

ABSTRACT

Proposal for an Adjusted Numerical Approximation to the Integrated Likelihood Function

Timothy Ruel

The primary focus of this prospectus is to motivate and explain an adapted version of numerical likelihood integration as a method for eliminating nuisance parameters from a statistical model.

Table of Contents

ABSTRACT	1
Table of Contents	2
Chapter 1. Introduction and Motivation	4
Chapter 2. Experiments and Statistical Models	5
Chapter 3. The Likelihood Function	8
3.1. Definition	8
3.2. Transformations	9
3.3. Regularity Conditions	12
3.4. Maximum Likelihood Estimation	13
3.5. The Bartlett Identities	15
3.6. One-Index Asymptotics	18
Chapter 4. Pseudolikelihood Functions	21
4.1. Model Parameter Decomposition	21
4.2. Types of Pseudolikelihoods	23
4.3. The Bartlett Identities Revisited	24
Chapter 5. Approximating the Integrated Likelihood Function	25
5.1. The Zero-Score Expectation Parameter	25
5.2. Two-Index Asymptotics	27
Chapter 6. Applications	32

	3
6.1. Multinomial Distribution	32
6.2. Standardized Mean Difference	32
References	33
Appendix A. Chapter 3	35
A.1. Definitions and Notation	35
A.2. Taylor's Theorem	35
A.3. Jensen's Inequality	36
Appendix B. Chapter 5	37
B.1. Desirable Properties of the Integrated Likelihood	37
B.2. Laplace's Method	37

CHAPTER 1

Introduction and Motivation

CHAPTER 2

Experiments and 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* θ 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 θ , say θ_1 , we can refer to its corresponding distribution in \mathcal{P} with the notation P_{θ_1} , and therefore \mathcal{P} itself may be written as $\mathcal{P} = \{P_\theta | \theta \in \Theta\}$. Θ is called the *parameter space* and represents the set of all possible values θ can take on.

We will restrict our attention in this paper to distributions that are absolutely continuous with respect to some σ -finite measure μ , so that they admit a probability density function by the Radon-Nikodym Theorem. Let p_θ denote the density function associated with the distribution P_θ . 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 θ . 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_θ is always defined with respect to some dominating measure μ .

In general, a model’s parameterization is not unique, and for a given parameter θ , we are free to choose any invertible function of θ as a new parameter. Once we have made our choice of parameterization, we will assume that Θ does contain a singular true parameter value, which we will denote by θ_0 . The conventional interpretation of θ_0 is as a fixed but unknown constant that represents the value of the parameter corresponding to the true density function p_{θ_0} in \mathcal{P} . Conversely, θ represents an arbitrary parameter value that is allowed to range over all possible elements of Θ , including θ_0 . In other words, θ 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 θ_0 is like trying to figure out the particular value (θ_0) to which a population’s dial (θ) has been set. Note that it is possible for the value

of θ_0 itself to change over time as well, depending on the population. In such cases, any estimate of θ_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.

Crucially, it must always be possible to identify the parameter in our model on the basis of the data we observe. A model is considered *identifiable* if having perfect knowledge of the population would enable us to determine θ_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 θ 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 Θ 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 3

The Likelihood Function**3.1. 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 θ_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 θ_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 θ_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 θ .

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

$$(3.1.1) \quad L(\theta) = L(\theta; \mathbf{x}_n) = p(\mathbf{x}_n; \theta), \quad \theta \in \Theta.$$

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 θ , the reverse is actually true for $L(\theta; \mathbf{x}_n)$. Indeed, we regard the likelihood as being a function of the parameter θ for fixed \mathbf{x}_n . The reversal of the order of the arguments θ 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 θ is the true parameter value. Crucially, this is *not* equivalent to the inverse probability that θ is the true parameter value given $X = x$. The likelihood does not directly tell us anything about the probability that θ 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 θ may still be defined as it is in Equation 3.1.1. 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 neighborhood of 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 θ 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 θ . Despite our use of probability in its definition, the likelihood itself is *not* a probability density function for the parameter θ and is subject to neither the same rules nor interpretations as one.

3.2. 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:

$$(3.2.1) \quad \ell(\theta) = \ell(\theta; \mathbf{x}_n) = \log L(\theta; \mathbf{x}_n), \quad \theta \in \Theta.$$

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 θ 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 θ 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. Consequently, we will use gradient notation (i.e. ∇ , ∇^2 , etc.) to emphasize the vector nature of these results.

