

# 4

## Loss functions and estimation

In the previous chapter we reviewed some results of linear (least squares) models without making the loss function explicit. In this chapter we remedy this and extend the discussion to various other (sometimes referred to as "robust") approaches. That the loss function determines the properties of estimators is common to classical and Bayesian statistics (whether made explicit or not). We'll review a few loss functions and the associated expected loss minimizing estimators. Then we briefly review maximum likelihood estimation (*MLE*) and nonlinear regression.

### 4.1 Loss functions

Let the loss function associated with the estimator  $\hat{\theta}$  for  $\theta$  be  $C(\hat{\theta}, \theta)$  and the posterior distribution function be  $f(\theta | y)$ ,<sup>1</sup> then minimum expected loss is

$$\min_{\hat{\theta}} E [C(\hat{\theta}, \theta)] = \int C(\hat{\theta}, \theta) f(\theta | y) d\theta$$

Briefly, a symmetric quadratic loss function results in an estimator equal to the posterior mean, a linear loss function results in an estimator equal to a quantile of the posterior distribution  $f(\theta | y)$ , and an all or nothing loss function results in an estimator for  $\theta$  equal to the posterior mode.

---

<sup>1</sup>A source of controversy is whether the focus is the posterior distribution  $f(\theta | y)$  or the likelihood function  $f(y | \theta)$ ; see Poirier [1995]. We initially focus on the posterior distribution then review *MLE*.

### 4.1.1 Quadratic loss

The quadratic loss function is

$$C(\hat{\theta}, \theta) = \begin{cases} c_1 (\hat{\theta} - \theta)^2 & \hat{\theta} \leq \theta \\ c_2 (\hat{\theta} - \theta)^2 & \hat{\theta} > \theta \end{cases}$$

First order conditions are

$$\frac{d}{d\hat{\theta}} \left\{ \begin{array}{l} \int_{\hat{\theta}}^{\infty} c_1 (\hat{\theta} - \theta)^2 f(\theta | y) d\theta \\ + \int_{-\infty}^{\hat{\theta}} c_2 (\hat{\theta} - \theta)^2 f(\theta | y) d\theta \end{array} \right\} = 0$$

Rearrangement produces

$$\frac{d}{d\hat{\theta}} \left\{ \begin{array}{l} c_1 (1 - F(\hat{\theta})) \hat{\theta}^2 - 2c_1 \hat{\theta} \int_{\hat{\theta}}^{\infty} \theta f(\theta | y) d\theta \\ + c_1 \int_{\hat{\theta}}^{\infty} \theta^2 f(\theta | y) d\theta + c_2 F(\hat{\theta}) \hat{\theta}^2 \\ - 2c_2 \int_{-\infty}^{\hat{\theta}} \theta f(\theta | y) d\theta + c_2 \int_{-\infty}^{\hat{\theta}} \theta^2 f(\theta | y) d\theta \end{array} \right\} = 0$$

where  $F(\hat{\theta})$  is the cumulative posterior distribution function for  $\theta$  given the data  $y$  evaluated at  $\hat{\theta}$ . Differentiation reveals

$$\left\{ \begin{array}{l} c_1 \left[ \begin{array}{l} 2\hat{\theta} (1 - F(\hat{\theta})) - \hat{\theta}^2 f(\hat{\theta}) \\ - 2 \int_{\hat{\theta}}^{\infty} \theta f(\theta | y) d\theta + 2\hat{\theta}^2 f(\hat{\theta}) - \hat{\theta}^2 f(\hat{\theta}) \end{array} \right] \\ + c_2 \left[ \begin{array}{l} 2\hat{\theta} F(\hat{\theta}) + \hat{\theta}^2 f(\hat{\theta}) \\ - 2 \int_{-\infty}^{\hat{\theta}} \theta f(\theta | y) d\theta - 2\hat{\theta}^2 f(\hat{\theta}) + \hat{\theta}^2 f(\hat{\theta}) \end{array} \right] \end{array} \right\} = 0$$

Simplification yields

$$\begin{aligned} & \hat{\theta} [c_1 (1 - F(\hat{\theta})) + c_2 F(\hat{\theta})] \\ &= c_1 (1 - F(\hat{\theta})) E[\theta | y, \hat{\theta} \leq \theta] + c_2 F(\hat{\theta}) E[\theta | y, \hat{\theta} > \theta] \end{aligned}$$

Or,

$$\hat{\theta} = \frac{c_1 (1 - F(\hat{\theta})) E[\theta | y, \theta \geq \hat{\theta}] + c_2 F(\hat{\theta}) E[\theta | y, \theta < \hat{\theta}]}{c_1 (1 - F(\hat{\theta})) + c_2 F(\hat{\theta})}$$

In other words, the quadratic expected loss minimizing estimator for  $\theta$  is a cost-weighted average of truncated means of the posterior distribution. If  $c_1 = c_2$  (symmetric loss), then  $\hat{\theta} = E[\theta | y]$ , the mean of the posterior distribution.

### 4.1.2 Linear loss

The linear loss function is

$$C(\hat{\theta}, \theta) = \begin{cases} c_1 |\hat{\theta} - \theta| & \hat{\theta} \leq \theta \\ c_2 |\hat{\theta} - \theta| & \hat{\theta} > \theta \end{cases}$$

First order conditions are

$$\frac{d}{d\hat{\theta}} \left\{ \begin{array}{l} -\int_{\hat{\theta}}^{\infty} c_1 (\hat{\theta} - \theta) f(\theta | y) d\theta \\ +\int_{-\infty}^{\hat{\theta}} c_2 (\hat{\theta} - \theta) f(\theta | y) d\theta \end{array} \right\} = 0$$

Rearranging yields

$$\frac{d}{d\hat{\theta}} \left\{ \begin{array}{l} -c_1 \hat{\theta} (1 - F(\hat{\theta})) + c_1 \int_{\hat{\theta}}^{\infty} \theta f(\theta | y) d\theta \\ +c_2 \hat{\theta} F(\hat{\theta}) - c_2 \int_{-\infty}^{\hat{\theta}} \theta f(\theta | y) d\theta \end{array} \right\} = 0$$

Differentiation produces

$$0 = c_1 \left[ -\left(1 - F(\hat{\theta})\right) + \hat{\theta} f(\hat{\theta}) - \hat{\theta} f(\hat{\theta}) \right] \\ + c_2 \left[ F(\hat{\theta}) + \hat{\theta} f(\hat{\theta}) - \hat{\theta} f(\hat{\theta}) \right]$$

Simplification reveals

$$c_1 (1 - F(\hat{\theta})) = c_2 F(\hat{\theta})$$

Or

$$F(\hat{\theta}) = \frac{c_1}{c_1 + c_2}$$

The expected loss minimizing estimator is the quantile that corresponds to the relative cost  $\frac{c_1}{c_1 + c_2}$ . If  $c_1 = c_2$ , then the estimator is the median of the posterior distribution.

