

Research Article

Linear Models with Response Functions Based on the Laplace Distribution: Statistical Formulae and Their Application to Epigenomics

C. Z. W. Hassell Sweatman,^{1,2} G. C. Wake,^{1,2,3} A. B. Pleasants,^{2,3,4} C. A. McLean,² and A. M. Sheppard^{2,3}

¹ INMS, Massey University, Albany Campus, Private Bag 102-904, North Shore Mail Centre, Auckland 0745, New Zealand

² Liggins Institute, The University of Auckland, Private Bag 92019, Victoria Street West, Auckland 1142, New Zealand

³ Gravida, National Centre for Growth and Development, The University of Auckland, Private Bag 92019,

Victoria Street West, Auckland 1142, New Zealand

⁴ Ruakara Research Centre, 10 Bisley Road, Private Bag 3115, Hamilton 3240, New Zealand

Correspondence should be addressed to G. C. Wake; g.c.wake@massey.ac.nz

Received 20 June 2013; Accepted 24 July 2013

Academic Editors: A. A. Ding and J. López-Fidalgo

Copyright © 2013 C. Z. W. Hassell Sweatman et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The statistical application considered here arose in epigenomics, linking the DNA methylation proportions measured at specific genomic sites to characteristics such as phenotype or birth order. It was found that the distribution of errors in the proportions of chemical modification (methylation) on DNA, measured at CpG sites, may be successfully modelled by a Laplace distribution which is perturbed by a Hermite polynomial. We use a linear model with such a response function. Hence, the response function is known, or assumed well estimated, but fails to be differentiable in the classical sense due to the modulus function. Our problem was to estimate coefficients for the linear model and the corresponding covariance matrix and to compare models with varying numbers of coefficients. The linear model coefficients may be found using the (derivative-free) simplex method, as in quantile regression. However, this theory does not yield a simple expression for the covariance matrix of the coefficients of the linear model. Assuming response functions which are \mathscr{C}^2 except where the modulus function attains zero, we derive simple formulae for the covariance matrix and a log-likelihood ratio statistic, using generalized calculus. These original formulae enable a generalized analysis of variance and further model comparisons.

1. Introduction and Motivation

This work arose in a biological context, in epigenomics, namely, the modelling of the distribution of errors in the proportions of chemical modification (methylation) on DNA, measured at specific genomic sites (CpG sites). It was observed that this error distribution may be suitably modelled by a truncated Laplace distribution perturbed by a Hermite polynomial.

This error distribution was first noticed in Sequenom measurements but has wider application. A survey of data generated by measurements on the Infinium, Illumina, Affymetrix, and MeDIP2 machines showed similar characteristics to that of the Sequenom, where such an amended Laplace distribution was required to properly describe the probability density function. It is thought that variation in the scattering angle of light in the measurement processes common to all of these platforms is responsible for the frequencies in the tails of the measurement distributions not conforming to a simple Laplace density and requiring our proposed amendment. Without the amendment the Laplace density gives tail probabilities for the deviations that are too high, potentially leading to an incorrect failure to reject a null hypothesis. Because the observed frequency distribution of epigenomic and gene expression measurements appears to be a common feature of molecular biology, it is important that the process of estimation and inference under the amended Laplace probability density be studied. The paper reports results from a study of estimation and inference under the amended Laplace density.

We extend the theory of linear models as given in [1] to deal with response variables with distributions more general than the exponential family. We consider the Laplace distribution, and amendments thereof, with probability density functions which have abrupt changes in gradient due to the presence of the modulus function. The theory in this paper corresponds to least absolute error (LAE) (or least absolute deviation (LAD)) regression [2–4], also called median regression [5], when the response function is the Laplace distribution without modification. We focus on coefficient estimation for a linear model with a response variable distribution assumed to be a truncated and/or perturbed version of the Laplace distribution, estimating the standard errors of these coefficients and understanding the asymptotic theory.

The theory of generalized linear models as described in [1] covers the case of distributions from the exponential family. These distributions have probability density functions which are twice continuously differentiable (\mathscr{C}^2), everywhere on their support. The usual expressions for the standard errors of the model coefficients for the generalized linear models in [1] are derived using Taylor series and assume distributions with probability density functions which are \mathscr{C}^2 everywhere on their support. They cannot be applied to our model due to the presence of the modulus function.

The modified Laplace probability density functions considered here have a sharp peak at the maximum. Maximum likelihood estimation (MLE) of coefficients may be done by non-gradient methods, such as the simplex method. However, the usual classical expressions for the standard errors of the coefficients, the information matrix, and the loglikelihood ratio statistic do not apply due to lack of differentiability. We derive expressions for generalized versions of these quantities using generalized functions. Consequently, we show that our MLE is asymptotically normal.

The method we present to estimate these statistics could in principle be applied to other probability density functions exhibiting abrupt changes in gradient. Response function parameters are assumed known or previously estimated. The theory is applied to find the standard errors for coefficients of a linear model, assuming the response function has a truncated Laplace distribution with added kurtosis, due to perturbation by a Hermite polynomial. To illustrate the application, we show how birth order can be linked to methylation status at two CpG sites in the promotor of the H19 gene.

2. The Model

2.1. The Expectation Is Linear. Let

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^T \in \mathbb{R}^n$$
(1)

be a vector of response variables; let

$$\mathbf{X} = \mathbf{X}_{n,m} = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1m} \\ 1 & x_{22} & \cdots & x_{2m} \\ & & \vdots \\ 1 & x_{n2} & \cdots & x_{nm} \end{pmatrix}$$
(2)

be an $n \times m$ matrix of explanatory variables (real-valued). The subscripts denote the dimensions and will be omitted when these are assumed fixed (in Sections 2 and 3). Let

$$\boldsymbol{\beta} = \left(\beta_1, \beta_2, \dots, \beta_m\right)^T \in \mathbb{R}^m \tag{3}$$

be a vector of coefficients for our linear model and assume that

$$E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}.\tag{4}$$

Then each component of the deviation (or error) vector

$$\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta} \tag{5}$$

has expectation zero. The explanatory variables may be continuous or discrete. We assume $n \ge m$. Let r_X denote the rank of **X**; then $r_X \le m$. In practice, we usually have *n*, the number of data points, much larger than *m*, the number of coefficients. Our goal is estimating the components of β by ML principles, and determining their standard errors, given a set of response variables **y**, explanatory variables **X**, and a response variable distribution based on the Laplace distribution as described below. In terms of generalized linear models, the link function is assumed to be the identity.

2.2. The Distribution of the Deviations—A Modified Laplace Distribution

Example 1. Let $f : \mathbb{R} \to \mathbb{R}$ be defined by

$$f(z;p) = \left(\frac{p}{Q(p)}\right)e^{(-p|z|)},\tag{6}$$

where p > 0 is a real-valued parameter and Q(p) is a real-valued normalizing function defined so that

$$\int_{-\infty}^{\infty} f(z; p) dz = 1.$$
 (7)

Then f is the probability density function for the Laplace distribution, with scale parameter p, centred at the origin, and with unbounded support. It is not differentiable at the origin in the classical sense.

The method of MLE for the response function (6) corresponds to least absolute error (LAE) regression [2–4]. However, the theory of LAE regression is not sufficiently general for our epigenomic modelling problem. We next describe the more general response functions we require.

Example 2. Now consider the case of bounded support. For finite B > 0, define $f : [-B, B] \rightarrow \mathbb{R}$ by

$$f(z;p) = \left(\frac{p}{Q(p;B)}\right)e^{(-p|z|)},\tag{8}$$

where p > 0 is a real-valued scale parameter and Q(p; B) is a real-valued normalizing function defined so that

$$\int_{-B}^{B} f(z; p) dz = 1.$$
 (9)

Then f is the probability density function for the truncated Laplace distribution with scale parameter p, centred at the origin, and with bounded support [-B, B].

Example 3. More generally, consider perturbations of the truncated Laplace probability density function of the following form. Let

$$f(z; p, \mathbf{q}) = \left(\frac{p}{Q(p, g, \mathbf{q}; B)}\right) e^{(-p|z|)} g(|z|; \mathbf{q}), \qquad (10)$$

where real-valued $g(z; \mathbf{q})$ is equal to the constant map

$$g_1(z) = 1 \tag{11}$$

plus a perturbation, **q** is a vector of parameters for *g*, and parameter vector $\mathbf{p} = (p, \mathbf{q}) \in \mathbb{R}^r$, $r \ge 1$. If r = 1, *g* has no parameters. Here, $Q(p, g, \mathbf{q}; B)$ is a real-valued normalizing function defined so that

$$\int_{-B}^{B} f(z; \mathbf{p}) dz = 1.$$
(12)

We assume that there exists some $\epsilon > 0$ such that $g(z; \mathbf{q})$ is \mathscr{C}^3 in z on $(-\epsilon, B + \epsilon)$ and that $g(z; \mathbf{q}) > 0$ on [0, B], for fixed parameter vector \mathbf{q} . Note that, as a consequence of using the modulus function, $g(|z|; \mathbf{q})$ will not be differentiable with respect to z, at z = 0, in general. As in Example 1, f is not differentiable at the origin due to the use of the modulus function.

Example 4. We could allow unbounded support if $\int_{-\infty}^{\infty} e^{(-p|z|)} g(|z|; \mathbf{q}) dz$ is finite.

Example 5. Now consider our motivating example, a truncated Laplace distribution with bounded support [-1, 1], perturbed by adding kurtosis. Such a distribution is used to model the deviations in the proportions of methylation measured at gene promoter CpG sites. Specifically, to fit with observations, kurtosis is added to a Laplace probability density function with bounded support by adding a third-order Hermite polynomial to give an amended version.

Consider

$$f(z;\mathbf{p}) = \left(\frac{p}{Q(p,g_2,q;B)}\right)e^{(-p|z|)}g_2(|z|;q).$$
(13)

Here, $\mathbf{p} = (p, q), q \ge 0$ is small, B = 1,

$$g_2(z;q) = 1 + qH_3(z),$$
 (14)

and $H_3(z) = z^3 - 3z$ is the third-order Hermite polynomial. Solving for Q yields

$$f(z; p, q) = \frac{p^4 e^{(-p|z|)} [1 + qH_3(|z|)]}{2 [(p^3 - 3qp^2 + 6q) - e^{-p} (p^3 (1 - 2q) + 6pq + 6q)]}.$$
(15)

3

Example 6. The functions $g_3(z;q) = 1 - qz$ and $g_4(z;q) = e^{-qz^2}$, assuming small positive *q* and bounded support, could be used in (10) to model distributions similar to the Laplace but with thinner tails.

We restrict to symmetric distributions satisfying f(z) = f(-z).

3. Maximum Likelihood Estimation

3.1. The Log-Likelihood Function. Let

$$f\left(z,\mathbf{p}\right) \tag{16}$$

be a probability density function with parameter vector \mathbf{p} , as described in Section 2. Let z_1, \ldots, z_n be a sequence of independent and identically distributed deviations with joint probability density function

$$f(z_1, \dots, z_n; \mathbf{p}) = f(z_1(\boldsymbol{\beta}), \dots, z_n(\boldsymbol{\beta}); \mathbf{p})$$

= $\Pi_{i=1}^n f(z_i(\boldsymbol{\beta}); \mathbf{p}).$ (17)

This is also the likelihood function

$$L_{\mathbf{z}}(\mathbf{z};\mathbf{p}) = L_{\boldsymbol{\beta}}(\boldsymbol{\beta};\mathbf{p};\mathbf{X},\mathbf{y})$$

= $L(\mathbf{z}(\boldsymbol{\beta});\mathbf{p}) = f(\mathbf{z}(\boldsymbol{\beta});\mathbf{p}),$ (18)

which may be regarded as a function of z or β ; here, the subscript reflects our point of view. We use the log-likelihood function l in the estimation of β , where, using various notation,

$$l_{\mathbf{z}}(\mathbf{z};\mathbf{p}) = \log_{e} \left(L_{\mathbf{z}}(\mathbf{z};\mathbf{p}) \right)$$

= $l_{\boldsymbol{\beta}} \left(\boldsymbol{\beta};\mathbf{p};\mathbf{X},\mathbf{y} \right) = l \left(\mathbf{z} \left(\boldsymbol{\beta} \right);\mathbf{p} \right).$ (19)

Substituting measured values of y_i and known inputs x_{ij} into l_β , we obtain a function of β and \mathbf{p} . The parameters \mathbf{p} are assumed known, but if not, may be estimated separately. In our biological application, they are estimated by MLE and are assumed fixed for a particular measuring process. Hence, we have a function of β , the coefficients of our linear model.

Our aim is to find a maximum likelihood estimator (MLE) denoted $\hat{\beta}_n \in \mathbb{R}^m$, that is, some point at which l attains its maximum value. The subscript n corresponds to the number of deviations. Now l is a continuous function. If B is finite, it has compact support in \mathbb{R}^m . Since a continuous function on a compact set attains its supremum, the existence of a MLE for l_β is guaranteed. Even if $B = \infty$, we may consider truncations with finite bounds $B_k = k, k = 1, 2, \ldots$ Since l is maximized when the z_i are small, truncating f to [-k, k] for k large enough will not affect the set of points at which l attains its maximum. We show in Section 3.3 that a MLE $\hat{\beta}_n$ is not necessarily unique.

Case 1. If $g(z; \mathbf{q}) = g_1(z) = 1$ and so f is the Laplace probability density function with parameter p as in (6) or (8), then, for $\boldsymbol{\beta} \in \mathbb{R}^m$ such that every $z_i(\boldsymbol{\beta})$ is in the support of f,

$$l_{\boldsymbol{\beta}} \left(\boldsymbol{\beta}; p; \mathbf{X}, \mathbf{y} \right) = l_{\boldsymbol{\beta}, n, m} \left(\boldsymbol{\beta}; p; \mathbf{X}_{n, m}, \mathbf{y} \right)$$

$$= -n \log \left(Q \right) + n \log \left(p \right) + \sum_{i=1}^{n} \left(-p \left| z_{i} \left(\boldsymbol{\beta} \right) \right| \right)$$

$$= -n \log \left(Q \right) + n \log \left(p \right)$$

$$+ \sum_{i=1}^{n} \left(-p \left| y_{i} - \left(\mathbf{X} \boldsymbol{\beta} \right)_{i} \right| \right),$$

(20)

where Q = 2 if the support of f is \mathbb{R} . If the support of f is [-B, B] then

$$Q = Q(p, g_1; B) = 2(1 - e^{-pB}).$$
(21)

The theory of LAE regression (corresponding to MLE using Laplace distributions without modification as response functions) may be found in various texts, for example, [3]. Here, it is proved that there exists a MLE $\hat{\beta}_n$ corresponding to at least r_x zero errors. We are concerned with the extension of these ideas to the case of perturbed and truncated Laplace response functions. For the truncated Laplace distribution, we prove that there exists a MLE $\hat{\beta}_n$ corresponding to at least r_x zero errors. Consider the following more general case.

Case 2. If *f* is a perturbed Laplace probability density function with perturbing function $g(z; \mathbf{q})$ and bounded support as in (10), then for $\boldsymbol{\beta} \in \mathbb{R}^m$ such that $|z_i(\boldsymbol{\beta})| \leq B, i = 1, 2, ..., n$,

$$l_{\beta}(\boldsymbol{\beta}; \mathbf{p}; \mathbf{X}, \mathbf{y}) = l_{\beta,n,m}(\boldsymbol{\beta}; p; \mathbf{X}_{n,m}, \mathbf{y})$$

$$= -n \log \left(Q\left(p, g, \mathbf{q}; B\right) \right) + n \log \left(p\right)$$

$$+ \sum_{i=1}^{n} \left(-p \left|z_{i}\left(\boldsymbol{\beta}\right)\right|\right) + \sum_{i=1}^{n} \log \left(g\left(\left|z_{i}\left(\boldsymbol{\beta}\right)\right|; \mathbf{q}\right)\right)$$

$$= -n \log \left(Q\left(p, g, \mathbf{q}; B\right)\right) + n \log \left(p\right)$$

$$+ \sum_{i=1}^{n} \left(-p \left|y_{i} - \left(\mathbf{X}\boldsymbol{\beta}\right)_{i}\right|\right)$$

$$+ \sum_{i=1}^{n} \log \left(g\left(\left|y_{i} - \left(\mathbf{X}\boldsymbol{\beta}\right)_{i}\right|; \mathbf{q}\right)\right).$$
(22)

Note since $g(z_i; \mathbf{q})$ is strictly positive on [0, B], so is $f(z_i; \mathbf{p})$ and so $\log(f(z_i; \mathbf{p}))$ is well-defined on [-B, B], i = 1, 2, ..., n.

In Section 3.3, we show that if the perturbing function g is such that $\log g(z; \mathbf{q})$ is convex and non-increasing on [0, B], there exists a MLE corresponding to at least $r_{\mathbf{X}}$ data points. We give an upper bound on $|d \log g(z; \mathbf{q})/dz|$ on [0, B] which, if not exceeded, ensures that there exists a MLE corresponding to at least one data point. We apply these results to our motivating example, the Laplace distribution with added kurtosis, described below.

