

C. Radhakrishna Rao
Helge Toutenburg
Shalabh
Christian Heumann

Linear Models and Generalizations

Least Squares and Alternatives

With Contributions by Michael Schomaker

With 35 Illustrations
Third Extended Edition

Springer

Contents

Preface to the First Edition	v
Preface to the Second Edition	vii
Preface to the Third Edition	ix
1 Introduction	1
1.1 Linear Models and Regression Analysis	1
1.2 Plan of the Book	3
2 The Simple Linear Regression Model	7
2.1 The Linear Model	7
2.2 Least Squares Estimation	8
2.3 Direct Regression Method	10
2.4 Properties of the Direct Regression Estimators	12
2.5 Centered Model	14
2.6 No Intercept Term Model	15
2.7 Maximum Likelihood Estimation	15
2.8 Testing of Hypotheses and Confidence Interval Estimation	17
2.9 Analysis of Variance	20
2.10 Goodness of Fit of Regression	23
2.11 Reverse Regression Method	24
2.12 Orthogonal Regression Method	24
2.13 Reduced Major Axis Regression Method	27

2.14	Least Absolute Deviation Regression Method	29
2.15	Estimation of Parameters when X Is Stochastic	30
3	The Multiple Linear Regression Model and Its Extensions	33
3.1	The Linear Model	33
3.2	The Principle of Ordinary Least Squares (OLS)	35
3.3	Geometric Properties of OLS	36
3.4	Best Linear Unbiased Estimation	38
3.4.1	Basic Theorems	38
3.4.2	Linear Estimators	43
3.4.3	Mean Dispersion Error	44
3.5	Estimation (Prediction) of the Error Term ϵ and σ^2	45
3.6	Classical Regression under Normal Errors	46
3.6.1	The Maximum-Likelihood (ML) Principle	47
3.6.2	Maximum Likelihood Estimation in Classical Normal Regression	47
3.7	Consistency of Estimators	49
3.8	Testing Linear Hypotheses	51
3.9	Analysis of Variance	57
3.10	Goodness of Fit	59
3.11	Checking the Adequacy of Regression Analysis	61
3.11.1	Univariate Regression	61
3.11.2	Multiple Regression	61
3.11.3	A Complex Example	65
3.11.4	Graphical Presentation	69
3.12	Linear Regression with Stochastic Regressors	70
3.12.1	Regression and Multiple Correlation Coefficient	70
3.12.2	Heterogenous Linear Estimation without Normality	72
3.12.3	Heterogeneous Linear Estimation under Normality	73
3.13	The Canonical Form	76
3.14	Identification and Quantification of Multicollinearity	77
3.14.1	Principal Components Regression	77
3.14.2	Ridge Estimation	79
3.14.3	Shrinkage Estimates	83
3.14.4	Partial Least Squares	84
3.15	Tests of Parameter Constancy	87
3.15.1	The Chow Forecast Test	88
3.15.2	The Hansen Test	91
3.15.3	Tests with Recursive Estimation	92
3.15.4	Test for Structural Change	93
3.16	Total Least Squares	96
3.17	Minimax Estimation	98
3.17.1	Inequality Restrictions	98

3.17.2	The Minimax Principle	101
3.18	Censored Regression	105
3.18.1	Overview	105
3.18.2	LAD Estimators and Asymptotic Normality . . .	107
3.18.3	Tests of Linear Hypotheses	108
3.19	Simultaneous Confidence Intervals	110
3.20	Confidence Interval for the Ratio of Two Linear Parametric Functions	112
3.21	Nonparametric Regression	112
3.21.1	Estimation of the Regression Function	114
3.22	Classification and Regression Trees (CART)	117
3.23	Boosting and Bagging	121
3.24	Projection Pursuit Regression	124
3.25	Neural Networks and Nonparametric Regression	126
3.26	Logistic Regression and Neural Networks	127
3.27	Functional Data Analysis (FDA)	127
3.28	Restricted Regression	130
3.28.1	Problem of Selection	130
3.28.2	Theory of Restricted Regression	130
3.28.3	Efficiency of Selection	132
3.28.4	Explicit Solution in Special Cases	133
3.29	LINEX Loss Function	135
3.30	Balanced Loss Function	137
3.31	Complements	138
3.31.1	Linear Models without Moments: Exercise	138
3.31.2	Nonlinear Improvement of OLSE for Nonnormal Disturbances	139
3.31.3	A Characterization of the Least Squares Estimator	139
3.31.4	A Characterization of the Least Squares Estimator: A Lemma	140
3.32	Exercises	140
4	The Generalized Linear Regression Model	143
4.1	Optimal Linear Estimation of β	144
4.1.1	R_1 -Optimal Estimators	145
4.1.2	R_2 -Optimal Estimators	149
4.1.3	R_3 -Optimal Estimators	150
4.2	The Aitken Estimator	151
4.3	Misspecification of the Dispersion Matrix	153
4.4	Heteroscedasticity and Autoregression	156
4.5	Mixed Effects Model: Unified Theory of Linear Estimation	164
4.5.1	Mixed Effects Model	164
4.5.2	A Basic Lemma	164
4.5.3	Estimation of $X\beta$ (the Fixed Effect)	166

4.5.4	Prediction of $U\xi$ (the Random Effect)	166
4.5.5	Estimation of ϵ	167
4.6	Linear Mixed Models with Normal Errors and Random Effects	168
4.6.1	Maximum Likelihood Estimation of Linear Mixed Models	171
4.6.2	Restricted Maximum Likelihood Estimation of Linear Mixed Models	174
4.6.3	Inference for Linear Mixed Models	178
4.7	Regression-Like Equations in Econometrics	183
4.7.1	Econometric Models	186
4.7.2	The Reduced Form	190
4.7.3	The Multivariate Regression Model	192
4.7.4	The Classical Multivariate Linear Regression Model	195
4.7.5	Stochastic Regression	196
4.7.6	Instrumental Variable Estimator	197
4.7.7	Seemingly Unrelated Regressions	198
4.7.8	Measurement Error Models	199
4.8	Simultaneous Parameter Estimation by Empirical Bayes Solutions	209
4.8.1	Overview	209
4.8.2	Estimation of Parameters from Different Linear Models	211
4.9	Supplements	215
4.10	Gauss-Markov, Aitken and Rao Least Squares Estimators	216
4.10.1	Gauss-Markov Least Squares	216
4.10.2	Aitken Least Squares	217
4.10.3	Rao Least Squares	218
4.11	Exercises	220
5	Exact and Stochastic Linear Restrictions	223
5.1	Use of Prior Information	223
5.2	The Restricted Least-Squares Estimator	225
5.3	Maximum Likelihood Estimation under Exact Restrictions	227
5.4	Stepwise Inclusion of Exact Linear Restrictions	228
5.5	Biased Linear Restrictions and MDE Comparison with the OLSE	233
5.6	MDE Matrix Comparisons of Two Biased Estimators	236
5.7	MDE Matrix Comparison of Two Linear Biased Estimators	242
5.8	MDE Comparison of Two (Biased) Restricted Estimators	243
5.9	Stein-Rule Estimators under Exact Restrictions	251
5.10	Stochastic Linear Restrictions	252
5.10.1	Mixed Estimator	252
5.10.2	Assumptions about the Dispersion Matrix	254