### 4.1.3 All or nothing loss

The all or nothing loss function is

$$C(\hat{\theta}, \theta) = \begin{cases} c_1 & \hat{\theta} < \theta \\ 0 & \hat{\theta} = \theta \\ c_2 & \hat{\theta} > \theta \end{cases}$$

If  $c_1 > c_2$ , then we want to choose  $\hat{\theta} > \theta$ , so  $\hat{\theta}$  is the upper limit of support for  $f(\theta | y)$ . If  $c_1 < c_2$ , then we want to choose  $\hat{\theta} < \theta$ , so  $\hat{\theta}$  is the lower limit of

support for  $f(\theta | y)$ . If  $c_1 = c_2$ , then we want to choose  $\hat{\theta}$  to maximize  $f(\theta | y)$ , so  $\hat{\theta}$  is the mode of the posterior distribution.<sup>2</sup>

## 4.2 Nonlinear Regression

Many accounting and business settings call for analysis of data involving limited dependent variables (such as discrete choice models discussed in the next chapter).<sup>3</sup> Nonlinear regression frequently complements our understanding of standard maximum likelihood procedures employed for estimating such models as well as providing a means for addressing alternative functional forms. Here we review some basics of nonlinear least squares including Newton's method of optimization, Gauss-Newton regression (*GNR*), and artificial regressions.

Our discussion revolves around minimizing a smooth, twice continuously differentiable function,  $Q(\beta)$ . It's convenient to think  $Q(\beta)$  equals  $SSR(\beta)$ , the residual sum of squares, but  $-Q(\beta)$  might also refer to maximization of the log-likelihood.

### 4.2.1 Newton's method

A second order Taylor series approximation of  $Q(\beta)$  around some initial values for  $\beta$ , say  $\beta_{(0)}$  yields

$$Q^*(\beta) = Q(\beta_{(0)}) + g_{(0)}^T (\beta - \beta_{(0)}) + \frac{1}{2} (\beta - \beta_{(0)})^T H_{(0)} (\beta - \beta_{(0)})$$

where  $g(\beta)$  is the  $k \times 1$  gradient of  $Q(\beta)$  with typical element  $\frac{\partial Q(\beta)}{\partial \beta_j}$ ,  $H(\beta)$  is the  $k \times k$  Hessian of  $Q(\beta)$  with typical element  $\frac{\partial^2 Q(\beta)}{\partial \beta_j \partial \beta_i}$ , and for notational simplicity,  $g_{(0)} \equiv g(\beta_{(0)})$  and  $H_{(0)} \equiv H(\beta_{(0)})$ . The first order conditions for a minimum of  $Q^*(\beta)$  with respect to  $\beta$  are

$$g_{(0)} + H_{(0)} (\beta - \beta_{(0)}) = 0$$

Solving for  $\beta$  yields a new value

$$\beta_{(1)} = \beta_{(0)} - H_{(0)}^{-1} g_{(0)}$$

This is the core of Newton's method. Successive values  $\beta_{(1)}, \beta_{(2)}, \dots$  lead to an approximation of the global minimum of  $Q(\beta)$  at  $\hat{\beta}$ . If  $Q(\beta)$  is approximately

---

<sup>2</sup>For a discrete probability mass distribution, the optimal estimator may be either the limit of support or the mode depending on the difference in cost. Clearly, large cost differentials are aligned with the limits and small cost differences are aligned with the mode.

<sup>3</sup>This section draws heavily from Davidson and MacKinnon [1993].

quadratic, as applies to sums of squares when sufficiently close to their minima, Newton's method usually converges quickly.<sup>4</sup>

### 4.2.2 Gauss-Newton regression

When minimizing a sum of squares function it is convenient to write the criterion as

$$Q(\beta) = \frac{1}{n} SSR(\beta) = \frac{1}{n} \sum_{t=1}^n (y_t - x_t(\beta))^2$$

Now, explicit expressions for the gradient and Hessian can be found. The gradient for the  $i^{\text{th}}$  element is

$$g_i(\beta) = -\frac{2}{n} \sum_{t=1}^n X_{ti}(\beta) (y_t - x_t(\beta))$$

where  $X_{ti}(\beta)$  is the partial derivative of  $x_t(\beta)$  with respect to  $\beta_i$ . The more compact matrix notation is

$$\mathbf{g}(\beta) = -\frac{2}{n} \mathbf{X}^T(\beta) (\mathbf{y} - \mathbf{x}(\beta))$$

The Hessian  $H(\beta)$  has typical element

$$H_{ij}(\beta) = -\frac{2}{n} \sum_{t=1}^n (y_t - x_t(\beta)) \frac{\partial X_{ti}(\beta)}{\partial \beta_j} - X_{ti}(\beta) X_{tj}(\beta)$$

Evaluated at  $\beta_0$ , this expression is asymptotically equivalent to<sup>5</sup>

$$\frac{2}{n} \sum_{t=1}^n X_{ti}(\beta) X_{tj}(\beta)$$

In matrix notation this is

$$D(\beta) = \frac{2}{n} \mathbf{X}^T(\beta) \mathbf{X}(\beta)$$

and  $D(\beta)$  is positive definite when  $\mathbf{X}(\beta)$  is full rank. Now, writing Newton's method as

$$\beta_{(j+1)} = \beta_{(j)} - D_{(j)}^{-1} g_{(j)}$$

---

<sup>4</sup>If  $Q^*(\beta)$  is strictly convex, as it is if and only if the Hessian is positive definite, then  $\beta_{(1)}$  is the global minimum of  $Q^*(\beta)$ . Please consult other sources, such as Davidson and MacKinnon [2003, ch. 6] and references therein, for additional discussion of Newton's method including search direction, step size, and stopping rules.

