

Chapter 8

Non-likelihood methods

8.1 Loss functions

Up until now our main focus has been on parameter estimating via the maximum likelihood. However, the negative maximum likelihood is simply one member of loss criterions. Loss functions are usually distances, such as the ℓ_1 and ℓ_2 distance. Typically we estimate a parameter by minimising the loss function, and using as the estimator the parameter which minimises the loss. Usually (but not always) the way to solve the loss function is to differentiate it and equate it to zero. Below we give examples of loss functions whose formal derivative does not exist.

8.1.1 L_1 -loss functions

The Laplacian

Consider the Laplacian (also known as the double exponential), which is defined as

$$f(y; \theta, \rho) = \frac{1}{2\rho} \exp\left(-\frac{|y - \theta|}{\rho}\right) = \begin{cases} \frac{1}{2\rho} \exp\left(\frac{y - \theta}{\rho}\right) & y < \theta \\ \frac{1}{2\rho} \exp\left(\frac{\theta - y}{\rho}\right) & y \geq \theta. \end{cases}$$

We observe $\{Y_i\}$ and our objective is to estimate the location parameter θ , for now the scale parameter ρ is not of interest. The log-likelihood is

$$\mathcal{L}_n(\theta, \rho) = -n \log 2\rho - \rho^{-1} \underbrace{\frac{1}{2} \sum_{i=1}^n |Y_i - \theta|}_{=L_n(\theta)}.$$

$$X_1, X_2, X_3 = 1, 3, 4$$

$$L_n(\theta) = \frac{1}{2} \{ |1-\theta| + |3-\theta| + |4-\theta| \}.$$

$$= \frac{1}{2} \begin{cases} (1-\theta) + (3-\theta) + (4-\theta) = 8-3\theta & \text{if } \theta < 1 \\ (\theta-1) + (3-\theta) + (4-\theta) = 6-\theta & \text{if } 1 \leq \theta < 3 \\ (\theta-1) + (\theta-3) + (4-\theta) = \theta & \text{if } 3 \leq \theta < 4 \\ (\theta-1) + (\theta-3) + (\theta-4) = 3\theta-8 & \text{if } \theta > 4 \end{cases}$$

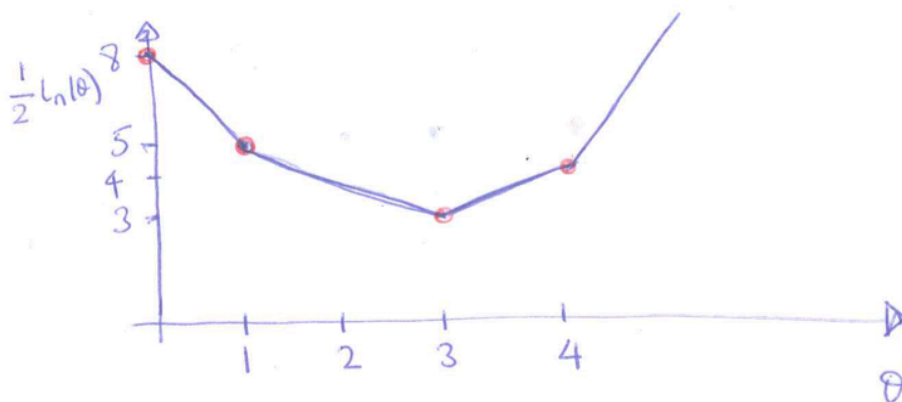


Figure 8.1: Plot of L_1 -norm

Since the θ which maximises the above does not depend on ρ we can simply focus on the component which maximises θ . We see that this is equivalent to minimising the loss function

$$L_n(\theta) = \frac{1}{2} \sum_{i=1}^n |Y_i - \theta| = \sum_{Y_{(i)} > \theta} \frac{1}{2} (Y_{(i)} - \theta) + \sum_{Y_{(i)} \leq \theta} \frac{1}{2} (\theta - Y_{(i)}).$$

If we make a plot of L_n over θ , and consider how L_n behaves at the ordered observations $\{Y_{(i)}\}$, we see that it is piecewise continuous (that is it is a piecewise continuous function, with joints at $Y_{(i)}$). On closer inspection (if n is odd) we see that L_n has its minimum at $\theta = Y_{(n/2)}$, which is the sample median (see Figure 8.1 for an illustration).

In summary, the normal distribution gives rise to the ℓ_2 -loss function and the sample mean. In contrast the Laplacian gives rise to the ℓ_1 -loss function and the sample median.

The asymmetric Laplacian

Consider the generalisation of the Laplacian, usually called the asymmetric Laplacian, which is defined as

$$f(y; \theta, \rho) = \begin{cases} \frac{p}{\rho} \exp\left(\frac{p(y-\theta)}{\rho}\right) & y < \theta \\ \frac{(1-p)}{\rho} \exp\left(-\frac{(1-p)(y-\theta)}{\rho}\right) & y \geq \theta. \end{cases}$$

where $0 < p < 1$. The corresponding negative likelihood to estimate θ is

$$L_n(\theta) = \sum_{Y_{(i)} > \theta} (1-p)(Y_i - \theta) + \sum_{Y_{(i)} \leq \theta} p(\theta - Y_i).$$

Using similar arguments to those in part (i), it can be shown that the minimum of L_n is approximately the p th quantile.

8.2 Estimating Functions

8.2.1 Motivation

Estimating functions are a unification and generalisation of the maximum likelihood methods and the method of moments. It should be noted that it is a close cousin of the *generalised method of moments* and *generalised estimating equation*. We first consider a few examples and will later describe a feature common to all these examples.

Example 8.2.1 (i) Let us suppose that $\{Y_i\}$ are iid random variables with $Y_i \sim \mathcal{N}(\mu, \sigma^2)$.