5.10.3	Biased Stochastic Restrictions	257
5.11	Stein-Rule Estimators under Stochastic Restrictions . . .	261
5.12	Weakened Linear Restrictions	262
5.12.1	Weakly (R, r) -Unbiasedness	262
5.12.2	Optimal Weakly (R, r) -Unbiased Estimators . . .	262
5.12.3	Feasible Estimators—Optimal Substitution of β in $\hat{\beta}_1(\beta, A)$	266
5.12.4	RLSE instead of the Mixed Estimator	268
5.13	Exercises	269
6	Prediction in the Generalized Regression Model	271
6.1	Introduction	271
6.2	Some Simple Linear Models	271
6.2.1	The Constant Mean Model	271
6.2.2	The Linear Trend Model	272
6.2.3	Polynomial Models	273
6.3	The Prediction Model	274
6.4	Optimal Heterogeneous Prediction	275
6.5	Optimal Homogeneous Prediction	277
6.6	MDE Matrix Comparisons between Optimal and Classical Predictors	280
6.6.1	Comparison of Classical and Optimal Prediction with Respect to the y_* Superiority . .	283
6.6.2	Comparison of Classical and Optimal Predictors with Respect to the $X_*\beta$ Superiority .	285
6.7	Prediction Regions	287
6.7.1	Concepts and Definitions	287
6.7.2	On q -Prediction Intervals	289
6.7.3	On q -Intervals in Regression Analysis	291
6.7.4	On (p, q) -Prediction Intervals	292
6.7.5	Linear Utility Functions	294
6.7.6	Normally Distributed Populations - Two-Sided Symmetric Intervals	296
6.7.7	Onesided Infinite Intervals	298
6.7.8	Utility and Length of Intervals	298
6.7.9	Utility and coverage	300
6.7.10	Maximal Utility and Optimal Tests	300
6.7.11	Prediction Ellipsoids Based on the GLSE	302
6.7.12	Comparing the Efficiency of Prediction Ellipsoids	305
6.8	Simultaneous Prediction of Actual and Average Values of y	306
6.8.1	Specification of Target Function	307
6.8.2	Exact Linear Restrictions	308
6.8.3	MDEP Using Ordinary Least Squares Estimator	309
6.8.4	MDEP Using Restricted Estimator	309
6.8.5	MDEP Matrix Comparison	310

6.8.6	Stein-Rule Predictor	310
6.8.7	Outside Sample Predictions	311
6.9	Kalman Filter	314
6.9.1	Dynamical and Observational Equations	314
6.9.2	Some Theorems	314
6.9.3	Kalman Model	317
6.10	Exercises	318
7	Sensitivity Analysis	321
7.1	Introduction	321
7.2	Prediction Matrix	321
7.3	Effect of Single Observation on Estimation of Parameters	327
7.3.1	Measures Based on Residuals	328
7.3.2	Algebraic Consequences of Omitting an Observation	329
7.3.3	Detection of Outliers	330
7.4	Diagnostic Plots for Testing the Model Assumptions . . .	334
7.5	Measures Based on the Confidence Ellipsoid	335
7.6	Partial Regression Plots	341
7.7	Regression Diagnostics for Removing an Observation with Graphics	343
7.8	Model Selection Criteria	350
7.8.1	Akaike's Information Criterion	351
7.8.2	Bayesian Information Criterion	353
7.8.3	Mallows C_p	353
7.8.4	Example	355
7.9	Exercises	356
8	Analysis of Incomplete Data Sets	357
8.1	Statistical Methods with Missing Data	358
8.1.1	Complete Case Analysis	358
8.1.2	Available Case Analysis	358
8.1.3	Filling in the Missing Values	359
8.1.4	Model-Based Procedures	359
8.2	Missing-Data Mechanisms	360
8.2.1	Missing Indicator Matrix	360
8.2.2	Missing Completely at Random	360
8.2.3	Missing at Random	360
8.2.4	Nonignorable Nonresponse	360
8.3	Missing Pattern	360
8.4	Missing Data in the Response	361
8.4.1	Least-Squares Analysis for Filled-up Data—Yates Procedure	362
8.4.2	Analysis of Covariance—Bartlett's Method . . .	363
8.5	Shrinkage Estimation by Yates Procedure	364

8.5.1	Shrinkage Estimators	364
8.5.2	Efficiency Properties	365
8.6	Missing Values in the X -Matrix	367
8.6.1	General Model	367
8.6.2	Missing Values and Loss in Efficiency	368
8.7	Methods for Incomplete X -Matrices	371
8.7.1	Complete Case Analysis	371
8.7.2	Available Case Analysis	371
8.7.3	Maximum-Likelihood Methods	372
8.8	Imputation Methods for Incomplete X -Matrices	373
8.8.1	Maximum-Likelihood Estimates of Missing Values	373
8.8.2	Zero-Order Regression	374
8.8.3	First-Order Regression	375
8.8.4	Multiple Imputation	377
8.8.5	Weighted Mixed Regression	378
8.8.6	The Two-Stage WMRE	382
8.9	Assumptions about the Missing Mechanism	384
8.10	Regression Diagnostics to Identify Non-MCAR Processes	384
8.10.1	Comparison of the Means	384
8.10.2	Comparing the Variance-Covariance Matrices	385
8.10.3	Diagnostic Measures from Sensitivity Analysis	385
8.10.4	Distribution of the Measures and Test Procedure	385
8.11	Treatment of Nonignorable Nonresponse	386
8.11.1	Joint Distribution of (X, Y) with Missing Values Only in Y	386
8.11.2	Conditional Distribution of Y Given X with Missing Values Only in Y	388
8.11.3	Conditional Distribution of Y Given X with Missing Values Only in X	389
8.11.4	Other Approaches	390
8.12	Further Literature	391
8.13	Exercises	391
9	Robust Regression	393
9.1	Overview	393
9.2	Least Absolute Deviation Estimators — Univariate Case	394
9.3	M-Estimates: Univariate Case	398
9.4	Asymptotic Distributions of LAD Estimators	401
9.4.1	Univariate Case	401
9.4.2	Multivariate Case	402
9.5	General M-Estimates	403
9.6	Tests of Significance	407