<sup>5</sup>Since  $y_t = x_t(\beta_0) + u_t$ , the first term becomes  $-\frac{2}{n} \sum_{t=1}^n \frac{\partial X_{ti}(\beta)}{\partial \beta_j} u_t$ . By the law of large numbers this term tends to 0 as  $n \rightarrow \infty$ .

and substituting the above results we have the classic Gauss-Newton result

$$\begin{aligned}\beta_{(j+1)} &= \beta_{(j)} - \left( \frac{2}{n} \mathbf{X}_{(j)}^T \mathbf{X}_{(j)} \right)^{-1} \left( -\frac{2}{n} \mathbf{X}_{(j)}^T (\mathbf{y} - \mathbf{x}_{(j)}) \right) \\ &= \beta_{(j)} + \left( \mathbf{X}_{(j)}^T \mathbf{X}_{(j)} \right)^{-1} \mathbf{X}_{(j)}^T (\mathbf{y} - \mathbf{x}_{(j)})\end{aligned}$$

Artificial regression

The second term can be readily estimated by an artificial regression. It's called an artificial regression because functions of the variables and model parameters are employed. This artificial regression is referred to as a Gauss-Newton regression (*GNR*)

$$\mathbf{y} - \mathbf{x}(\beta) = \mathbf{X}(\beta) \mathbf{b} + \text{residuals}$$

To be clear, Gaussian projection (*OLS*) produces the following estimate

$$\hat{\mathbf{b}} = \left( \mathbf{X}^T(\beta) \mathbf{X}(\beta) \right)^{-1} \mathbf{X}^T(\beta) (\mathbf{y} - \mathbf{x}(\beta))$$

To appreciate the *GNR*, consider a linear regression where  $X$  is the matrix of regressors. Then  $\mathbf{X}(\beta)$  is simply replaced by  $X$ , the *GNR* is

$$\mathbf{y} - \mathbf{X}\beta_{(0)} = \mathbf{X}\mathbf{b} + \text{residuals}$$

and the artificial parameter estimates are

$$\hat{\mathbf{b}} = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta_{(0)}) = \hat{\beta} - \beta_{(0)}$$

where  $\hat{\beta}$  is the *OLS* estimate. Rearranging we see that the Gauss-Newton estimate replicates *OLS*,  $\beta_{(1)} = \beta_{(0)} + \hat{\mathbf{b}} = \beta_{(0)} + \hat{\beta} - \beta_{(0)} = \hat{\beta}$ , as expected.

Covariance matrices

Return to the *GNR* above and substitute the nonlinear parameter estimates

$$\mathbf{y} - \mathbf{x}(\hat{\beta}) = \mathbf{X}(\hat{\beta}) \mathbf{b} + \text{residuals}$$

The artificial regression estimate is

$$\hat{\mathbf{b}} = \left( \mathbf{X}^T(\hat{\beta}) \mathbf{X}(\hat{\beta}) \right)^{-1} \mathbf{X}^T(\hat{\beta}) (\mathbf{y} - \mathbf{x}(\hat{\beta}))$$

Since the first order or moment conditions require

$$\mathbf{X}^T(\hat{\beta}) (\mathbf{y} - \mathbf{x}(\hat{\beta})) = 0$$

this regression cannot have any explanatory power,  $\hat{\mathbf{b}} = \mathbf{0}$ . Though this may not seem very interesting, it serves two useful functions. First, it provides a check on

the consistency of the nonlinear optimization routine. Second, as it is the *GNR* variance estimate, it provides a quick estimator of the covariance matrix for the parameter estimates

$$\widehat{Var}[\hat{b}] = s^2 \left( \mathbf{X}^T (\hat{\beta}) \mathbf{X} (\hat{\beta}) \right)^{-1}$$

and it is readily available from the artificial regression.

Further, this same *GNR* readily supplies a heteroskedastic-consistent covariance matrix estimator. If  $E[uu^T] = \Omega$ , then a heteroskedastic-consistent covariance matrix estimator is

$$\widehat{Var}[\hat{\mathbf{b}}] = \left( \mathbf{X}^T (\hat{\beta}) \mathbf{X} (\hat{\beta}) \right)^{-1} \mathbf{X}^T (\hat{\beta}) \hat{\Omega} \mathbf{X} (\hat{\beta}) \left( \mathbf{X}^T (\hat{\beta}) \mathbf{X} (\hat{\beta}) \right)^{-1}$$

where  $\hat{\Omega}$  is a diagonal matrix with  $t^{th}$  element equal to the squared residual  $u_t^2$ . Next, we turn to maximum likelihood estimation and exploit some insights gained from nonlinear regression as they relate to typical *MLE* settings.

### 4.3 Maximum likelihood estimation (*MLE*)

Maximum likelihood estimation (*MLE*) applies to a wide variety of problems.<sup>6</sup> Since it is the most common method for estimating discrete choice models and discrete choice models are central to the discussion of accounting choice, we focus the discussion of *MLE* around discrete choice models.

#### 4.3.1 Parameter estimation

The most common method for estimating the parameters of discrete choice models is maximum likelihood. Recall the likelihood is defined as the joint density for the parameters of interest  $\beta$  conditional on the data  $X_t$ . For binary choice models and  $Y_t = 1$  the contribution to the likelihood is  $F(X_t\beta)$ , and for  $Y_t = 0$  the contribution to the likelihood is  $1 - F(X_t\beta)$  where these are combined as binomial draws. Hence,

$$L(\beta|X) = \prod_{t=1}^n F(X_t\beta)^{Y_t} [1 - F(X_t\beta)]^{1-Y_t}$$

The log-likelihood is

$$\ell(\beta|X) \equiv \log L(\beta|X) = \sum_{t=1}^n Y_t \log(F(X_t\beta)) + (1 - Y_t) \log(1 - F(X_t\beta))$$

---

<sup>6</sup>This section draws heavily from Davidson and MacKinnon [1993], chapter 8.