Example 7. If $g(z; \mathbf{q}) = g_2(z; q) = 1 + qH_3(z)$ and so f is a Laplace probability density function with scale parameter p, with added kurtosis and bounded support as in (13), then for $\boldsymbol{\beta} \in \mathbb{R}^m$ such that $|z_i(\boldsymbol{\beta})| \leq 1, i = 1, 2, ..., n$,

$$l_{\beta}(\beta; p, q; \mathbf{X}, \mathbf{y}) = l_{\beta,n,m}(\beta; p; \mathbf{X}_{n,m}, \mathbf{y})$$

$$= -n \log (Q(p, g_{2}, q; 1)) + n \log (p)$$

$$+ \sum_{i=1}^{n} (-p |z_{i}(\beta)|)$$

$$+ \sum_{i=1}^{n} \log (1 + qH_{3}(|z_{i}(\beta)|))$$

$$= -n \log (Q(p, g_{2}, q; 1)) + n \log (p)$$

$$+ \sum_{i=1}^{n} (-p |y_{i} - (\mathbf{X}\beta)_{i}|)$$

$$+ \sum_{i=1}^{n} \log (1 + qH_{3}(|y_{i} - (\mathbf{X}\beta)_{i}|)).$$
(23)

Now l_{β} is not differentiable in the classical sense with respect to the linear model coefficients β_j when any $z_i =$ 0. Hence, we cannot assume that l_{β} is differentiable at a MLE. This paper addresses this issue firstly by proposing a non-gradient method of coefficient estimation and secondly (in Section 4) by using generalized functions to calculate statistical estimates including estimates of standard errors. In Section 5, we discuss the asymptotic theory of our MLE.

3.2. Coefficient Estimation Dealing with Abrupt Changes in Gradient. Although L_{β} and l_{β} are continuous functions of β , their first derivatives are not. Consider the geometry of the coefficient space \mathbb{R}^m , where $\beta \in \mathbb{R}^m$. For each index *i*, since the response function distribution is defined in terms of absolute values, we can find a hyperplane H_i^0 in \mathbb{R}^m on which $L = L_{\beta}$ and $l = l_{\beta}$ are not differentiable, defined by setting $z_i = 0$. Let $\mathbf{x}_i^T = (x_{i1}, x_{i2}, \dots, x_{im})$, the *i*th row of \mathbf{X} , it is never the zero vector since $x_{i1} = 1, i = 1, 2, \dots, n$. Choose $\mathbf{w} \in \mathbb{R}^m$ so that $\mathbf{x}_i^T \mathbf{w} = 0$. Then

$$\boldsymbol{\beta} = \mathbf{w} + \frac{y_i \mathbf{x}_i}{\mathbf{x}_i^T \mathbf{x}_i} \tag{24}$$

yields $z_i = 0$. Let H_i^0 be the set of all such β . For example, for m = 2, for each error term there is a line in \mathbb{R}^2 on which L_β and l_β have a sharp ridge. By inspection of the geometry, we would expect the values of β which maximize l_β to be either on the union of the hyperplanes or very close to intersections of the hyperplanes H_i^0 , i = 1, 2, ..., n. Imagine searching in β -space near the hyperplanes H_i^0 . Even if L_β has a local maximum near but not on the union of the hyperplanes, it would be difficult to use a method based on the gradient of either L_β or l_β , since the gradient changes sharply whenever we cross one of the H_i^0 . The simplex method of coefficient

estimation, which does not require any partial derivatives, suits this geometry.

3.3. The Maximum Likelihood Estimator Corresponds to a Data Point

3.3.1. Convex and Non-Increasing Perturbations of the Laplace Probability Density Function. Consider the probability density function

$$f(z; p, \mathbf{q}) = \left(\frac{p}{Q(p, g, \mathbf{q}; B)}\right) e^{(-p|z|)} g(|z|; \mathbf{q}), \qquad (25)$$

with support [-B, B], for some finite B > 0 as described in Section 2.2 (Example 3). Recall we assume that there exists some $\epsilon > 0$ such that $g(z; \mathbf{q})$ is \mathcal{C}^3 in z on $(-\epsilon, B + \epsilon)$ and that $g(z; \mathbf{q}) > 0$ on [0, B], for fixed parameter \mathbf{q} . Let

$$\Omega_B = \left\{ \mathbf{z} \in \mathbb{R}^n : \left| z_i \right| \le B, i = 1, \dots, n \right\} = \left[-B, B \right]^n.$$
(26)

We consider the log-likelihood function

$$l_{\mathbf{z}}(\mathbf{z};\mathbf{p}):\Omega_{B}\longrightarrow\mathbb{R}$$
(27)

(conditional on **p**) with its domain restricted to

$$\mathcal{A} = \{ \mathbf{z}(\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^m \} \cap \Omega_B$$

= $\{ \mathbf{y} - \mathbf{X}\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^m \} \cap \Omega_B,$ (28)

that is, constrained to ${\mathscr A}$ or equivalently, consider the log-likelihood function

$$l_{\boldsymbol{\beta}}\left(\boldsymbol{\beta};\mathbf{p};\mathbf{X},\mathbf{y}\right):\mathscr{B}\longrightarrow\mathbb{R},$$
(29)

constrained to \mathscr{B} , where $\mathscr{B} \subset \mathbb{R}^m$ is defined as $\{\beta \in \mathbb{R}^m : z(\beta) \in \Omega_B\}$. Note that $\mathscr{A} = \mathscr{A}(\mathbf{X}, \mathbf{y}, B)$ and similarly $\mathscr{B} = \mathscr{B}(\mathbf{X}, \mathbf{y}, B)$. Also, if the function *Z* is defined by

$$Z: \mathbb{R}^{m} \longrightarrow \mathbb{R}^{n}, \boldsymbol{\beta}$$
$$\longmapsto \mathbf{z}(\boldsymbol{\beta}) = \mathbf{y} - \mathbf{X}\boldsymbol{\beta},$$
(30)

then $Z(\mathscr{B}) = \mathscr{A}$.

Lemma 8. If $\log(g) : (-\epsilon, B + \epsilon) \to \mathbb{R}$ is convex and nonincreasing (i.e., $d \log g(z; \mathbf{q})/dz \leq 0$ on [0, B], then there exists a maximum of $l_{\beta} : \mathcal{B} \to \mathbb{R}$ corresponding to at least $r_{\mathbf{X}}$ data points. That is, there exists $\hat{\beta}_n \in \mathcal{B} \subset \mathbb{R}^m$ such that the constrained l_{β} attains its maximum at $\hat{\beta}_n$, and there exists at least $r_{\mathbf{X}}$ indices $i_j \in \{1, 2, ..., n\}$ such that $z_{i_j}(\hat{\beta}_n) = 0$, $j = 1, 2, ..., r_{\mathbf{X}}$.

Corollary 9. Let f(z; p) be the Laplace probability density function (8) with support [-B, B] and parameter p. Then there exists a maximum of $l_{\beta} : \mathcal{B} \to \mathbb{R}$ corresponding to at least $r_{\mathbf{X}}$ data points.

Proof of Corollary 9. Let *g* be the constant map $g_1(z) = 1$; then $\log(g(z)) = 0$, and hence $\log(g(z))$ is convex and non-increasing on [0, B]. Corollary 9 follows from Lemma 8.

Proof of Lemma 8. To begin, assume that **X** has full rank *m*, recall $n \ge m$ and that *f* has bounded support. Then \mathscr{A} is a compact convex subset of an *m*-dimensional affine subspace of \mathbb{R}^n . Since the mapping *Z* is linear and has full rank *m*, the inverse image $\mathscr{B} = Z^{-1}(\mathscr{A})$ is a compact convex subset of \mathbb{R}^m . We partition $\mathscr{B} \subset \mathbb{R}^m$, which is the support of L_β , into a finite collection of compact convex sets, so that, on each subset, l_β is convex.

various applications (see [6], a survey paper, and also [7]).

The convex analysis results we require are in Appendix A.

Let H_i^{δ} be the hyperplane in \mathbb{R}^m corresponding to the error term $z_i(\beta) = \delta$. Then H_i^{-B} and H_i^B are the hyperplanes in \mathbb{R}^m corresponding to errors $z_i = -B$ and $z_i = B$, respectively. It follows that the log-likelihood function

$$l_{\boldsymbol{\beta}}(\boldsymbol{\beta}; \mathbf{p}; \mathbf{X}, \mathbf{y}) = -n \log \left(Q\left(p, g, \mathbf{q}; B\right) \right)$$

$$+ n \log \left(p\right) + \sum_{i=1}^{n} \left(-p \left|z_{i}\left(\boldsymbol{\beta}\right)\right|\right)$$

$$+ \sum_{i=1}^{n} \log \left(g\left(\left|z_{i}\left(\boldsymbol{\beta}\right)\right|; \mathbf{q}\right)\right)$$

$$= -n \log \left(Q\left(p, g, \mathbf{q}; B\right)\right) + n \log \left(p\right)$$

$$+ \sum_{i=1}^{n} \left(-p \left|y_{i} - \left(\mathbf{X}\boldsymbol{\beta}\right)_{i}\right|\right)$$

$$+ \sum_{i=1}^{n} \log \left(g\left(\left|y_{i} - \left(\mathbf{X}\boldsymbol{\beta}\right)_{i}\right|; \mathbf{q}\right)\right)$$
(31)

is a convex function in between the the hyperplanes H_i^{-B} , H_i^0 , and H_i^B , i = 1, ..., n. These 3n hyperplanes divide the domain \mathscr{B} in \mathbb{R}^m into at most 2^n open sets bounded by (but not intersecting) the hyperplanes. Each such open set (and hence its closure) may be labelled by a set of n signs. For any $\boldsymbol{\beta} \in \mathbb{R}^m$ such that $0 < |z_i(\boldsymbol{\beta})| < B, i = 1, ..., n$; (sgn $(z_1(\boldsymbol{\beta}))$), sgn $(z_2(\boldsymbol{\beta})), \ldots$, sgn $(z_n(\boldsymbol{\beta}))$) labels the open set containing $\boldsymbol{\beta}$.

Next, consider \mathbb{R}^n as the union of its orthants, which we denote \mathcal{O}_k , $k = 1, ..., 2^n$. We assume the orthants are closed sets. For example, the non-negative orthant is $\{z \in \mathbb{R}^n : z_i \ge 0, i = 1, ..., n\}$. In the interior of any \mathcal{O}_k , the sign of z_i does not change, i = 1, ..., n. Relabel the open subsets $\mathcal{B}_k = \mathcal{B}_k(\mathbf{X}, \mathbf{y}, B)$, where $k \in \{1, 2, ..., 2^n\}$, so that $Z(\mathcal{B}_k) \subset int(\mathcal{O}_k)$. Let $\mathcal{A}_k = Z(\mathcal{B}_k) = \mathcal{A}_k(\mathbf{X}, \mathbf{y}, B)$.

Now, $\mathscr{B} = \bigcup_k \operatorname{cl}(\mathscr{B}_k)$, where $\operatorname{cl}(\mathscr{B}_k)$ denotes the closure of the set. Since $\operatorname{cl}(\mathscr{B}_k)$ is bounded by hyperplanes, it is convex. It is closed and bounded and hence compact. Choose $k \in \{1, \ldots, 2^n\}$, such that \mathscr{B}_k is non-empty. Since continuous functions are bounded on compact sets, the supremum of l_β , when restricted to $\operatorname{cl}(\mathscr{B}_k)$, must be attained at one or more points in $\operatorname{cl}(\mathscr{B}_k)$. By Corollary A.3, the supremum (in our case the maximum) of l_β on $\operatorname{cl}(\mathscr{B}_k)$ is attained on the whole set or on a union of faces of dimension less than m or at a vertex. Since there are a finite number of sets to consider, the maximum of l_{β} must occur at a vertex but might occur, for example, on the whole of a set or on a union of faces. This is important to consider when using search algorithms such as the simplex method; as repeated application with different starting points may give a set of solutions which, for example, lie on a line segment. Note the following points.

- (i) Assuming that $r_{\mathbf{X}} = m$, any vertex of the set in \mathbb{R}^m at which l_{β} attains its maximum must correspond to m data points, possibly more (degeneracy). This is due to the fact that in \mathbb{R}^n , the gradient $\nabla l_{\mathbf{z}}(\mathbf{z}; \mathbf{p})$ points in the direction of the boundary of the corresponding orthant and away from the boundary of Ω_B .
- (ii) A MLE is not necessarily unique.
- (iii) If $r_X < m$, then we may apply the same reasoning to a subspace of \mathbb{R}^m of dimension r_X on which the error mapping has full rank r_X .
- (iv) Since at the MLE the absolute values of the deviations $|z_i(\hat{\beta}_n)|$ will all be small, this proof for finite *B* may be extended to $B = \infty$.

Lemma 8 is useful but we need to know what happens for more general perturbing functions g. First, we consider the Laplace distribution without perturbation.

3.3.2. The Truncated Laplace Probability Density Function. For n = 1, let

$$f(z;p) = \left(\frac{p}{Q(p;1)}\right)e^{(-p|z|)}$$
(32)

be the Laplace probability density function (8) with scale parameter p and with support [-1, 1]. Then

$$Q(p;1) = 2(1 - e^{-p}),$$

$$f(z;p) = \frac{p}{2(1 - e^{-p})}e^{(-p|z|)},$$
(33)

and

$$\log(f(z; p)) = \log\left(\frac{p}{2(1 - e^{-p})}\right) - p|z|.$$
(34)

The latter is a piecewise affine function in *z*. It has a maximum value of $\log(p/(2(1 - e^{-p})))$ when z = 0. More generally, for n > 1 and independent deviations (error terms) $z_1, z_2, ..., z_n$,

$$l_{z}(z;p) = \sum_{i=1}^{n} \log (f(z_{i};p))$$

$$= \sum_{i=1}^{n} (-p |z_{i}|) + n \log \left(\frac{p}{2(1-e^{-p})}\right).$$
(35)

Hence, the log-likelihood function is, up to a constant term, a piecewise linear function in the error terms, with a maximum

attained when all the errors are zero. However, we must restrict our domain to $\mathscr{A} \subset \mathbb{R}^n$, or equivalently to $\mathscr{B} \subset \mathbb{R}^m$. The log-likelihood function

$$l_{\boldsymbol{\beta}}(\boldsymbol{\beta}; p) = -p \sum_{i=1}^{n} \left| (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})_{i} \right|$$

$$+ n \log \left(\frac{p}{2 \left(1 - e^{-p} \right)} \right)$$
(36)

is a piecewise linear function, up to a constant term, in between the hyperplanes H_i^{-1} , H_i^0 , and H_i^{+1} , i = 1, ..., n. Let \mathbf{X}_i denote the *j*th column of \mathbf{X} , j = 1, 2..., m. Let

$$\operatorname{Sp} \{\mathbf{X}\} = \operatorname{Sp} \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m\} \subset \mathbb{R}^n, \quad (37)$$

denote the span of the columns of X.

Now $l_{\mathbf{z}} : \operatorname{cl}(\mathscr{A}_k) \mapsto \mathbb{R}$ has a critical point at $\mathbf{z}(\boldsymbol{\beta}) \in \mathscr{A}_k$ if and only if the gradient $\nabla l_{\mathbf{z}} = (\partial l/\partial z_1, \dots, \partial l/\partial z_n)^T$ (evaluated at $\mathbf{z}(\boldsymbol{\beta})$) is orthogonal to $\operatorname{Sp}\{\mathbf{X}\}$. The columns of \mathbf{X} are tangent vectors to \mathscr{A} at this point. The gradient $\nabla l_{\mathbf{z}} = -p(\operatorname{sgn}(z_1(\boldsymbol{\beta})), \operatorname{sgn}(z_2(\boldsymbol{\beta})), \dots, \operatorname{sgn}(z_n(\boldsymbol{\beta})))^T$ is constant in the interior of any orthant. If we travel along a straight line path in any orthant, $l_{\mathbf{z}}$ either always increases, always decreases, or remains constant. Hence, we will not find an isolated local maximum or minimum in \mathscr{A}_k , an open set, for the constrained $l_{\mathbf{z}}$.

We need to be aware of the case where \mathscr{A}_k lies in or very close to a level set of l_z . We might need to test for this. This happens when the sign vector $(\operatorname{sgn}(z_1), \ldots, \operatorname{sgn}(z_n))^T$ is orthogonal to $\operatorname{Sp}\{\mathbf{X}\}$, or nearly so. In the former case, the constrained l is constant on $\operatorname{cl}(\mathscr{A}_k)$. In the latter case, the constrained l will differ very little around the maximum on $\operatorname{cl}(\mathscr{A}_k)$. If this is the case for all the \mathscr{A}_k , then the ML values for the coefficients β_j will not be sharply defined (will have large variance).

3.3.3. More General Perturbations of the Laplace Probability Density Function. The question is, given a nontrivial perturbing function $g(z; \mathbf{q})$, does the log-likelihood function attain its maximum at a data point? We have given conditions on g in Lemma 8 which are sufficient to ensure the maximum is attained at a data point. We give a more general criterion in Lemma 10.

