

### 8.2.2 Maximum Likelihood Inference

It turns out that the parametric bootstrap agrees with least squares in the previous example because the model (8.5) has additive Gaussian errors. In general, the parametric bootstrap agrees not with least squares but with maximum likelihood, which we now review.

We begin by specifying a probability density or probability mass function for our observations

$$z_i \sim g_\theta(z). \quad (8.8)$$

In this expression  $\theta$  represents one or more unknown parameters that govern the distribution of  $Z$ . This is called a *parametric model* for  $Z$ . As an example, if  $Z$  has a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , then

$$\theta = (\mu, \sigma^2), \quad (8.9)$$

and

$$g_\theta(z) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(z-\mu)^2/\sigma^2}. \quad (8.10)$$

Maximum likelihood is based on the *likelihood function*, given by

$$L(\theta; \mathbf{Z}) = \prod_{i=1}^N g_\theta(z_i), \quad (8.11)$$

the probability of the observed data under the model  $g_\theta$ . The likelihood is defined only up to a positive multiplier, which we have taken to be one. We think of  $L(\theta; \mathbf{Z})$  as a function of  $\theta$ , with our data  $\mathbf{Z}$  fixed.

Denote the logarithm of  $L(\theta; \mathbf{Z})$  by

$$\begin{aligned} \ell(\theta; \mathbf{Z}) &= \sum_{i=1}^N \ell(\theta; z_i) \\ &= \sum_{i=1}^N \log g_\theta(z_i), \end{aligned} \quad (8.12)$$

which we will sometimes abbreviate as  $\ell(\theta)$ . This expression is called the log-likelihood, and each value  $\ell(\theta; z_i) = \log g_\theta(z_i)$  is called a log-likelihood component. The method of maximum likelihood chooses the value  $\theta = \hat{\theta}$  to maximize  $\ell(\theta; \mathbf{Z})$ .

The likelihood function can be used to assess the precision of  $\hat{\theta}$ . We need a few more definitions. The *score function* is defined by

$$\dot{\ell}(\theta; \mathbf{Z}) = \sum_{i=1}^N \dot{\ell}(\theta; z_i), \quad (8.13)$$

where  $\dot{\ell}(\theta; z_i) = \partial \ell(\theta; z_i) / \partial \theta$ . Assuming that the likelihood takes its maximum in the interior of the parameter space,  $\dot{\ell}(\hat{\theta}; \mathbf{Z}) = 0$ . The *information matrix* is

$$\mathbf{I}(\theta) = - \sum_{i=1}^N \frac{\partial^2 \ell(\theta; z_i)}{\partial \theta \partial \theta^T}. \quad (8.14)$$

When  $\mathbf{I}(\theta)$  is evaluated at  $\theta = \hat{\theta}$ , it is often called the *observed information*. The *Fisher information* (or expected information) is

$$\mathbf{i}(\theta) = E_{\theta}[\mathbf{I}(\theta)]. \quad (8.15)$$

Finally, let  $\theta_0$  denote the true value of  $\theta$ .

A standard result says that the sampling distribution of the maximum likelihood estimator has a limiting normal distribution

$$\hat{\theta} \rightarrow N(\theta_0, \mathbf{i}(\theta_0)^{-1}), \quad (8.16)$$

as  $N \rightarrow \infty$ . Here we are independently sampling from  $g_{\theta_0}(z)$ . This suggests that the sampling distribution of  $\hat{\theta}$  may be approximated by

$$N(\hat{\theta}, \mathbf{i}(\hat{\theta})^{-1}) \text{ or } N(\hat{\theta}, \mathbf{I}(\hat{\theta})^{-1}), \quad (8.17)$$

where  $\hat{\theta}$  represents the maximum likelihood estimate from the observed data.

The corresponding estimates for the standard errors of  $\hat{\theta}_j$  are obtained from

$$\sqrt{\mathbf{i}(\hat{\theta})_{jj}^{-1}} \quad \text{and} \quad \sqrt{\mathbf{I}(\hat{\theta})_{jj}^{-1}}. \quad (8.18)$$

Confidence points for  $\theta_j$  can be constructed from either approximation in (8.17). Such a confidence point has the form

$$\hat{\theta}_j - z^{(1-\alpha)} \cdot \sqrt{\mathbf{i}(\hat{\theta})_{jj}^{-1}} \quad \text{or} \quad \hat{\theta}_j - z^{(1-\alpha)} \cdot \sqrt{\mathbf{I}(\hat{\theta})_{jj}^{-1}},$$

respectively, where  $z^{(1-\alpha)}$  is the  $1 - \alpha$  percentile of the standard normal distribution. More accurate confidence intervals can be derived from the likelihood function, by using the chi-squared approximation

$$2[\ell(\hat{\theta}) - \ell(\theta_0)] \sim \chi_p^2, \quad (8.19)$$

where  $p$  is the number of components in  $\theta$ . The resulting  $1 - 2\alpha$  confidence interval is the set of all  $\theta_0$  such that  $2[\ell(\hat{\theta}) - \ell(\theta_0)] \leq \chi_p^{2(1-2\alpha)}$ , where  $\chi_p^{2(1-2\alpha)}$  is the  $1 - 2\alpha$  percentile of the chi-squared distribution with  $p$  degrees of freedom.

Let's return to our smoothing example to see what maximum likelihood yields. The parameters are  $\theta = (\beta, \sigma^2)$ . The log-likelihood is

$$\ell(\theta) = -\frac{N}{2} \log \sigma^2 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - h(x_i)^T \beta)^2. \quad (8.20)$$

The maximum likelihood estimate is obtained by setting  $\partial\ell/\partial\beta = 0$  and  $\partial\ell/\partial\sigma^2 = 0$ , giving

$$\begin{aligned} \hat{\beta} &= (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}, \\ \hat{\sigma}^2 &= \frac{1}{N} \sum (y_i - \hat{\mu}(x_i))^2, \end{aligned} \quad (8.21)$$

which are the same as the usual estimates given in (8.2) and below (8.3).

The information matrix for  $\theta = (\beta, \sigma^2)$  is block-diagonal, and the block corresponding to  $\beta$  is

$$\mathbf{I}(\beta) = (\mathbf{H}^T \mathbf{H}) / \sigma^2, \quad (8.22)$$

so that the estimated variance  $(\mathbf{H}^T \mathbf{H})^{-1} \hat{\sigma}^2$  agrees with the least squares estimate (8.3).

### 8.2.3 Bootstrap versus Maximum Likelihood

In essence the bootstrap is a computer implementation of nonparametric or parametric maximum likelihood. The advantage of the bootstrap over the maximum likelihood formula is that it allows us to compute maximum likelihood estimates of standard errors and other quantities in settings where no formulas are available.

In our example, suppose that we adaptively choose by cross-validation the number and position of the knots that define the  $B$ -splines, rather than fix them in advance. Denote by  $\lambda$  the collection of knots and their positions. Then the standard errors and confidence bands should account for the adaptive choice of  $\lambda$ , but there is no way to do this analytically. With the bootstrap, we compute the  $B$ -spline smooth with an adaptive choice of knots for each bootstrap sample. The percentiles of the resulting curves capture the variability from both the noise in the targets as well as that from  $\hat{\lambda}$ . In this particular example the confidence bands (not shown) don't look much different than the fixed  $\lambda$  bands. But in other problems, where more adaptation is used, this can be an important effect to capture.

## 8.3 Bayesian Methods

In the Bayesian approach to inference, we specify a sampling model  $\Pr(\mathbf{Z}|\theta)$  (density or probability mass function) for our data given the parameters,