The gradient of ℓ with respect to θ 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

$$(3.2.2) \quad \mathcal{S}(\theta) = \mathcal{S}(\theta; \mathbf{x}_n) = \nabla \ell(\theta; \mathbf{x}_n) = \begin{pmatrix} \frac{\partial \ell}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell}{\partial \theta_d} \end{pmatrix}, \quad \theta \in \Theta.$$

Similarly, the negative Hessian matrix of the log-likelihood function with respect to θ is called the *observed information*, or just the *information*, and is denoted by

$$(3.2.3) \quad \mathcal{J}(\theta) = \mathcal{J}(\theta; \mathbf{x}_n) = -\nabla^2 \ell(\theta; \mathbf{x}_n) = - \begin{pmatrix} \frac{\partial^2 \ell}{\partial \theta_1^2} & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_d} \\ \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ell}{\partial \theta_d \partial \theta_1} & \frac{\partial^2 \ell}{\partial \theta_d \partial \theta_2} & \cdots & \frac{\partial^2 \ell}{\partial \theta_d^2} \end{pmatrix}, \quad \theta \in \Theta.$$

Note that this is also equal to the negative gradient of the score function. The use of the term “information” here derives from the fact that the second partial derivatives of ℓ with respect to the components of θ are all related to the curvature of ℓ near its maximum - the sharper the curve, the less uncertainty and therefore more information we have about θ .

Recall that $L(\theta; \mathbf{x}_n)$ is defined as a function of θ 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 $L(\theta; \mathbf{x}_n)$ as a random variable in the following sense: for a given θ , the value of $L(\theta; \mathbf{x}_n)$ depends entirely

on the values of X_1, \dots, X_n that we happened to observe, and so $L(\theta; \mathbf{X}_n)$ 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, $\mathcal{S}(\theta; \mathbf{X}_n)$ is a random variable for which we have observed the value $\mathcal{S}(\theta; \mathbf{x}_n)$.

The random nature of these likelihood-based quantities further implies that finding their expectations and variances with respect to $p_\theta(\mathbf{x}_n)$ 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 θ_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

$$(3.2.4) \quad \mathcal{J}_{\mathbf{X}_n}(\theta) = \text{Var}_\theta[\mathcal{S}(\theta; \mathbf{X}_n)], \quad \theta \in \Theta.$$

Since we are working in the more general framework in which $\mathcal{S}(\theta)$ 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 $\mathcal{S}(\theta)$. Hence, we have

$$\begin{aligned} \mathcal{J}_{\mathbf{X}_n}(\theta) &= \text{Var}_\theta[\mathcal{S}(\theta; \mathbf{X}_n)] && \text{(by Eq. 2.2.4)} \\ &= \text{Cov}_\theta \left[\left(\frac{\partial \ell}{\partial \theta_1}, \dots, \frac{\partial \ell}{\partial \theta_d} \right)^T \right] && \text{(by Eq. 2.2.2)} \\ (3.2.5) \quad &= \begin{pmatrix} \text{Var}_\theta \left(\frac{\partial \ell}{\partial \theta_1} \right) & \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_2} \right) & \dots & \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_1}, \frac{\partial \ell}{\partial \theta_d} \right) \\ \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_2}, \frac{\partial \ell}{\partial \theta_1} \right) & \text{Var}_\theta \left(\frac{\partial \ell}{\partial \theta_2} \right) & \dots & \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_2}, \frac{\partial \ell}{\partial \theta_d} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_d}, \frac{\partial \ell}{\partial \theta_1} \right) & \text{Cov}_\theta \left(\frac{\partial \ell}{\partial \theta_d}, \frac{\partial \ell}{\partial \theta_2} \right) & \dots & \text{Var}_\theta \left(\frac{\partial \ell}{\partial \theta_d} \right) \end{pmatrix}. \end{aligned}$$

Note that 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,

$$(3.2.6) \quad \mathcal{J}_{\mathbf{X}_n}(\theta) = \sum_{i=1}^n \mathcal{J}_{X_i}(\theta).$$

If the observations are also identically distributed according to the distribution of some random variable X , then $\mathcal{I}_{X_i}(\theta) = \mathcal{I}_X(\theta)$ 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 :

$$(3.2.7) \quad \mathcal{I}_{\mathbf{X}_n}(\theta) = n\mathcal{I}_X(\theta).$$

3.3. Regularity Conditions

As a consequence of the random variable interpretation of likelihood-based quantities such as the score function and the MLE for θ_0 , 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 maximum likelihood estimate converges as n tends toward infinity, if indeed one exists. 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 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) \mathcal{P} is identifiable, i.e., $P_{\theta_1} = P_{\theta_2} \implies \theta_1 = \theta_2$ for all $\theta_1, \theta_2 \in \Theta$.
- RC3) The distributions in \mathcal{P} have a common support \mathcal{X} not depending on θ .
- RC4) The parameter space Θ is a convex set.
- RC5) The Fisher information matrix $\mathcal{I}(\theta)$ is positive definite and has only finite components for all $\theta \in \Theta$.
- RC6) There exists an open set $\Theta^* \subseteq \Theta$ of which θ_0 is an interior point.
- RC7) For almost all x and for all $\theta \in \Theta$, $p(x; \theta)$ is twice continuously differentiable with respect to θ .
- RC8) The integral $\int_{-\infty}^{\infty} p(x; \theta) dx$ can be differentiated twice under the integral sign with respect to $\theta \in \Theta$.