Since this function for binary response models like probit and logit is globally concave, numerical maximization is straightforward. The first order conditions for a maximum are

$$\sum_{t=1}^n \frac{Y_t f(X_t \beta) X_{ti}}{F(X_t \beta)} - \frac{(1-Y_t) f(X_t \beta) X_{ti}}{1-F(X_t \beta)} = 0 \quad i = 1, \dots, k$$

where  $f(\cdot)$  is the density function. Simplifying yields

$$\sum_{t=1}^n \frac{[Y_t - F(X_t \beta)] f(X_t \beta) X_{ti}}{F(X_t \beta) [1 - F(X_t \beta)]} = 0 \quad i = 1, \dots, k$$

For the logit model the first order conditions simplify to

$$\sum_{t=1}^n [Y_t - \Lambda(X_{ti})] X_{ti} = 0 \quad i = 1, \dots, k$$

since the logit density is  $\lambda(X_{ti}) = \Lambda(X_{ti}) [1 - \Lambda(X_{ti})]$  where  $\Lambda(\cdot)$  is the logit (cumulative) distribution function.

Notice the above first order conditions look like the first order conditions for weighted nonlinear least squares with weights given by  $[F(1-F)]^{-1/2}$ . This is sensible because the error term in the nonlinear regression

$$Y_t = F(X_t \beta) + \varepsilon_t$$

has mean zero and variance

$$\begin{aligned} E[\varepsilon_t^2] &= E[\{Y_t - F(X_t \beta)\}^2] \\ &= Pr(Y_t = 1) [1 - F(X_t \beta)]^2 + Pr(Y_t = 0) [0 - F(X_t \beta)]^2 \\ &= F(X_t \beta) [1 - F(X_t \beta)]^2 + [1 - F(X_t \beta)] F(X_t \beta)^2 \\ &= F(X_t \beta) [1 - F(X_t \beta)] \end{aligned}$$

As *ML* is equivalent to weighted nonlinear least squares for binary response models, the asymptotic covariance matrix for  $n^{1/2}(\hat{\beta} - \beta)$  is  $(n^{-1} X^T \Psi X)^{-1}$  where  $\Psi$  is a diagonal matrix with elements  $\frac{f(X_t \beta)^2}{F(X_t \beta) [1 - F(X_t \beta)]}$ . In the logit case,  $\Psi$  simplifies to  $\lambda$  (see Davidson and MacKinnon, p. 517-518).

### 4.3.2 Estimated asymptotic covariance for MLE of $\hat{\theta}$

There are (at least) three common estimators for the variance of  $\hat{\theta}_{MLE}$ :<sup>7</sup>

- (i)  $\left[-H(\hat{\theta})\right]^{-1}$  the negative inverse of Hessian evaluated at  $\hat{\theta}_{MLE}$ ,
- (ii)  $\left[g(\hat{\theta}) g(\hat{\theta})^T\right]^{-1}$  the outer product of gradient (*OPG*) or Berndt, Hall, Hall, and Hausman (*BHHH*) estimator,

<sup>7</sup>This section draws heavily from Davidson and MacKinnon [1993], pp. 260-267.

(iii)  $[\mathfrak{I}(\hat{\theta})]^{-1}$  inverse of information matrix or negative expected value of Hessian, where the following definitions apply:

- *MLE* is defined as the solution to the first order conditions (*FOC*):  $g(Y, \hat{\theta}) = 0$  where gradient or score vector  $g$  is defined by  $g^T(Y, \theta) = D_{\theta}\ell(Y, \theta)$  (since  $D_{\theta}\ell$  is row vector,  $g$  is column vector of partial derivatives of with respect to  $\theta$ ).
- Define  $G(g, \theta)$  as the matrix of contributions to the gradient (*CG* matrix) with typical element  $G_{ti}(g, \theta) \equiv \frac{\partial \ell_t(Y, \theta)}{\partial \theta_i}$ .
- $H(Y, \theta)$  is the Hessian matrix for the log-likelihood with typical element  $H_{ij}(Y, \theta) \equiv \frac{\partial^2 \ell_t(Y, \theta)}{\partial \theta_i \partial \theta_j}$ .
- Define the expected average Hessian for sample of size  $n$  as  $H_n(\theta) \equiv E_{\theta} [n^{-1}H(Y, \theta)]$ .
- The limiting Hessian or asymptotic Hessian (if it exists) is  $H(\theta) \equiv \lim_{n \rightarrow \infty} H_n(\theta)$  (the matrix is negative semidefinite).
- Define the information in observation  $t$  as  $\mathfrak{I}_t(\theta)$  a  $k \times k$  matrix with typical element  $(\mathfrak{I}_t(\theta))_{ij} \equiv E_{\theta} [G_{ti}(\theta) G_{tj}(\theta)]$  (the information matrix is positive semidefinite).
- The average information matrix is  $\mathfrak{I}_n(\theta) \equiv n^{-1} \sum_{t=1}^n \mathfrak{I}_t(\theta) = n^{-1} \mathfrak{I}_n$  and the limiting information matrix or asymptotic information matrix (if it exists) is  $\mathfrak{I}(\theta) \equiv \lim_{n \rightarrow \infty} \mathfrak{I}_n(\theta)$ .

The short explanation for these variance estimators is that *ML* estimators (under suitable regularity conditions) achieve the Cramer-Rao lower bound for consistent

estimators.<sup>8</sup> That is,

$$\text{Asy.Var} [\hat{\theta}] = \left\{ -E \left[ \frac{\partial^2 \ell(Y, \theta)}{\partial \theta \partial \theta^T} \right] \right\}^{-1} = \left\{ E \left[ \left( \frac{\partial \ell(Y, \theta)}{\partial \theta} \right) \left( \frac{\partial \ell(Y, \theta)}{\partial \theta^T} \right) \right] \right\}^{-1}$$

The expected outer product of the gradient (*OPG*) is an estimator of the inverse of the variance matrix for the gradient. Roughly speaking, the inverse of the gradient function yields *MLE* (type 2) parameter estimates and the inverse of expected *OPG* estimates the parameter variance matrix (see Berndt, Hall, Hall, and Hausman [1974]). Also, the expected value of the Hessian equals the negative of the information matrix.<sup>9</sup> In turn, the inverse of the information matrix is an estimator for the estimated parameter variance matrix.