Assume that log(g) is non-linear in any orthant. Otherwise, we can write g in the form of a scaled Laplace distribution and apply Lemma 8. Then l is the sum of an affine function and a non-linear function in any orthant. This affine function is

$$l_{p,B} = -n \log (Q(p, g_1; B)) + n \log (p) + \sum_{i=1}^{n} (-p |z_i|),$$
(38)

the log-likelihood function corresponding to the Laplace distribution. The non-linear function is

$$l_{\text{nlin}} = -n \log (Q(p, g, \mathbf{q}; B)) + n \log (Q(p, g_1; B)) + \sum_{i=1}^{n} \log g(|z_i|; \mathbf{q}).$$
(39)

Then $l = l_{p,B} + l_{nlin}$, and so $\nabla l = \nabla l_{p,B} + \nabla l_{nlin}$, where

$$\nabla l_{p,B} = -p(\operatorname{sgn}(z_1), \operatorname{sgn}(z_2), \dots, \operatorname{sgn}(z_n))^T,$$

$$\nabla l_{\operatorname{nlin}} = \left(\frac{d\log(g(z_1; \mathbf{q}))}{dz_1}, \frac{d\log(g(z_2; \mathbf{q}))}{dz_2}\right)$$

$$(40)$$

$$(40)$$

It is possible that there exist orthants \mathcal{O}_k in which the set $\mathcal{A}_k(\mathbf{X}, \mathbf{y}, B)$ is orthogonal to the gradient $\nabla l_{p,B}$. It is possible that $l_{p,B}$ attains its global maximum (with respect to $\mathcal{A}(\mathbf{X}, \mathbf{y}, B)$) on the whole of such an $\mathcal{A}_k(\mathbf{X}, \mathbf{y}, B)$. Hence, it is important to consider the behaviour of l_{nlin} in such orthants.

Lemma 10. Assume bounded support and let

$$\gamma = \sup_{0 \le z \le B} \left\{ \left| \frac{dg(z; \mathbf{q}) / dz}{g(z; \mathbf{q})} \right| \right\}$$

$$= \sup_{0 \le z \le B} \left\{ \left| \frac{d \log \left(g(z; \mathbf{q}) \right)}{dz} \right| \right\},$$
(41)

where sup denotes supremum. Then, if $\gamma < p$, the supremum or maximum of *l* is attained at a data point. In the special case that $Sp\{X\}$ is orthogonal to $(-sgn(z_1), -sgn(z_2), \ldots, -sgn(z_n))^T$ in any orthant, it may be that the supremum is also attained elsewhere.

Proof of Lemma 10. Choose $\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B) = \mathscr{A}(\mathbf{X}, \mathbf{y}, B) \cap \mathscr{O}_k$, where $k \in \{1, 2..., 2^n\}$, such that the set $\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$ is non-empty, choose $\mathbf{z} \in \mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$ (an open set relative to \mathscr{A}), and let $\mathbf{w}(k, \mathbf{X})$ be the projection of

$$\frac{\nabla l_{p,B}}{p} = \left(-\operatorname{sgn}\left(z_{1}\right), -\operatorname{sgn}\left(z_{2}\right), \dots, -\operatorname{sgn}\left(z_{n}\right)\right)^{T}$$
(42)

onto Sp{**X**}. Consider the case **w**(k, **X**) is non-zero. Then the affine function $l_{p,B}$ is not constant on $\mathcal{A}_k(\mathbf{X}, \mathbf{y}, B)$. Given any **v** in $\mathcal{A}_k(\mathbf{X}, \mathbf{y}, B)$, there exists some $\epsilon > 0$ and a map

$$\begin{aligned} h_{p,B} : (-\epsilon, \epsilon) &\longrightarrow \mathbb{R} \\ t &\longmapsto l_{p,B} \left(\mathbf{v} + t \mathbf{w} \left(k, \mathbf{X} \right) \right) \end{aligned}$$
 (43)

which is increasing at t = 0, since

$$\left(\frac{dh_{p,B}}{dt}\right)_{t=0} = \nabla l_{p,B}^{T} \left(\frac{d\left(\mathbf{v} + t\mathbf{w}\left(k,\mathbf{X}\right)\right)}{dt}\right)_{t=0}$$
$$= \nabla l_{p,B}^{T} \mathbf{w}\left(k,\mathbf{X}\right)$$
$$= \left(\mathbf{w}\left(k,\mathbf{X}\right)\right)^{T} \mathbf{w}\left(k,\mathbf{X}\right).$$
(44)

Now consider

$$h: (-\epsilon, \epsilon) \longrightarrow \mathbb{R}$$

$$t \longmapsto l \left(\mathbf{v} + t \mathbf{w} \left(k, \mathbf{X} \right) \right).$$
(45)

Since $g(z; \mathbf{q})$ is continuous and strictly positive on [0, B], γ is finite. We show that if $\gamma < p$, dh/dt > 0 at t = 0, that is, at **v**. Specifically,

$$\frac{dh}{dt} = -p \sum_{i=1}^{n} \left(\operatorname{sgn} \left(v_{i} + tw_{i} \left(k, \mathbf{X} \right) \right) w_{i} \left(k, \mathbf{X} \right) \right) \\
+ \sum_{i=1}^{n} \frac{\left(dg \left(\left| z_{i} \right| ; \mathbf{q} \right) / d \left| z_{i} \right| \right)_{h(t)}}{g \left(\left| v_{i} + tw_{i} \left(k, \mathbf{X} \right) \right| ; \mathbf{q} \right)} \\
\times \left(\operatorname{sgn} \left(v_{i} + tw_{i} \left(k, \mathbf{X} \right) \right) \right) w_{i} \left(k, \mathbf{X} \right) \\
\geq -p \sum_{i=1}^{n} \left(\operatorname{sgn} \left(v_{i} + tw_{i} \left(k, \mathbf{X} \right) \right) w_{i} \left(k, \mathbf{X} \right) \right) \\
+ \gamma \sum_{i=1}^{n} \left(\operatorname{sgn} \left(v_{i} + tw_{i} \left(k, \mathbf{X} \right) \right) w_{i} \left(k, \mathbf{X} \right) \right) \\
= \left(1 - \frac{\gamma}{p} \right) \left(\mathbf{w} \left(k, \mathbf{X} \right) \right)^{T} \mathbf{w} \left(k, \mathbf{X} \right) \\
> 0$$

if $\gamma < p$ and $\mathbf{w}(k, \mathbf{X})$ is non-zero. Then the supremum of l, constrained to $cl(\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B))$, must be attained on the boundary of $\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$, relative to \mathscr{A} , and this must be at a data point due to the direction of ∇l .

In the limiting case, where $\mathbf{w}(k, \mathbf{X}) = \mathbf{0}$, but we still have $\gamma < p$, we find that the supremum of *l*, constrained to $cl(\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B))$ (relative to \mathscr{A}), must be attained at a data point but may be attained at other points in $\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$ as well. Assume this is not the case, that is, there exists no data point at which *l* attains its supremum when constrained to $cl(\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B))$. We find a contradiction as follows. Assume there exists some $\mathbf{v} \in \mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$, an open set (relative to \mathscr{A}), such that

$$l_{\mathbf{z}}(\mathbf{v}) - \sup \left\{ l_{\mathbf{z}}(\mathbf{z}) : \mathbf{z} \in \mathrm{bd}\left(\mathcal{O}_{k}\right) \cap \mathrm{cl}\left(\mathscr{A}_{k}\left(\mathbf{X}, \mathbf{y}, B\right)\right) \right\}$$

= $\tau > 0,$ (47)

where bd denotes boundary. Imagine perturbing the columns of **X** slightly (and continuously) to obtain **X**' so that l_{β} , a continuous function, conditional on **X**, changes by at most $\tau/3$, that is,

$$\left| l_{\beta} \left(\boldsymbol{\beta}; \mathbf{p}; \mathbf{X}, \mathbf{y} \right) - l_{\beta} \left(\boldsymbol{\beta}; \mathbf{p}; \mathbf{X}', \mathbf{y} \right) \right| < \frac{\tau}{3}$$
(48)

for all $\beta \in \mathscr{B}_k(\mathbf{X}, \mathbf{y}, B) \cap \mathscr{B}_k(\mathbf{X}', \mathbf{y}, B)$, and so that now $\mathbf{w}(k, \mathbf{X}') \neq \mathbf{0}$. We also require the perturbation small enough that

$$\begin{aligned} \left| \sup \left\{ l_{\mathbf{z}} \left(\mathbf{z}; \mathbf{p} \right) : \mathbf{z} \in \mathrm{bd} \left(\mathcal{O}_{k} \right) \cap \mathrm{cl} \left(\mathcal{A}_{k} \left(\mathbf{X}, \mathbf{y}, B \right) \right) \right\} \\ &- \sup \left\{ l_{\mathbf{z}} \left(\mathbf{z}; \mathbf{p} \right) : \mathbf{z} \in \mathrm{bd} \left(\mathcal{O}_{k} \right) \cap \mathrm{cl} \left(\mathcal{A}_{k} \left(\mathbf{X}', \mathbf{y}, B \right) \right) \right\} \right| \\ &< \frac{\tau}{3}. \end{aligned}$$

$$\tag{49}$$

Now $\mathbf{v} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}_{\mathbf{v}}$ for some unique $\boldsymbol{\beta}_{\mathbf{v}}$. Let $\mathbf{v}' = \mathbf{y} - \mathbf{X}'\boldsymbol{\beta}_{\mathbf{v}}$. If the perturbation of **X** is small enough, then $\boldsymbol{\beta}_{\mathbf{v}} \in \mathcal{B}_k(\mathbf{X}', \mathbf{y}, B)$ and so $\mathbf{v}' \in \mathcal{A}_k(\mathbf{X}', \mathbf{y}, B)$.

Now,
$$|l_{\mathbf{z}}(\mathbf{v}') - l_{\mathbf{z}}(\mathbf{v})| < \tau/3$$
. Hence,
 $l_{\mathbf{z}}(\mathbf{v}') - \sup \{l_{\mathbf{z}}(\mathbf{z}) : \mathbf{z} \in \mathrm{bd}(\mathcal{O}_k) \cap \mathrm{cl}(\mathcal{A}_k(\mathbf{X}', \mathbf{y}, B))\}$
 $> \frac{\tau}{3} > 0,$
(50)

which is not possible since $\mathbf{w}(k, \mathbf{X}') \neq \mathbf{0}$. Hence, we have a contradiction. Now, for completeness, assume there exists some $\mathbf{r} \in \mathrm{bd}(\mathscr{A}_k(\mathbf{X}, \mathbf{y}, B))$, \mathbf{r} not a data point (that is some $r_i = \pm B$) such that

$$l_{\mathbf{z}}(\mathbf{r}) - \sup \left\{ l_{\mathbf{z}}(\mathbf{z}) : \mathbf{z} \in \mathrm{bd}\left(\mathcal{O}_{k}\right) \cap \mathrm{cl}\left(\mathscr{A}_{k}\left(\mathbf{X}, \mathbf{y}, B\right)\right) \right\}$$

= $\tau > 0.$ (51)

Then there exists $\mathbf{v} \in \mathscr{A}_k(\mathbf{X}, \mathbf{y}, B)$, such that

$$l_{\mathbf{z}}\left(\mathbf{r}\right) - l_{\mathbf{z}}\left(\mathbf{v}\right) < \frac{\tau}{6}.$$
(52)

Hence,

$$l_{\mathbf{z}}(\mathbf{v}) - \sup \left\{ l_{\mathbf{z}}(\mathbf{z}) : \mathbf{z} \in \mathrm{bd}\left(\mathcal{O}_{k}\right) \cap \mathrm{cl}\left(\mathscr{A}_{k}\left(\mathbf{X}, \mathbf{y}, B\right)\right) \right\}$$

$$> \frac{5\tau}{6} > 0.$$
(53)

The contradiction follows as above. Hence, we have proved the lemma. $\hfill \Box$

3.3.4. An Amended Laplace Probability Density Function with Added Kurtosis. We may apply Lemma 10 to show that, for the amended Laplace probability density function with added kurtosis (Example 5) and for realistic values of p and q, the maximum of the log-likelihood function must be attained at a data point. Here, $g(z;q) = g_2(z;q) = 1 + q(z^3 - 3z) = 1 + qH_3(z)$ on [0, 1], that is, $g_2(|z|;q) = 1 + q(H_3(|z|))$ on [-1, 1]. In n dimensions,

$$l(\boldsymbol{\beta}; \mathbf{p}; \mathbf{X}, \mathbf{y}) = -n \log (Q(p, g_2, q; 1)) + n \log (p) + \sum_{i=1}^{n} (-p | y_i - (\mathbf{X}\boldsymbol{\beta})_i |) + \sum_{i=1}^{n} \log (1 + qH_3 (| y_i - (\mathbf{X}\boldsymbol{\beta})_i |)) = -n \log (Q(p, g_2, q; 1)) + n \log (p) + \sum_{i=1}^{n} (-p | y_i - (\mathbf{X}\boldsymbol{\beta})_i |) + \sum_{i=1}^{n} \log (1 + q (| y_i - (\mathbf{X}\boldsymbol{\beta})_i |^3 -3 | y_i - (\mathbf{X}\boldsymbol{\beta})_i |)).$$

(54)

Consider, for $z \in [0, 1]$,

$$\frac{d \log (g_2(z;q))}{dz} = \frac{d \left(\log \left(1 + q \left(z^3 - 3z \right) \right) \right)}{dz}$$

$$= \frac{3q \left(z^2 - 1 \right)}{1 + q \left(z^3 - 3z \right)},$$

$$\frac{d^2 \log \left(g_2(z;q) \right)}{dz^2} = \frac{d^2 \left(\log \left(1 + q \left(z^3 - 3z \right) \right) \right)}{dz^2}$$

$$= \frac{(3q) \left(2z - q \left(3 + z^4 \right) \right)}{\left(1 + q \left(z^3 - 3z \right) \right)^2}.$$
(55)

When 0 < q < 0.5, $d^2 \log(g_2(z;q))/dz^2$ is negative when z = 0 and positive when z = 1, and so the continuous monotonic function $\log(g_2)$ has a point of inflection, where $d^2 \log(g_2(z;q))/dz^2 = 0$ in (0, 1), at say $\omega(q)$. Hence, $\log(g_2)$ and l are both concave and convex on [0, 1]. When q is small, the concavity occurs close to the origin and the functions are convex on most of (0, 1). Moreover, $d \log(g_2(z;q))/dz$ is always negative on (0, 1), and its limit as $z \to \infty$ is also negative, so we will not get isolated local maxima in one dimension. The function $\log(g_2(z;q))$ and its first and second derivatives are plotted in Figures 1, 2, and 3, respectively; setting q = 0.025, a typical value. Reading from Figure 2, the upper bound γ is approximately 0.075, when q = 0.025. Since typically $p \ge 3$, the criterion ($\gamma < p$) for Lemma 10 is satisfied.

3.3.5. Non-Increasing Perturbations Both Concave and Convex. In certain situations we might find the criteria for both Lemmas 8 and 10 are not satisfied but that $\log g(z; \mathbf{q})$ is non-increasing on [0, B] and convex on [w(q), B], where w(q) is small. Since a sum of convex functions is convex (see [8]), l will be convex wherever $\log g(|z_i|; \mathbf{q})$ is convex for i = 1, 2, ..., n. For example, for $g = g_2$ and B = 1, we can prove a partial result as follows. By restricting to the domain

$$\mathscr{B}(\mathbf{X}, \mathbf{y}, [\omega(q), 1])$$

= $\bigcap_{i=1}^{n} \{ \boldsymbol{\beta} \in \mathbb{R}^{m} : \omega(q) \le |z_{i}(\boldsymbol{\beta})| \le 1 \} \subset \mathscr{B}(\mathbf{X}, \mathbf{y}, 1)$
(56)

or equivalently by substituting for each set $\mathscr{A}_k(\mathbf{X},\mathbf{y},1)$ the subset

$$\{\mathbf{z} \in \cap \mathscr{A}_{k} (\mathbf{X}, \mathbf{y}, 1) : \omega(q) < |z_{i}| < 1, i = 1, \dots, n\}$$

$$\subset \mathscr{A}_{k} (\mathbf{X}, \mathbf{y}, 1),$$
(57)

and applying convex function theory as for Lemma 8 we can prove that there exists $\beta_1 \in \mathscr{B}(\mathbf{X}, \mathbf{y}, [\omega(q), 1])$ at which lattains its maximum, and there exists at least one index isuch that $z_i(\beta_1) = \omega(q)$. Hence, there exists some $\beta_2 \in \mathscr{B}(\mathbf{X}, \mathbf{y}, 1) = \mathscr{B}(\mathbf{X}, \mathbf{y}, [0, 1])$ at which l attains its maximum, and there exists at least one index i such that $z_i(\beta_2) \leq \omega(q)$. In other words, there exists some point at which l has a



FIGURE 1: The non-linear part of the log-likelihood function, n = 1; $\log(g(z;q)), q = 0.025$.



FIGURE 2: The first derivative of $\log(g(z;q))$, q = 0.025.

maximum at which at least one of the errors is small. It makes sense to search for the maxima of *l* near or at vertices.