RC9) There exists a random function $M(x)$ (that does not depend on θ) satisfying $E_{\theta_0}[M(X)] < \infty$ and a radius $r > 0$ such that for all $\theta \in N_r(\theta_0)$ and all integers $1 \leq i, j, k \leq d$, $|\nabla^3 \ell(\theta; x)_{ijk}| \leq M(x)$.

The phrase “almost all x ” in RC6) means that there may exist a set of points \mathcal{X} satisfying $\int_{\mathcal{X}} p_{\theta}(x) dx = 0$ for which $p(x; \theta)$ is not differentiable with respect to θ at any point $x \in \mathcal{X}$.

RC1) is included merely to simplify calculations and frame our discussion in the context of a standard case in which likelihood theory holds. RC2-3) guarantee that it is always possible to solve the likelihood equation for $\hat{\theta}$. RC4-6) ensure that there is no issue with defining derivatives of the log-likelihood function. Finally, by Leibniz’s integral rule RC7) implies that , and that the remainder term in the first-order Taylor series expansion of the score function centered at θ_0 is negligible.

3.4. Maximum Likelihood Estimation

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, θ_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 θ_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, it would seem that a sensible choice for an estimate of θ_0 is the value that maximizes the likelihood function based on an observed dataset \mathbf{x}_n . To make this argument rigorous, it suffices to show that with probability tending to 1 as the sample size tends toward infinity, the likelihood will be strictly larger at θ_0 than for any other $\theta \in \Theta$. We start

by observing that RC1) implies

$$(3.4.1) \quad L(\theta; \mathbf{x}_n) = p(\mathbf{x}_n; \theta) = \prod_{i=1}^n p(x_i; \theta)$$

and

$$(3.4.2) \quad \ell(\theta; \mathbf{x}_n) = \log p(\mathbf{x}_n; \theta) = \sum_{i=1}^n \log p(x_i; \theta).$$

It follows that

$$(3.4.3) \quad \begin{aligned} L(\theta; \mathbf{x}_n) < L(\theta_0; \mathbf{x}_n) &\iff \ell(\theta; \mathbf{x}_n) < \ell(\theta_0; \mathbf{x}_n) \\ &\iff \sum_{i=1}^n \log p(x_i; \theta) - \sum_{i=1}^n \log p(x_i; \theta_0) < 0 \\ &\iff \sum_{i=1}^n [\log p(x_i; \theta) - \log p(x_i; \theta_0)] < 0 \\ &\iff \sum_{i=1}^n \log \frac{p(x_i; \theta)}{p(x_i; \theta_0)} < 0 \\ &\iff \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i; \theta)}{p(x_i; \theta_0)} < 0. \end{aligned}$$

Note that RC3) guarantees that the ratio $p(x; \theta)/p(x; \theta_0)$ is well-defined and finite for all $x \in \mathcal{X}$, the region of common support. Then by the Weak Law of Large Numbers,

$$(3.4.4) \quad \frac{1}{n} \sum_{i=1}^n \log \frac{p(X_i; \theta)}{p(X_i; \theta_0)} \rightarrow E_{\theta_0} \left[\log \frac{p(X; \theta)}{p(X; \theta_0)} \right]$$

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

$$(3.4.5) \quad E_{\theta_0} \left[\frac{p(X; \theta)}{p(X; \theta_0)} \right] = \int_{\mathcal{X}} \left[\frac{p(x; \theta)}{p(x; \theta_0)} \right] p(x; \theta_0) dx = \int_{\mathcal{X}} p(x; \theta) dx = 1.$$

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

$$(3.4.6) \quad E_{\theta_0} \left[\log \frac{p(X; \theta)}{p(X; \theta_0)} \right] < \log E_{\theta_0} \left[\frac{p(X; \theta)}{p(X; \theta_0)} \right] = \log 1 = 0.$$

Hence, the quantity on the left-hand side of Equation 3.4.4 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 3.4.3, it therefore follows that

$$(3.4.7) \quad \lim_{n \rightarrow \infty} P_{\theta_0} [L(\theta; \mathbf{x}_n) < L(\theta_0; \mathbf{x}_n)] = \lim_{n \rightarrow \infty} P_{\theta_0} \left[\frac{1}{n} \sum_{i=1}^n \log \frac{p(X_i; \theta)}{p(X_i; \theta_0)} < 0 \right] = 1,$$

and the result is proved.