The log-likelihood is proportional to

$$\mathcal{L}_n(\mu, \sigma^2) = -\frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2.$$

We know that to estimate μ and σ^2 we use the μ and σ^2 which are the solution of

$$\frac{-1}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \mu)^2 = 0 \quad \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = 0. \quad (8.1)$$

(ii) In general suppose $\{Y_i\}$ are iid random variables with $Y_i \sim f(\cdot; \theta)$. The log-likelihood is $\mathcal{L}_n(\theta) = \sum_{i=1}^n \log f(\theta; Y_i)$. If the regularity conditions are satisfied then to estimate θ we use the solution of

$$\frac{\partial \mathcal{L}_n(\theta)}{\partial \theta} = 0. \quad (8.2)$$

(iii) Let us suppose that $\{X_i\}$ are iid random variables with a Weibull distribution $f(x; \theta) = (\frac{x}{\phi})^\alpha \exp(-(x/\phi)^\alpha)$, where $\alpha, \phi > 0$.

We know that $E(X) = \phi\Gamma(1 + \alpha^{-1})$ and $E(X^2) = \phi^2\Gamma(1 + 2\alpha^{-1})$. Therefore $E(X) - \phi\Gamma(1 + \alpha^{-1}) = 0$ and $E(X^2) - \phi^2\Gamma(1 + 2\alpha^{-1}) = 0$. Hence by solving

$$\frac{1}{n} \sum_{i=1}^n X_i - \phi\Gamma(1 + \alpha^{-1}) = 0 \quad \frac{1}{n} \sum_{i=1}^n X_i^2 - \phi^2\Gamma(1 + 2\alpha^{-1}) = 0, \quad (8.3)$$

we obtain estimators of α and Γ . This is essentially a method of moments estimator of the parameters in a Weibull distribution.

(iv) We can generalise the above. It can be shown that $E(X^r) = \phi^r\Gamma(1 + r\alpha^{-1})$. Therefore, for any distinct s and r we can estimate α and Γ using the solution of

$$\frac{1}{n} \sum_{i=1}^n X_i^r - \phi^r\Gamma(1 + r\alpha^{-1}) = 0 \quad \frac{1}{n} \sum_{i=1}^n X_i^s - \phi^s\Gamma(1 + s\alpha^{-1}) = 0. \quad (8.4)$$

(v) Consider the simple linear regression $Y_i = \alpha x_i + \varepsilon_i$, with $E(\varepsilon_i) = 0$ and $\text{var}(\varepsilon_i) = 1$, the least squares estimator of α is the solution of

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \alpha x_i)x_i = 0. \quad (8.5)$$

We observe that all the above estimators can be written as the solution of a homogeneous equations - see equations (8.1), (8.2), (8.3), (8.4) and (8.5). In other words, for each case we can define a random function $G_n(\theta)$, such that the above estimators are the solutions of $G_n(\tilde{\theta}_n) = 0$. In the case that $\{Y_i\}$ are iid then $G_n(\theta) = \sum_{i=1}^n g(Y_i; \theta)$, for some function $g(Y_i; \theta)$. The function $G_n(\tilde{\theta})$ is called an estimating function. All the function G_n , defined above, satisfy the unbiased property which we define below.

Definition 8.2.1 (Estimating function) An estimating function G_n is called unbiased if at the true parameter θ_0 $G_n(\cdot)$ satisfies

$$E[G_n(\theta_0)] = 0.$$

If there are p unknown parameters and p estimating equations, the estimation equation estimator is the θ which solves $G_n(\theta) = 0$.

Hence the estimating function is an alternative way of viewing parameter estimating. Until now, parameter estimators have been defined in terms of the maximum of the likelihood. However, an alternative method for defining an estimator is as the solution of a function. For example, suppose that $\{Y_i\}$ are random variables, whose distribution depends in some way on the parameter θ_0 . We want to estimate θ_0 , and we know that there exists a function such that $G(\theta_0) = 0$. Therefore using the data $\{Y_i\}$ we can define a random function, G_n where $E(G_n(\theta)) = G(\theta)$ and use the parameter $\tilde{\theta}_n$, which satisfies $G_n(\tilde{\theta}) = 0$, as an estimator of θ . We observe that such estimators include most maximum likelihood estimators and method of moment estimators.

Example 8.2.2 *Based on the examples above we see that*

(i) *The estimating function is*

$$G_n(\mu, \sigma) = \left(\begin{array}{c} \frac{-1}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \mu)^2 \\ \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) \end{array} \right).$$

(ii) *The estimating function is* $G_n(\theta) = \frac{\partial \mathcal{L}_n(\theta)}{\partial \theta}$.

(iii) *The estimating function is*

$$G_n(\alpha, \phi) = \left(\begin{array}{c} \frac{1}{n} \sum_{i=1}^n X_i - \phi \Gamma(1 + \alpha^{-1}) \\ \frac{1}{n} \sum_{i=1}^n X_i^2 - \phi^2 \Gamma(1 + 2\alpha^{-1}) \end{array} \right).$$

(iv) *The estimating function is*

$$G_n(\alpha, \phi) = \left(\begin{array}{c} \frac{1}{n} \sum_{i=1}^n X_i^s - \phi^s \Gamma(1 + s\alpha^{-1}) \\ \frac{1}{n} \sum_{i=1}^n X_i^r - \phi^r \Gamma(1 + r\alpha^{-1}) \end{array} \right).$$

(v) *The estimating function is*

$$G_n(a) = \frac{1}{n} \sum_{i=1}^n (Y_i - ax_i)x_i.$$

Observe that regardless of the distribution of the errors (or dependency between $\{Y_i\}$)

$$E\left(\frac{1}{n} \sum_{i=1}^n (Y_i - ax_i)x_i\right) = 0, \quad (8.6)$$

*is true regardless of the distribution of Y_i ($\{\varepsilon_i\}$) and is also true if there $\{Y_i\}$ are dependent random variables (see Rao (1973), *Linear Statistical Inference and its applications*).*

The advantage of this approach is that sometimes the solution of an estimating equation will have a smaller finite sample variance than the MLE. Even though asymptotically (under certain conditions) the MLE will asymptotically attain the Cramer-Rao bound (which is the smallest variance). Moreover, MLE estimators are based on the assumption that the distribution is known (else the estimator is misspecified - see Section 5.1.1), however sometimes an estimating equation can be free of such assumptions.