For $g = g_2$, we may apply Lemma 10 and so do not need this partial result, but for a perturbation similar in shape to g_2 , non-increasing everywhere on [0, 1], with $\log g(z; \mathbf{q})$ convex everywhere on $(\omega(g; \mathbf{q}), 1]$, concave everywhere on $[0, \omega(g; \mathbf{q}))$, for some small positive $\omega(g; \mathbf{q})$, and with steep slope at the point of inflection $(z = \omega(g; \mathbf{q}))$, such analysis would be useful.

4. Statistics for Linear Model Coefficients Assuming Perturbed and Truncated Laplace Response Functions

4.1. Dealing with Abrupt Changes in Gradient. The inclusion of the modulus (absolute value) function in the Laplace probability density function (6) (and variations thereof) is the cause of abrupt changes in the gradient of the log-likelihood

FIGURE 3: The second derivative of log(g(z; q)), q = 0.025.

function l_{β} (see Section 4.2). This section is devoted to dealing with the problems which are associated with these abrupt changes, encountered when deriving statistical formulae, for example, for standard errors.

The fact that l_{β} is not differentiable in the classical sense at a local maximum means that the assumptions made in the derivation of the usual classical formulae for the information matrix, the expected value of the Hessian of the log-likelihood function and the variance-covariance matrix for the model coefficients β_j , j = 1, 2, ..., m, are not met. For \mathscr{C}^2 probability density functions (and \mathscr{C}^2 log-likelihood functions), these formulae are derived using Taylor series. We find alternative expressions for these quantities assuming the truncated and/or perturbed Laplace response functions (as defined in Section 2) which are \mathscr{C}^3 where the modulus function is nonzero. In Section 5 these expressions will be used to prove the asymptotic convergence of our MLE to a random variable with a normal distribution.

4.2. Differentiation in a Generalized Sense. The following generalized functions are required to determine the first and second partial derivatives of the log-likelihood function l_{β} , with respect to the coefficients β_j . These derivatives are needed for the calculation of the standard errors. We require

$$\operatorname{sgn}(z) = \begin{cases} 1, & z > 0, \\ 0, & z = 0, \\ -1, & z < 0 \end{cases}$$
(58)

and $\delta(z)$ which is the delta function, that is, $\delta(z) = 0$ except at z = 0, and $\int_{-\infty}^{\infty} \delta(z) dz = 1$. These expressions and the modulus function are connected by

$$\frac{d|z|}{dz} = \operatorname{sgn}(z), \qquad (59)$$

$$\frac{d\operatorname{sgn}(z)}{dz} = 2\delta(z), \qquad (60)$$

where the differentiation is taken in the generalized sense (see [9, 10]). Hence, for $z_i \in [-B, B]$, the generalized derivative

$$\frac{dg\left(\left|z_{i}\right|;\mathbf{q}\right)}{dz_{i}} = \operatorname{sgn}\left(z_{i}\right)\frac{dg\left(\left|z_{i}\right|;\mathbf{q}\right)}{d\left|z_{i}\right|}.$$
(61)

Also, the derivative of the delta function may be defined via integration by parts, assuming $h : \mathbb{R} \to \mathbb{R}$ is \mathcal{C}^1 , we have

$$\int_{-\infty}^{\infty} \delta'(t) h(t) dt = -\int_{-\infty}^{\infty} \delta(t) h'(t) dt = -h'(0). \quad (62)$$

In Section 5, we investigate the behaviour of our model as $n \rightarrow \infty$, and so use subscripts to clarify the variables under consideration (**z** or β) and/or the dimension of the space(s) under consideration. Using (10) and (59), it follows that

$$\frac{\partial l}{\partial z_i} = \frac{\partial l_{\mathbf{z},n}}{\partial z_i} = -p \operatorname{sgn}\left(z_i\right) + \frac{d}{dz_i} \log\left(g\left(|z_i|;\mathbf{q}\right)\right)$$
$$= -p \operatorname{sgn}\left(z_i\right) + \frac{\operatorname{sgn}\left(z_i\right)}{g\left(|z_i|;\mathbf{q}\right)} \left(\frac{dg\left(|z_i|;\mathbf{q}\right)}{d|z_i|}\right).$$
(63)

Since $\mathbf{z} = \mathbf{y} - \mathbf{X}_{n,m} \boldsymbol{\beta}$ (see (5)),

$$\frac{\partial l}{\partial \beta_j} = \frac{\partial l_{\beta,n,m}}{\partial \beta_j} = p \sum_{i=1}^n x_{ij} \operatorname{sgn}(z_i) - \sum_{i=1}^n x_{ij} \frac{\operatorname{sgn}(z_i)}{g(|z_i|;\mathbf{q})} \frac{dg(|z_i|;\mathbf{q})}{d|z_i|}.$$
(64)

Letting $\nabla l_{\beta,n,m}$ denote the gradient of $l_{\beta,n,m}$ and letting $\nabla l_{z,n}$ denote the gradient of $l_{z,n}$, we have

$$\nabla l_{\boldsymbol{\beta},n,m} = -\mathbf{X}_{n,m}^T \nabla l_{\mathbf{z},n}.$$
(65)

In addition, using (60) and omitting the dependence of g upon its parameters for brevity,

$$\begin{aligned} \frac{\partial^2 l_{\mathbf{z},n}}{\partial z_i^2} &= -2p\delta\left(z_i\right) + 2\delta\left(z_i\right) \frac{1}{g\left(|z_i|\right)} \left(\frac{dg\left(|z_i|\right)}{d\left|z_i\right|}\right) \\ &+ \operatorname{sgn}\left(z_i\right) \frac{d}{dz_i} \left(\frac{1}{g\left(|z_i|\right)} \frac{dg\left(|z_i|\right)}{d\left|z_i\right|}\right) \\ &= -2p\delta\left(z_i\right) + 2\delta\left(z_i\right) \frac{1}{g\left(|z_i|\right)} \left(\frac{dg\left(|z_i|\right)}{d\left|z_i\right|}\right) \\ &+ \left(\operatorname{sgn}\left(z_i\right)\right)^2 \frac{d}{d\left|z_i\right|} \left(\frac{1}{g\left(|z_i|\right)} \frac{dg\left(|z_i|\right)}{d\left|z_i\right|}\right) \\ &= -2p\delta\left(z_i\right) + 2\delta\left(z_i\right) \frac{1}{g\left(|z_i|\right)} \left(\frac{dg\left(|z_i|\right)}{d\left|z_i\right|}\right) \end{aligned}$$

$$+ (\operatorname{sgn}(z_{i}))^{2} \times \left[\frac{1}{g(|z_{i}|)} \frac{d^{2}g(|z_{i}|)}{d|z_{i}|^{2}} - \left(\frac{1}{g(|z_{i}|)^{2}}\right) \left(\frac{dg(|z_{i}|)}{d|z_{i}|}\right)^{2} \right]$$

$$= -2p\delta(z_{i}) + 2\delta(z_{i}) \frac{1}{g(|z_{i}|)} \left(\frac{dg(|z_{i}|)}{d|z_{i}|}\right)$$

$$+ \frac{(\operatorname{sgn}(z_{i}))^{2}}{(g(|z_{i}|))^{2}} \left[g(|z_{i}|) \frac{d^{2}g(|z_{i}|)}{d|z_{i}|^{2}} - \left(\frac{dg(|z_{i}|)}{d|z_{i}|}\right)^{2} \right],$$

$$(66)$$

())2

and, if $i \neq j$,

$$\frac{\partial^2 l_{\mathbf{z},n}}{\partial z_i \partial z_j} = 0. \tag{67}$$

Let $\mathscr{H}_{\mathbf{z},n}$ denote the generalized Hessian of $l_{\mathbf{z},n}$, where

$$\mathscr{H}_{\mathbf{z},n} = \begin{pmatrix} \frac{\partial^2 l}{\partial z_1^2} & \frac{\partial^2 l}{\partial z_1 \partial z_2} & \cdots & \frac{\partial^2 l}{\partial z_1 \partial z_n} \\ \frac{\partial^2 l}{\partial z_2 \partial z_1} & \frac{\partial^2 l}{\partial z_2^2} & \cdots & \frac{\partial^2 l}{\partial z_2 \partial z_n} \\ & \vdots \\ \frac{\partial^2 l}{\partial z_n \partial z_1} & \frac{\partial^2 l}{\partial z_n \partial z_2} & \cdots & \frac{\partial^2 l}{\partial z_n^2} \end{pmatrix}, \quad (68)$$

and let $E(\mathcal{H}_{z,n})$ denote its expected value. Then $\mathcal{H}_{z,n}$ and $E(\mathcal{H}_{z,n})$ are diagonal matrices. Since the diagonal elements of $E(\mathcal{H}_{z,n})$ are all equal (see Section 4.6), this matrix is a multiple of the identity. We have

$$E\left(\mathscr{H}_{\mathbf{z},n}\right) = E\left(\frac{\partial^2 l_{\mathbf{z},n}}{\partial z_1^2}\right) I_n,\tag{69}$$

where I_n denotes the $n \times n$ identity matrix. Let $\mathcal{H}_{\beta,n,m}$ denote the generalized Hessian of $l_{\beta,n,m}$, where

$$\mathscr{H}_{\beta,n,m} = \begin{pmatrix} \frac{\partial^2 l}{\partial \beta_1^2} & \frac{\partial^2 l}{\partial \beta_1 \partial \beta_2} & \cdots & \frac{\partial^2 l}{\partial \beta_1 \partial \beta_m} \\ \frac{\partial^2 l}{\partial \beta_2 \partial \beta_1} & \frac{\partial^2 l}{\partial \beta_2^2} & \cdots & \frac{\partial^2 l}{\partial \beta_2 \partial \beta_m} \\ & \vdots \\ \frac{\partial^2 l}{\partial \beta_m \partial \beta_1} & \frac{\partial^2 l}{\partial \beta_m \partial \beta_2} & \cdots & \frac{\partial^2 l}{\partial \beta_m^2} \end{pmatrix}, \quad (70)$$

and let $E(\mathcal{H}_{\beta,n,m})$ denote its expected value. Then

$$\mathcal{H}_{\boldsymbol{\beta},n,m} = \mathbf{X}_{n,m}^T \mathcal{H}_{\mathbf{z},n} \mathbf{X}_{n,m}$$
(71)

and

$$E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right) = \mathbf{X}_{n,m}^{T} E\left(\mathscr{H}_{\mathbf{z},n}\right) \mathbf{X}_{n,m}$$
$$= E\left(\frac{\partial^{2} l_{\mathbf{z},n}}{\partial z_{1}^{2}}\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}.$$
(72)

If any $z_i = 0$, then the *i*th diagonal element of $\mathcal{H}_{z,n}$ is infinite. In this case $\mathcal{H}_{\beta,n,m}$ has infinite components. We prove in Section 4.6 that $E(\partial^2 l_{z,n}/\partial z_1^2)$ is finite.

4.3. A Classical Relation to Be Generalized. Let $\mathcal{J}_{\beta,n,m}$ denote the $m \times m$ Fisher information matrix, where

$$\left(\mathscr{F}_{\boldsymbol{\beta},n,m}\right)_{jk} = E\left[\left(\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \beta_j} - E\left(\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \beta_j}\right)\right)\left(\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \beta_k} - E\left(\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \beta_k}\right)\right)\right].$$
(73)

The components of our MLE $\hat{\beta}_n$ depend on the errors z_i and have a distribution whose variance-covariance matrix is denoted $\mathcal{V}_{\beta,n,m}$, where

$$\left(\mathcal{V}_{\boldsymbol{\beta},n,m}\right)_{jk} = E\left[\left(\boldsymbol{\beta}_{j} - E\left(\boldsymbol{\beta}_{j}\right)\right)\left(\boldsymbol{\beta}_{k} - E\left(\boldsymbol{\beta}_{k}\right)\right)\right], \quad (74)$$

j = 1, 2, ..., m, and k = 1, 2, ..., m. If the log-likelihood function l was sufficiently smooth around the region of interest (i.e., around its maximum value), then Taylor series expansions could be used to derive a relationship between $E(\mathcal{H}_{\beta,n,m}), \mathcal{J}_{\beta,n,m}$, and $\mathcal{V}_{\beta,n,m}$, namely,

$$\mathscr{V}_{\boldsymbol{\beta},n,m} = \mathscr{J}_{\boldsymbol{\beta},n,m}^{-1} = \left(-E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right)\right)^{-1},\tag{75}$$

(see [1]). However, our l is not sufficiently smooth, and so we cannot make use of this relationship without further justification. In general, (75) does not hold, assuming a truncated (and possibly perturbed) Laplace distribution.

In Sections 4.5 and 4.6 we calculate the expected values of the first and second partial derivatives of the log-likelihood function $l_{\beta,n,m}$ using generalized functions; this enables us to derive a generalized Taylor series expansion for the log-likelihood function about a maximum even when the maximum is, for example, on a ridge or at a vertex. Also, this enables us to derive an expression for the generalized variance-covariance matrix for the MLEs of the model coefficients and an expression for the generalized log-likelihood ratio statistic. These formulae differ from the standard formulae for the case of smooth log-likelihood functions, although their form is similar. Specifically, in our case, we prove that $\mathcal{F}_{\beta,n,m}$ is a multiple of $E(\mathcal{H}_{\beta,n,m})$, but that the multiple is not -1, rather, it is a negative real number that depends on p, the perturbation g, its parameters \mathbf{q} , and the bound B. We prove that our generalized $\mathcal{V}_{\beta,n,m}$ is a multiple of $\mathcal{J}_{\beta,n,m}^{-1}$, where the multiple is a positive real number depending on p, g, q, and B. We assume independent error distributions.

4.4. The Mean and Variance of the Partial Derivatives of the Log-Likelihood Function. The mean and the variance of $\partial l/\partial z_i$ are required in the calculation of $\mathcal{J}_{\beta,n,m}$. Recall $f(z_1, \ldots, z_n; \mathbf{p}) = \prod_{i=1}^n f(z_i; \mathbf{p})$ is the joint probability density function for the independent deviations (errors), and that

$$l = l_{\mathbf{z},n} = \sum_{i=1}^{n} \log\left(f\left(z_{i};\mathbf{p}\right)\right),\tag{76}$$

so

$$\frac{\partial l}{\partial z_i} = \frac{\partial l_{\mathbf{z},n}}{\partial z_i} = \frac{\partial \log\left(f\left(z_i;\mathbf{p}\right)\right)}{\partial z_i}.$$
(77)

Let

$$\mu_{i}(p, g, \mathbf{q}; B) = E\left(\frac{\partial l_{\mathbf{z},n}}{\partial z_{i}}\right)$$
$$= \int_{\Omega_{B}} \left(\frac{\partial \log\left(f\left(z_{i}; \mathbf{p}\right)\right)}{\partial z_{i}}\right) (f\left(\mathbf{z}; \mathbf{p}\right)) dz_{1} \cdots dz_{n}$$
$$= \int_{-B}^{B} \left(\frac{\partial \log\left(f\left(z_{i}; \mathbf{p}\right)\right)}{\partial z_{i}}\right) (f\left(z_{i}; \mathbf{p}\right)) dz_{i},$$
(78)

i = 1, ..., n; then $\mu(p, g, \mathbf{q}; B) = \mu_i(p, g, \mathbf{q}; B)$ is independent of index *i*. Since $f(z_i; \mathbf{p}) = f(-z_i; \mathbf{p})$, i = 1, ..., n, *L* and *l* are symmetric about the origin, and so

$$\mu = \mu \left(p, g, \mathbf{q}; B \right) = 0 \tag{79}$$

for any choice of *p*, *g*, **q**, and *B*. Let

$$\nu_{i}(p, g, \mathbf{q}; B) = \operatorname{var}\left(\frac{\partial l_{\mathbf{z}, n}}{\partial z_{i}}\right)$$
$$= \int_{\Omega_{B}} \left(\frac{\partial \log\left(f\left(z_{i}; \mathbf{p}\right)\right)}{\partial z_{i}}\right)^{2} (f\left(\mathbf{z}; \mathbf{p}\right)) dz_{1} \cdots dz_{n}$$
$$= \int_{-B}^{B} \left(\frac{\partial \log\left(f\left(z_{i}; \mathbf{p}\right)\right)}{\partial z_{i}}\right)^{2} (f\left(z_{i}; \mathbf{p}\right)) dz_{i},$$
(80)

where var denotes variance. Using (63) for $\partial l_{z,n}/\partial z_i$,

$$\nu_i(p, g, \mathbf{q}; B) = 2 \int_0^B \left(F(z_i; \mathbf{p}) \right)^2 f(z_i; \mathbf{p}) dz_i, \qquad (81)$$

where, for $z \in (0, B]$,

$$F(z;\mathbf{p}) = F(z;p,g,\mathbf{q}) = -p + \frac{1}{g(z)} \frac{dg(z;\mathbf{q})}{dz}$$
(82)

and $F(0; \mathbf{p}) = 0$. Since $v_i(p, g, \mathbf{q}; B)$ is independent of index *i*, we omit the subscript. Hence, the information matrix $\mathcal{J}_{\mathbf{z},n}$, defined with respect to $l_{\mathbf{z},n}$, is

$$\mathcal{J}_{\mathbf{z},n} = \nu I_n. \tag{83}$$

If $g(z; \mathbf{q}) = g_1(z) = 1$, then $F(z; p) = F(z; p, g_1) = -p$, and so

$$\nu\left(p,g_1;B\right) = p^2. \tag{84}$$

For the nontrivial perturbing function g_2 , for fixed p and q, one can show that ν depends on B by direct calculation.