Let $\hat{\theta} = \hat{\theta}(\mathbf{x}_n)$ be the value that maximizes the likelihood at the observed $\mathbf{X}_n = \mathbf{x}_n$, i.e.,

$$(3.4.8) \quad \hat{\theta}(\mathbf{x}_n) = \arg \max_{\theta \in \Theta} L(\theta; \mathbf{x}_n).$$

If $\hat{\theta}(\mathbf{X}_n)$ is unique, it is called the *maximum likelihood estimator* (MLE) of θ_0 . For an arbitrarily chosen model, there is no guarantee that $\hat{\theta}(\mathbf{x}_n)$ will exist, and even if it does, it will not necessarily be unique. We say that the MLE is undefined in these cases. However, it can be shown that for parametric models satisfying RC1-3), the MLE exists, is unique with probability tending to 1, and is consistent for θ_0 . We will explore more of the asymptotic properties of the MLE in Section 3.6.

3.5. 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 θ will satisfy the Bartlett identities. It is guaranteed, however, that the log-likelihood function associated with X and θ will satisfy them, provided that the model is regular. Thus, 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 θ is a scalar. For a single observation $X = x$, the expectation of $\frac{\partial}{\partial \theta} \ell(\theta; X)$ gives

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

Equation 3.5.1 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 θ may be rewritten as

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

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

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

Therefore,

$$(3.5.2) \quad \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \ell(\theta; X) \right] + \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \ell(\theta; X) \right)^2 \right] = 0.$$

Equation 3.5.2 is called the second Bartlett identity. It follows that, for models satisfying RC1-3) the expected information can be rewritten as

$$(3.5.3) \quad \mathbb{E}_\theta \left[-\frac{\partial^2}{\partial \theta^2} \ell(\theta; X) \right] = \mathbb{V}[\ell_\theta(\theta; x); \theta].$$

Any function that satisfies the second Bartlett identity is said to be *information-unbiased*.

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

3.6. One-Index Asymptotics

The one-index asymptotics framework describes the behavior of likelihood-based statistics as the sample size (n) grows to infinity while the dimension of the nuisance parameter (q) remains fixed. The aim of this section is to present a basic overview of the theory's results so that we will have a readily available baseline against which to compare the results of the following section discussing the two-index asymptotics framework, in which q is allowed to increase with n .

Assume the regularity conditions of the previous section apply to our model and let $\hat{\theta}$ denote the MLE for θ_0 . The traditional method for analyzing the asymptotic behavior of $\hat{\theta}$ is to use a first-order Taylor series expansion of the score function.¹

To do this, first note that the likelihood function for θ based on $\mathbf{x}_n = (x_1, \dots, x_n)$ as

$$L(\theta; \mathbf{x}_n) = p_\theta(\mathbf{x}_n) = \prod_{i=1}^n p_\theta(x_i).$$

¹See Appendix A for a review of Taylor's theorem and the conditions under which it is satisfied.

We may therefore write the log-likelihood function as

$$\begin{aligned}
 \ell(\theta; \mathbf{x}_n) &= \log L(\theta; \mathbf{x}_n) \\
 &= \log p_\theta(\mathbf{x}_n) \\
 &= \log \left[\prod_{i=1}^n p_\theta(x_i) \right] \\
 &= \sum_{i=1}^n \log p_\theta(x_i) \\
 &= \sum_{i=1}^n \ell(\theta; x_i).
 \end{aligned}$$

It follows that the score function is equal to

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

In other words, the score function for the parameter θ based on data x_1, \dots, x_n can be written as the sum of individual contributions $\mathcal{S}(\theta; x_i)$, $i = 1, \dots, n$, where each $\mathcal{S}(\theta; x_i)$ can be thought of as the score function for θ had we drawn only one observation, x_i . Note that these individual contributions are all independent from one another, a consequence of the independence assumption in RC1).

Equation 3.6.1 implies that a Taylor series expansion of $\mathcal{S}(\theta; \mathbf{x}_n)$ will be equal to the sum of the Taylor series expansions of its individual contributions, plus a remainder term that grows with n . We have assumed each x_i is identically distributed, so it suffices to consider the expansion for an arbitrary contribution, $\mathcal{S}(\theta; x_i)$.

Since θ is a $d \times 1$ vector, the k -th derivative of the log-likelihood with respect to θ will be a k -dimensional array having d entries along each of its k indices. In particular, the score function (i.e. the

first derivative of the log-likelihood) will also be a $d \times 1$ vector,

$$\mathcal{S}(\theta; x_i) = \begin{pmatrix} \mathcal{S}_1(\theta; x_i) \\ \vdots \\ \mathcal{S}_d(\theta; x_i) \end{pmatrix},$$

where each component is a function $\mathcal{S}_j : \Theta \rightarrow \mathbb{R}$. Similarly, the first and second derivatives of the score function will be a $d \times d$ matrix and a three-dimensional $d \times d \times d$ array, respectively. To simplify notation, we will perform a component-wise first-order Taylor series expansion of the score function around the point $\theta = \theta_0$ wherein each component $\mathcal{S}_j(\theta; x_i)$ of $\mathcal{S}(\theta; x_i)$ is expanded separately.