Example 8.2.3 *In many statistical situations it is relatively straightforward to find a suitable estimating function rather than find the likelihood. Consider the time series $\{X_t\}$ which is “stationary” (moments are invariant to shift i.e $E[X_t X_{t+r}] = E[X_0 X_r]$) which satisfies*

$$X_t = a_1 X_{t-1} + a_2 X_{t-2} + \sigma \varepsilon_t,$$

where $\{\varepsilon_t\}$ are iid zero mean random variables (zero mean ensures that $E[X_t] = 0$). We do not know the distribution of ε_t , but under certain conditions on a_1 and a_2 (causality conditions) ε_t is independent of X_{t-1} and X_{t-2} . Thus by multiplying the above equation by X_{t-1} or X_{t-2} and taking expectations we have

$$\begin{aligned} E(X_t X_{t-1}) &= a_1 E(X_{t-1}^2) + a_2 E(X_{t-1} X_{t-2}) \\ E(X_t X_{t-2}) &= a_1 E(X_{t-1} X_{t-2}) + a_2 E(X_{t-2}^2). \end{aligned}$$

Since the above time series is ‘stationary’ (we have not formally defined this - but basically it means the properties of $\{X_t\}$ do not “evolve” over time), the above reduces to

$$\begin{aligned} c(1) &= a_1 c(0) + a_2 c(1) \\ c(2) &= a_1 c(1) + a_2 c(0), \end{aligned}$$

where $E[X_t X_{t+r}] = c(r)$. Given $\{X_t\}_{t=1}^n$, it can be shown that $\hat{c}_n(r) = n^{-1} \sum_{t=|r|+1}^n X_t X_{t-|r|}$ is an estimator of $c(r)$ and that for small r $E[\hat{c}_n(r)] \approx c(r)$ (and is consistent). Hence replacing the above with its estimators we obtain the estimating equations

$$G_1(a_1, a_2) = \begin{pmatrix} \hat{c}_n(1) - a_1 \hat{c}_n(0) - a_2 \hat{c}_n(1) \\ \hat{c}_n(2) - a_1 \hat{c}_n(1) - a_2 \hat{c}_n(0) \end{pmatrix}$$

8.2.2 The sampling properties

We now show that under certain conditions $\tilde{\theta}_n$ is a consistent estimator of θ .

Theorem 8.2.1 *Suppose that $G_n(\theta)$ is an unbiased estimating function, where $G_n(\tilde{\theta}_n) = 0$ and $E(G_n(\theta_0)) = 0$.*

(i) *If θ is a scalar, for every n $G_n(\theta)$ is a continuous monotonically decreasing function in θ and for all θ $G_n(\theta) \xrightarrow{P} E(G_n(\theta))$ (notice that we do require an equicontinuous assumption), then we have $\tilde{\theta}_n \xrightarrow{P} \theta_0$.*

(ii) *If we can show that $\sup_{\theta} |G_n(\theta) - E(G_n(\theta))| \xrightarrow{P} 0$ and $E(G_n(\theta))$ is uniquely zero at θ_0 then we have $\tilde{\theta}_n \xrightarrow{P} \theta_0$.*

PROOF. The proof of case (i) is relatively straightforward (see also page 318 in Davison (2002)). The idea is to exploit the monotonicity property of $G_n(\cdot)$ to show for every $\varepsilon > 0$ $P(\tilde{\theta}_n < \theta_0 - \varepsilon \text{ or } \tilde{\theta}_n > \theta_0 + \varepsilon) \rightarrow 0$ as $n \rightarrow \infty$. The proof is best understood by making a plot of $G_n(\theta)$ with $\tilde{\theta}_n < \theta_0 - \varepsilon < \theta_0$ (see Figure 8.2). We first note that since $E[G_n(\theta_0)] = 0$, then for any fixed $\varepsilon > 0$

$$G_n(\theta_0 - \varepsilon) \xrightarrow{P} E[G_n(\theta_0 - \varepsilon)] > 0, \quad (8.7)$$

since G_n is monotonically decreasing for all n . Now, since $G_n(\theta)$ is monotonically decreasing we see that $\tilde{\theta}_n < (\theta_0 - \varepsilon)$ implies $G_n(\tilde{\theta}_n) - G_n(\theta_0 - \varepsilon) > 0$ (and visa-versa) hence

$$P(\tilde{\theta}_n - (\theta_0 - \varepsilon) \leq 0) = P(G_n(\tilde{\theta}_n) - G_n(\theta_0 - \varepsilon) > 0).$$

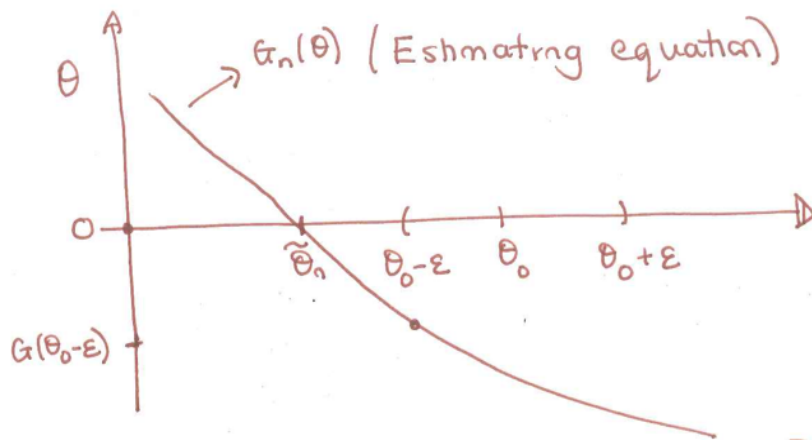
But we have from (8.7) that $E(G_n(\theta_0 - \varepsilon)) \xrightarrow{P} E(G_n(\theta_0 - \varepsilon)) > 0$. Thus $P(G_n(\tilde{\theta}_n) - G_n(\theta_0 - \varepsilon) > 0) \xrightarrow{P} 0$ and

$$P(\tilde{\theta}_n - (\theta_0 - \varepsilon) < 0) \xrightarrow{P} 0 \quad \text{as } n \rightarrow \infty.$$