4.5. The Information Matrix. We calculate the information matrix $\mathcal{F}_{\beta,n,m}$ (conditional on p, g, \mathbf{q} , and B). We are trying to quantify the steepness of the slope of $l_{\beta,n,m}$ around a maximum, in the directions represented by the coefficients β_j . If $l_{\beta,n,m}$ is very flat in one direction, then the model coefficient representing that direction is not well defined (will have large variance). When calculating $\mathcal{F}_{\beta,n,m}$, we are taking into account the behaviour of the gradient of $l_{\beta,n,m}$ on a whole neighbourhood of the MLE (how it differs from the expected value) and discontinuities on sets of measure zero can be accommodated. Recall (64), for j = 1, 2, ..., m,

$$\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \boldsymbol{\beta}_{j}} = p \sum_{i=1}^{n} x_{ij} \operatorname{sgn}(z_{i}) \\ - \sum_{i=1}^{n} x_{ij} \frac{\operatorname{sgn}(z_{i})}{g(|z_{i}|;\mathbf{q})} \left(\frac{dg(|z_{i}|;\mathbf{q})}{d|z_{i}|}\right).$$
(85)

Hence, omitting some subscripts on *l* for brevity, we obtain

$$E\left(\frac{\partial l_{\boldsymbol{\beta},n,m}}{\partial \boldsymbol{\beta}_{j}}\right)$$

$$= \int_{\Omega_{B}} \frac{\partial l}{\partial \beta_{j}} (z_{1},...,z_{n}) f(z_{1},...,z_{n};\mathbf{p}) dz_{1}\cdots dz_{n}$$

$$= \int_{\Omega_{B}} \left(\sum_{i=1}^{n} \left(\frac{\partial l}{\partial z_{i}} \frac{\partial z_{i}}{\partial \beta_{j}}\right)\right)$$

$$\times f(z_{1};\mathbf{p})\cdots f(z_{n};\mathbf{p}) dz_{1}\cdots dz_{n} \qquad (86)$$

$$= \sum_{i=1}^{n} (-x_{ij}) \int_{-B}^{B} \left(\frac{\partial l}{\partial z_{i}}\right) f(z_{i};\mathbf{p}) dz_{i}$$

$$= \sum_{i=1}^{n} (-x_{ij}) E\left(\frac{\partial l}{\partial z_{i}}\right)$$

$$=\sum_{i=1}^{n}\left(-x_{ij}\right)\mu=0$$

Since the expectations of the partial derivatives of $l_{\beta,n,m}$ are zero, the diagonal elements of $\mathcal{J}_{\beta,n,m}$ are an indication of the steepness of the gradient around the maximum likelihood estimate. Now,

$$\left(\mathscr{F}_{\boldsymbol{\beta},n,m}\right)_{jk} = E\left[\left(\frac{\partial l}{\partial \beta_j} - E\left(\frac{\partial l}{\partial \beta_j}\right)\right)\left(\frac{\partial l}{\partial \beta_k} - E\left(\frac{\partial l}{\partial \beta_k}\right)\right)\right]$$
$$= E\left[\left(\frac{\partial l}{\partial \beta_j}\right)\left(\frac{\partial l}{\partial \beta_k}\right)\right]$$
$$= \int_{\Omega_B} \left(\frac{\partial l}{\partial \beta_j}\right)\left(\frac{\partial l}{\partial \beta_k}\right)f(\mathbf{z};\mathbf{p})\,dz_1\cdots dz_n$$
$$= \int_{\Omega_B} \left(\sum_{i=1}^n \left(\frac{\partial l}{\partial z_i}\frac{\partial z_i}{\partial \beta_j}\right)\right)$$
$$\times \left(\sum_{t=1}^n \left(\frac{\partial l}{\partial z_t}\frac{\partial z_t}{\partial \beta_k}\right)\right)f(\mathbf{z};\mathbf{p})\,dz_1\cdots dz_n$$

$$= \int_{\Omega_{B}} \left(\sum_{i=1}^{n} \left(-x_{ij} \frac{\partial l}{\partial z_{i}} \right) \right) \\ \times \left(\sum_{t=1}^{n} \left(-x_{tk} \frac{\partial l}{\partial z_{t}} \right) \right) \\ \times f\left(z_{1}; \mathbf{p} \right) \cdots f\left(z_{n}; \mathbf{p} \right) dz_{1} \cdots dz_{n} \\ = \int_{\Omega_{B}} \left(\sum_{i=1}^{n} \left(x_{ij} \frac{\partial l}{\partial z_{i}} \right) \left(x_{ik} \frac{\partial l}{\partial z_{i}} \right) \right) \\ \times f\left(z_{1}; \mathbf{p} \right) \cdots f\left(z_{n}; \mathbf{p} \right) dz_{1} \cdots dz_{n},$$
(87)

since the cross-terms indexed by $i \neq t$ equal zero by symmetry, that is,

$$\int_{z_{i}=-B, z_{t}=-B}^{z_{i}=B, z_{t}=-B} \left(x_{ij} \frac{\partial l}{\partial z_{i}} \right) \left(x_{tk} \frac{\partial l}{\partial z_{t}} \right) f(z_{i}; \mathbf{p})$$

$$\times f(z_{t}; \mathbf{p}) dz_{i} dz_{t} = x_{ij} x_{tk} \mu^{2} = 0.$$
(88)

So

$$\left(\mathscr{F}_{\boldsymbol{\beta},n,m}\right)_{jk} = \int_{\Omega_{B}} \left(\sum_{i=1}^{n} \left(x_{ij}x_{ik}\left(\frac{\partial l}{\partial z_{i}}\right)^{2}\right)\right) \\ \times f\left(z_{1};\mathbf{p}\right)\cdots f\left(z_{n};\mathbf{p}\right)dz_{1}\cdots dz_{n} \\ = \sum_{i=1}^{n} \left(x_{ij}x_{ik}\right)\int_{-B}^{B} \left(\frac{\partial l}{\partial z_{i}}\right)^{2} f\left(z_{i};\mathbf{p}\right)dz_{i} \\ = \nu\left(p,g,\mathbf{q};B\right)\sum_{i=1}^{n} \left(x_{ij}x_{ik}\right).$$

$$(89)$$

Hence,

$$\mathcal{J}_{\boldsymbol{\beta},n,m} = \mathbf{X}_{n,m}^T \mathcal{J}_{\mathbf{z},n} \mathbf{X}_{n,m} = \nu \left(p, g, \mathbf{q}; B \right) \mathbf{X}_{n,m}^T \mathbf{X}_{n,m}.$$
(90)
If $g(z; \mathbf{q}) = g_1(z) = 1$, then $\mathcal{J}_{\boldsymbol{\beta},n,m} = p^2 \mathbf{X}_{n,m}^T \mathbf{X}_{n,m}.$

4.6. *The Expected Value of the Generalized Hessian*. In order to calculate the expected value of the generalized Hessian we require

$$E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right) = \mathbf{X}_{n,m}^{T} E\left(\mathscr{H}_{\mathbf{z},n}\right) \mathbf{X}_{n,m} = E\left(\frac{\partial^{2} l_{\mathbf{z},n}}{\partial z_{1}^{2}}\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m},$$
(91)

and

$$\frac{\partial^{2} l_{\mathbf{z},n}}{\partial z_{i}^{2}} = -2p\delta\left(z_{i}\right) + 2\delta\left(z_{i}\right)\frac{1}{g\left(\left|z_{i}\right|;\mathbf{q}\right)}\left(\frac{dg\left(\left|z_{i}\right|;\mathbf{q}\right)}{d\left|z_{i}\right|}\right) + \frac{\left(\operatorname{sgn}\left(z_{i}\right)\right)^{2}}{\left(g\left(\left|z_{i}\right|;\mathbf{q}\right)\right)^{2}} \times \left[g\left(\left|z_{i}\right|;\mathbf{q}\right)\frac{d^{2}g\left(\left|z_{i}\right|;\mathbf{q}\right)}{d\left|z_{i}\right|^{2}} - \left(\frac{dg\left(\left|z_{i}\right|;\mathbf{q}\right)}{d\left|z_{i}\right|}\right)^{2}\right].$$
(92)

Let

$$\begin{aligned} \zeta_{i}\left(p,g,\mathbf{q};B\right) &= E\left(\frac{\partial^{2}l_{\mathbf{z},n}}{\partial z_{i}^{2}}\right) = \int_{\Omega_{B}} \frac{\partial^{2}l}{\partial z_{i}^{2}} f\left(\mathbf{z};\mathbf{p}\right) dz_{1} \cdots z_{n} \\ &= \int_{-B}^{B} \frac{\partial^{2}l}{\partial z_{i}^{2}} f\left(z_{i};\mathbf{p}\right) dz_{i} \\ &= -2pf\left(0;\mathbf{p}\right) + 2f\left(0;\mathbf{p}\right) \frac{1}{g\left(0;\mathbf{q}\right)} \\ &\times \left(\frac{dg\left(|z_{i}|;\mathbf{q}\right)}{d|z_{i}|}\right)_{z_{i}=0} \\ &+ \int_{-B}^{B} \frac{1}{g\left(|z_{i}|;\mathbf{q}\right)} \frac{d^{2}g\left(|z_{i}|;\mathbf{q}\right)}{d|z_{i}|^{2}} f\left(z_{i};\mathbf{p}\right) dz_{i} \\ &- \int_{-B}^{B} \frac{1}{\left(g\left(|z_{i}|;\mathbf{q}\right)\right)^{2}} \left(\frac{dg\left(|z_{i}|;\mathbf{q}\right)}{d|z_{i}|}\right)^{2} \\ &\times f\left(z_{i};\mathbf{p}\right) dz_{i}, \end{aligned}$$

$$(93)$$

since $(dg(|z_i|; \mathbf{q})/d|z_i|)_{z_i=0} = (dg(u; \mathbf{q})/du)_{u=0}$, where $u = |z_i|$. The quantity $\zeta_i(p, g, \mathbf{q}; B)$ does not depend on *i*, so we omit the index. Hence,

$$E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) = \zeta\left(\boldsymbol{p},\boldsymbol{g},\mathbf{q};\boldsymbol{B}\right)\mathbf{X}_{n,m}^{T}\mathbf{X}_{n,m}.$$
(94)

If $g(z; \mathbf{q}) = g_1(z) = 1$, then

$$\zeta(p, g_1; B) = \frac{-2p^2}{Q(p, g_1; B)} = \frac{-p^2}{(1 - e^{-pB})},$$
(95)

and so although the generalised Hessian $\mathcal{H}_{\beta,n,m}$ has infinite elements, its expected value is negative definite. By continuity, if $g(z; \mathbf{q})$ is close enough to the constant map g_1 , the expected value $E(\mathcal{H}_{\beta,n,m})$ will still be negative definite. Note that if g has negative slope at the origin, the peak of $l_{z,n}$ at the origin becomes sharper, compared to that for the case $g = g_1$. If g has positive slope at the origin, we see the opposite effect.

4.7. The Generalized Variance-Covariance Matrix for the Model Coefficients. We use a generalized Taylor series expansion (in the coefficients β_j) to approximate $l_{\beta,n,m}$ by a negative definite quadratic function about a local maximum. Although we know that, for finite *n*, this approximation is not exact, we show in Section 5 that we would expect it to become more accurate as $n \to \infty$. Assuming $\mathbf{X}_{n,m}$ and $\mathbf{y} \in \mathbb{R}^n$ are given, conditional on *p* and **q**, we could write a Taylor series expansion for $l_{\beta,n,m}$ about a ML estimator $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}_n \in \mathcal{B}$ as follows:

$$l_{\boldsymbol{\beta},n,m}\left(\boldsymbol{\beta}\right) \approx l_{\boldsymbol{\beta},n,m}\left(\widehat{\boldsymbol{\beta}}_{n}\right) + \left(\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_{n}\right)^{T} \left(\nabla l_{\boldsymbol{\beta}}\right)_{\widehat{\boldsymbol{\beta}}_{n}} + \left(\frac{1}{2}\right) \left(\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_{n}\right)^{T} \mathscr{H}_{\widehat{\boldsymbol{\beta}}_{n}}\left(\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}_{n}\right) + \cdots,$$
(96)

if *f* and hence $L_{\beta,n,m}$ and $l_{\beta,n,m}$ were sufficiently smooth. Now our probability density function *f* and hence $l_{\beta,n,m}$ are not

sufficiently smooth, but we can replace the second derivative of $l_{\beta,n,m}$ by its expected value using generalized functions. Let $\Delta \beta = \beta - \hat{\beta}_n$. This yields the approximation

$$l_{\boldsymbol{\beta},n,m}\left(\boldsymbol{\beta}\right) \approx l_{\boldsymbol{\beta},n,m}\left(\widehat{\boldsymbol{\beta}}_{n}\right) + \left(\frac{1}{2}\right)\left(\Delta\boldsymbol{\beta}\right)^{T} E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) \Delta\boldsymbol{\beta}.$$
(97)

Assuming $E(\mathcal{H}_{\beta,n,m})$ is nonsingular, which is true when $g = g_1$, we ignore higher order terms. Equation (97) provides an indication of the behaviour of $l_{\beta,n,m}$ about a maximum since, for example, if $g = g_1$, then $E(\mathcal{H}_{\beta,n,m})$ is negative definite.

Next we consider the score function $\nabla l_{\beta,n,m}$ and use a Taylor series approximation incorporating generalized functions (about a local maximum $\hat{\beta}_n$) to derive a relationship between the expected value of the generalized Hessian $E(\mathcal{H}_{\beta,n,m})$, the information matrix $\mathcal{F}_{\beta,n,m}$, and the generalized variance-covariance matrix $\mathcal{V}_{\beta,n,m}$. This approximation is

$$\left(\nabla l_{\boldsymbol{\beta},n,m} \right)_{\boldsymbol{\beta}} \approx \left(\nabla l_{\boldsymbol{\beta},n,m} \right)_{\widehat{\boldsymbol{\beta}}_{n}} + E\left(\mathscr{H}_{\boldsymbol{\beta},n,m} \right) \Delta \boldsymbol{\beta}$$

$$= E\left(\mathscr{H}_{\boldsymbol{\beta},n,m} \right) \Delta \boldsymbol{\beta}.$$

$$(98)$$

Since $E(\nabla l_{\beta,n,m}) = 0$ and $E(\mathcal{H}_{\beta,n,m})$ has full rank, $E(\widehat{\beta}_n) = \beta$. We multiply each side by its own transpose and take expected values to obtain

$$E\left[\left(\nabla l_{\boldsymbol{\beta},n,m}\right)_{\boldsymbol{\beta}}\left(\nabla l_{\boldsymbol{\beta},n,m}\right)_{\boldsymbol{\beta}}^{T}\right]$$

$$=\left(E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right)\right)E\left[\Delta\boldsymbol{\beta}\left(\Delta\boldsymbol{\beta}\right)^{T}\right]\left(E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right)\right)^{T}.$$
(99)

Hence,

$$\mathcal{J}_{\boldsymbol{\beta},n,m} = \left(E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) \right) \mathcal{V}_{\boldsymbol{\beta},n,m} \left(E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) \right)^{T}$$

$$= \left(E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) \right) \mathcal{V}_{\boldsymbol{\beta},n,m} \left(E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) \right),$$
(100)

so

$$\mathcal{V}_{\boldsymbol{\beta},n,m} = \mathcal{V}_{\boldsymbol{\beta},n,m} \left(p, g, \mathbf{q}, \mathbf{X}_{n,m}; B \right)$$
$$= \left(E \left(\mathcal{H}_{\boldsymbol{\beta},n,m} \right) \right)^{-1} \mathcal{F}_{\boldsymbol{\beta},n,m} \left(E \left(\mathcal{H}_{\boldsymbol{\beta},n,m} \right) \right)^{-1}$$
$$= \frac{\nu \left(p, g, \mathbf{q}; B \right)}{\left(\zeta \left(p, g, \mathbf{q}; B \right) \right)^2} \left(\mathbf{X}_{n,m}^T \mathbf{X}_{n,m} \right)^{-1}.$$
(101)

Here, we require $\mathbf{X}_{n,m}$ to have full rank and that $\zeta(p, g, \mathbf{q}; B) \neq 0$, which is certainly true when $g = g_1$.