10 Models for Categorical Response Variables	411
10.1 Generalized Linear Models	411
10.1.1 Extension of the Regression Model	411
10.1.2 Structure of the Generalized Linear Model	413
10.1.3 Score Function and Information Matrix	416
10.1.4 Maximum-Likelihood Estimation	417
10.1.5 Testing of Hypotheses and Goodness of Fit	420
10.1.6 Overdispersion	421
10.1.7 Quasi Loglikelihood	423
10.2 Contingency Tables	425
10.2.1 Overview	425
10.2.2 Ways of Comparing Proportions	427
10.2.3 Sampling in Two-Way Contingency Tables	429
10.2.4 Likelihood Function and Maximum-Likelihood Estimates	430
10.2.5 Testing the Goodness of Fit	432
10.3 GLM for Binary Response	435
10.3.1 Logit Models and Logistic Regression	435
10.3.2 Testing the Model	437
10.3.3 Distribution Function as a Link Function	438
10.4 Logit Models for Categorical Data	439
10.5 Goodness of Fit—Likelihood-Ratio Test	440
10.6 Loglinear Models for Categorical Variables	441
10.6.1 Two-Way Contingency Tables	441
10.6.2 Three-Way Contingency Tables	444
10.7 The Special Case of Binary Response	448
10.8 Coding of Categorical Explanatory Variables	450
10.8.1 Dummy and Effect Coding	450
10.8.2 Coding of Response Models	453
10.8.3 Coding of Models for the Hazard Rate	455
10.9 Extensions to Dependent Binary Variables	457
10.9.1 Overview	458
10.9.2 Modeling Approaches for Correlated Response	460
10.9.3 Quasi-Likelihood Approach for Correlated Binary Response	460
10.9.4 The GEE Method by Liang and Zeger	462
10.9.5 Properties of the GEE Estimate $\hat{\beta}_G$	463
10.9.6 Efficiency of the GEE and IEE Methods	465
10.9.7 Choice of the Quasi-Correlation Matrix $R_t(\alpha)$	465
10.9.8 Bivariate Binary Correlated Response Variables	466
10.9.9 The GEE Method	467
10.9.10 The IEE Method	468
10.9.11 An Example from the Field of Dentistry	469
10.9.12 Full Likelihood Approach for Marginal Models	474

10.10 Exercises	486
A Matrix Algebra	489
A.1 Overview	489
A.2 Trace of a Matrix	491
A.3 Determinant of a Matrix	492
A.4 Inverse of a Matrix	494
A.5 Orthogonal Matrices	495
A.6 Rank of a Matrix	495
A.7 Range and Null Space	496
A.8 Eigenvalues and Eigenvectors	496
A.9 Decomposition of Matrices	498
A.10 Definite Matrices and Quadratic Forms	501
A.11 Idempotent Matrices	507
A.12 Generalized Inverse	508
A.13 Projectors	516
A.14 Functions of Normally Distributed Variables	517
A.15 Differentiation of Scalar Functions of Matrices	520
A.16 Miscellaneous Results, Stochastic Convergence	523
B Tables	527
C Software for Linear Regression Models	531
C.1 Software	531
C.2 Special-Purpose Software	536
C.3 Resources	537
References	539
Index	563

$$E(y) = \beta_0 + \beta_1 X$$

and

$$\text{var}(y) = \sigma^2 .$$

Sometimes X can also be a random variable. Such an aspect is explained later in Section 2.15. In such a case, instead of simple mean and simple variance of y , we consider the conditional mean of y given $X = x$ as

$$E(y|x) = \beta_0 + \beta_1 x$$

and the conditional variance of y given $X = x$ as

$$\text{var}(y|x) = \sigma^2 .$$

The parameters β_0 , β_1 and σ^2 are generally unknown and e is unobserved. The determination of the statistical model (2.1) depends on the determination (*i.e.*, estimation) of β_0 , β_1 and σ^2 .

Only T pairs of observations (x_t, y_t) ($t = 1, \dots, T$) on (X, y) are observed which are used to determine the unknown parameters.

Different methods of estimation can be used to determine the estimates of the parameters. Among them, the least squares and maximum likelihood principles are the most popular methods of estimation.

2.2 Least Squares Estimation

We observe a sample of T sets of observations (x_t, y_t) ($t = 1, \dots, T$) and in view of (2.1), we can write

$$y_t = \beta_0 + \beta_1 x_t + e_t \quad (t = 1, \dots, T) . \quad (2.2)$$

The principle of least squares aims at estimating β_0 and β_1 so that the sum of squares of difference between the observations and the line in the scatter diagram is minimum. Such an idea is viewed from different perspectives. When the vertical difference between the observations and the line in the scatter diagram (see Fig. 2.1(a)) is considered and its sum of squares is minimized to obtain the estimates of β_0 and β_1 , the method is known as *direct regression*.

Another approach is to minimize the sum of squares of difference between the observations and the line in horizontal direction in the scatter diagram (see Fig. 2.1(b)) to obtain the estimates of β_0 and β_1 . This is known as *reverse (or inverse) regression* method.

Alternatively, the sum of squares of perpendicular distance between the observations and the line in the scatter diagram (see Fig. 2.1(c)) is minimized to obtain the estimates of β_0 and β_1 . This is known as *orthogonal regression* or *major axis regression* method.

The *least absolute deviation regression* method considers the sum of the absolute deviation of the observations from the line in the vertical direction in the scatter diagram (see Fig. 2.1(a)) to obtain the estimates of β_0 and β_1 .

The *reduced major axis regression* method proposes to minimize the sum of the areas of rectangles defined between the observed data points and the nearest point on the line in the scatter diagram to obtain the estimates of the regression coefficients (see Fig. 2.1(d)).