A similar argument can be used to show that that $P(\tilde{\theta}_n - (\theta_0 + \varepsilon) > 0) \xrightarrow{P} 0$ as $n \rightarrow \infty$. As the above is true for all ε , together they imply that $\tilde{\theta}_n \xrightarrow{P} \theta_0$ as $n \rightarrow \infty$.

The proof of (ii) is more involved, but essentially follows the lines of the proof of Theorem 2.6.1. □

We now show normality, which will give us the variance of the limiting distribution of $\tilde{\theta}_n$.



$$P\{\tilde{\theta}_n < \theta_0 - \epsilon\} = P\{G_n(\tilde{\theta}_n) < G_n(\theta_0 - \epsilon)\}$$
 thanks to the assumption of monotonicity (decreasing) of G_n

Figure 8.2: Plot of $G_n(\cdot)$

Theorem 8.2.2 *Let us suppose that $\{Y_i\}$ are iid random variables, where $E[g(Y_i, \theta)] = 0$. Define the estimating equation $G_n(\theta) = \frac{1}{n} \sum_{i=1}^n g(Y_i; \theta)$ and suppose $G_n(\tilde{\theta}_n) = 0$.*

Suppose that $\tilde{\theta}_n \xrightarrow{\mathcal{P}} \theta_0$ and the first and second derivatives of $g(Y_i, \cdot)$ have a finite expectation (we will assume that θ is a scalar to simplify notation). Then we have

$$\sqrt{n}(\tilde{\theta}_n - \theta_0) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{\text{var}(g(Y_i; \theta_0))}{\left[E\left(\frac{\partial g(Y_i; \theta)}{\partial \theta}\right)\Big|_{\theta_0}\right]^2}\right),$$

as $n \rightarrow \infty$.

Suppose $\{Y_i\}$ are independent but not identically distributed random variables, where for all i $E[g_i(Y_i, \theta)] = 0$. Define the estimating equation $G_n(\theta) = \frac{1}{n} \sum_{i=1}^n g(Y_i; \theta)$ and suppose $G_n(\tilde{\theta}_n) = 0$. Suppose that $\tilde{\theta}_n \xrightarrow{\mathcal{P}} \theta_0$ and the first and second derivatives of $g_i(Y_i, \cdot)$ have a finite, uniformly bounded expectation (we will assume that θ is a scalar to simplify notation). Then we have

$$\sqrt{n}(\tilde{\theta}_n - \theta_0) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \frac{n^{-1} \sum_{i=1}^n \text{var}(g_i(Y_i; \theta_0))}{\left[E\left(n^{-1} \sum_{i=1}^n \frac{\partial g(Y_i; \theta)}{\partial \theta}\right)\Big|_{\theta_0}\right]^2}\right), \quad (8.8)$$

as $n \rightarrow \infty$.

PROOF. We use the standard Taylor expansion to prove the result (which you should be expert in by now). Using a Taylor expansion and that $\tilde{\theta}$

$$\begin{aligned} G_n(\tilde{\theta}_n) &= G_n(\theta_0) + (\tilde{\theta}_n - \theta_0) \frac{\partial G_n(\theta)}{\partial \theta} \Big|_{\tilde{\theta}_n} \\ \Rightarrow (\tilde{\theta}_n - \theta_0) &= \left(E\left(-\frac{\partial g_n(\theta)}{\partial \theta}\right) \Big|_{\theta_0} \right)^{-1} G_n(\theta_0), \end{aligned} \quad (8.9)$$

where the above is due to $\frac{\partial G_n(\theta)}{\partial \theta} \Big|_{\tilde{\theta}_n} \xrightarrow{\mathcal{P}} E\left(\frac{\partial g_n(\theta)}{\partial \theta}\right) \Big|_{\theta_0}$ as $n \rightarrow \infty$. Now, since $G_n(\theta_0) = \frac{1}{n} \sum_i g(Y_i; \theta)$ is the sum of iid random variables we have

$$\sqrt{n}G_n(\theta_0) \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \underbrace{\text{var}(G_n(\theta_0))}_{=\text{var}[g(Y_i; \theta_0)]}\right), \quad (8.10)$$

(since $E(g(Y_i; \theta_0)) = 0$). Therefore (8.9) and (8.10) together give

$$\sqrt{n}(\tilde{\theta}_n - \theta_0) \xrightarrow{\mathcal{P}} \mathcal{N}\left(0, \frac{\text{var}(g(Y_i; \theta_0))}{\left[E\left(\frac{-\partial g(Y_i; \theta)}{\partial \theta}\right)\Big|_{\theta_0}\right]^2}\right),$$

as required. □

In most cases $\frac{\text{var}(g(Y_i; \theta_0))}{\left[E\left(\frac{-\partial g(Y_i; \theta)}{\partial \theta}\right)\Big|_{\theta_0}\right]^2} \geq I(\theta_0)^{-1}$ (where $I(\theta)$ is the Fisher information).

Example 8.2.4 (The Huber estimator) We describe the Huber estimator which is a well known estimator of the mean which is robust to outliers. The estimator can be written as an estimating function.