Example: Consider the *MLE* of a standard linear regression model with *DGP*:  $Y = X\beta + \varepsilon$  where  $\varepsilon \sim N(0, \sigma^2 I)$  and  $E[X^T \varepsilon] = 0$ . Of course, the *MLE* for  $\beta$  is  $b = (X^T X)^{-1} X^T Y$  as

$$g(\beta) \equiv \frac{\partial \ell(Y, \beta)}{\partial \beta} = -\frac{1}{\sigma^2} \begin{bmatrix} X_1^T (Y - X\beta) \\ \vdots \\ X_p^T (Y - X\beta) \end{bmatrix}$$

<sup>8</sup>See Theil [1971], pp. 384-385 and Amemiya [1985], pp. 14-17.

$$E \left[ \frac{\partial^2 \ell}{\partial \theta \partial \theta^T} \right] = E \left[ \frac{\partial}{\partial \theta} \left( \frac{1}{L} \frac{\partial L}{\partial \theta^T} \right) \right]$$

by the chain rule

$$\begin{aligned} &= E \left[ -\frac{1}{L^2} \frac{\partial L}{\partial \theta} \frac{\partial L}{\partial \theta^T} + \frac{1}{L} \frac{\partial^2 L}{\partial \theta \partial \theta^T} \right] \\ &= E \left[ -\left( \frac{1}{L} \frac{\partial L}{\partial \theta} \right) \left( \frac{\partial L}{\partial \theta^T} \frac{1}{L} \right) \right] + \int \frac{1}{L} \frac{\partial^2 L}{\partial \theta \partial \theta^T} L dx \\ &= E \left[ -\left( \frac{1}{L} \frac{\partial L}{\partial \theta} \right) \left( \frac{\partial L}{\partial \theta^T} \frac{1}{L} \right) \right] + \int \frac{\partial^2 L}{\partial \theta \partial \theta^T} dx \\ &= -E \left[ \frac{\partial \ell}{\partial \theta} \frac{\partial \ell}{\partial \theta^T} \right] + \int \frac{\partial^2 L}{\partial \theta \partial \theta^T} dx \end{aligned}$$

since the regulatory conditions essentially make the order of integration and differentiation interchangeable the last term can be rewritten

$$\int \frac{\partial^2 L}{\partial \theta \partial \theta^T} dx = \frac{\partial}{\partial \theta} \int \frac{\partial L}{\partial \theta^T} dx = \frac{\partial}{\partial \theta} \frac{\partial}{\partial \theta^T} \int L dx = 0$$

Now we have

$$E \left[ \frac{\partial^2 \ell}{\partial \theta \partial \theta^T} \right] = -E \left[ \frac{\partial \ell}{\partial \theta} \frac{\partial \ell}{\partial \theta^T} \right]$$

<sup>9</sup>This is motivated by the fact that  $\text{plim} \frac{1}{n} \sum_{i=1}^n g(y_i) = E[g(y)]$  for a random sample provided the first two moments of  $g(y)$  are finite (see Greene [1997], ch. 4).

where  $X_j$  refers to column  $j$  of  $X$ . Substituting  $X\beta + \varepsilon$  for  $Y$  produces

$$\left(\frac{\partial \ell(Y, \beta)}{\partial \beta}\right) \left(\frac{\partial \ell(Y, \beta)}{\partial \beta^T}\right) = \left(\frac{1}{\sigma^2}\right)^2 \begin{bmatrix} X_1^T \varepsilon \varepsilon^T X_1 & \cdots & X_1^T \varepsilon \varepsilon^T X_p \\ \vdots & \ddots & \vdots \\ X_p^T \varepsilon \varepsilon^T X_1 & \cdots & X_p^T \varepsilon \varepsilon^T X_p \end{bmatrix}$$

Now,

$$E \left[ \left(\frac{\partial \ell(Y, \beta)}{\partial \beta}\right) \left(\frac{\partial \ell(Y, \beta)}{\partial \beta^T}\right) \right] = \left(\frac{1}{\sigma^2}\right) \begin{bmatrix} X_1^T X_1 & \cdots & X_1^T X_p \\ \vdots & \ddots & \vdots \\ X_p^T X_1 & \cdots & X_p^T X_p \end{bmatrix}$$

Since

$$H(\beta) \equiv \frac{\partial^2 \ell(Y, \beta)}{\partial \beta \partial \beta^T} = - \left(\frac{1}{\sigma^2}\right) \begin{bmatrix} X_1^T X_1 & \cdots & X_1^T X_p \\ \vdots & \ddots & \vdots \\ X_p^T X_1 & \cdots & X_p^T X_p \end{bmatrix}$$

we have

$$E \left[ \left(\frac{\partial \ell(Y, \beta)}{\partial \beta}\right) \left(\frac{\partial \ell(Y, \beta)}{\partial \beta^T}\right) \right] = -E \left[ \frac{\partial^2 \ell(Y, \beta)}{\partial \beta \partial \beta^T} \right]$$

and the demonstration is complete as

$$\begin{aligned} \text{Asy.Var}[b] &= \left\{ E \left[ \left(\frac{\partial \ell(Y, \beta)}{\partial \beta}\right) \left(\frac{\partial \ell(Y, \beta)}{\partial \beta^T}\right) \right] \right\}^{-1} \\ &= - \left\{ E \left[ \frac{\partial^2 \ell(Y, \beta)}{\partial \beta \partial \beta^T} \right] \right\}^{-1} \\ &= \sigma^2 (X^T X)^{-1} \end{aligned}$$

A more complete explanation (utilizing results and notation developed in the appendix) starts with the *MLE* first order condition (*FOC*)  $g(\hat{\theta}) = 0$ . Now, a Taylor series expansion of the likelihood *FOC* around  $\theta$  yields  $0 = g(\hat{\theta}) \approx g(\theta) + H(\bar{\theta})(\hat{\theta} - \theta)$  where  $\bar{\theta}$  is convex combination (perhaps different for each row) of  $\theta$  and  $\hat{\theta}$ . Solve for  $(\hat{\theta} - \theta)$  and rewrite so every term is  $O(1)$

$$n^{1/2}(\hat{\theta} - \theta) = - [n^{-1}H(\bar{\theta})]^{-1} [n^{-1/2}g(\theta)]$$

By *WULLN* (weak uniform law of large numbers), the first term is asymptotically nonstochastic, by *CLT* (the central limit theorem) the second term is asymptotically normal, so  $n^{1/2}(\hat{\theta} - \theta)$  is asymptotically normal. Hence, the asymptotic variance of  $n^{1/2}(\hat{\theta} - \theta)$  is the asymptotic expectation of  $n(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T$ .