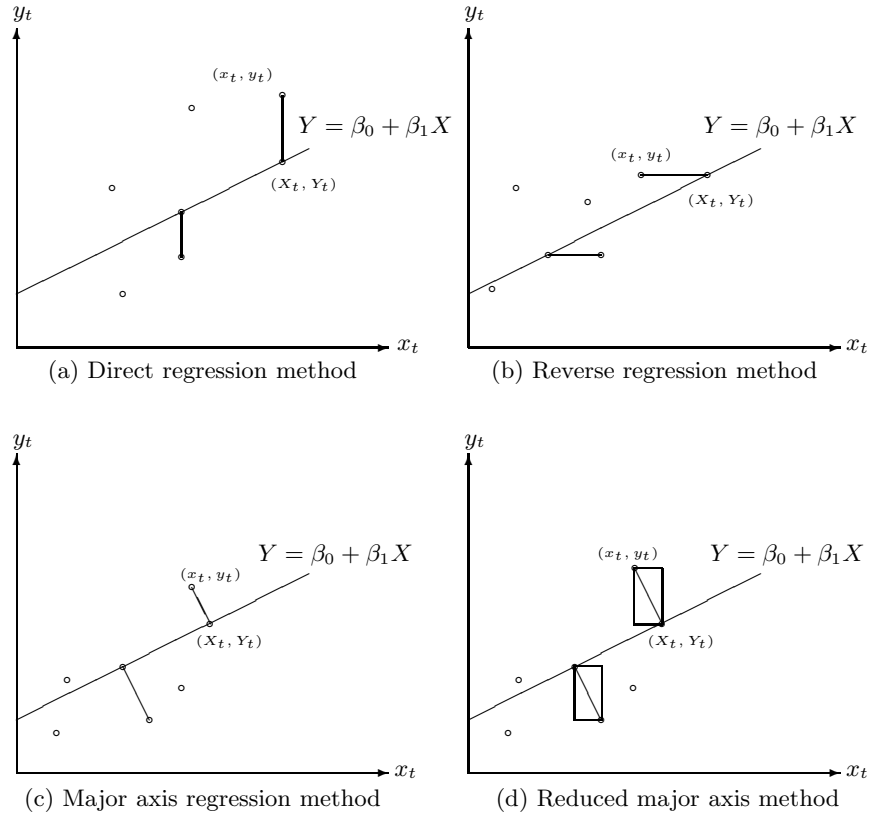


FIGURE 2.1. Scatter diagrams of different methods of regression

One may note that the principle of least squares does not require any assumption about the form of probability distribution of e_t in deriving the least squares estimates. For the purpose of deriving the statistical inferences only, we assume that e_t 's are observed as random variable ϵ_t with $E(\epsilon_t) = 0$, $\text{var}(\epsilon_t) = \sigma^2$ and $\text{cov}(\epsilon_t, \epsilon_{t^*}) = 0$ for all $t \neq t^*$ ($t, t^* = 1, \dots, T$). This

assumption is needed to find the mean and variance of the least squares estimates. The assumption that ϵ_t 's are normally distributed is utilized while constructing the tests of hypotheses and confidence intervals of the parameters.

Based on these approaches, different estimates of β_0 and β_1 are obtained which have different statistical properties. Among them the direct regression approach is more popular. Generally, the direct regression estimates are referred as the least squares estimates. We will consider here the direct regression approach in more detail. Other approaches are also discussed.

2.3 Direct Regression Method

This method is also known as the *ordinary least squares estimation*. The regression models (2.1) and (2.2) can be viewed as the regression models for population and sample, respectively. The direct regression approach minimizes the sum of squares

$$S(\beta_0, \beta_1) = \sum_{t=1}^T e_t^2 = \sum_{t=1}^T (y_t - \beta_0 - \beta_1 x_t)^2 \quad (2.3)$$

with respect to β_0 and β_1 .

The partial derivatives of (2.3) with respect to β_0 and β_1 are

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} = -2 \sum_{t=1}^T (y_t - \beta_0 - \beta_1 x_t) \quad (2.4)$$

and

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} = -2 \sum_{t=1}^T (y_t - \beta_0 - \beta_1 x_t) x_t, \quad (2.5)$$

respectively. The solution of β_0 and β_1 is obtained by setting (2.4) and (2.5) equal to zero. Thus obtained solutions are called the direct regression estimators, or usually called as the Ordinary Least Squares (OLS) estimators of β_0 and β_1 .

This gives the ordinary least squares estimates b_0 of β_0 and b_1 of β_1 as

$$b_0 = \bar{y} - b_1 \bar{x} \quad (2.6)$$

$$b_1 = \frac{SXY}{SXX} \quad (2.7)$$

where

$$SXY = \frac{1}{T} \sum_{t=1}^T (x_t - \bar{x})(y_t - \bar{y}), \quad SXX = \frac{1}{T} \sum_{t=1}^T (x_t - \bar{x})^2, \quad \bar{x} = \frac{1}{T} \sum_{t=1}^T x_t$$

$$\text{and } \bar{y} = \frac{1}{T} \sum_{t=1}^T y_t.$$

Further, using (2.4) and (2.5), we have

$$\begin{aligned}\frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_0^2} &= -2 \sum_{t=1}^T (-1) = 2T, \\ \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_1^2} &= 2 \sum_{t=1}^T x_t^2, \\ \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_0 \partial \beta_1} &= 2 \sum_{t=1}^T x_t = 2T\bar{x}.\end{aligned}$$

Thus we get the Hessian matrix which is the matrix of second order partial derivatives as

$$\begin{aligned}H &= \begin{pmatrix} \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_0^2} & \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_0 \partial \beta_1} & \frac{\partial^2 S(\beta_0, \beta_1)}{\partial \beta_1^2} \end{pmatrix} \\ &= 2 \begin{pmatrix} T & T\bar{x} \\ T\bar{x} & \sum_{t=1}^T x_t^2 \end{pmatrix} \\ &= 2 \begin{pmatrix} \mathbf{1}' \\ x' \end{pmatrix} (\mathbf{1}, x)\end{aligned}\quad (2.8)$$

where $\mathbf{1} = (1, \dots, 1)'$ is a T -vector of elements unity and $x = (x_1, \dots, x_T)'$ is a T -vector of observations on X . The matrix (2.8) is positive definite if its determinant and the element in the first row and column of H are positive. The determinant of H is

$$\begin{aligned}|H| &= 2 \left(T \sum_{t=1}^T x_t^2 - T^2 \bar{x}^2 \right) \\ &= 2T \sum_{t=1}^n (x_t - \bar{x})^2 \\ &\geq 0.\end{aligned}\quad (2.9)$$

The case when $\sum_{t=1}^T (x_t - \bar{x})^2 = 0$ is not interesting because then all the observations are identical, *i.e.*, $x_t = c$ (some constant). In such a case there is no relationship between x and y in the context of regression analysis. Since $\sum_{t=1}^T (x_t - \bar{x})^2 > 0$, therefore $|H| > 0$. So H is positive definite for any (β_0, β_1) ; therefore $S(\beta_0, \beta_1)$ has a global minimum at (b_0, b_1) .