Let us suppose that $\{Y_i\}$ are iid random variables with mean θ , and density function which is symmetric about the mean θ . So that outliers do not effect the mean, a robust method of estimation is to truncate the outliers and define the function

$$g_{(c)}(Y_i; \theta) = \begin{cases} -c & Y_i < -c + \theta \\ Y_i - c & -c + \theta \leq Y_i \leq c + \theta \\ c & Y_i > c + \theta \end{cases} .$$

The estimating equation is

$$G_{c,n}(\theta) = \sum_{i=1}^n g_{(c)}(Y_i; \theta).$$

And we use as an estimator of θ , the $\tilde{\theta}_n$ which solves $G_{c,n}(\tilde{\theta}_n) = 0$.

(i) In the case that $c = \infty$, then we observe that $G_{\infty,n}(\theta) = \sum_{i=1}^n (Y_i - \theta)$, and the estimator is $\tilde{\theta}_n = \bar{Y}$. Hence without truncation, the estimator of the mean is the sample mean.

(ii) In the case that c is small, then we have truncated many observations.

Definition 8.2.2 (Generalized method of moments) We observe from Example 8.2.1(iii,iv) that there are several estimating equations which can be used to estimate a finite number of parameters (number of estimating equations is more than the number of parameters). In this case, we can use M estimating equations to construct the estimator by minimising the L_2 criterion

$$L_n(\alpha, \phi) = \sum_{r=1}^M \left(\frac{1}{n} \sum_{i=1}^n X_i^r - \phi^r \Gamma(1 + r\alpha^{-1}) \right)^2 .$$

This is an example of the generalized method of moments, which generalizes the ideas of solving estimating equations to obtain parameter estimators.

8.2.3 A worked problem

- (1) Let us suppose we observe the response Y_i and regressor X_i . We assume they satisfy the random coefficient regression model

$$Y_i = (\phi + \xi_i) X_i + \varepsilon_i,$$

where $\{\xi_i\}_i$ and $\{\varepsilon_i\}_i$ are zero mean iid random variables which are independent of each other, with $\sigma_\xi^2 = \text{var}[\xi_i]$ and $\sigma_\varepsilon^2 = \text{var}[\varepsilon_i]$. In this question we will consider how to estimate ϕ , ξ_i and ε_i based on the observations $\{Y_i, X_i\}$.

- (a) What is the Expectation of Y_i given (conditioned on) X_i ?
- (b) What is the variance of Y_i given (conditioned on) X_i ?
- (c) Use your answer in part (a) and least squares to obtain an explicit expression for estimating ϕ .
- (d) Use your answer in part (c) to define the ‘residual’.
- (e) Use your answer in part (b) and (d) and least squares to obtain an explicit expression for estimating σ_ξ^2 and σ_ε^2 .
- (f) By conditioning on the regressors $\{X_i\}_{i=1}^n$, obtain the negative log-likelihood of $\{Y_i\}_{i=1}^n$ under the assumption of Gaussianity of ξ_i and ε_i . Explain the role that (c) and (e) plays in your maximisation algorithm.
- (g) Assume that the regressors, $\{X_i\}$, are iid random variables that are independent of ε_i and ξ_i .

Show that the expectation of the *negative* log-likelihood is minimised at the true parameters ϕ , σ_ξ^2 and σ_ε^2 *even* when ξ_i and ε_i are not Gaussian.

Hint: You may need to use that $-\log x + x$ is minimum at $x = 1$.

Solution:

- (a) *What is the Expectation of Y_i given (conditioned on) X_i ?*

$$E[Y_i|X_i] = \phi X_i.$$

- (b) *What is the variance of Y_i given (conditioned on) X_i ?*

$$\text{var}[Y_i|X_i] = E[(\xi_i X_i + \varepsilon_i)^2|X_i] = \sigma_\xi^2 X_i^2 + \sigma_\varepsilon^2$$

- (c) Use your answer in part (a) and least squares to obtain an explicit expression for estimating ϕ .

$$\hat{\phi} = \arg \min_{\phi} \sum_{i=1}^n (Y_i - \phi X_i)^2 = \frac{\sum_{i=1}^n Y_i X_i}{\sum_{i=1}^n X_i^2}$$

- (d) Use your answer in part (c) to define the ‘residual’.

$$\text{For } 1 \leq i \leq n, \hat{r}_i = Y_i - \hat{\phi} X_i$$

- (e) Use your answer in part (b) and (d) and least squares to obtain an explicit expression for estimating σ_{ξ}^2 and σ_{ε}^2 .

Let

$$r_i = Y_i - \mathbb{E}[Y_i] = Y_i - \phi X_i = \xi_i X_i + \varepsilon_i.$$

From (b) it is clear that $\mathbb{E}[r_i|X_i] = 0$ and $\mathbb{E}[r_i^2|X_i] = \sigma_{\xi}^2 X_i^2 + \sigma_{\varepsilon}^2$, thus we can write

$$r_i^2 = \sigma_{\xi}^2 X_i^2 + \sigma_{\varepsilon}^2 + \epsilon_i$$

where $\epsilon_i = r_i^2 - \mathbb{E}[r_i^2|X_i]$ hence $\mathbb{E}[\epsilon_i] = 0$, resembles a simple linear equation (with hetero errors). Since \hat{r}_i is an estimator of r_i we can use least squares to estimate σ_{ξ}^2 and σ_{ε}^2 , where we replace r_i with \hat{r}_i and minimise

$$\sum_{i=1}^n (\hat{r}_i^2 - \sigma_{\xi}^2 X_i^2 - \sigma_{\varepsilon}^2)^2$$

with respect to σ_{ξ}^2 and σ_{ε}^2 . These gives use explicit estimators.