Since  $n^{1/2}(\hat{\theta} - \theta) \stackrel{a}{=} -[n^{-1}H(\theta)]^{-1}[n^{-1/2}g(\theta)]$ , the asymptotic variance is  $(-H^{-1}(\theta))(n^{-1}E_{\theta}[g(\theta)g^T(\theta)])(-H^{-1}(\theta))$ . Simplifying yields

$$\text{Asym.Var}[n^{1/2}(\hat{\theta} - \theta)] = H^{-1}(\theta)\mathfrak{S}(\theta)H^{-1}(\theta)$$

This can be simplified since  $H(\theta) = -\mathfrak{S}(\theta)$  by *LLN*. Hence,

$$\text{Asy.Var}[n^{1/2}(\hat{\theta} - \theta)] = -H^{-1}(\theta) = \mathfrak{S}^{-1}(\theta)$$

And the statistic relies on estimation of  $H^{-1}(\theta)$  or  $\mathfrak{S}^{-1}(\theta)$ .

- A common estimator of the empirical Hessian is

$$\hat{H} \equiv H_n(Y, \hat{\theta}) = n^{-1}D_{\theta\theta}^2\ell_t(Y, \hat{\theta})$$

(*LLN* and consistency of  $\hat{\theta}$  guarantee consistency of  $\hat{H}$  for  $H(\theta)$ ).

- The *OPG* or *BHHH* estimator is

$$\mathfrak{S}_{OPG} \equiv n^{-1}\sum_{t=1}^n D_{\theta}^T\ell_t(Y, \hat{\theta})D_{\theta}\ell_t(Y, \hat{\theta}) = n^{-1}G^T(\hat{\theta})G(\hat{\theta})$$

(consistency is guaranteed by *CLT* and *LLN* for the sum).

- The third estimator evaluates the expected values of the second derivatives of the log-likelihood at  $\hat{\theta}$ . Since this form is not always known, this estimator may not be available. However, as this estimator does not depend on the realization of  $Y$  it is less noisy than the other estimators.

We round out this discussion of *MLE* by reviewing a surprising case where *MLE* is not the most efficient estimator. Next, we discuss James-Stein shrinkage estimators.

## 4.4 James-Stein shrinkage estimators

Stein [1955] showed that when estimating  $K$  parameters from independent normal observations with (for simplicity) unit variance, we can uniformly improve on the conventional maximum likelihood estimator in terms of expected squared error loss for  $K > 2$ . James and Stein [1961] determined such a shrinkage estimator can be written as a function of the maximum likelihood estimator  $\hat{\theta}$

$$\theta^* = \hat{\theta} \left( 1 - \frac{a}{\hat{\theta}^T \hat{\theta}} \right)$$

where  $0 \leq a \leq 2(K-2)$ . The expected squared error loss of the James-Stein estimator  $\theta^*$  is

$$\begin{aligned}
\rho(\theta, \theta^*) &= E \left[ (\theta - \theta^*)^T (\theta - \theta^*) \right] \\
&= E \left[ \left\{ (\hat{\theta} - \theta) - \frac{a\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right\}^T \left\{ (\hat{\theta} - \theta) - \frac{a\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right\} \right] \\
&= E \left[ (\hat{\theta} - \theta)^T (\hat{\theta} - \theta) \right] - 2aE \left[ (\hat{\theta} - \theta)^T \frac{\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] \\
&\quad + a^2 E \left[ \frac{\hat{\theta}^T \hat{\theta}}{(\hat{\theta}^T \hat{\theta})^2} \right] \\
&= E \left[ (\hat{\theta} - \theta)^T (\hat{\theta} - \theta) \right] - 2aE \left[ \frac{\hat{\theta}^T \hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] + 2a\theta^T E \left[ \frac{\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] \\
&\quad + a^2 E \left[ \frac{\hat{\theta}^T \hat{\theta}}{(\hat{\theta}^T \hat{\theta})^2} \right] \\
&= E \left[ (\hat{\theta} - \theta)^T (\hat{\theta} - \theta) \right] - 2a + 2a\theta^T E \left[ \frac{\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] + a^2 E \left[ \frac{1}{\hat{\theta}^T \hat{\theta}} \right]
\end{aligned}$$

This can be further simplified by exploiting the following theorems; we conclude this section with Judge and Bock's [1978, p. 322-3] proof following discussion of the James-Stein shrinkage estimator.

**Theorem 4.1**  $E \left[ \frac{\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] = \theta E \left[ \frac{1}{\chi^2_{(K+2, \lambda)}} \right]$  where  $\hat{\theta}^T \hat{\theta} \sim \chi^2_{(K, \lambda)}$  and  $\lambda = \theta^T \theta$  is a noncentrality parameter.<sup>10</sup>

Using

$$\begin{aligned}
E \left[ (\hat{\theta} - \theta)^T (\hat{\theta} - \theta) \right] &= E \left[ \hat{\theta}^T \hat{\theta} \right] - 2\theta^T E \left[ \hat{\theta} \right] + \theta^T \theta \\
&= K + \lambda - 2\lambda + \lambda = K
\end{aligned}$$

for the first term, a convenient substitution for one in the second term, and the above theorem for the third term, we rewrite the squared error loss (from above)

$$\rho(\theta, \theta^*) = E \left[ (\hat{\theta} - \theta)^T (\hat{\theta} - \theta) \right] - 2a + 2a\theta^T E \left[ \frac{\hat{\theta}}{\hat{\theta}^T \hat{\theta}} \right] + a^2 E \left[ \frac{1}{\hat{\theta}^T \hat{\theta}} \right]$$

<sup>10</sup>We adopt the convention the noncentrality parameter is the sum of squared means  $\theta^T \theta$ ; others, including Judge and Bock [1978], employ  $\frac{\theta^T \theta}{2}$ .

as

$$\rho(\theta, \theta^*) = K - 2aE \left[ \frac{\chi_{(K-2, \lambda)}^2}{\chi_{(K-2, \lambda)}^2} \right] + 2a\theta^T \theta E \left[ \frac{1}{\chi_{(K+2, \lambda)}^2} \right] + a^2 E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right]$$

**Theorem 4.2** For any real-valued function  $f$  and positive definite matrix  $A$ ,

$$\begin{aligned} E \left[ f \left( \widehat{\theta}^T \widehat{\theta} \right) \left( \widehat{\theta}^T A \widehat{\theta} \right) \right] &= E \left[ f \left( \chi_{(K+2, \lambda)}^2 \right) \text{tr}(A) \right] \\ &\quad + E \left[ f \left( \chi_{(K+4, \lambda)}^2 \right) \right] \left( \theta^T A \theta \right) \end{aligned}$$

where  $\text{tr}(A)$  is trace of  $A$ .