For any $\theta \in N_r(\theta_0)$, there exists $\bar{\theta}_j$ on the line segment connecting θ and θ_0 such that

$$\begin{aligned} \mathcal{S}_j(\theta; x_i) &= \mathcal{S}_j(\theta_0; x_i) + \nabla \mathcal{S}_j(\theta_0; x_i)^T (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^T \nabla^2 \mathcal{S}_j(\bar{\theta}_j; x_i) (\theta - \theta_0) \\ &= \mathcal{S}_j(\theta_0; x_i) + [\nabla \mathcal{S}_j(\theta_0; x_i) + \frac{1}{2} \nabla^2 \mathcal{S}_j(\bar{\theta}_j; x_i) (\theta - \theta_0)]^T (\theta - \theta_0) \\ &= \mathcal{S}_j(\theta_0; x_i) + [\nabla \mathcal{S}_j(\theta_0; x_i) + M(x_i) O(\|\theta - \theta_0\|)]^T (\theta - \theta_0), \end{aligned}$$

where the last equality follows as a result of RC7). Summing over each of the individual contributions to the score function yields

$$\begin{aligned} \mathcal{S}_j(\theta; \mathbf{x}_n) &= \sum_{i=1}^n \mathcal{S}_j(\theta; x_i) \\ &= \sum_{i=1}^n [\mathcal{S}_j(\theta_0; x_i) + [\nabla \mathcal{S}_j(\theta_0; x_i) + M(x_i) O(\|\theta - \theta_0\|)]^T (\theta - \theta_0)] \\ &= \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left[\nabla \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left\{ \sum_{i=1}^n M(x_i) \right\} O(\|\theta - \theta_0\|) \right]^T (\theta - \theta_0) \\ &= \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left[\frac{1}{n} \nabla \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O(\|\theta - \theta_0\|) \right]^T n(\theta - \theta_0). \end{aligned}$$

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

$$\frac{1}{\sqrt{n}} \mathcal{S}_j(\theta; \mathbf{x}_n) = \frac{1}{\sqrt{n}} \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left[\frac{1}{n} \nabla \mathcal{S}_j(\theta_0; \mathbf{x}_n) + \left\{ \frac{1}{n} \sum_{i=1}^n M(x_i) \right\} O(\|\theta - \theta_0\|) \right]^T \sqrt{n}(\theta - \theta_0).$$

CHAPTER 4

Pseudolikelihood Functions

4.1. 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 ψ taking values in a set $\Psi \subseteq \mathbb{R}^p$, where $p < d$. In such an event, we refer to ψ as the *parameter of interest*. Crucially, as we will see, ψ can always be expressed as a function of θ .

Since ψ is of lower dimension than θ , it necessarily follows that there is another parameter λ , 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 θ . We refer to λ 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.¹ 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 ψ 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 μ_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 ψ . We will examine some of the issues posed by high-dimensional nuisance parameters in greater detail in the next chapter.

¹Note that nuisance parameters are not always uniquely defined. Depending on the choice of parameter of interest, there may be multiple or even infinite ways to define a nuisance parameter.

4.1.1. 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 ψ to be explicit and λ to be implicit, or vice versa.

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

$$(4.1.1) \quad \psi = (\theta_{I_1}, \dots, \theta_{I_p}).$$

It immediately follows that λ is the sub-vector of all components of θ that are not part of ψ . 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

$$(4.1.2) \quad \lambda = (\theta_{J_1}, \dots, \theta_{J_q}).$$

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

4.1.2. Implicit Parameters

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

$$(4.1.3) \quad \psi = \varphi(\theta).$$

As before, Ψ 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 λ in terms of the original components of θ need not exist. For this reason, implicit nuisance parameters are in general more difficult to eliminate compared to their explicit counterparts.

Note that when the parameter of interest and nuisance parameter are explicit, it is always possible to define a function φ such that

$$(4.1.4) \quad \varphi(\theta) = (\theta_{I_1}, \dots, \theta_{I_p}) \equiv \psi,$$

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 $\psi = \varphi(\theta)$. With this understanding in mind, we will use the notation $\psi = \varphi(\theta)$ 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.

4.2. Types of Pseudolikelihoods

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

Any summary of the values in \mathcal{L}_ψ that does not depend on λ theoretically constitutes a pseudolikelihood function for ψ . 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. We will explore each of these methods in more detail in the following sections.

4.2.1. The Profile Likelihood

The profile likelihood is the most straightforward method for eliminating a nuisance parameter from a likelihood function.