- (f) By conditioning on the regressors $\{X_i\}_{i=1}^n$, obtain the negative log-likelihood of $\{Y_i\}_{i=1}^n$ under the assumption of Gaussianity of ξ_t and ε_t . Explain the role that (c) and (e) plays in your maximisation algorithm.

The log-likelihood is equal to

$$\sum_{i=1}^n \log f(Y_i|X_i; \theta).$$

We recall from (a) and (b) that $\mathbb{E}[Y_i|X_i] = \phi X_i$ and $\text{var}[Y_i|X_i] = \sigma_{\varepsilon}^2 + \sigma_{\xi}^2 X_i^2$. Therefore $Y_i|X_i \sim \mathcal{N}(\phi X_i, \sigma_{\varepsilon}^2 + \sigma_{\xi}^2 X_i^2)$. Thus the *negative* log likelihood is proportional to

$$L(\theta; \underline{Y}_n) = \sum_{i=1}^n \left(\log[\sigma_{\varepsilon}^2 + \sigma_{\xi}^2 X_i^2] + \frac{(Y_i - \phi X_i)^2}{\sigma_{\varepsilon}^2 + \sigma_{\xi}^2 X_i^2} \right).$$

We choose the parameters which *minimise* $L(\theta; \underline{Y}_n)$. We note that this means we need to take the derivative of $L(\theta; \underline{Y}_n)$ with respect to the three parameters and solve using the Newton Raphson scheme. However, the estimators obtained in (c) and (d) can be used as initial values in the scheme.

- (g) *Let us assume that the regressors are iid random variables. Show that the expectation of the negative log-likelihood is minimised at the true parameters ϕ , σ_ξ^2 and σ_ε^2 even when ξ_t and ε_t are not Gaussian.*

Hint: You may need to use that $-\log x + x$ is minimum at $x = 1$.

Since $\{X_i\}$ are iid random variables, $\{Y_i\}$ are iid random variables the expectation of $\frac{1}{n}L(\theta; \underline{Y}_n)$ is

$$L(\theta) = \mathbb{E} \left(\frac{1}{n} L(\theta; \underline{Y}_n) \right) = \frac{1}{n} \sum_{i=1}^n L_i(\theta)$$

where

$$\begin{aligned} L_i(\theta) &= \mathbb{E} \log[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2] + \mathbb{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] \\ &= \log[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2] + \frac{1}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \mathbb{E} [(Y_i - \phi X_i)^2] \end{aligned}$$

Let θ_0 denote the true parameter in the model. Our aim is to show that

$$L(\theta) - L(\theta_0) = \frac{1}{n} \sum_{i=1}^n (L_i(\theta) - L_i(\theta_0)) \geq 0,$$

where equality to zero arises when $\theta = \theta_0$. Taking differences we have

$$\begin{aligned} &L_i(\theta) - L_i(\theta_0) \\ &= \log \frac{[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2]}{[\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2]} + \mathbb{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - \left[\frac{(Y_i - \phi_0 X_i)^2}{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2} \right] \\ &= -\log \frac{[\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2]}{[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2]} + \mathbb{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - \mathbb{E} \left[\frac{(Y_i - \phi_0 X_i)^2}{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2} \right] \end{aligned}$$

We will show that $L_i(\theta) - L_i(\theta_0)$ is non-negative for all θ and zero when $\theta = \theta_0$. This immediately implies that θ_0 minimises the negative pseudo (pseudo because we do not assume Gaussianity) likelihood.

Our aim is to place the difference in the form $-\log x + x$ plus an additional positive term (it is similar in idea to completing the square), but requires a lot of algebraic manipulation. Let

$$L_i(\theta) - L_i(\theta_0) = A_i(\theta) + B_i(\theta)$$

where

$$\begin{aligned} A_i(\theta) &= - \left(\log \frac{[\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2]}{[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2]} \right) \\ B_i(\theta) &= \text{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - \text{E} \left[\frac{(Y_i - \phi_0 X_i)^2}{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2} \right]. \end{aligned}$$

First consider the difference

$$\begin{aligned} B_i(\theta) &= \text{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - \underbrace{\text{E} \left[\frac{(Y_i - \phi_0 X_i)^2}{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2} \right]}_{=(\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2)^{-1} \text{var}(Y_i) = 1} \\ &= \text{E} \left[\frac{(Y_i - \phi X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - 1. \end{aligned}$$

Now replace ϕ by ϕ_0

$$\begin{aligned} B_i(\theta) &= \text{E} \left[\frac{(Y_i - \phi_0 X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] + \text{E} \left[\frac{(Y_i - \phi X_i)^2 - (Y_i - \phi_0 X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - 1 \\ &= \text{E} \left[\frac{(\varepsilon_t + \xi_t X_i)^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] + \text{E} \left[\frac{2(\phi - \phi_0)(Y_i - \phi_0 X_i)X_i}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] + \\ &\quad \text{E} \left[\frac{(\phi - \phi_0)^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] - 1 \\ &= \frac{\text{E}[(\varepsilon_t + \xi_t X_i)^2]}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} + \text{E} \left[\frac{2(\phi - \phi_0)(Y_i - \phi_0 X_i)X_i}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} \right] + \\ &\quad \frac{(\phi - \phi_0)^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} - 1 \\ &= \frac{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} + \frac{(\phi - \phi_0)^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} - 1. \end{aligned}$$

Therefore, substituting this into $L_i(\theta) - L_i(\theta_0)$ we have

$$\begin{aligned} & L_i(\theta) - L_i(\theta_0) \\ = & -\log \frac{[\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2]}{[\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2]} + \frac{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} + (\phi - \phi_0)^2 \frac{X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2} - 1. \end{aligned}$$

Let

$$x = \frac{\sigma_{0,\varepsilon}^2 + \sigma_{0,\xi}^2 X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2}.$$

Hence

$$L_i(\theta) - L_i(\theta_0) = -\log x + x - 1 + (\phi - \phi_0)^2 \frac{X_i^2}{\sigma_\varepsilon^2 + \sigma_\xi^2 X_i^2}.$$

Since $-\log x + x$ is minimum at $x = 1$ where it is 1, we can see that $L_i(\theta) - L_i(\theta_0)$ is non-negative and zero at $\theta = \theta_0$. As this is true for all i we have that

$$L(\theta) - L(\theta_0) = \frac{1}{n} \sum_{i=1}^n (L_i(\theta) - L_i(\theta_0)) \geq 0,$$

where equality to zero arises when $\theta = \theta_0$.