Letting  $f \left( \widehat{\theta}^T \widehat{\theta} \right) = \frac{1}{\chi_{(K-2, \lambda)}^2}$  and  $A = I$  with rank  $K - 2$ ,

$$-2aE \left[ \frac{\chi_{(K-2, \lambda)}^2}{\chi_{(K-2, \lambda)}^2} \right] = -2aE \left[ \frac{K-2}{\chi_{(K, \lambda)}^2} \right] - 2a\theta^T \theta E \left[ \frac{1}{\chi_{(K+2, \lambda)}^2} \right]$$

and

$$\rho(\theta, \theta^*) = K - a[2(K-2) - a] E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right]$$

Hence,  $\rho(\theta, \theta^*) = K - a[2(K-2) - a] E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right] \leq \rho(\theta, \widehat{\theta}) = K$  for all  $\theta$  if  $0 < a < 2(K-2)$  with strict inequality for some  $\theta^T \theta$ .

Now, we can find the optimal James-Stein shrinkage estimator. Solving the first order condition

$$\begin{aligned} \frac{\partial \rho(\theta, \theta^*)}{\partial a} &= 0 \\ (-2(K-2) - a + 2a) E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right] &= 0 \end{aligned}$$

leads to  $a^* = K - 2$ ; hence,  $\theta^* = \widehat{\theta} \left( 1 - \frac{K-2}{\widehat{\theta}^T \widehat{\theta}} \right)$ . As  $E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right] = \frac{1}{K-2}$ , the James-Stein estimator has minimum expected squared error loss when  $\theta = 0$ ,

$$\begin{aligned} \rho(\theta, \theta^*) &= K - (K-2)^2 E \left[ \frac{1}{\chi_{(K, \lambda)}^2} \right] \\ &= K - (K-2) = 2 \end{aligned}$$

and its  $MSE$  approaches that for the  $MLE$  as  $\lambda = \theta^T \theta$  approaches infinity. Next, we sketch proofs of the theorems.

Stein [1966] identified a key idea used in the proofs. Suppose a  $J \times 1$  random vector  $w$  is distributed as  $N(\theta, I)$ , then its quadratic form  $w^T w$  has a noncentral

$\chi_{(J,\lambda)}^2$  where  $\lambda = \theta^T \theta$ . This quadratic form can be regarded as having a central  $\chi_{(J+2H)}^2$  where  $H$  is a Poisson random variable with parameter  $\frac{\lambda}{2}$ . Hence,

$$\begin{aligned} E \left[ f \left( \chi_{(J,\lambda)}^2 \right) \right] &= E_H \left[ H E \left[ f \left( \chi_{(J+2H)}^2 \right) \right] \right] \\ &= \sum_{t=0}^{\infty} \left( \frac{\lambda}{2} \right)^t \frac{\exp \left[ -\frac{\lambda}{2} \right]}{t!} E \left[ f \left( \chi_{(J+2t)}^2 \right) \right] \end{aligned}$$

Now, we proceed with proofs to the above theorems.

**Theorem 4.3**  $E \left[ \frac{\hat{\theta}}{\theta^T \hat{\theta}} \right] = \theta E \left[ \frac{1}{\chi_{(K+2,\lambda)}^2} \right]$ .

**Proof.** Write

$$\begin{aligned} E \left[ f \left( w^2 \right) w \right] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f \left( w^2 \right) w \exp \left[ -\frac{(w-\theta)^2}{2} \right] dw \\ &= \exp \left[ -\frac{\theta^2}{2} \right] \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f \left( w^2 \right) w \exp \left[ -\frac{w^2}{2} + \theta w \right] dw \end{aligned}$$

Rewrite as

$$\exp \left[ -\frac{\theta^2}{2} \right] \frac{1}{\sqrt{2\pi}} \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} f \left( w^2 \right) \exp \left[ -\frac{w^2}{2} + \theta w \right] dw$$

complete the square and write

$$\exp \left[ -\frac{\theta^2}{2} \right] \frac{\partial}{\partial \theta} \left\{ E \left[ f \left( w^2 \right) \right] \exp \left[ \frac{\theta^2}{2} \right] \right\}$$

Since  $w \sim N(\theta, 1)$ ,  $w^2 \sim \chi_{(1,\theta^2)}^2$ . Now, apply Stein's observation

$$\begin{aligned} &\exp \left[ -\frac{\theta^2}{2} \right] \frac{\partial}{\partial \theta} \left\{ E \left[ f \left( w^2 \right) \right] \exp \left[ \frac{\theta^2}{2} \right] \right\} \\ &= \exp \left[ -\frac{\theta^2}{2} \right] \frac{\partial}{\partial \theta} \left\{ \exp \left[ \frac{\theta^2}{2} \right] \sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{\exp \left[ -\frac{\theta^2}{2} \right]}{j!} E \left[ f \left( \chi_{(1+2j)}^2 \right) \right] \right\} \\ &= \exp \left[ -\frac{\theta^2}{2} \right] \frac{\partial}{\partial \theta} \left\{ \sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{1}{j!} E \left[ f \left( \chi_{(1+2j)}^2 \right) \right] \right\} \end{aligned}$$

Taking the partial derivative yields

$$\exp \left[ -\frac{\theta^2}{2} \right] \theta \left\{ \sum_{j=1}^{\infty} \left( \frac{\theta^2}{2} \right)^{j-1} \frac{1}{(j-1)!} E \left[ f \left( \chi_{(3+2(j-1))}^2 \right) \right] \right\}$$

or

$$E [f (w^2) w] = \theta E \left[ f \left( \chi_{(3, \theta^2)}^2 \right) \right]$$

