Taylor's Theorem

Suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is twice continuously differentiable in a neighborhood  $N_{\mathbf{x}_0}$  of some point  $\mathbf{x}_0 \in \mathbb{R}^d$ . Then for any  $\mathbf{x} \in N_{\mathbf{x}_0}$ , there exists  $\tilde{\mathbf{x}} \in \text{LS}(\mathbf{x}, \mathbf{x}_0)$  such that

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H}(f(\tilde{\mathbf{x}}))(\mathbf{x} - \mathbf{x}_0),$$

where  $\mathbf{H}(f(\tilde{\mathbf{x}}))$  denotes the Hessian matrix of  $f(\mathbf{x})$  evaluated at  $\mathbf{x} = \tilde{\mathbf{x}}$ .

This is not Taylor's theorem in its full generality, but rather a particular case of it that is well-suited for use in this paper. More accurately, we could refer to it as a second-order Taylor series expansion of a multivariable scalar-valued function with the Lagrange form of the remainder.

### A.2.4. Slutsky's Theorem

Let  $X_n, Y_n$  be sequences of scalar, vector, or matrix random variables. If  $X_n \xrightarrow{d} X$ , where  $X$  is another random variable, and  $Y_n \xrightarrow{p} c$ , where  $c$  is a constant, then

- i)  $X_n + Y_n \xrightarrow{d} X + c$ ;
- ii)  $X_n Y_n \xrightarrow{d} Xc$ ;
- iii)  $X_n/Y_n \xrightarrow{d} X/c$ , assuming  $c$  is invertible.

As an extension, suppose  $\mathbf{X}_n \xrightarrow{d} \mathbf{X}$ , where  $\mathbf{X}$  is a  $d \times 1$  vector,  $\mathbf{A}_n \xrightarrow{p} \mathbf{A}$ , where  $\mathbf{A}$  is a positive definite matrix, and  $\mathbf{X}_n = \mathbf{A}_n \mathbf{Y}_n$ . Then

$$\mathbf{Y}_n \xrightarrow{d} \mathbf{A}^{-1} \mathbf{X}.$$

## APPENDIX B

### Chapter 4

#### B.1. Desirable Properties of the Integrated Likelihood

Severini (2007) considered the following four statements to be essential properties that any integrated likelihood function must satisfy if it is to be useful for the purpose of non-Bayesian inference.

##### B.1.1. Property 1

Suppose the likelihood function for a parameter  $\boldsymbol{\theta} = (\psi, \lambda)$  can be decomposed as the product  $L(\boldsymbol{\theta}) = L_1(\psi)L_2(\lambda)$ . Then the integrated likelihood for  $\psi$  should satisfy

$$\bar{L}(\psi) = L_1(\psi).$$

##### B.1.2. Property 2

The frequency properties of an integrated likelihood function should mirror that of a genuine likelihood function as the sample size tends to infinity. In particular,  $\bar{L}(\psi)$  should be approximately score- and information-unbiased, i.e., asymptotically satisfy the first and second Bartlett identities, respectively.

##### B.1.3. Property 3

The integrated likelihood function should be insensitive to the choice of prior density.

##### B.1.4. Property 4

The integrated likelihood function should be invariant with respect to reparameterizations of the model that leave the parameter of interest unchanged.



## B.2. Laplace's Method

Let  $\theta$  be a scalar parameter taking values in  $\mathbb{R}$  and consider an integral of the form

$$I = \int_{-\infty}^{\infty} f(\theta) \exp[-nh(\theta)] d\theta.$$

Suppose the function  $-h(\theta)$  is smooth, bounded, and unimodal so that it attains a maximum at a point  $\hat{\theta}$ . Then Laplace's method states that an approximation for  $I$  is given by

$$\hat{I} = f(\hat{\theta}) \sqrt{\frac{2\pi}{n}} \sigma \exp[-nh(\hat{\theta})],$$

where

$$\sigma = \left[ \frac{\partial^2 h}{\partial \theta^2} \bigg|_{\theta=\hat{\theta}} \right]^{-1/2}.$$

It can further be shown that

$$I = \hat{I} \left\{ 1 + O\left(\frac{1}{n}\right) \right\},$$

where the term  $n$  may be interpreted as the sample size.

## APPENDIX C

## Chapter 5

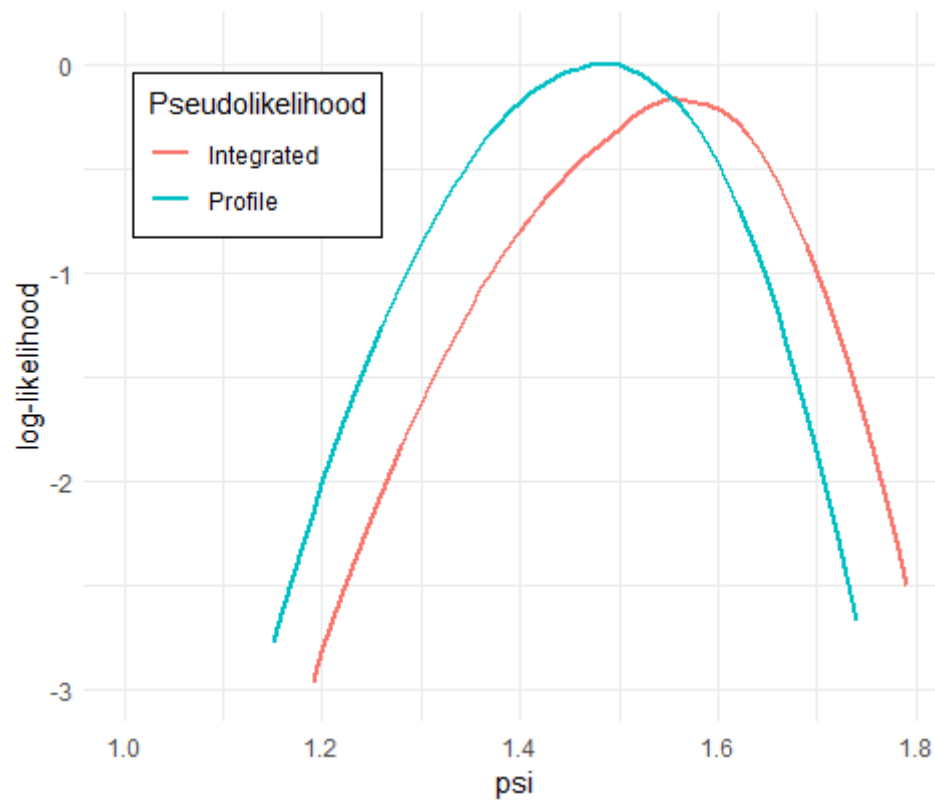


Figure C.1. Comparison between pseudolikelihoods for entropy of multinomial distribution