For example, suppose we are interested in estimating the mean of a random variable Y , where $Y \sim N(\mu, \sigma^2)$. The full model parameter is $\theta = (\mu, \sigma^2)$ but since we are only interested in estimating the mean, the parameter of interest is $\psi = \mu$ and the nuisance parameter is $\lambda = \sigma^2$.

4.2.2. The Conditional Likelihood

4.2.3. The Marginal Likelihood

4.2.4. The Integrated Likelihood

4.3. The Bartlett Identities Revisited

The Bartlett identities offer an alternative way of characterizing the difference between likelihood and pseudolikelihood functions. A genuine likelihood function of θ can be characterized as any nonnegative random function of θ for which all of the Bartlett identities hold. Similarly, we can think of a pseudolikelihood function of θ as being any nonnegative random function of θ 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.

CHAPTER 5

Approximating the Integrated Likelihood Function

5.1. The Zero-Score Expectation Parameter

Let $\psi = \varphi(\theta)$ and λ denote the parameter of interest and nuisance parameter, respectively, for some statistical model $(\mathcal{S}, \mathcal{P}_\theta)$. Then the general expression to obtain an integrated likelihood for ψ may be written as

$$(5.1.1) \quad \bar{L}(\psi) = \int_{\Lambda} L(\psi, \lambda) \pi(\lambda|\psi) d\lambda,$$

where $\pi(\lambda|\psi)$ is a conditional prior density for λ given ψ .

Severini (2007) considered the problem of selecting $\pi(\lambda|\psi)$ such that when the likelihood function is integrated with respect to this density, the result is useful for non-Bayesian inference. To do this, he outlined four properties (see Appendix B) that an integrated likelihood function must satisfy if it is to be of any use. He went on to prove that an integrated likelihood satisfying these properties could be obtained by first constructing a new nuisance parameter $\phi \in \Phi$ that is unrelated to the parameter of interest (in the sense that its maximum likelihood estimator remains roughly constant for all values of ψ) and then choosing a prior density $\pi(\phi)$ that is independent of ψ . Once chosen, the desired integrated likelihood function for ψ is given by

$$(5.1.2) \quad \bar{L}(\psi) = \int_{\Phi} \tilde{L}(\psi, \phi) \pi(\phi) d\phi,$$

where $\tilde{L}(\psi, \phi)$ is the likelihood function for the model after it has been reparameterized in terms of ϕ . It is important to note that the exact choice of prior density for ϕ is not particularly important; the only restriction we place upon it is that it must not depend on ψ .

Suppose that we have an explicit parameter of interest and nuisance parameter, so that $\theta = (\psi, \lambda)$. Then Severini (2007) defines this new nuisance parameter ϕ as the solution to the equation

$$(5.1.3) \quad \mathbb{E}(\ell_\lambda(\psi, \lambda); \psi_0, \lambda_0) \Big|_{(\psi_0, \lambda_0) = (\hat{\psi}, \phi)} = 0,$$

where $\ell_\lambda(\psi, \lambda) = \frac{\partial \ell(\psi, \lambda)}{\partial \lambda}$, ψ_0 and λ_0 denote the true values of ψ and λ , and $\hat{\psi}$ is the MLE for ψ_0 . In other words, for a particular value of $(\psi, \lambda, \hat{\psi})$, we can find the corresponding value of ϕ by solving for it in Equation 5.1.3. ϕ is called the *zero-score expectation* (ZSE) parameter because it is defined as the value that makes the expectation of the score function (where the derivative is taken with respect to λ) evaluated at the point $(\hat{\psi}, \phi)$ equal to zero. Note that ϕ is a function of the data through $\hat{\psi}$. Normally we 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 ϕ that depends on the data to parameterize it.

For a given value of $(\psi, \phi, \hat{\psi})$, it is also possible to solve Equation 5.1.3 for λ . This allows us to write Equation 5.1.2 in terms of $L(\psi, \lambda)$:

$$(5.1.4) \quad \bar{L}(\psi) = \int_{\Phi} L(\psi, \lambda(\psi, \phi)) \pi(\phi) d\phi.$$

Severini (2018) proved that reparameterizing the nuisance parameter in terms of the ZSE parameter yields the same desirable properties in the subsequent integrated likelihood when ψ and λ are implicit. Suppose $\psi = \varphi(\theta)$, for some function $\varphi : \Theta \rightarrow \Psi$, and consider the set of all values of θ satisfying $\varphi(\theta) = \hat{\psi}$. Call this set $\Omega_{\hat{\psi}}$ so that

$$(5.1.5) \quad \Omega_{\hat{\psi}} = \left\{ \omega \in \Theta : \varphi(\omega) = \hat{\psi} \right\}.$$

Elements of $\Omega_{\hat{\psi}}$ take the form $(\hat{\psi}, \phi)$, where $\phi \in \Lambda$.

