

where  $a_0$  and  $a_j$  are model coefficients,  $NV$  is the number of explanatory variables, and  $X_j$  is an explanatory variable<sup>2</sup>. Equation (4) is then exponentiated to yield an estimate of instantaneous load:

$$\hat{L}_{RC} = \exp \left( a_0 + \sum_{j=1}^{NV} a_j X_j \right) \quad (5)$$

where  $\hat{L}_{RC}$  is a “rating curve” estimate of instantaneous load. Development of load estimates using equations 4 and 5 is thus a 3-step process:

(1) **Model Formulation.** The form of the linear model (the right-hand side of equation 4) is determined based on the user’s knowledge of the hydrologic and biogeochemical system. Each explanatory variable ( $X_j$ ) is a function of a data variable (streamflow or time, for example) that is thought to influence instantaneous load. The number and form of explanatory variables is highly dependent on the system under study and the constituent of interest. A simple model with a single explanatory variable (log streamflow) is often sufficient for prediction of suspended-sediment load (Crawford, 1991), whereas a model with six explanatory variables based on various functions of streamflow and time is often applicable to nutrients (Cohn and others, 1992a). Additional guidance on model formulation is provided elsewhere (Judge and others, 1988; Draper and Smith, 1998; Helsel and Hirsch, 2002).

(2) **Model Calibration.** Given the form of the regression model, a time series of constituent load and the explanatory variables is used to develop the model coefficients ( $a_0$  and  $a_j$ , equ. 4) by using ordinary least squares (OLS) regression. The regression equation then is used to calculate estimates of log load [ $\ln(\hat{L})$ ] for each observation in the time series (the calibration data set). Residual error for each observation is equal to the difference between observed and estimated values of log load [ $\ln(L) - \ln(\hat{L})$ ].

(3) **Load Estimation.** Estimates of the instantaneous load are obtained using the retransformed version of the regression model (equ. 5) and a time series of explanatory variables (the estimation data set). Individual estimates of instantaneous load then are used to determine the total (equ. 2) or mean (equ. 3) load.

As outlined above, estimation of constituent loads using the regression approach is theoretically straightforward. Several statistical complications arise, however, when dealing with real-world data. Load calculations within LOADEST are therefore more complex than the calculations described above. Three of these complicating factors (retransformation bias, data censoring, and nonnormality) are described below, where the three load estimation methods used within LOADEST are detailed. Additional issues that are germane to all three methods are described in Sections 2.3 and 2.4.

## 2.2 Load Estimation Methods used within LOADEST

The load estimation process is complicated by retransformation bias, data censoring, and nonnormality. As noted by Ferguson (1986), rating curve estimates (equ. 5) of instantaneous load are biased; estimates may underestimate the true load by as much as 50 percent. This retransformation bias is addressed by introducing bias correction factors for the calculation of instantaneous load. Data censoring occurs when one or more observations used in the calibration step have constituent concentrations that are less than the laboratory detection limit (Gilbert, 1987). Although substitution (setting  $C$  equal to one-half the detection limit, for example) appears to be a simple remedy for the replacement of less-than values, none of the substitution methods commonly used yield adequate results (Helsel and Cohn, 1988). A more rigorous treatment of censored data is therefore required. A final complication is the assumption of OLS regression that the model residuals are normally distributed. Alternate methods for estimating model coefficients are applicable when model residuals do not follow a normal

<sup>2</sup>The  $i$  subscript is omitted from  $L$  in equation 4 and all subsequent equations.

distribution. Because of these complications, LOADEST provides three methods for load estimation; each method is described below.

### 2.2.1 Maximum Likelihood Estimation (MLE)

As an alternative to OLS regression, model coefficients ( $a_0$  and  $a_j$ , equ. 4) may be calculated using the method of maximum likelihood (MLE). When the calibration data set includes censored data, implementation of MLE also is known as tobit regression (Helsel and Hirsch, 2002). As with OLS, tobit regression assumes that model residuals are normally distributed with constant variance.

Given the model coefficients provided by regression, estimates of instantaneous load may be obtained by retransforming equation 4. When the calibration data set is uncensored, the bias correction factor of Bradu and Mundlak (1970) provides a minimum variance unbiased estimate (MVUE) of instantaneous load (Cohn and others, 1989):

$$\hat{L}_{MVUE} = \exp \left( a_0 + \sum_{j=1}^{NV} a_j X_j \right) g_m(m, s^2, V) \quad (6)$$

where  $\hat{L}_{MVUE}$  is the MLE estimate of instantaneous load,  $m$  is the number of degrees of freedom,  $s^2$  is the residual variance, and  $V$  is a function of the explanatory variables (Cohn and others, 1989). The model coefficients in equation 6 ( $a_0$  and  $a_j$ ) are estimated by maximum likelihood; the bias correction factor [ $g_m(m, s^2, V)$ ] is an approximation of the infinite series given in Finney (1941). Within LOADEST,  $g_m(m, s^2, V)$  is replaced by a similar function, phi (Likes, 1980).

Under the MLE method, estimates of instantaneous load are developed for all of the observations in the estimation data set using equation 6. Mean load estimates for various time periods then are calculated using equation 3 (where  $\hat{L}_\tau = \hat{L}_{MVUE}$ ). Standard errors reflecting the uncertainty in each estimate of mean load are calculated by using the method described by Likes (1980) and Gilroy and others (1990) (for specifics, see equations 9–25 in Gilroy and others, 1990).