The fitted line or the fitted linear regression model is

$$y = b_0 + b_1 X \quad (2.10)$$

and the predicted values are

$$\hat{y}_t = b_0 + b_1 x_t \quad (t = 1, \dots, T). \quad (2.11)$$

3.20 Confidence Interval for the Ratio of Two Linear Parametric Functions

Let $\theta_1 = P_1'\beta$ and $\theta_2 = P_2'\beta$ be two linear parametric functions and we wish to find a confidence interval of $\lambda = \frac{\theta_1}{\theta_2}$.

The least squares estimators of θ_1 and θ_2 are

$$\hat{\theta}_1 = P_1'\hat{\beta} \quad \text{and} \quad \hat{\theta}_2 = P_2'\hat{\beta}$$

with the variance-covariance matrix

$$\sigma^2 \begin{pmatrix} P_1'HP_1 & P_1'HP_2 \\ P_2'HP_1 & P_2'HP_2 \end{pmatrix} = \sigma^2 \begin{pmatrix} a & b \\ b' & c \end{pmatrix}, \text{ say.}$$

Then

$$E(\hat{\theta}_1 - \lambda\hat{\theta}_2) = 0, \quad \text{var}(\hat{\theta}_1 - \lambda\hat{\theta}_2) = \sigma^2(a - 2\lambda b + \lambda^2 c).$$

Hence

$$F = \frac{(\hat{\theta}_1 - \lambda\hat{\theta}_2)^2}{s^2(a - 2\lambda b + \lambda^2 c)} \sim F_{1, T-K}$$

and

$$P \left\{ (\hat{\theta}_1 - \lambda\hat{\theta}_2)^2 - F_{1-\alpha} s^2 (a - 2\lambda b + \lambda^2 c) \leq 0 \right\} = 1 - \alpha. \quad (3.377)$$

The inequality within the brackets in (3.377) provides a $(1 - \alpha)$ confidence region for λ . Because the expression in (3.377) is quadratic in λ , the confidence region is the interval between the roots of the quadratic equation or outside the interval, depending on the nature of the coefficients of the quadratic equation.

3.21 Nonparametric Regression

The nonparametric regression model describes the dependence of study variable on explanatory variables without specifying the function that relates them. The general nonparametric regression model is expressed as

$$y = \psi(X) + \epsilon \quad (3.378)$$

where y is a study variable, X is a vector of explanatory variables, ϵ is disturbance term and $\psi(x)$ is the unspecified real valued function of X at some fixed value x given by

$$\psi(x) = E(y|X = x). \quad (3.379)$$

The first derivative of $\psi(x)$ indicates the response or regression coefficient of y with respect to x and the second derivative of $\psi(x)$ indicates the curvature of $\psi(x)$.

We assume that both y and X_1, \dots, X_K are stochastic and related as

$$y = \psi(X_1, \dots, X_K) + \epsilon \quad (3.380)$$

with $E(\epsilon|X = x) = 0$ and $V(\epsilon|X = x) = \sigma^2 I$. We observe T identically and independently distributed observations $(y_t, x_{t1}, \dots, x_{tK})$, $t = 1, \dots, T$ from an absolutely continuous $(K + 1)$ -variate distribution with density $f(y, X_1, \dots, X_K) = f(y, X)$. If $E(|y|) < \infty$, then the conditional mean of y given $X = x$ exists as in (3.379).

The regression coefficient related to x_j ($j = 1, \dots, K$) is

$$\beta_j(x) = \beta(x) = \frac{\partial \psi(x)}{\partial x_j} = \lim_{h \rightarrow 0} \frac{\psi(x + h) - \psi(x - h)}{2h} \quad (3.381)$$

where $\psi(x - h) = \psi(x_1, \dots, x_j - h, \dots, x_K)$. When $\psi(x)$ is linear, then $\beta_j(x)$ is the j^{th} regression coefficient and is fixed for all x . When $\psi(x)$ is non-linear, then $\beta_j(x)$ depends on x and $\beta_j(x)$ is a varying regression coefficient. The fixed regression coefficient can be defined as $\beta(\bar{x})$, *i.e.*, $\beta(x)$ evaluated at $x = \bar{x} = (\bar{x}_1, \dots, \bar{x}_K)$. Similarly the second order partial derivative is

$$\beta^{(2)}(x) = \frac{\partial^2}{\partial x_j^2} \psi(x) = \lim_{h \rightarrow 0} \frac{\psi(x + 2h) - 2\psi(x) - \psi(x - 2h)}{(2h)^2}, \quad (3.382)$$

and, in general, the p^{th} order partial derivative ($p = 1, 2, \dots$) is

$$\begin{aligned} \beta^{(p)}(x) &= \frac{\partial^p}{\partial x_j^p} \psi(x) \\ &= \lim_{h \rightarrow 0} \left[\left(\frac{1}{2h} \right)^p \sum_{m=0}^p (-1)^m \binom{p}{m} \psi(x + (p - 2m)h) \right] \end{aligned} \quad (3.383)$$

and the cross partial derivative is

$$\begin{aligned} \frac{\partial^{p_1 + \dots + p_r}}{\partial x_{j_1}^{p_1}, \dots, \partial x_{j_r}^{p_r}} &= \\ \lim_{h \rightarrow 0} \left[\left(\frac{1}{2h} \right)^{p_1 + \dots + p_r} \sum_{m_1=0}^{p_1} \dots \sum_{m_r=0}^{p_r} (-1)^{m_1 + \dots + m_r} \right. \\ &\quad \times \binom{p_1}{m_1} \dots \binom{p_r}{m_r} \psi(x + (p_1 - 2m_1)h, \dots, x + (p_r - 2m_r)h) \left. \right] \end{aligned} \quad (3.384)$$

respectively, where each of $j_1, \dots, j_r = 1, \dots, K$ ($j_1 \neq \dots \neq j_r$), $x + (p_1 - 2m_1)h = (x_1, \dots, x_{j_1} + (p_1 - 2m_1)h, \dots, x_K)$, $x + (p_r - 2m_r)h = (x_1, \dots, x_{j_r} + (p_r - 2m_r)h, \dots, x_K)$.