5.1.1. Weight Functions

5.2. 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 ψ 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 ψ 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 ψ and the individual components of λ 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

$$(5.2.1) \quad X_{ij} \sim p_{ij}(x_{ij}; \psi, \lambda_i),$$

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

¹It 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.

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 5.2.1 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

$$(5.2.2) \quad L^{(i)}(\psi, \lambda_i) = p_i(\mathbf{x}_i; \psi, \lambda_i),$$

and

$$(5.2.3) \quad \ell^{(i)}(\psi, \lambda_i) = \log L^{(i)}(\psi, \lambda_i),$$

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

$$(5.2.4) \quad \bar{L}^{(i)}(\psi) = \int_{\Lambda} L^{(i)}(\psi, \lambda_i) g(\lambda_i; \psi) d\lambda_i.$$

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 5.2.4 as

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

One consequence of the regularity conditions we have assumed is that $L^{(i)}(\psi, \lambda_i)$ is that an MLE for θ_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 ψ - 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 λ_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

$$(5.2.5) \quad \hat{\bar{L}}^{(i)}(\psi) = f(\hat{\lambda}_{i\psi}) \sqrt{\frac{2\pi}{m}} \sigma \exp[-mh(\hat{\lambda}_{i\psi})],$$

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 λ_i only (i.e. the negative second partial derivative of the log-likelihood with respect to λ_i) evaluated at $\hat{\lambda}_{i\psi}$. Plugging in the appropriate quantities for f , h , and σ into Equation 5.2.5, we arrive at

$$\begin{aligned} (5.2.6) \quad \hat{\bar{L}}^{(i)}(\psi) &= f(\hat{\lambda}_{i\psi}) \sqrt{\frac{2\pi}{m}} \sigma \exp[-mh(\hat{\lambda}_{i\psi})] \\ &= g(\hat{\lambda}_{i\psi}; \psi) \sqrt{\frac{2\pi}{m}} \left[\frac{1}{m} \mathcal{J}(\hat{\lambda}_{i\psi}) \right]^{-1/2} \exp \left\{ -m \cdot -\frac{1}{m} \ell^{(i)}(\psi, \hat{\lambda}_{i\psi}) \right\} \\ &= \frac{\sqrt{2\pi}}{m} L^{(i)}(\psi, \hat{\lambda}_{i\psi}) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2} \\ &= \frac{\sqrt{2\pi}}{m} L_P^{(i)}(\psi) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2}, \end{aligned}$$

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

$$(5.2.7) \quad \bar{L}^{(i)}(\psi) = \hat{\bar{L}}^{(i)}(\psi) \left\{ 1 + O\left(\frac{1}{m}\right) \right\} \quad \text{as } m \rightarrow \infty.$$

Let

$$(5.2.8) \quad \bar{\ell}^{(i)}(\psi) = \log \bar{L}^{(i)}(\psi)$$

denote the integrated log-likelihood for ψ . Putting the results in Equation 5.2.6, Equation 5.2.7, and Equation 5.2.8 together, we have

$$\begin{aligned}
\bar{\ell}^{(i)}(\psi) &= \log \bar{L}^{(i)}(\psi) && \text{(by Equation 5.2.8)} \\
&= \log \left(\hat{\bar{L}}^{(i)}(\psi) \left\{ 1 + O\left(\frac{1}{m}\right) \right\} \right) && \text{(by Equation 5.2.7)} \\
&= \log \hat{\bar{L}}^{(i)}(\psi) + \log \left\{ 1 + O\left(\frac{1}{m}\right) \right\} && \text{(by Equation 5.2.6)} \\
&= \log \left\{ \frac{\sqrt{2\pi}}{m} L_P^{(i)}(\psi) g(\hat{\lambda}_{i\psi}; \psi) [\mathcal{J}(\hat{\lambda}_{i\psi})]^{-1/2} \right\} + O\left(\frac{1}{m}\right) \\
&= \frac{1}{2} \log(2\pi) - \log(m) + \log L_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{1}{m}\right) \\
&= \ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}) + \frac{1}{2} \log(2\pi) - \log(m) + O\left(\frac{1}{m}\right) \quad \text{as } m \rightarrow \infty,
\end{aligned}$$

where $\ell_P^{(i)}(\psi) = \ell^{(i)}(\psi, \hat{\lambda}_{i\psi})$ is the profile log-likelihood for ψ . 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 :

$$(5.2.9) \quad \hat{\bar{\ell}}^{(i)}(\psi) = \ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \log \mathcal{J}(\hat{\lambda}_{i\psi}).$$

The error in this approximation is given by

$$(5.2.10) \quad \bar{\ell}^{(i)}(\psi) = \hat{\bar{\ell}}^{(i)}(\psi) + O\left(\frac{1}{m}\right).$$