4.8. Generalized Statistical Expressions and Relations. We have shown that the expected value of the generalized Hessian

 $E(\mathcal{H}_{\beta,n,m})$ and the information matrix $\mathcal{J}_{\beta,n,m}$ are multiples of $\mathbf{X}_{n,m}^T \mathbf{X}_{n,m}$. Specifically,

$$E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) = \mathbf{X}_{n,m}^{T} E\left(\mathcal{H}_{\mathbf{z},n}\right) \mathbf{X}_{n,m}$$

$$= E\left(\frac{\partial^{2}l\left(z_{1}; p, g, \mathbf{q}\right)}{\partial z_{1}^{2}}\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}$$

$$= \zeta\left(p, g, \mathbf{q}; B\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}, \qquad (102)$$

$$\mathcal{J}_{\boldsymbol{\beta},n,m} = E\left(\left(\frac{\partial l\left(z_{1}; p, g, \mathbf{q}\right)}{\partial z_{1}}\right)^{2}\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}$$

$$= \nu\left(p, g, \mathbf{q}; B\right) \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}$$

(see (94) and (90)). Hence, if $v(p, g, \mathbf{q}; B) \neq 0$ (true if $g = g_1$), then

$$E\left(\mathscr{H}_{\boldsymbol{\beta},n,m}\right) = \frac{E\left(\partial^{2}l\left(z_{1}; p, g, \mathbf{q}\right)/\partial z_{1}^{2}\right)}{E\left(\left(\partial l\left(z_{1}; p, g, \mathbf{q}\right)/\partial z_{1}\right)^{2}\right)}\mathscr{F}_{\boldsymbol{\beta},n,m}$$

$$= \frac{\zeta\left(p, g, \mathbf{q}; B\right)}{\nu\left(p, g, \mathbf{q}; B\right)}\mathscr{F}_{\boldsymbol{\beta},n,m}.$$
(103)

In addition, assuming that the scalar $\zeta(p, g, \mathbf{q}; B)$ is also nonzero (true if $g = g_1$) and that $\mathbf{X}_{n,m}$ has full rank *m*, we have proved the following relations:

$$\left(\mathscr{V}_{\boldsymbol{\beta},n,m}\left(\boldsymbol{p},\boldsymbol{g},\mathbf{q},\mathbf{X}_{n,m};\boldsymbol{B}\right)\right)^{-1} = \frac{\left(\zeta\left(\boldsymbol{p},\boldsymbol{g},\mathbf{q};\boldsymbol{B}\right)\right)^{2}}{\nu\left(\boldsymbol{p},\boldsymbol{g},\mathbf{q};\boldsymbol{B}\right)} \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m}$$
(104)

$$= \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} E\left(\mathscr{H}_{\beta,n,m}\right)$$
(105)
$$= \frac{\left(\zeta(p, g, \mathbf{q}; B)\right)^2}{\left(\nu(p, g, \mathbf{q}; B)\right)^2} \mathscr{J}_{\beta,n,m},$$
(106)

to be used in the derivation of the generalized log-likelihood ratio statistic (Section 4.11).

4.9. Statistical Relations for the Laplace Distribution. If $g(z, \mathbf{q}) = g_1(z) = 1$, then using (21), (84), and (95):

$$Q(p, g_{1}; B) = 2(1 - e^{-pB}),$$

$$\nu(p, g_{1}; B) = p^{2},$$

$$(p, g_{1}; B) = \frac{-2p^{2}}{Q(p, g_{1}; B)} = \frac{-p^{2}}{(1 - e^{-pB})}.$$
(107)

Hence,

ζ

$$\mathcal{J}_{\boldsymbol{\beta},n,m} = p^{2} \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m},$$

$$E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) = \frac{\left(-p^{2}\right)}{\left(1-e^{-pB}\right)} \mathbf{X}_{n,m}^{T} \mathbf{X}_{n,m},$$
(108)

SO

$$\mathscr{V}_{\boldsymbol{\beta},n,m} = \frac{\left(1 - e^{-pB}\right)^2}{p^2} \left(\mathbf{X}_{n,m}^T \mathbf{X}_{n,m}\right)^{-1}.$$
 (109)

Here,

$$\mathcal{F}_{\boldsymbol{\beta},n,m} = -\left(1 - e^{-pB}\right) E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right),\tag{110}$$

$$\mathscr{V}_{\boldsymbol{\beta},n,m} = \left(1 - e^{-pB}\right)^2 \mathscr{J}_{\boldsymbol{\beta},n,m}^{-1},\tag{111}$$

where p > 0, B > 0, and $\mathbf{X}_{n,m}$ has full rank m. Hence, (75), derived for \mathscr{C}^2 distributions, does not hold for the Laplace distribution with bounded support. However, (75) describes the limiting behaviour, as $B \to \infty$.

4.10. Statistical Relations for a Laplace Distribution with Added Kurtosis. Now consider our motivating example, a Laplace distribution with bounded support [-1, 1] amended by adding kurtosis. In this case $g(z, \mathbf{q}) = g_2(z; q) = 1+qH_3(z)$ (see (14)), B = 1,

$$Q(p,q) = \frac{2\left[\left(p^{3} - 3qp^{2} + 6q\right) - e^{-p}\left(p^{3}\left(1 - 2q\right) + 6pq + 6q\right)\right]}{p^{3}},$$

f(z; p,q)

$$=\frac{p^4 e^{(-p|z|)} \left[1+q H_3\left(|z|\right)\right]}{2 \left[\left(p^3-3q p^2+6q\right)-e^{-p} \left(p^3 \left(1-2q\right)+6 p q+6q\right)\right]}.$$
(112)

Also, using (82) and (81), we obtain

$$F(z; p, q) = F(z; p, g_2, q) = -p + \frac{3q(z^2 - 1)}{1 + q(z^3 - 3z)},$$

$$\nu(p, g_2, q; 1) = 2 \int_0^1 (F(z; p, q))^2 f(z; p, q) dz,$$
(113)

and using (93), we obtain

$$\begin{split} \zeta\left(p,g_{2},q;1\right) &= -2pf\left(0;p,q\right) + 2f\left(0;p,q\right) \\ &\times \frac{1}{g\left(0;q\right)} \left(\frac{dg\left(|z|;q\right)}{d|z|}\right)_{z=0} \\ &+ \int_{-1}^{1} \frac{1}{g\left(|z|;q\right)} \frac{d^{2}g\left(|z|;q\right)}{d|z|^{2}} f\left(z;p,q\right) dz \\ &- \int_{-1}^{1} \frac{1}{\left(g\left(|z|;q\right)\right)^{2}} \left(\frac{dg\left(|z|;q\right)}{d|z|}\right)^{2} \\ &\times f\left(z;p,q\right) dz \\ &= \frac{-2p^{2}}{Q\left(p,q\right)} + \frac{-6pq}{Q\left(p,q\right)} \\ &+ 2\int_{0}^{1} \frac{1}{g\left(z;q\right)} \frac{d^{2}g\left(z;q\right)}{dz^{2}} f\left(z;p,q\right) dz \\ &- 2\int_{0}^{1} \frac{1}{\left(g\left(z;q\right)\right)^{2}} \left(\frac{dg\left(z;q\right)}{dz}\right)^{2} \\ &\times f\left(z;p,q\right) dz \end{split}$$

$$= \frac{-2p^{2} - 6pq}{Q(p,q)} + 2\int_{0}^{1} \frac{6qz}{(1+q(z^{3} - 3z))} f(z; p,q) dz - 2\int_{0}^{1} \frac{1}{(1+q(z^{3} - 3z)^{2})} \times (3q(z^{2} - 1))^{2} f(z; p,q) dz.$$
(114)

For typical values p = 5.254 and q = 0.025, numerical integration gives v = 28.3561 and $\zeta = -28.4957$, to four decimal places. In this case,

$$\mathcal{F}_{\boldsymbol{\beta},n,m} = 28.3561 \mathbf{X}_{n,m}^T \mathbf{X}_{n,m},$$

$$E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right) = -28.4957 \mathbf{X}_{n,m}^T \mathbf{X}_{n,m},$$
(115)

and so

$$\mathscr{V}_{\boldsymbol{\beta},n,m} = \frac{28.3561}{\left(-28.4957\right)^2} \left(\mathbf{X}_{n,m}^T \mathbf{X}_{n,m}\right)^{-1}.$$
 (116)

Here

$$\mathcal{F}_{\boldsymbol{\beta},n,m} = \frac{28.3561}{-28.4957} E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right),$$

$$\mathcal{V}_{\boldsymbol{\beta},n,m} = \frac{\left(28.3561\right)^2}{\left(-28.4957\right)^2} \mathcal{F}_{\boldsymbol{\beta},n,m}^{-1}.$$
 (117)

In this example, (75) could be used as an approximation but does not exactly describe the relationship between $E(\mathcal{H}_{\beta,n,m})$, $\mathcal{J}_{\beta,n,m}$, and $\mathcal{V}_{\beta,n,m}$.

4.11. The Generalized Log-Likelihood Ratio Statistic. The loglikelihood ratio statistic enables us to assess the adequacy of a model. It enables us to compare a model with M coefficients (parameters) with a model of interest which differs only in that it has fewer coefficients, say P, with $M > P \ge 1$, see [1], for example. We wish to compare a linear model with mcoefficients, the β_j , for j = 1, 2, ..., m with a lesser linear model with fewer coefficients. The aim is to decide whether or not the excluded coefficients are useful. Our comparison is conditional on the parameters p and \mathbf{q} .

Let $L_{\rho,n,M}(\rho_1,\ldots,\rho_M; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y})$ denote the likelihood function for the model with M parameters, and let $L_{\psi,n,P}(\psi_1,\ldots,\psi_P; p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y})$ denote the likelihood function for the lesser model. Let $\hat{\boldsymbol{\rho}}$ (or $\hat{\boldsymbol{\rho}}_n$) be the maximum likelihood estimator of $\boldsymbol{\rho} = (\rho_1,\ldots,\rho_M)^T$, and let $\hat{\boldsymbol{\psi}}$ (or $\hat{\boldsymbol{\psi}}_n$) be the maximum likelihood estimator of $\boldsymbol{\psi} = (\psi_1,\ldots,\psi_P)^T$. Then the likelihood ratio

$$\lambda = \frac{L_{\rho,n,M}\left(\hat{\boldsymbol{\rho}}; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y}\right)}{L_{\psi,n,P}\left(\hat{\psi}; p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}\right)}$$
(118)

15

is a ratio of two probabilities and will be greater than one since the model with *M* parameters provides the more complete description of the model. In our application, M = m, $\rho = \beta$, and usually (but not necessarily) P = m - 1. We show that a multiple (conditional on *p* and **q**) of

$$\log (\lambda) = \log \left(L_{\boldsymbol{\rho},n,M} \left(\widehat{\boldsymbol{\rho}}; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y} \right) \right) - \log \left(L_{\boldsymbol{\psi},n,P} \left(\widehat{\boldsymbol{\psi}}; p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y} \right) \right)$$
(119)

has a chi-squared distribution as follows.

The derivation of the log-likelihood ratio statistic (for smooth functions) may be found in the textbooks. For example, for generalized linear models (see [1]) where the log-likelihood function is smooth in a neighbourhood of its maximum, $D = 2\log(\lambda)$ is distributed approximately as $\chi^2(M-P, \delta_D)$. Here δ_D is a noncentrality parameter, a positive constant which will be near zero if the lesser model fits the data almost as well as the model with more coefficients. Consider our situation in which the log-likelihood function is continuous at a maximum, but where this maximum occurs at a vertex or possibly on a ridge in (β , l)-space, and generalized calculus is required to consider Taylor series expansions. Then, as in (97), around $\hat{\rho}$, replacing the Hessian of $\log(L_M) = l_M$ by its expected value, we obtain the following approximation:

$$\log \left(L_{\rho,n,M} \left(E\left(\widehat{\rho} \right) ; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y} \right) \right)$$

- $\log \left(L_{\rho,n,M} \left(\widehat{\rho} ; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y} \right) \right)$
= $\log \left(L_{\rho,n,M} \left(E\left(\widehat{\rho} \right) \right) \right) - \log \left(L_{\rho,n,M} \left(\widehat{\rho} \right) \right)$
 $\approx \frac{1}{2} \left(E\left(\widehat{\rho} \right) - \widehat{\rho} \right)^{T} E\left(\mathscr{H}_{\rho,n,M} \right) \left(E\left(\widehat{\rho} \right) - \widehat{\rho} \right).$ (120)

So, by (105) assuming that the scalar factors $v(p, g, \mathbf{q}; B)$ and $\zeta(p, g, \mathbf{q}; B)$ are both non-zero and that $\mathbf{X}_{n,M}$ has full rank

$$2\frac{E\left(\left(\partial^{2}l\left(z_{1}; p, g, \mathbf{q}\right)\right)/\partial z_{1}^{2}\right)}{E\left(\left(\left(\partial l\left(z_{1}; p, g, \mathbf{q}\right)\right)/\partial z_{1}\right)^{2}\right)}\left(l_{\rho,n,M}\left(E\left(\widehat{\rho}\right)\right) - l_{\rho,n,M}\left(\widehat{\rho}\right)\right)$$

$$\approx \left(E\left(\widehat{\rho}\right) - \widehat{\rho}\right)^{T}\frac{\zeta\left(p, g, \mathbf{q}; B\right)}{\nu\left(p, g, \mathbf{q}; B\right)}E\left(\mathscr{H}_{\rho,n,M}\right)\left(E\left(\widehat{\rho}\right) - \widehat{\rho}\right)$$

$$= \left(E\left(\widehat{\rho}\right) - \widehat{\rho}\right)^{T}\mathscr{V}_{\rho,n,M}^{-1}\left(E\left(\widehat{\rho}\right) - \widehat{\rho}\right),$$
(121)

which has the distribution $\chi^2(M)$ if the MLE has a normal distribution. We show in Section 5 that the distribution of the

MLE is asymptotically normal. Hence, if n is large enough, this will be a good approximation. Similarly

$$2\frac{E\left(\left(\partial^{2}l\left(z_{1}; p, g, \mathbf{q}\right)\right)/\partial z_{1}^{2}\right)}{E\left(\left(\left(\partial l\left(z_{1}; p, g, \mathbf{q}\right)\right)/\partial z_{1}\right)^{2}\right)}\left(l_{\psi,n,P}\left(E\left(\widehat{\psi}\right)\right) - l_{\psi,n,P}\left(\widehat{\psi}\right)\right)$$

$$\approx \left(E\left(\widehat{\psi}\right) - \widehat{\psi}\right)^{T}\frac{\zeta\left(p, g, \mathbf{q}; B\right)}{\nu\left(p, g, \mathbf{q}; B\right)}E\left(\mathscr{H}_{\psi,n,P}\right)\left(E\left(\widehat{\psi}\right) - \widehat{\psi}\right)$$

$$= \left(E\left(\widehat{\psi}\right) - \widehat{\psi}\right)^{T}\mathscr{V}_{\psi,n,P}^{-1}\left(E\left(\widehat{\psi}\right) - \widehat{\psi}\right),$$
(122)

which has the distribution $\chi^2(P)$, approximately. Noting that

$$\frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} = \frac{-1}{(1 - e^{-pB})} < 0,$$
(123)

when $g = g_1$, let

$$D_{\text{gen}} = -2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \log(\lambda).$$
(124)

Then D_{gen} , the log-likelihood ratio statistic calculated with generalized functions, a positive number, may be expressed as

$$D_{\text{gen}} = -2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \\ \times \left(l_{\rho,n,M}(\hat{\rho}; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y}) - l_{\psi,n,P}(\widehat{\psi}; p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}) \right) \\ = -2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \\ \times \left(l_{\rho,n,M}(\hat{\rho}; p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y}) - l_{\rho,n,M}(E(\hat{\rho}); p, \mathbf{q}; \mathbf{X}_{n,M}, \mathbf{y}) \right) \\ + 2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \\ \times \left(l_{\psi,n,P}(\widehat{\psi}; p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}) - l_{\psi,n,P}(E(\widehat{\psi}); p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}) \right) \\ - 2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \\ \times \left(l_{\rho,n,M}(E(\hat{\rho}); p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}) \right) \\ - 2 \frac{\zeta(p, g, \mathbf{q}; B)}{\nu(p, g, \mathbf{q}; B)} \\ \times \left(l_{\rho,n,M}(E(\hat{\rho}); p, \mathbf{q}; \mathbf{X}_{n,P}, \mathbf{y}) \right) .$$
(125)

Hence, D_{gen} is the sum of three terms; the first (positive) has the distribution $\chi^2(M)$ (approximately). The second (negative) has the distribution $\chi^2(P)$ (approximately). The third is a positive constant (say δ_{gen}) that depends on p and \mathbf{q} . Hence, D_{gen} is distributed approximately as $\chi^2(M - P, \delta_{\text{gen}})$. If the lesser model gives a good description of the data, then δ_{gen} will be small. The generalized log-likelihood ratio statistic D_{gen} is easily calculated and is hence a potentially useful statistic for assessing our linear model for which the log-likelihood function is not \mathscr{C}^2 . Note if $g = g_1$, $D_{\text{gen}} = (2 \log(\lambda))/(1 - e^{-pB})$, and so $D_{\text{gen}} \rightarrow 2 \log(\lambda)$, as $B \rightarrow \infty$.

5. The Maximum Likelihood Estimator Is Consistent and Asymptotically Normal

The generalized expressions derived in Section 4 will be used to prove the asymptotic convergence of our MLE to a random variable with a normal distribution. Recall the following assumptions.

- (i) Our model is linear (see (4)).
- (ii) The response function *f* is a Laplace probability density function, generally perturbed and/or truncated, as given in (10).