Now we consider the nonparametric estimation of partial derivatives of $\psi(x)$ without specifying its form. We first consider a nonparametric estimator of $\psi(x)$ and then take its partial derivatives.

3.21.1 Estimation of the Regression Function

Most of the estimation methods of nonparametric regression assume that the regression function is smooth in some sense.

Since $\psi(x)$ depends on unknown densities, so we use the data of $(K + 1)$ variables as $z_t = (x_t, y_t)$, $t = 1, \dots, T$ and note that

$$g_T(z) = \frac{1}{Th^{K+1}} \sum_{t=1}^T \mathcal{K}\left(\frac{z_t - z}{h}\right), \quad (3.385)$$

$$g_{1T}(x) = \int g_T(z) dy = \frac{1}{Th^K} \sum_{t=1}^T \mathcal{K}_1\left(\frac{x_t - x}{h}\right), \quad (3.386)$$

where h is the window width (also called as band width or smoothing parameter) which is a positive function of T that goes to zero as $T \rightarrow \infty$, \mathcal{K} is a kernel or a weight function such that $\int \mathcal{K}(z) dz = 1$ and $\mathcal{K}_1(x) = \int \mathcal{K}(z) dy$. The kernel \mathcal{K} determines the shape of the curve and h determines their width. See, Prakasa-Rao (1983), Silverman (1986), Ullah and Vinod (1988) and Pagan and Ullah (1999) for the details on kernel density estimation. Substituting (3.385) and (3.386) in (3.379), we have

$$\psi_T(x) = \psi_T = \int y \frac{g_T(z)}{g_{1T}(x)} dy = \sum_{t=1}^T y_t w_t(x) \quad (3.387)$$

where

$$w_t(x) = \frac{\mathcal{K}_1\left(\frac{x_t - x}{h}\right)}{\sum_{t=1}^T \mathcal{K}_1\left(\frac{x_t - x}{h}\right)}. \quad (3.388)$$

The estimator ψ_T of ψ in (3.387) is known as Nadaraya–Watson type estimator due to Nadaraya (1964) and Watson (1964) and is a kernel nonparametric regression estimate. Note that (3.387) is a weighted average of the observed values y_t where the weight of t^{th} observation depends on the distance x_t to x through the kernel \mathcal{K} . The fitted nonparametric regression model that is obtained by without making any assumption about the functional form of $\psi(x)$ is

$$y = \psi_T(x) + \hat{\epsilon} \quad (3.389)$$

where $\hat{\epsilon}$ is the nonparametric residual.

This estimator (3.387) is also the weighted least squares estimator of $\psi(x)$ because $\psi_T(x)$ is the value of $\psi(x)$ for which the weighted squared error

$$\sum_{t=1}^T \mathcal{K}\left(\frac{x_t - x}{h}\right) (y_t - \psi(x))^2$$

is minimum. The method of moments also yields the same estimator of $\psi(x)$ as in (3.387).

When the window width h is not the same for all data points, then some alternative estimators of $\psi(x)$ are suggested. The recursive regression estimator of $\psi(x)$ in such a case is

$$\hat{\psi}_T(x) = \frac{\sum_{t=1}^T \frac{y_t}{h_t^K} \mathcal{K}\left(\frac{x_t - x}{h_t}\right)}{\sum_{t=1}^T \frac{1}{h_t^K} \mathcal{K}\left(\frac{x_t - x}{h_t}\right)} \quad (3.390)$$

where h_t denotes a sequence of positive numbers, assumed to satisfy $\sum h_t^K \rightarrow \infty$ as $T \rightarrow \infty$. An alternative estimator is

$$\tilde{\psi}_T(x) = \frac{\sum_{t=1}^T y_t \mathcal{K}\left(\frac{x_t - x}{h_t}\right)}{\sum_{t=1}^T \mathcal{K}\left(\frac{x_t - x}{h_t}\right)}. \quad (3.391)$$

Both (3.390) and (3.391) are recursive as

$$\hat{\psi}_T(x) = \hat{\psi}_{T-1}(x) + \frac{y_T - \hat{\psi}_{T-1}(x)}{1 + (T-1) \frac{\hat{f}_{T-1}(x)}{h_T^K} \mathcal{K}\left(\frac{x_T - x}{h_T}\right)} \quad (3.392)$$

$$\tilde{\psi}_T(x) = \tilde{\psi}_{T-1}(x) + \vartheta_T^{-1} \left[y_T - \tilde{\psi}_{T-1}(x) \mathcal{K}\left(\frac{x_T - x}{h_T}\right) \right] \quad (3.393)$$

where

$$\vartheta_T = \vartheta_{T-1} + \mathcal{K}\left(\frac{x_T - x}{h_T}\right), \quad \vartheta_0 = 0.$$

Both (3.392) and (3.393) can be updated as additional data points are available.

When ϵ 's are such that $V(\epsilon) = \Sigma (\neq \sigma^2 I)$, a $T \times T$ positive definite matrix, then the generalized least squares estimator of $\psi(x)$ is obtained by minimizing $\epsilon' \mathcal{K}^{1/2} \Sigma^{-1} \mathcal{K}^{1/2} \epsilon$ with respect to $\psi(x)$ as

$$\psi_T^* = (\mathbf{1}' \mathcal{K}^{\frac{1}{2}} \Sigma^{-1} \mathcal{K}^{\frac{1}{2}} \mathbf{1})^{-1} \mathbf{1}' \mathcal{K}^{\frac{1}{2}} \Sigma^{-1} y \quad (3.394)$$

where $\mathbf{1} = (1, \dots, 1)'$, $\mathcal{K} = \text{diag}(\mathcal{K}_1, \dots, \mathcal{K}_T)$ is a diagonal matrix with $\mathcal{K}_T = \mathcal{K}\left(\frac{x_T - x}{h}\right)$.

An operational version of a consistent estimator of $\beta(x)$ in (3.381) is

$$b_T(x) = \frac{\psi_T(x+h) - \psi_T(x-h)}{2h} \quad (3.395)$$

where $\psi_T(x)$ is given by (3.387). Similarly, the estimators of the p^{th} order partial derivative and cross partial derivatives can be obtained by replacing $\psi(\cdot)$ with $\psi_T(\cdot)$ in (3.383) and (3.384), respectively.