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

$$(5.2.11) \quad L(\psi, \lambda) = \prod_{i=1}^q L^{(i)}(\psi, \lambda_i)$$

and

$$(5.2.12) \quad \ell(\psi, \lambda) = \sum_{i=1}^q \ell^{(i)}(\psi, \lambda_i),$$

respectively. If we define the weight function

$$(5.2.13) \quad G(\lambda; \psi) \equiv \prod_{i=1}^q g(\lambda_i; \psi)$$

then the integrated likelihood function for ψ becomes separable. That is,

$$(5.2.14) \quad \begin{aligned} \bar{L}(\psi) &= \int_{\Lambda^q} L(\psi, \lambda) G(\lambda; \psi) d\lambda \\ &= \int_{\Lambda} \cdots \int_{\Lambda} \left[\prod_{i=1}^q L^{(i)}(\psi, \lambda_i) \right] \left[\prod_{i=1}^q g(\lambda_i; \psi) \right] d\lambda_1 \cdots d\lambda_q \\ &= \prod_{i=1}^q \int_{\Lambda} L^{(i)}(\psi, \lambda_i) g(\lambda_i; \psi) d\lambda_i \\ &= \prod_{i=1}^q \bar{L}^{(i)}(\psi). \end{aligned}$$

Let $\bar{\ell}(\psi) = \log \bar{L}(\psi)$ denote the integrated log-likelihood function for ψ . Taking the logarithm of both sides in Equation 5.2.14, we have

$$(5.2.15) \quad \bar{\ell}(\psi) = \sum_{i=1}^q \bar{\ell}^{(i)}(\psi).$$

Plugging in our approximation to $\bar{\ell}^{(i)}(\psi)$ and its error term in Equation 5.2.9 and Equation 5.2.10, respectively, yields

$$(5.2.16) \quad \begin{aligned} \bar{\ell}(\psi) &= \sum_{i=1}^q \left[\ell_P^{(i)}(\psi) + \log g(\hat{\lambda}_{i\psi}; \psi) - \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{1}{m}\right) \right] \\ &= \ell_P(\psi) + \sum_{i=1}^q \log g(\hat{\lambda}_{i\psi}; \psi) - \frac{1}{2} \sum_{i=1}^q \log \mathcal{J}(\hat{\lambda}_{i\psi}) + O\left(\frac{q}{m}\right) \quad \text{as } m \rightarrow \infty. \end{aligned}$$

CHAPTER 6

Applications**6.1. Multinomial Distribution****6.2. Standardized Mean Difference**

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>.
- 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>.
- Kalbfleisch, J. D., and D. A. Sprott. 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>.
- 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>.

APPENDIX A

Chapter 3

A.1. Definitions and Notation

A.1.1. Neighborhoods

A *neighborhood* of a point $\mathbf{p} \in \mathbb{R}^d$ is the set of all points $\mathbf{x} \in \mathbb{R}^d$ such that the Euclidean distance between \mathbf{p} and \mathbf{x} is less than some radius $r > 0$. We use the following notation to refer to such neighborhoods:

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

A.1.2. Line Segments

For two points $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$, a third point $\bar{\mathbf{x}}$ 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 $\bar{\mathbf{x}} = \omega\mathbf{x}_1 + (1 - \omega)\mathbf{x}_2$. We use the following notation to refer to such line segments:

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

A.2. Taylor's Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function that is $(k + 1)$ -times continuously differentiable in a neighborhood $N_r(\mathbf{x}_0)$ of some point $\mathbf{x}_0 \in \mathbb{R}^d$ and suppose there exists M satisfying $|D^\alpha f| \leq M$ for all $x \in N_r(x_0)$ and all α such that $|\alpha| = k + 1$. Then

$$f(\mathbf{x}) = \sum_{0 \leq |\alpha| \leq k} \frac{D^\alpha f(\mathbf{x}_0)}{\alpha!} (\mathbf{x} - \mathbf{x}_0)^\alpha + R_k(\mathbf{x}),$$

where the remainder term $R_k(\mathbf{x})$ satisfies

$$|R_k(\mathbf{x})| \leq \frac{M}{\boldsymbol{\alpha}!} |\mathbf{x} - \mathbf{x}_0|^\alpha$$

for all $\boldsymbol{\alpha}$ such that $|\boldsymbol{\alpha}| = k + 1$

A.3. Jensen's Inequality

Jensen's inequality states that for a real-valued random variable X with finite expectation and a strictly concave function φ ,

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

APPENDIX B

Chapter 5

B.1. Desirable Properties of the Integrated Likelihood**B.1.1. Property 1**

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

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

B.1.2. Property 2**B.1.3. Property 3****B.1.4. Property 4****B.2. Laplace's Method**

Let θ 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.