This example, illustrates the versatility of the models based on the assumption of Gaussianity. Even if the Gaussian assumption does not hold, often we can obtain reasonable (consistent) estimators of the known parameters by treating the errors as if they were Gaussian.

8.3 Optimal estimating functions

As illustrated in Example 8.2.2(iii,iv) there are several different estimators of the same parameters. But which estimator does one use?

Suppose that $\{Y_i\}$ are independent random variables with mean $\{\mu_i(\theta_0)\}$ and variance $\{V_i(\theta_0)\}$, where the parametric form of $\{\mu_i(\cdot)\}$ and $\{V_i(\cdot)\}$ are known, but θ_0 is unknown. One possible estimating equation is

$$G_{1,n}(\theta) = \sum_{i=1}^n [Y_i - \mu_i(\theta)],$$

which is motivated by the observation $E(G_{1,n}(\theta_0)) = 0$. Another estimating equation comes from the least squares criterion

$$\sum_{i=1}^n [Y_i - \mu_i(\theta)]^2,$$

which leads to the estimating equation

$$G_{2,n}(\theta) = \sum_{i=1}^n \frac{\mu_i(\theta)}{\partial \theta} [Y_i - \mu_i(\theta)],$$

again it can be seen that $E(G_{2,n}(\theta_0)) = 0$. Based on the above examples, we see that by simply weighting $[Y_i - \mu_i(\theta)]$ we obtain a valid estimating equation

$$G_n^{(W)}(\theta) = \sum_{i=1}^n w_i(\theta) [Y_i - \mu_i(\theta)].$$

We observe that $E(G_n^{(W)}(\theta_0)) = 0$, thus giving a valid estimating equation. But we need to select the weights $w_i(\theta)$. It seems reasonable to select the weights which minimise the asymptotic “variance”

$$\text{var}(\tilde{\theta}_n) \approx \frac{\sum_{i=1}^n \text{var}(g_i(Y_i; \theta_0))}{[E(\sum_{i=1}^n \frac{\partial g(Y_i; \theta)}{\partial \theta} \Big|_{\theta_0})]^2}. \quad (8.11)$$

Note the above comes from (8.8) (observe the n^{-1} has been removed, since we have not standardized $\tilde{\theta}_n$). Since $\{Y_i\}$ are independent we observe that

$$\begin{aligned} \text{var}(G_n^{(W)}(\theta_0)) &= n^{-1} \sum_{i=1}^n \text{var}(g_i(Y_i; \theta_0)) = \sum_{i=1}^n w_i(\theta_0)^2 V_i(\theta_0) \\ E\left(\frac{\partial G_n^{(W)}(\theta)}{\partial \theta} \Big|_{\theta_0}\right) &= E\left(\sum_{i=1}^n \frac{\partial g(Y_i; \theta)}{\partial \theta} \Big|_{\theta_0}\right) \\ &= E\left(\sum_{i=1}^n w_i'(\theta_0) [Y_i - \mu_i(\theta_0)] - \sum_{i=1}^n w_i(\theta_0) \mu_i'(\theta_0)\right) = -\sum_{i=1}^n w_i(\theta_0) \mu_i'(\theta_0). \end{aligned}$$

Substituting the above into (8.11) gives

$$\text{var}(\tilde{\theta}_n) \approx \frac{\sum_{i=1}^n w_i(\theta_0)^2 V_i(\theta_0)}{(\sum_{i=1}^n w_i(\theta_0) \mu_i'(\theta_0))^2}.$$

Now we want to choose the weights, thus the estimation function, which has the smallest variance. Therefore we look for weights which minimise the above. Since the above is a

ratio, and we observe that a small $w_i(\theta)$ leads to a large denominator but a small numerator. To resolve this, we include a Lagrangian multiplier (this, essentially, minimises the numerator by controlling the magnitude of the denominator). We constrain the numerator to equal one; $(\sum_{i=1}^n w_i(\theta)\mu'_i(\theta))^2 = 1$ and minimise under this constraint

$$\sum_{i=1}^n w_i(\theta)^2 V_i(\theta) + \lambda \left[\sum_{i=1}^n w_i(\theta)\mu'_i(\theta) - 1 \right],$$

with respect to $\{w_i(\theta)\}$ and λ . Partially differentiating the above with respect to $\{w_i(\theta)\}$ and λ and setting to zero gives for all i

$$2w_i(\theta)V_i(\theta) + \mu'_i(\theta) = 0 \text{ subject to } \sum_{i=1}^n w_i(\theta)\mu'_i(\theta) = 1.$$

Thus we choose

$$w_i(\theta) = -\frac{\mu'_i(\theta)}{2V_i(\theta)}$$

but standardize to ensure $\sum_{i=1}^n w_i(\theta)\mu'_i(\theta) = 1$;

$$w_i(\theta) = \left(\sum_{j=1}^n V_j(\theta)^{-1}\mu'_j(\theta) \right)^{-1} \frac{\mu'_i(\theta)}{V_i(\theta)}.$$

Since $\left(\sum_{j=1}^n V_j(\theta)^{-1}\mu'_j(\theta) \right)^{-1}$ is common for all weights $w_i(\theta)$ it can be ignored, thus leading to the optimal estimating function is

$$G_n^{(\mu'V^{-1})}(\theta) = \sum_{i=1}^n \frac{\mu'_i(\theta)}{V_i(\theta)} (Y_i - \mu_i(\theta)). \quad (8.12)$$

The interesting point about the optimal estimating equation, is that even if the *variance* has been misspecified, the estimating equation can still be used to consistently estimate θ (it just will not be optimal).