Ullah and Vinod (1988) analytically derived the estimator of $\beta(x) = \partial\psi(x)/\partial x$ as

$$\hat{\beta}_T(x) = \frac{\partial}{\partial x} \psi_T(x) = \sum_{t=1}^T y_t (\omega_{1t} - \omega_{2t})$$

where

$$\omega_{1t} = \frac{\mathcal{K}'\left(\frac{x_t - x}{h}\right)}{\sum_{t=1}^T \mathcal{K}\left(\frac{x_t - x}{h}\right)}$$

and $\omega_{2t} = \omega_t(x) \sum \omega_{1t}$; $\omega_t(x)$ is as in (3.388) and

$$\mathcal{K}'\left(\frac{x_t - x}{h}\right) = \frac{\partial}{\partial x_j} \mathcal{K}\left(\frac{x_t - x}{h}\right).$$

Alternatively, $\hat{\beta}_T(x)$ and its generalization for p^{th} order derivatives of $\psi(x)$, $\hat{\beta}^{(p)}(x)$ can be obtained as a solution of

$$\sum_{m=0}^p \binom{p}{m} \beta_T^{(m)}(x) f_T^{(p-1)}(x) = g_T^{(p)}(x), \quad (p = 1, 2, \dots)$$

where $g^{(p)}(x)$ is the p^{th} order partial derivative of $g(x) = \int y f(y, x) dy$ with respect to x_j .

The restricted least squares estimator of $\psi(x)$ under the exact linear restrictions $R\beta(x) = r$ is

$$\hat{\beta}_T(x) = b_T(x) - R'(RR')^{-1}[Rb_T(x) - r] \quad (3.396)$$

where $b_T(x)$ is given by (3.395).

The regression function can also be estimated by using various nonparametric procedures like nearest neighbor kernel estimation, local polynomial regression, and smoothing splines.

The method of nearest neighborhood kernel estimation is based on defining a symmetric unimodal weight function $W(x)$ which is centered on the focal observation and goes to zero at the boundaries of the neighborhood around the focal value. Let x_{fo} be a focal x -value at which $\psi(x)$ is to be estimated. Now find ν nearest x -neighbors of x_{fo} where ν/x_{fo} is the span of the kernel smoother. The larger the span, smoother is the estimated regression function. Using the weights defined by $W(x)$, calculate the weighted average of y and obtain the fitted value

$$\hat{y}_{fo} = \hat{f}(x_{fo}) = \frac{\sum_{t=1}^T y_t W(x_t)}{\sum_{t=1}^T W(x_t)}. \quad (3.397)$$

Repetition of this procedure at a range of x -values spanning the data and connecting the fitted values produces an estimate of the regression function.

In local polynomial regression, the fitted values are produced by locally weighted regression rather than by locally weighted averaging. Another method of nonparametric regression is smoothing splines which are the solution to the penalized regression problem. Additive regression models are an alternate to nonparametric regression with several explanatory variables.

The readers are referred to Prakasa-Rao (1983), Silverman (1986), Ullah (1989a), Ullah (1989b), Härdle (1990) and Pagan and Ullah (1999) for the asymptotic properties of the estimators, related testing of hypothesis and other aspects on nonparametric regression.

3.22 Classification and Regression Trees (CART)

Nonparametric regression with multiple explanatory variables suffers from the problem of *curse of dimensionality*. This means that if the number of explanatory variables is high, then it may be difficult to catch the relevant features of the problem in hand, e.g. the influence of interactions of explanatory variables on the study variable may be difficult to study. We have only a finite sample available, but there may be big volumes in the space of explanatory variables where there may be no observation or only a few observations are obtained (*sparseness problem*). Therefore a reliable statistical estimation is not possible in these volumes. Parametric models, such as simple linear models, or additive models as proposed by Hastie and Tibshirani (1990) try to catch at least the *main effects* of the explanatory variables and discard any global or local interactions of the explanatory variables on the study variable. Furthermore, the results from the approaches like Projection Pursuit Regression (see Section 3.24) or Neural Networks (see Section 3.25) may be hard to interpret. In such situations, CART is more useful. CART tries to catch the relevant interactions of the explanatory variables in their influence on the study variable and present the results in a simple way.

Consider a general regression setup in which the study variable y is either real-valued or categorical and X_1, \dots, X_K are the explanatory variables. In the usual nonparametric regression setup, we assume

$$y = \psi(X_1, \dots, X_K) + \epsilon, \quad (3.398)$$

with $E(\epsilon|X) = 0$. If the function ψ is unknown and not parameterized by a finite dimensional parameter, Breiman, H., Olshen and Stone (1984) suggested a recursive partitioning algorithm of the covariate space which results in a tree structure. If y is real-valued, as in (3.398), then the resulting

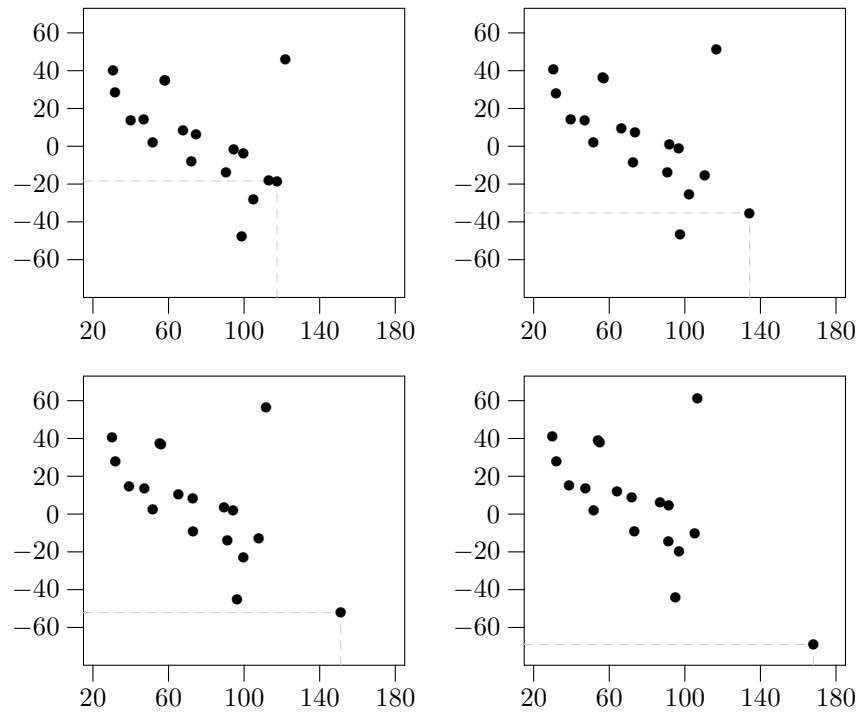


FIGURE 7.8. Four frames $\lambda = 0$ (a), $\lambda = \frac{1}{3}$ (b), $\lambda = \frac{2}{3}$ (c) and $\lambda = 1$ (d) (left to right, top down) of an animated plot of $\hat{e}(\lambda)$ versus $\hat{y}(\lambda)$ for data in Example 7.4 when removing the 17th observation (marked by dotted lines).