We make the following further assumptions.

- (i) There exists a unique true vector of coefficients $\beta_0 \in \mathbb{R}^m$ whose value we are trying to estimate.
- (ii) The matrix $\mathbf{X}_{n,m}$ has full rank *m*, so that $\mathbf{X}_{n,m}^T \mathbf{X}_{n,m}$ has full rank *m*.
- (iii) The $\lim_{n\to\infty} (1/n) \mathbf{X}_{n,m}^T \mathbf{X}_{n,m} = \mathbf{W}_m$, a positive definite matrix.
- (iv) Assuming fixed *m*, for $n \ge m$, denote by $\hat{\beta}_n$ a (not necessarily unique) MLE of the true value β_0 , corresponding to the explanatory variables in $\mathbf{X}_{n,m}$.

Lemma 11. The ML estimates β_n exist, for $n \ge m$.

Proof of Lemma 11. Since a continuous function on a compact set attains it maximum, the existence of a maximum of the log-likelihood function $l_{\beta,n,m}$ is guaranteed for finite bound *B*. Even if the domain is not bounded, we can work with finite bound *B*, and then let $B \to \infty$.

Theorem 12. The random variable $\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)$ converges in distribution to an *m*-dimensional normally distributed random vector with mean **0** and covariance matrix $(\nu/\zeta^2)\mathbf{W}_m^{-1}$, that is, as $n \to \infty$,

$$\sqrt{n}\left(\widehat{\boldsymbol{\beta}}_{n}-\boldsymbol{\beta}_{0}\right) \xrightarrow{D} N\left(\mathbf{0},\left(\frac{\nu}{\zeta^{2}}\right)\mathbf{W}_{m}^{-1}\right).$$
 (126)

Lemma 13. The random variable $(1/\sqrt{n})\nabla l_{\beta,n,m}$ converges in distribution to an *m*-dimensional normally distributed random vector with mean **0** and covariance matrix $v \mathbf{W}_m$, that is, as $n \to \infty$,

$$\left(\frac{1}{\sqrt{n}}\right) \nabla l_{\boldsymbol{\beta},n,m} \xrightarrow{D} N\left(\mathbf{0}, \nu \mathbf{W}_{m}\right).$$
(127)

Hence, $(1/n)\nabla l_{\beta,n,m}$ converges in probability to **0**.

Proof of Lemma 13. Consider the random variable $\nabla l_{\beta,n,m} = -\mathbf{X}_{n,m}^T \nabla l_{z,n}$. If we repeat the sampling of $n \ge m$ data points, using $\mathbf{X}_{n,m}$ a total of *t* times, we may write

$$\nabla l_{\boldsymbol{\beta},nt,m} = -\mathbf{X}_{nt,m}^{T} \nabla l_{z,nt},$$

$$= -t\mathbf{X}_{n,m}^{T} \overline{\nabla l_{z,n}},$$
(128)

where $\nabla l_{z,n}$ denotes the average of *t* samples of the random vector $\nabla l_{z,n}$. Now $\nabla l_{z,n}$ has mean **0** and covariance matrix νI_n , and so by the multivariate Central Limit Theorem, as $t \to \infty$,

$$\sqrt{t} \ \overline{\nabla l_{z,n}} \xrightarrow{D} N(\mathbf{0}, \nu I_n),$$

$$\mathbf{X}_{n,m}^T \sqrt{t} \ \overline{\nabla l_{z,n}} \xrightarrow{D} N(\mathbf{0}, \mathbf{X}_{n,m}^T \nu \mathbf{X}_{n,m}),$$

$$\left(\frac{1}{\sqrt{t}}\right) \nabla l_{\boldsymbol{\beta},nt,m} \xrightarrow{D} N(\mathbf{0}, \mathbf{X}_{n,m}^T \nu \mathbf{X}_{n,m}),$$

$$\left(\frac{1}{\sqrt{nt}}\right) \nabla l_{\boldsymbol{\beta},nt,m} \xrightarrow{D} N(\mathbf{0}, \frac{\nu}{n} \mathbf{X}_{n,m}^T \mathbf{X}_{n,m}).$$
(129)

Hence, as $n \to \infty$ and $t \to \infty$,

$$\left(\frac{1}{\sqrt{nt}}\right) \nabla l_{\boldsymbol{\beta}, nt, m} \xrightarrow{D} N\left(\mathbf{0}, \nu \mathbf{W}_{m}\right).$$
(130)

Hence, as $n \to \infty$,

$$\left(\frac{1}{\sqrt{n}}\right) \nabla l_{\boldsymbol{\beta},n,m} \xrightarrow{D} N\left(\mathbf{0}, \boldsymbol{\nu} \mathbf{W}_{m}\right).$$
(131)

Lemma 14. Although for finite n, the log-likelihood function $l_{\beta,n,m}$ is not differentiable at a maximum; $(1/n)l_{\beta,n,m}$ converges in distribution to a negative definite quadratic function centred at β_0 . Hence, the MLEs $\hat{\beta}_n$ are consistent, that is, they converge in probability to the true value β_0 .

Proof of Lemma 14. Using a generalized Taylor series expansion about β_0 and noting that by Lemma 13, $(1/n)\nabla l_{\beta,n,m}$ converges in probability to **0**, we can write

$$\lim_{n \to \infty} \frac{l_{\boldsymbol{\beta},n,m} (\boldsymbol{\beta}) - l_{\boldsymbol{\beta},n,m} (\boldsymbol{\beta}_0)}{n}$$

$$= \lim_{n \to \infty} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T \frac{E(\mathcal{H}_{\boldsymbol{\beta},n,m})}{n} (\boldsymbol{\beta} - \boldsymbol{\beta}_0)$$

$$= (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T \left(\lim_{n \to \infty} \frac{\zeta \mathbf{X}_{n,m}^T \mathbf{X}_{n,m}}{n}\right) (\boldsymbol{\beta} - \boldsymbol{\beta}_0)$$

$$= \zeta (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T \mathbf{W}_m (\boldsymbol{\beta} - \boldsymbol{\beta}_0).$$
(132)

This shows that, in the limit as $n \to \infty$, the log-likelihood function $l_{\beta,n,m}$ has an isolated maximum at β_0 , and so a sequence of MLEs $\hat{\beta}_n$ must converge to β_0 in probability (consistency).

In the Proof of Lemma 14 we ignored the third partial derivatives of $l_{\beta,n,m}$ in the generalized Taylor series expansion. The justification is as follows. Since g is assumed to be \mathscr{C}^3 on an open interval which contains [0, B], $|d^3g/dz^3|$ must be bounded above by some positive real number G_3 on [0, B]. Hence, the absolute values of the third partial derivatives of l must be bounded above by some positive real number except at points \mathbf{z} , where some $z_i = 0$ (data points). Note that $E(\partial^3 l/\partial z_i^3) = 0$, i = 1, 2, ..., n. This follows from (66), using the symmetry introduced by the modulus function and (62). Hence, the third partial derivative terms will be small

Proof of Theorem 12. Consider the first degree approximation

compared to the second partial derivative terms near a critical point, and so we may ignore them in our generalized Taylor series expansion for *l*, when the expected value of the Hessian

$$\nabla l_{\boldsymbol{\beta},n,m}\left(\boldsymbol{\beta}\right) \sim E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right)\left(\boldsymbol{\beta}-\boldsymbol{\beta}_{0}\right). \tag{133}$$

Since $E(\mathcal{H}_{\beta,n,m})$ has full rank,

$$(\boldsymbol{\beta} - \boldsymbol{\beta}_0) \sim \left(E\left(\mathscr{H}_{\boldsymbol{\beta},n,m} \right) \right)^{-1} \nabla l_{\boldsymbol{\beta},n,m} \left(\boldsymbol{\beta} \right), \qquad (134)$$

so

has full rank.

$$\frac{n}{\sqrt{n}} \left(\boldsymbol{\beta} - \boldsymbol{\beta}_0\right) \sim \left(\frac{E\left(\mathcal{H}_{\boldsymbol{\beta},n,m}\right)}{n}\right)^{-1} \frac{\nabla l_{\boldsymbol{\beta},n,m}\left(\boldsymbol{\beta}\right)}{\sqrt{n}}.$$
 (135)

By Lemma 13, as $n \to \infty$,

$$\left(\frac{1}{\sqrt{n}}\right) \nabla l_{\boldsymbol{\beta},n,m} \xrightarrow{D} N\left(\mathbf{0}, \boldsymbol{\nu} \mathbf{W}_{m}\right).$$
(136)

Hence, as $n \to \infty$,

$$\sqrt{n} \left(\boldsymbol{\beta} - \boldsymbol{\beta}_{0}\right) \xrightarrow{D} N\left(\mathbf{0}, \left(\zeta \mathbf{W}_{m}\right)^{-1} \nu \mathbf{W}_{m} \left(\zeta \mathbf{W}_{m}\right)^{-T}\right).$$
(137)

Hence, as $n \to \infty$,

$$\sqrt{n} \left(\boldsymbol{\beta} - \boldsymbol{\beta}_0 \right) \xrightarrow{D} N \left(\mathbf{0}, \left(\frac{\nu}{\zeta^2} \right) \mathbf{W}_m^{-1} \right).$$
 (138)

An alternative Proof of Theorem 12 for the simplest case, that is, in the absence of truncation or perturbation, (i.e., $g = g_1$ and $B = \infty$) is in [2]. This case is also discussed in [11, page 451].

6. Real and Simulated Data Illustrations

6.1. Empirical Distribution of Methylation Proportion Deviations. Quantitative analysis of DNA methylation at specific genomic sites (known as CpG sites) was carried out with the Sequenom MassARRAY Compact System (http:// www.sequenom.com/). Briefly, this involves accurate determination and comparison of the mass of transcription cleavage products following chemical modification of the DNA which is dependent upon the *a priori* methylation status,

FIGURE 4: Deviations in methylation proportion.

using MALDI-TOF mass spectrometry (Bruker-Sequenom) (see [12]). Quantitative CpG methylation was assessed using proprietary EpiTyper software v1.0.5 (http://www.sequenom .com/). Sequenom measurements of 1440 CpG sites in each of 41 human umbilical cord tissue samples were performed in duplicate, and the difference between the measurements was recorded. This difference represents the deviation in the measurement of CpG methylation due to sample nanodispensing and MALDI-TOF mass spectrometry detection. Figure 4 is a histogram of the deviations of these 1440 repeated measurements. Although not obvious from the histogram, about 1.04% of values are greater than 0.2 in absolute value. The histogram has heavy tails. The data was shown not to conform to a normal distribution using the test proposed by [13] (*P* value <0.00001).

6.2. Simulated Data Example Using Methylation Proportion Deviations. In order to illustrate the application of the theory developed, a sample of 40 deviations was chosen at random from the total pool of 1440 available CpG methylation proportion deviations. A constant value of 0.48 was added to 20 of these samples and designated treatment H, while a constant value of 0.45 was added to the other 20 samples and designated treatment L. A uniform random variable sampled between -0.01 and 0.01 was added to each value to simulate the additional differences expected to occur between individuals.

We analysed the data using our amended Laplace distribution (13), with $g = g_2$, p = 53.41, q = 0.0314 (machine characteristics), and B = 1. The parameter MLEs $p = 53.41 \pm 3.45$ and $q = 0.0314 \pm 0.0097$ (estimates plus or minus standard errors) were given by Buckland's algorithm [14, 15], with the *q* being highly significantly different from zero, (with a log-likelihood ratio statistic of 135.6 distributed as $\chi^2(1)$). Assuming q = 0.0 (LAE regression with truncation) yields $p = 36.22 \pm 1.167$. Buckland's algorithm does not take into account the truncation, but since *p* is large and $(1 - e^{-pB})$ is so close to 1 in this example, the effect of truncation on the standard errors and D_{gen} (for *P* values) is negligible. If

the machine characteristic p was smaller, say $1 \le p \le 3$, assuming B = 1, this effect would become more significant.

We coded a low value treatment (L) by setting $x_{i2} = -1$, i = 1, 2, ..., 20 and a high value treatment (H) by setting $x_{i2} = 1$, i = 21, 22, ..., 40. Estimates for β_1 and β_2 were calculated by maximum likelihood estimation, using the simplex method. The standard errors of the coefficient estimates were calculated using our generalized \mathcal{V} . Hence, the estimated treatment means ($\beta_1 \pm \beta_2$), and their standard errors were calculated. A *P* value was obtained using D_{gen} , which is assumed distributed $\chi^2(1)$. This simulation was performed twice, and the results are given in Table 1. Note for each simulation, the *P* value is less than 0.01, indicating a significant difference between the means.

For both simulations $v(p, g_2, q; 1) = 2862.70$ and $\zeta(p, g_2, q; 1) = -2862.71$ (see Section 4.10). Also,

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 40 & 0\\ 0 & 40 \end{pmatrix},\tag{139}$$

rounding to five significant figures,

$$E(\mathscr{H}) = \begin{pmatrix} -114510 & 0\\ 0 & -114510 \end{pmatrix},$$
$$\mathscr{J} = \begin{pmatrix} 114510 & 0\\ 0 & 114510 \end{pmatrix},$$
$$(140)$$
$$\mathscr{V} = \begin{pmatrix} 8.7330e - 06 & 0\\ 0 & 8.7330e - 06 \end{pmatrix}.$$

Here,

$$E(\mathscr{H}) = -1.0000 \mathscr{J},$$

 $\mathscr{V}^{-1} = 1.0000 \mathscr{J},$ (141)
 $D_{\text{gen}} = 2.0000 \log(\lambda).$

These expressions correspond to the classical expressions for smooth functions, up to our numerical tolerance.

For comparison, the results of a standard analysis of variance, assuming the deviations have a normal distribution, are also given in Table 1, for each of the two simulations. Using the *P* values obtained, 4e-01 and 2e-01, we would not reject the null hypothesis (that the means are equal) at any usual level of significance, using a least squares approach. However, our algorithm based on the amended Laplace distribution correctly identified the structure of the simulated data set, separating two means which were fairly close. The standard least squares algorithm failed to do this. The two data sets analysed are given in Tables 3 and 4 in Appendix B.

For comparison we also included LAE regression, estimating the model coefficients assuming the response function is the Laplace probability density function truncated to [-1, 1], using p = 36.22, given by Buckland's algorithm. We see that including the perturbing Hermite polynomial improves (decreases) the *P* values, meaning we can be even more confident, than when using LAE regression, that the means are different.

	Amended Laplace distribution p = 53.41, q = 0.0314			La	place distribut	ion	Normal distribution		
Т				p = 36.22, q = 0.0			MATLAB glmfit		
	Mean	Std. err.	P value	Mean	Std. err.	P value	Mean	Std. err.	P value
Н	0.4817	0.0042	2e - 05	0.4812	0.0062	4a - 04	0.4658	0.0094	4a = 01
L	0.4532	0.0042		0.4519	0.0062	40 - 04	0.4540	0.0094	40 - 01
Н	0.4803	0.0042	2e - 04	0.4809	0.0062	2a - 03	0.4887	0.0119	2a = 01
L	0.4592	0.0042		0.4586	0.0062	2e = 03	0.4673	0.0119	20 01

TABLE 1: Two simulation results, adding high (H) or low (L) treatments (T) to DNA methylation proportions.

TABLE 2: Primiparous (p) versus multiparous (m) effects on DNA methylation proportion at the promoter of the H19 gene.

Site	Parity	Amended Laplace distribution p = 75.53, q = 0.4999			N	Mann-Whitney		
		Mean	Std. err.	P value	Mean	Std. err.	P value	P value
CpG9	р	0.180	0.0029		0.300	0.059	1e = 01	5e - 02
CpG9	m	0.450	0.0029	<1.0e-09	0.441	0.059	10 01	50 02
CpG13	р	0.230	0.0029		0.326	0.061	4a = 02	4a - 02
CpG13	m	0.560	0.0029	<1.0e-09	0.523	0.061	HL = 02	-10 = 02

TABLE 3: First simulation data, treatments either H ($x_i = 1$) or L ($x_i = -1$) plus randomly sampled methylation deviances and randomly sampled uniformly distributed individual variation.

Т	y_i	Т	y_i	Т	y_i	Т	y_i
L	0.4579	L	0.4467	Η	0.4841	Η	0.4873
L	0.4243	L	0.4610	Η	0.4735	Η	0.4761
L	0.4993	L	0.4609	Η	0.4878	Η	0.2391
L	0.4131	L	0.4851	Η	0.4823	Η	0.4805
L	0.4463	L	0.4340	Η	0.4462	Η	0.4877
L	0.4317	L	0.4573	Η	0.4664	Η	0.4779
L	0.4473	L	0.4360	Η	0.4845	Η	0.4929
L	0.4347	L	0.4584	Η	0.4817	Η	0.4863
L	0.4760	L	0.4420	Η	0.4861	Η	0.4751
L	0.4776	L	0.4914	Η	0.4668	Η	0.4543

TABLE 4: Second simulation data, treatments either H ($x_i = 1$) or L ($x_i = -1$) plus randomly sampled methylation deviances and randomly sampled uniformly distributed individual variation.