### 2.2.2 Adjusted Maximum Likelihood Estimation (AMLE)

For the case of censored data, model coefficients estimated by tobit regression (MLE, Section 2.2.1) exhibit first-order bias. In addition, the Bradu-Mundlak bias correction factor ( $g_m$ , equation 6) results in biased estimates of instantaneous load. By using adjusted maximum likelihood estimation (AMLE, Cohn 1988; Cohn and others, 1992b), first order bias in the model coefficients is eliminated using the calculations given in Shenton and Bowman (1977). A “nearly unbiased” (Cohn 1988) estimate of instantaneous load then is given by:

$$\hat{L}_{AMLE} = \exp \left( a_0 + \sum_{j=1}^{NV} a_j X_j \right) H(a, b, s^2, \alpha, \kappa) \quad (7)$$

where  $\hat{L}_{AMLE}$  is the AMLE estimate of instantaneous load,  $a$  and  $b$  are functions of the explanatory variables (Cohn and others, 1992b),  $\alpha$  and  $\kappa$  are parameters of the gamma distribution, and  $s^2$  is the residual variance. The model coefficients in equation 7 ( $a_0$  and  $a_j$ ) are maximum likelihood estimates corrected for first-order bias; the bias correction factor [ $H(a, b, s^2, \alpha, \kappa)$ ] is an approximation of the infinite series given in Cohn and others (1992b).

Under AMLE, estimates of instantaneous load are developed for all of the observations in the estimation data set using equation 7. Mean load estimates for various time periods then are calculated using equation 3 (where  $\hat{L}_\tau = \hat{L}_{AMLE}$ ). The uncertainty associated with each estimate of mean load is expressed in terms of the standard error (SE) and the standard error of prediction (SEP). The SE for each mean load estimate (Cohn and others, 1992b; equ. 35) represents the variability that may

be attributed to the model calibration (parameter uncertainty). Calculation of the SEP begins with an estimate of parameter uncertainty (the SE) and adds the unexplained variability about the model (random error). Because SEP incorporates parameter uncertainty and random error, it is larger than SE and provides a better description of how closely estimated loads correspond to actual loads. The SEP is therefore the preferred method of describing uncertainty in loads and is used within LOADEST to develop 95 percent confidence intervals for each estimate of mean load.

### 2.2.3 Least Absolute Deviation (LAD)

All of the regression methods discussed thus far (OLS, MLE, AMLE) assume the model residuals are normally distributed with constant variance. When model residuals do not conform to the assumption, alternate techniques may be appropriate. One such technique, the least absolute deviation (LAD) method, is implemented within LOADEST. Model coefficients for LAD are developed using the regression method of Powell (1984), as implemented by Buchinsky (1994). Given the model coefficients, estimates of instantaneous load are developed using the “smearing” approach of Duan (1983):

$$\hat{L}_{LAD} = \exp\left(a_0 + \sum_{j=1}^{NV} a_j X_j\right) \frac{\sum_{k=1}^n \exp(e_k)}{n} \quad (8)$$

where  $\hat{L}_{LAD}$  is the LAD estimate of instantaneous load,  $a_0$  and  $a_j$  are model coefficients developed by the LAD regression,  $e$  is the residual error, and  $n$  is the number of uncensored observations in the calibration data set<sup>3</sup>.

LAD estimates of instantaneous load are developed for all of the observations in the estimation data set using equation 8. Mean load estimates for various time periods then are calculated using equation 3 (where  $\hat{L}_\tau = \hat{L}_{LAD}$ ). Standard errors reflecting the uncertainty in each estimate of mean load are calculated using the jackknife method described by Efron (1982).

### 2.2.4 Summary of MLE, AMLE, and LAD for Load Estimation

The primary load estimation method used within LOADEST is AMLE. AMLE has been shown to have negligible bias when the calibration data set is censored (Cohn and others, 1992b). For the special case where the calibration data set is uncensored, the AMLE method converges to MLE (Cohn and others, 1992b), resulting in a minimum variance unbiased estimate of constituent loads. MLE estimates are provided as a check on AMLE results and as a means of comparing LOADEST results with standard statistical packages that implement MLE.

AMLE and MLE results are contingent upon the assumption that model residuals are normally distributed. Following model formulation and calibration (Section 2.1), AMLE residuals should be examined to see if the normality assumption is valid. Checks for normality include calculation of the PPCC (probability plot correlation coefficient; Vogel, 1986) and Turnbull-Weiss likelihood ratio (Turnbull and Weiss, 1978) statistics, construction of a normal-probability plot (the graphical analog of the PPCC; Helsel and Hirsch, 2002), and examination of standardized residuals. If the residuals do not adhere to the assumption of normality, AMLE (and MLE) results for censored data may not be optimal. Load estimates from the LAD method should therefore be considered in lieu of AMLE, as the LAD load estimates are not dependent on the normality assumption.

---

<sup>3</sup>Because of a lack of published results for the verification of the numerical algorithm, the LOADEST implementation of LAD is limited to the case of uncensored data. LAD results are omitted when the calibration data set contains censored observations.