7.8 Model Selection Criteria

Consider a set of regression models with a study variable y and K explanatory variables X_1, \dots, X_K that can be fitted to a given set of data. The question then arises how to select a good model that has good agreement with the observed data. This also pertains to how to select the number of explanatory variables in the model which yield a ‘good’ model. A simple approach for model selection is to choose a model with smallest residual sum of squares (RSS). Another criterion is based on the coefficient of determination R^2 and adjusted R^2 which suggest to select the model with the higher coefficient of determination. Some other popular criteria of model selection are Akaike’s Information Criterion (AIC, Akaike (1973)), the Bayesian Information Criterion (BIC, Schwarz (1978)) and Mallows C_p (Mallows (1973)) which we introduce in the following subsections.

7.8.1 Akaike's Information Criterion

One approach to select a good model from a set of candidate models is to use Akaike's Information Criterion (AIC) due to Akaike (1973). Its concept is based on the relationship between information theory and likelihood theory.

Kullback and Leibler (1951) presented an approach to quantify the information and defined the Kullback-Leibler distance function as

$$\begin{aligned} I(f, g) &= \int f(x) \log \left(\frac{f(x)}{g(x)} \right) dx \\ &= \int f(x) \cdot \log(f(x)) dx - \int f(x) \cdot \log(g(x)) dx \end{aligned} \quad (7.98)$$

which can be considered as a measure of distance between the two functions f and g . These functions can be two different models from the selection point of view in the context of linear regression analysis. Let f denote the underlying but unknown true model and $g(x|\theta)$ denote any approximation to it which depends on the unknown parameters to be estimated. Then the Kullback-Leibler distance function is

$$\begin{aligned} I(f, g(x|\theta)) &= \int f(x) \cdot \log(f(x)) dx - \int f(x) \cdot \log(g(x|\theta)) dx \\ &= E_f[\log(f(x))] - E_f[\log(g(x|\theta))]. \end{aligned} \quad (7.99)$$

The model $g_a(x|\theta)$ is said to be better than $g_b(x|\theta)$ in the sense of Kullback-Leibler distance if $I(f, g_a(x|\theta)) < I(f, g_b(x|\theta))$. In such a case the distance between model $g_a(x|\theta)$ and f is smaller than the distance between $g_b(x|\theta)$ and f . Therefore the aim is to search a model $g(x|\theta)$ which minimizes the Kullback-Leibler distance $I(f, g(x|\theta))$. Using (7.99), we note that

$$\begin{aligned} I(f, g_a(x|\theta)) &< I(f, g_b(x|\theta)) \\ \Leftrightarrow E_f[\log(f(x))] - E_f[\log(g_a(x|\theta))] &< E_f[\log(f(x))] - E_f[\log(g_b(x|\theta))] \\ \Leftrightarrow -E_f[\log(g_a(x|\theta))] &< -E_f[\log(g_b(x|\theta))]. \end{aligned} \quad (7.100)$$

Thus (7.100) indicates that the term $E_f[\log(f(x))]$ can be treated as a constant for the comparison between two models and therefore $E_f[\log(g(x|\theta))]$ becomes the function of interest. Let $\hat{\theta}(y)$ denote an estimator of θ based on observations y . Since θ is usually unknown and has to be estimated from the data, so one can consider on minimizing the *expected* Kullback-Leibler distance $E_y \left[I(f, g(\cdot|\hat{\theta}(y))) \right]$ instead of (7.99). One good approach in such cases is to use the maximum likelihood estimate of θ .

So, using (7.100),

$$\begin{aligned} E_y \left[I(f, g(\cdot|\hat{\theta}(y))) \right] &= \int f(x) \log(f(x)) dx \\ &\quad - E_y \left[\int f(x) \log[g(x|\hat{\theta}(y))] dx \right] \\ &= C - E_y E_x \left[\log[g(x|\hat{\theta}(y))] \right] \end{aligned} \quad (7.101)$$

where C is a constant term. Akaike found a relationship between the second term of (7.101) and the log-likelihood function. He showed, that

$$E_y E_x \left[\log[g(x|\hat{\theta}(y))] \right] \approx L(\hat{\theta}|y) - K \quad (7.102)$$

where $L(\hat{\theta}|y)$ denotes the maximized log-likelihood for the model g and K is the the number of estimable parameters in the approximating model. The readers are referred to Burnham and Anderson (2002) for the derivation of (7.102). The AIC can further be expressed as

$$AIC = -2L(\hat{\theta}|y) + 2K. \quad (7.103)$$

Burnham and Anderson (2002) state that Akaike multiplied the bias-corrected log-likelihood by -2 for the reasons like e.g., that it is well known that -2 times the logarithm of the ratio of two maximized likelihood values is asymptotically chi-squared under certain conditions and assumptions. Two points frequently arise and we note this here. Firstly, the model associated with the minimum AIC remains unchanged if the bias corrected likelihood (*i.e.*, $\log L - K$) is multiplied by -0.17 , -34 , or -51.3 , or any other negative number. Thus the minimization is not changed by the multiplication of both the terms by a negative constant; Akaike merely chooses -2 . Secondly, some investigators have not realized the formal link between Kullback-Leibler information and AIC and believed that the number 2 in the second term in AIC was somehow arbitrary and other numbers should also be considered. This error has led to considerable confusion in the technical literature; clearly, K is the asymptotic bias correction and is not arbitrary. Akaike chooses to work with $-2 \log L$ rather than $\log L$; thus the term $+2K$ is theoretically correct for large sample size. As long as both the terms (the log-likelihood and the bias correction) are multiplied by the same negative constant, the minimization of AIC in (7.103) gives the same result and nothing is arbitrary.

The model with smallest AIC, which in turn is closest to the unknown true model, is said to be the best model under AIC criterion. In case of a linear regression model, AIC can be written as:

$$AIC_{LM} = T \cdot \ln(RSS) + 2K - T \cdot \ln(T). \quad (7.104)$$

The second terms in (7.103) and (7.104) are often termed as penalty term.