Example 8.3.1 (i) Consider the case where $\{Y_i\}$ is such that $E[Y_i] = \mu_i(\beta) = \exp(\beta'x_i)$ and $\text{var}(Y_i) = V_i(\beta) = \exp(\beta'x_i)$. Then, $\frac{d\mu(\beta'x_i)}{d\beta} = \exp(\beta'x_i)x_i$. Substituting this yields the optimal estimating equation

$$\sum_{i=1}^n (Y_i - e^{\beta'x_i})x_i = 0.$$

In general if $E[Y_i] = \text{var}[Y_i] = \mu(\beta'x_i)$, the optimal estimating equation is

$$\sum_{i=1}^n \frac{[Y_i - \mu(\beta'x_i)]}{\mu(\beta'x_i)} \mu'(\beta'x_i)x_i = 0,$$

where we use the notation $\mu'(\theta) = \frac{d\mu(\theta)}{d\theta}$. But it is interesting to note that when Y_i comes from a Poisson distribution (where the main feature is that the mean and variance are equal), the above estimating equation corresponds to the score of the likelihood.

(ii) Suppose $\{Y_i\}$ are independent random variables where $E[Y_i] = \mu_i(\beta)$ and $\text{var}[Y_i] = \mu_i(\beta)(1 - \mu_i(\beta))$ (thus $0 < \mu_i(\beta) < 1$). Then the optimal estimating equation corresponds to

$$\sum_{i=1}^n \frac{[Y_i - \mu(\beta'x_i)]}{\mu(\beta'x_i)[1 - \mu(\beta'x_i)]} \mu'(\beta'x_i)x_i = 0,$$

where we use the notation $\mu'(\theta) = \frac{d\mu(\theta)}{d\theta}$. This corresponds to the score function of binary random variables. More of this in the next chapter!

Example 8.3.2 Suppose that $Y_i = \sigma_i Z_i$ where σ_i and Z_i are positive, $\{Z_i\}$ are iid random variables and the regressors x_i influence σ_i through the relation $\sigma_i = \exp(\beta_0 + \beta_1'x_i)$. To estimate β_0 and β_1 we can simply take logarithms of Y_i

$$\log Y_i = \beta_0 + \beta_1'x_i + \log Z_i.$$

Least squares can be used to estimate β_0 and β_1 . However, care needs to be taken since in general $E[\log Z_i] \neq 0$, this will mean the least squares estimator of the intercept β_0 will be biased, as it estimates $\beta_0 + E[\log Z_i]$.

Examples where the above model can arise is $Y_i = \lambda_i Z_i$ where $\{Z_i\}$ are iid with exponential density $f(z) = \exp(-z)$. Observe this means that Y_i is also exponential with density $\lambda_i^{-1} \exp(-y/\lambda_i)$.

Remark 8.3.1 (Weighted least squares) Suppose that $E[Y_i] = \mu_i(\theta)$ and $\text{var}[Y_i] = V_i(\theta)$, motivated by the normal distribution, we can construct the weighted least squared criterion

$$\mathcal{L}_n(\theta) = \sum_{i=1}^n \left[\frac{1}{V_i(\theta)} (Y_i - \mu_i(\theta))^2 + \log V_i(\theta) \right].$$

Taking derivatives, we see that this corresponds to the estimating equation

$$\begin{aligned} G_n(\theta) &= \sum_{i=1}^n \left[-\frac{2}{V_i(\theta)} \{Y_i - \mu_i(\theta)\} \frac{d\mu_i(\theta)}{d\theta} - \frac{1}{V_i(\theta)^2} \{Y_i - \mu_i(\theta)\}^2 \frac{dV_i(\theta)}{d\theta} + \frac{1}{V_i(\theta)} \frac{dV_i(\theta)}{d\theta} \right] \\ &= G_{1,n}(\theta) + G_{2,n}(\theta) \end{aligned}$$

where

$$\begin{aligned} G_{1,n}(\theta) &= -2 \sum_{i=1}^n \frac{1}{V_i(\theta)} \{Y_i - \mu_i(\theta)\} \frac{d\mu_i(\theta)}{d\theta} \\ G_{2,n}(\theta) &= - \sum_{i=1}^n \left[\frac{1}{V_i(\theta)^2} \{Y_i - \mu_i(\theta)\}^2 \frac{dV_i(\theta)}{d\theta} - \frac{1}{V_i(\theta)} \frac{dV_i(\theta)}{d\theta} \right]. \end{aligned}$$

Observe that $E[G_{1,n}(\theta_0)] = 0$ and $E[G_{2,n}(\theta_0)] = 0$, which implies that $E[G_n(\theta_0)] = 0$. This proves that the true parameter θ_0 corresponds to either a local minimum or saddle point of the weighted least squares criterion $\mathcal{L}_n(\theta)$. To show that it is the global minimum one must use an argument similar to that given in Section 8.2.3.

Remark 8.3.2 We conclude this section by mentioning that one generalisation of estimating equations is the generalised method of moments. We observe the random vectors $\{Y_i\}$ and it is known that there exist a function $g(\cdot; \theta)$ such that $E(g(Y_i; \theta_0)) = 0$. To estimate θ_0 , rather than find the solution of $\frac{1}{n} \sum_{i=1}^n g(Y_i; \theta)$, a matrix M_n is defined and the parameter which mimimises

$$\left(\frac{1}{n} \sum_{i=1}^n g(Y_i; \theta) \right)' M_n \left(\frac{1}{n} \sum_{i=1}^n g(Y_i; \theta) \right)$$

is used as an estimator of θ .