For the multivariate case at hand, this implies

$$E \left[ \frac{\widehat{\theta}}{\widehat{\theta}^T \widehat{\theta}} \right] = \theta E \left[ \frac{1}{\chi_{(K+2, \lambda)}^2} \right]$$

■

**Lemma 4.1**  $E [f (w^2) w^2] = E [f (\chi_{(3, \theta^2)}^2)] + \theta^2 E [f (\chi_{(5, \theta^2)}^2)]$ .

**Proof.** Let  $z \sim \chi_{(1, \theta^2)}^2$ .

$$\begin{aligned} E [f (w^2) w^2] &= E [f (z) z] \\ &= \sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{1}{j!} \exp \left[ -\frac{\theta^2}{2} \right] E [f (\chi_{(1+2j)}^2)] \end{aligned}$$

Since

$$E [f (\chi_{(n)}^2) \chi_{(n)}^2] = \int_0^{\infty} \frac{f(s) s \exp \left[ -\frac{s}{2} \right] s^{\frac{n}{2}-1}}{\Gamma \left( \frac{n}{2} \right) 2^{\frac{n}{2}}} ds$$

combining terms involving powers of  $s$  and rewriting

$$\Gamma(t) = \int_0^{\infty} x^{t-1} \exp[-x] dx$$

for  $t > 0$ , as  $\Gamma(t+1) = t\Gamma(t)$ , leads to

$$\begin{aligned} &n \int_0^{\infty} \frac{f(s) \exp \left[ -\frac{s}{2} \right] s^{\frac{n+2}{2}-1}}{\Gamma \left( \frac{n+2}{2} \right) 2^{\frac{n+2}{2}}} ds \\ &= n E [f (\chi_{(n+2)}^2)] \end{aligned}$$

Now,

$$\begin{aligned} &\sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{1}{j!} \exp \left[ -\frac{\theta^2}{2} \right] E [f (\chi_{(1+2j)}^2)] \\ &= \sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{1}{j!} \exp \left[ -\frac{\theta^2}{2} \right] (1+2j) E [f (\chi_{(3+2j)}^2)] \end{aligned}$$

Rewrite as

$$\begin{aligned} &\sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^j \frac{1}{j!} \exp \left[ -\frac{\theta^2}{2} \right] E [f (\chi_{(3+2j)}^2)] \\ &+ 2 \frac{\theta^2}{2} \sum_{j=0}^{\infty} \left( \frac{\theta^2}{2} \right)^{j-1} \frac{1}{(j-1)!} \exp \left[ -\frac{\theta^2}{2} \right] E [f (\chi_{(5+2(j-1))}^2)] \end{aligned}$$

Again, apply Stein's observation to produce

$$E [f(w^2) w^2] = E [f(\chi_{(3, \theta^2)}^2)] + \theta^2 E [f(\chi_{(5, \theta^2)}^2)]$$

■

**Theorem 4.4**

$$\begin{aligned} E [f(\widehat{\theta}^T \widehat{\theta}) (\widehat{\theta}^T A \widehat{\theta})] &= E [f(\chi_{(K+2, \lambda)}^2) \text{tr}(A)] \\ &\quad + E [f(\chi_{(K+4, \lambda)}^2)] (\theta^T A \theta) \end{aligned}$$

**Proof.** Let  $P$  be an orthogonal matrix such that  $PAP^T = D$ , a diagonal matrix with eigenvalues of  $A$ ,  $d_j > 0$ , along the diagonal. Define vector  $\omega = Pw \sim N(P\theta, I)$ . Since

$$\omega^T \omega = w^T P^T P w = w^T w$$

and

$$\omega^T D \omega = \omega^T P^T A P \omega = w^T A w$$

$$E [f(\omega^T \omega) \omega^T D \omega] = \sum_{i=1}^J d_i E \left[ E \left[ f \left( \omega_i^2 + \sum_{j \neq i} \omega_j^2 \right) \omega_i^2 \mid \omega_j, i \neq j \right] \right]$$

Using the lemma, this can be expressed as

$$\sum_{i=1}^J d_i \left\{ \begin{aligned} &E \left[ f \left( \chi_{(3, (p_i^T \theta)^2)}^2 + \sum_{j \neq i} \omega_j^2 \right) \right] \\ &+ (p_i^T \theta)^2 E \left[ f \left( \chi_{(5, (p_i^T \theta)^2)}^2 + \sum_{j \neq i} \omega_j^2 \right) \right] \end{aligned} \right\}$$

where  $p_i^T$  is the  $i$ th row of  $P$ . Since  $\sum_{i=1}^J d_i (p_i^T \theta)^2 = \theta^T A \theta$  and  $\sum_{i=1}^J d_i = \text{tr}(A)$ ,

$$\begin{aligned} E [f(\widehat{\theta}^T \widehat{\theta}) (\widehat{\theta}^T A \widehat{\theta})] &= E [f(\chi_{(K+2, \lambda)}^2) \text{tr}(A)] \\ &\quad + E [f(\chi_{(K+4, \lambda)}^2)] (\theta^T A \theta) \end{aligned}$$

■

## 4.5 Summary

This chapter has briefly reviewed loss functions, nonlinear regression, maximum likelihood estimation, and some alternative estimation methods (including James-Stein shrinkage estimators). It is instructive to revisit nonlinear regression (especially, *GNR*) in the next chapter when we address specification and estimation of discrete choice models.

## 4.6 Additional reading

Poirier [1995] provides a nice discussion of loss functions. Conditional linear loss functions lead to quantile regression (see Koenker and Bassett [1978], Koenker [2005], and Koenker [2009] for an **R** computational package). Shugan and Mitra [2008] offer an intriguing discussion of when and why non-averaging statistics (e.g., maximum and variance) explain more variance than averaging metrics. Maximum likelihood estimation is discussed by a broad range of authors including Davidson and MacKinnon [1993], Greene [1997], Amemiya [1985], Rao [1973], and Theil [1971]. Stigler [2007] provides a fascinating account of the history of maximum likelihood estimation including the pioneering contributions of Gauss and Fisher as well as their detractors. The nonlinear regression section draws heavily from a favorite reference, Davidson and MacKinnon [2003]. Their chapter 6 and references therein provide a wealth of ideas related to estimation and specification of nonlinear models.