Т	y_i	Т	y_i	Т	y_i	Т	y_i
L	0.5287	L	0.4416	Η	0.4881	Η	0.4829
L	0.5224	L	0.4547	Η	0.4803	Η	0.4246
L	0.5162	L	0.4496	Η	0.4789	Η	0.4609
L	0.4564	L	0.4568	Η	0.4790	Η	0.4233
L	0.4628	L	0.4574	Η	0.4739	Η	0.5412
L	0.5230	L	0.4599	Η	0.4974	Η	0.5193
L	0.3731	L	0.5124	Η	0.7010	Η	0.4921
L	0.4389	L	0.4592	Η	0.4725	Η	0.4702
L	0.4519	L	0.4458	Η	0.4871	Η	0.5520
L	0.4675	L	0.4685	Η	0.4885	Η	0.3600

6.3. Primiparous versus Multiparous Effects on DNA Methylation Proportion at the Promoter of the H19 Gene. The CpG methylation at two CpG sites in the promoter of the H19

TABLE 5: CpG methylation measurements at sites 9 and 13 on the promoter of the H19 gene versus primiparous (p) or multiparous (m).

CpG9	p/m	CpG9	p/m	CpG13	p/m	CpG13	p/m
1.00	р	0.16	р	0.30	р	0.16	р
0.08	р	0.19	р	0.00	р	0.36	р
0.04	р	0.15	р	0.03	р	0.02	р
0.17	р	0.35	р	0.25	р	0.60	р
0.46	р	0.04	р	0.80	р	0.01	р
1.00	р	0.27	р	0.71	р	0.70	р
0.18	р	0.32	р	0.17	р	0.56	р
0.33	m	0.37	р	0.56	m	0.70	р
0.28	m	0.05	р	0.40	m	0.00	р
0.82	m	0.39	р	0.57	m	0.61	р
0.20	р	0.07	р	0.18	р	0.02	р
0.08	р	0.17	р	0.03	р	0.23	р
0.15	р	0.14	m	0.09	р	0.99	m
1.00	р	0.61	m	0.96	р	0.60	m
0.10	m	0.53	m	0.79	m	0.35	m
0.89	m	0.45	m	0.83	m	0.63	m
0.07	m	0.09	m	0.02	m	0.07	m
0.62	m	0.57	m	0.38	m	0.53	m
0.48	m	0.73	m	0.68	m	0.72	m
0.31	m	0.30	m	0.27	m	0.22	m
0.62	m			0.80	m		

gene was measured by Sequenom in umbilical cord tissues collected as part of an ongoing prospective birth cohort study. Phenotype variables in this population include birth order or parity, defined as first born child (primiparous) or later born (multiparous). We analysed the relationship between H19 gene methylation status and birth order in this study, using our amended Laplace distribution (13), with $g = g_2$ and

B = 1, and for comparison, the usual least squares algorithm. Here an amended Laplace distribution (13) has been used to model the measurement errors as opposed to the differences of errors in repeated measurements (the deviations in Section 6.2). In practice this works well, despite the observation that if the measurement errors have a distribution of the form (13), then the deviations would not be expected to have exactly this form of distribution. Modelling the deviations as the difference between two random variables both distributed Laplace also works well. We used the estimate of $p_d = 36.22$ for the deviations to estimate $p_e = 51.22 \sim 36.22 * \sqrt{2}$ for the errors. The factor $\sqrt{2}$ is chosen to halve the variance in the underlying Laplace distribution. For the error distribution, we estimated q_e to be as close to 0.5 as we allow, say $q_e =$ 0.4999. This was done by assuming the errors were distributed as an amended Laplace distribution (13) with parameters p_e and q_e , deriving a corresponding distribution for the deviations and using MLE with the deviation data to obtain parameter estimates. We also tried $p_e = 75.53 \sim 53.41 * \sqrt{2}$ with $q_e = 0.4999$. The latter (Sequenom) parameter estimates gave smaller standard errors and P values and were consistent with Infinium parameter estimates.

The problem was coded by substituting $x_{i2} = 1$ for primiparous and $x_{i2} = -1$ for multiparous. Estimates for β_1 and β_2 and their standard errors were calculated. The data used is given in Table 5 in Appendix B. The estimated means $(\beta_1 \pm \beta_2)$ and their standard errors are given in Table 2. The small *P* values associated with the amended Laplace distribution (13), calculated using D_{gen} , identify a difference between the mean methylation proportions (primiparous versus multiparous) at a given CpG site. This demonstrates the power of accounting properly for the distributional properties of the methylation errors and hence enables clearer inference of the epigenetic mechanisms underlying these biological phenomena.

In this example, setting q = 0.0 (LAE regression with truncation) yields similar means and variances to our MLE algorithm (the same to two decimal places) and also gives small *P* values. For the CpG13 example with p = 75.53 and q = 0.4999, $D_{gen} = 201.479$. For this example with p = 75.53 and q = 0.0, $D_{gen} = 197.889$. These values are very close as *pB* is large and so e^{-pB} is small. For smaller values of *pB*, the value of taking into account the truncation and perturbation increases. A Mann-Whitney 2-sample *U*-test was also done for comparison.

7. Discussion

7.1. A Comparison with LAE Regression. The original MLE theory and methods in this paper were developed assuming the response function is a modified version of the Laplace probability density function, that is, assuming nontrivial perturbation and/or truncation to compact support [-B, B]. Such response functions have been observed in measurement data generated by nearly all the analytical platforms currently used to assess DNA methylation, including the Sequenom EpiTyper, Infinium Mass Array and Restricted Representation Bisulphite Sequencing platforms [16].

In the absence of perturbation or truncation of the response function, the results in this paper correspond to the

theory of LAE (or median) regression as found in [2–4]. That is, if g(z) = 1 and $B = \infty$, then our MLE method corresponds with LAE regression. In this case, Q = 2, $\nu = p^2$, $\zeta = -p^2$, $f(z) = (p/2)e^{-p|z|}$ (in one dimension), and so

$$\frac{\nu}{\zeta^2} = \frac{p^2}{p^4} = \frac{1}{p^2} = \frac{1}{\left(2f\left(0\right)\right)^2},\tag{142}$$

which is the asymptotic variance of the ordinary sample median for f [2].

We present an original and practical method of obtaining the covariance matrix for the model coefficients. This involves evaluating $(\mathbf{X}_{n,m}^T \mathbf{X}_{n,m})^{-1}$, where although *n* might be large, *m* generally is not, and using generalized functions to numerically evaluate two one-dimensional integrals (to find ν and ζ). The calculation of ν and ζ takes into account the characteristics of the response function (truncation and perturbation) which would be ignored if we used median regression. This is possible when we know the response function parameters, or have fairly accurate estimates, as in our epigenetic application, modelling DNA methylation proportion deviations.

For LAE regression, other methods of determining approximations to this covariance matrix may be found in the literature. In particular, in the method of quantile regression [5] implemented in the statistical package R, the covariance matrix for the MLE is calculated by resampling techniques, by bootstrapping or by using hierarchical spline models [17].

We prove that, even for truncated and perturbed Laplace response functions, subject to certain restrictions, the maximum of the log-likelihood function occurs at a data point. This result is well-known in the case of LAE regression. A proof that the LAE estimator passes through at least r_x data points may be found in references [3, 4].

Three asymptotically equivalent test statistics for LAE regression may be found in [18], namely, a likelihood ratio test statistic, a Wald test statistic, and a Lagrange multiplier test statistic. Our likelihood ratio test statistic is an original modification of the former, applicable to our general case (not restricted to LAE regression), calculated using generalized functions. An *F* test statistic for LAE regression is found in [4].

When working with a model for which the response function is assumed to be a truncated Laplace probability density function, we could ignore the truncation to [-B, B]. However, taking the truncation into account reduces the variance in the model coefficient estimates by a factor $(1 - e^{-pB})^2$ and increases the log-likelihood ratio statistic by a factor $(1 - e^{-pB})^{-1}$. This effect is small if e^{-pB} is small but becomes more significant as pB decreases, that is, as more of the density function is truncated. Refer to Section 4.9, (111) and Section 4.11, (123), and (124) which show that, for example, when $g = g_1$,

$$D_{\text{gen}} = \frac{2}{\left(1 - e^{-pB}\right)} \log\left(\lambda\right). \tag{143}$$

Hence, by taking into account the truncation, we can be more confident of our coefficient estimates and the value of appropriate beta coefficients than the standard formulae for LAE regression indicate. This effect will also be seen for small perturbations of the density function.

7.2. Original Formulae Enable Model Comparison. The original formulae derived in Section 4 for the covariance matrix of the model coefficients and the generalized log-likelihood ratio statistic enable us to do a generalized analysis of variance, in our case comparing the means of two different groups of data. We see that in Section 6.3, the primiparous and multiparous means could be judged distinct with very high probability using either the Laplace or an amended Laplace response function. The *P*-values are slightly lower when using the Hermite modification. In the simulation in Section 6.2, using the Hermite modification decreases the *P*value by an order of magnitude. In a situation in which *p* is low or two means are less distinct, the Hermite modification may prove very useful.

Preliminary results indicate the use of an amended Laplace distribution enables a clearer separation of means than that given by other more standard procedures, for example, beta regression, in cases where independent evidence suggests that the means are different [16].

7.3. Summary. The Laplace distribution is the basis of many mathematical models (see [19]). Our focus has been modelling the distributions of errors in the proportions of DNA methylation measured at genomic CpG sites.

Molecular biology deals with complex interactions both in terms of the physiology of the processes of interest and in the instrumentation required to measure these effects. The non-linearity of these processes can result in frequency distributions that are far from normal, so that application of "standard" methods of statistical inference based on least squares may be inadequate. Methods which deal with the form of the frequency distribution directly such as maximum likelihood are necessary for adequate inference to be made.

The Laplace or double exponential distribution considered here has been observed in molecular biology studies, where a significant proportion of high deviations appears to occur regularly [16, 20]. The extension by Hermite polynomials considered here provides flexibility for describing the tails of the distribution. However, as noted, the use of the Laplace distribution as the "key" function introduces problems in finding maximum likelihood estimators and particularly their standard errors. This paper presents both a practical method for dealing with these problems and the underlying asymptotic theory.

Appendices

A. Useful Convex Analysis Results

We prove Lemma 8 by applying results from convex analysis (see [8]). The following definitions are taken from [8]. A face of a convex set $C \subset \mathbb{R}^n$ is a convex subset C' of C such that every (closed) line segment in C with a relative interior point in C' has both endpoints in C'. The empty set and C itself are faces of C. The zero-dimensional faces of C are called the extreme points of C. The relative interior of a convex set $C \subset$

 \mathbb{R}^n is defined as the interior which results when *C* is regarded as a subset of its affine hull. The affine hull of a set $S \subset \mathbb{R}^n$ is the unique smallest affine set containing *S*. An alternative definition of an extreme point of a convex set *C* is a point $z \in$ *C* that cannot be written as $z = \theta u + (1 - \theta)v$ with $0 < \theta < 1$, $u \in C$, $v \in C$, and $u \neq v$ [21, p686].

Theorem A.1. Let C be a compact convex set in \mathbb{R}^n , and let $f : \mathbb{R}^n \to \mathbb{R}$ be a linear function. The maximum and minimum of f are attained at extreme points of C.

Theorem A.1 (on the Maximum/Minimum Property (see [21])) is useful but does not give a complete characterisation of the set of points in *C* at which *f* has a maximum.

Rockafellar [8] gives a more general definition of a convex function than we require. It is enough for our purposes to say that if the domain of real-valued function f is a convex set in \mathbb{R}^n and if for any u and v in this domain, $f(\lambda u + (1 - \lambda)v) \le \lambda f(u) + (1 - \lambda)f(v)$ for all $\lambda \in [0, 1]$, then f is convex.

Theorem A.2 (see [8, Theorem 32.1]). Let f be a convex function, and let C be a convex set contained in the effective domain of f. If f attains its supremum relative to C at some point of the relative interior of C, then f is actually constant throughout C.

For our purposes, the effective domain of f is the domain of f since the functions we consider are finite-valued. (See [8] for definitions.)

Corollary A.3 (see [8, Corollary 32.1.1]). Let f be a convex function, and let C be a convex set contained in the effective domain of f. Let W be the set of points (if any), where the supremum of f relative to C is attained. Then W is a union of faces of C.

Corollary A.4 (see [8, Corollary 32.3.2]). Let *f* be a convex function, and let *C* be a non-empty closed bounded convex set contained in the relative interior of the effective domain of *f*. Then the supremum of *f* relative to *C* is finite, and it is attained at some extreme point of *C*.

Theorem A.5 (see [8, Theorem 5.4]). Let f be a twice continuously differentiable real-valued function on a open convex set C in \mathbb{R}^n . Then f is convex on C if and only if its Hessian matrix is positive semidefinite for every $z \in C$.

B. Data Sets for Section 5

The simulated high (H) and low (L) treatment data analysed in Section 6.2 are given in Tables 3 and 4. The CpG methylation measurements analysed in Section 6.3 are given in Table 5.

Acknowledgments

The authors wish to acknowledge funding support provided by the National Research Centre for Growth and Development, New Zealand (G. Wake, A. Pleasants, A. Sheppard), and the Foundation of Research Science and Technology, New Zealand (UOAX0808, A. Sheppard). Further, the authors acknowledge their collaborative link with the GUSTO birth cohort, led by Professors P. D. Gluckman, University of Auckland, and Yap-Seng Chong, National University of Singapore.

References

- A. J. Dobson and A. G. Barnett, An Introduction to Generalized Linear Models, Chapman and Hall/CRC Press, 3rd edition, 2008.
- [2] G. Bassett Jr. and R. Koenker, "Asymptotic theory of least absolute error regression," *Journal of the American Statistical Association*, vol. 73, no. 363, pp. 618–622, 1978.
- [3] D. Birkes and Y. Dodge, Alternative Methods of Regression, John Wiley & Sons, New York, NY, USA, 1993.
- [4] P. Bloomfield and W. L. Steiger, Least Absolute Deviations, Theory Applications and Algorithms, Birkhäuser, Boston, Mass, USA, 1983.
- [5] R. Koenker and G. Bassett, Jr., "Regression Quantiles," *Econometrica*, vol. 46, no. 1, pp. 33–50, 1978.
- [6] S. C. Narula and J. F. Wellington, "The minimum sum of absolute errors regression: a state of the art survey," *International Statistical Review*, vol. 50, no. 3, pp. 317–326, 1982.
- [7] R. M. Norton, "The double exponential distribution: using calculus to find a maximum likelihood estimator," *The American Statistician*, vol. 38, no. 2, pp. 135–136, 1984.
- [8] R. T. Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, NJ, USA, 10th edition, 1970.
- [9] M. J. Lighthill, Introduction to Fourier Analysis and Generalized Functions, Cambridge University Press, New York, NY, USA, 1958.
- [10] I. Stakgold, *Boundary Value Problems of Mathematical Physics, Vol. 1*, The Macmillan, New York, NY, USA, 1967.
- [11] E. L. Lehmann and G. Casella, *Theory of Point Estimation*, Springer, New York, NY, USA, 2nd edition, 1998.
- [12] M. Ehrich, M. R. Nelson, P. Stanssens et al., "Quantitative highthroughput analysis of DNA methylation patterns by basespecific cleavage and mass spectrometry," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 102, no. 44, pp. 15785–15790, 2005.
- [13] S. S. Shapiro and M. B. Wilk, "An analysis of variance test for normality: complete samples," *Biometrika*, vol. 52, pp. 591–611, 1965.
- [14] S. T. Buckland, "Fitting density functions with polynomials," *Journal of the Royal Statistical Society Series C (Applied Statistics)*, vol. 41, no. 1, pp. 63–76, 1992.
- [15] S. T. Buckland, "Maximum likelihood fitting of hermite and simple polynomial densities," *Journal of the Royal Statistical Society Series C (Applied Statistics)*, vol. 41, pp. 241–267, 1992.
- [16] Pleasants, unpublished observations.
- [17] W. Hendricks and R. Koenker, "Hierarchical spline models for conditional quantiles and the demand for electricity," *Journal of the American Statistical Association*, vol. 87, no. 417, pp. 58–68, 1978.
- [18] R. Koenker and G. Bassett, "Tests of linear hypotheses ond l₁ estimation," *Econometrica*, vol. 50, pp. 1157–1583, 1982.
- [19] S. Kotz, T. J. Kozubowski, and K. Podgórski, *The Laplace Distribution and Generalizations*, Birkhäuser, Boston, Mass, USA, 2001.

- [20] E. Purdom and S. P. Holmes, "Error distribution for gene expression data," *Statistical Applications in Genetics and Molecular Biology*, vol. 4, no. 1, article 16, 2005.
- [21] D. Kincaid and W. Cheney, *Mathematics of Scientific Computing*, Brooks/Cole, 3rd edition, 2002.

The Scientific World Journal

Decision Sciences

Journal of Probability and Statistics

Hindawi Submit your manuscripts at

International Journal of Differential Equations

International Journal of Combinatorics

Mathematical Problems in Engineering

Abstract and Applied Analysis

Discrete Dynamics in Nature and Society

Journal of Function Spaces

International Journal of Stochastic Analysis

Journal of Optimization