

Statistical Modelling For Heterogeneous Dataset

Sibnarayan Guria

Department of Statistics, West Bengal State University, Barasat, India.

ABSTRACT: *Statistical modelling is of prime importance in each and every sphere of data analysis. This paper reviews the justification of fitting linear model to the collected data. Inappropriateness of the fitted model may be due two reasons 1.wrong choice of the analytical form, 2. Suffers from the adverse effects of outliers and/or influential observations. The aim is to identify outliers using the deletion technique. In I extend the result of deletion diagnostics to the ex- changeable model and reviews some results of model analytical form checking and the technique illustrated through an example.*

KEYWORDS AND PHRASES: Parsimony, Auto-regressive model, Exchangeable model, outlier detection, deletion technique, Akaike Information Criterion.

AMS 2000 Subject Classification: 62J20

I. INTRODUCTION

In fitting a regression model to a given data set it is vital to know how good the fit is. Use of the model, and particularly predictions based on it, requires that the fitted model is compatible with the data. However, there may often be observations which are different from the others. These outliers frequently have an inordinate influence on least squares estimates. And quite often an outlier may neither show large residuals nor exhibit any influence on the regression line. But outliers need to be detected not only because they may affect the fit but also because they may lead to valuable information regarding the data. To this effect, diagnostics play an important role in regression analysis (see for example Belsley, Kuh and Welsch (1980), Chatterjee and Hadi (1988) or Sen and Srivastava (1990)). These diagnostics provide various ways of studying the residuals and assessing the impact of the respective observations on the regression line. Initially, diagnostic studies used the deletion of observation technique to assess this impact. In recent years Cook's (1986) method for assessing the local influence through model perturbation has also been used by several authors. However, most of the studies thus far conducted have been restricted to models having uncorrelated disturbances with constant variances. But as has been frequently observed, the dispersion matrix may not be spherical. In such cases it is necessary to apply the generalized least-squares method to estimate the parameters of the model. Hence the usual regression diagnostics need to be modified too. [1] Beach and MacKinnon (1978) circumvent this problem by developing a computationally efficient technique for maximizing the full likelihood function for an auto correlated linear regression model. In the generalized least-square context, Putterman (1988) studied the influence of the first transformed observation on the parameter estimates. However, Kim and Huggins (1998) claim that the deletion approach is inappropriate in studying the diagnostics in a regression model with auto correlated errors. They discuss the effects of simultaneous perturbation of the response vector on all the parameters as also the autocorrelation coefficient. Sharing their concern, Tsai and Wu (1992) used the profile likelihood function to examine the diagnostics through the effects of small perturbations. Schall and Dunne (1991) use a similar technique to study a regression model where the disturbances follow the ARMA model.

and devise some remedial procedure for estimation. Sen Roy & Guria (2009) We introduce the model and outline parameter estimation in Section 2, their justifications are discussed in Section 3. Diagnostics for exchangeable model are derived in Section 4. A numerical example is given in Section 5.

II. THE MODEL

Consider the regression model

$$y_t = \underline{x}_t \underline{\beta} + u_t, \quad t = 1, \dots, n. \quad (2.1)$$

where y_t is the response at time t , u_t the disturbance term, \underline{x}_t the vector of observations on the p explanatory variables and $\underline{\beta}$ the $p \times 1$ parameter vector. In general, it is assumed that u_t 's are independent, but in any real life situation u_t may not be independent. Data sets often clustered or otherwise correlated due to the way data were collected and intrinsic ecological pattern. If standard Gauss- Markov model is assumed, the likelihood of type I error is increased. Here I assumed that observations are correlated among themselves with correlation coefficient ρ .

i.e. u_t 's are:

$$E(u_t) = 0 \quad \forall \quad t$$

and $Cov(u_t, u_s) = \sigma^2 \quad \forall \quad t = s$ and $= \rho\sigma^2 \quad \forall t \neq s.$

(2.1) can be written in matrix form as

$$\mathbf{y} = \mathbf{X} \underline{\beta} + \underline{u}, \quad (2.2)$$

where \mathbf{y} and \underline{u} are $n \times 1$ vectors and \mathbf{X} is a $n \times p$ matrix.

Using the ordinary least-squares, the parameter $\underline{\beta}$ is estimated as

$$\underline{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}.$$

However, owing to conditions satisfied by the disturbance term in (2.2) above,

\underline{b} is not the best linear unbiased estimator (b.l.u.e.) of $\underline{\beta}$.

Since the dispersion of \underline{u} is

$$D(\underline{u}) = \sigma^2\Omega,$$

where

$$\Omega = \sigma^2 \begin{bmatrix} 1 & \rho & \rho & \cdots & \rho \\ \rho & 1 & \rho & \cdots & \rho \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ \rho & \rho & \rho & \cdots & 1 \end{bmatrix} \quad (2.3)$$

the b.l.u.e of $\underline{\beta}$ is given by

$$\underline{b}^* = (\mathbf{X}'\Omega^{-1}\mathbf{X})^{-1}\mathbf{X}'\Omega^{-1}\mathbf{y}. \quad (2.4)$$

In studying the regression diagnostics it is thus imperative that the residuals considered are not

$$\underline{\varepsilon} = \mathbf{y} - \mathbf{X}\underline{\mathbf{b}}$$

but $\underline{\varepsilon}^* = \mathbf{y} - \mathbf{X}\underline{\mathbf{b}}^*.$

Since Ω is a positive definite matrix, there exists a nonsingular matrix $P \ni \Omega^{-1} = P'P$. Therefore, defining $\mathbf{X}^* = P\mathbf{X}$, $\mathbf{y}^* = P\mathbf{y}$ and $\underline{u}^* = P \underline{u}$, (2.5) can be obtained from the transformed model

$$\mathbf{y}^* = \mathbf{X}^* \beta + \underline{u}^* \tag{2.5}$$

as

$$\begin{aligned} \underline{b}^* &= (\mathbf{X}' P' P \mathbf{X})^{-1} \mathbf{X}' P' P \mathbf{y} \\ &= (\mathbf{X}^{*'} \mathbf{X}^*)^{-1} \mathbf{X}^{*'} \mathbf{y}^*. \end{aligned} \tag{2.6}$$

Hence when the diagnostic relates to the deletion of one observation at a time and observing the discrepancy it produces, the corresponding row and column of Ω must be deleted. This would mean that the corresponding column of P is deleted. Here we study the effect of this on the model (2.1)-(2.2).

III. JUSTIFICATION

It common practice in any experimental science and social science fitting straight line to the collected data to understand what is happening in the system and guess future behaviour of the system. Suppose we have collected n data points on $p + 1$ variables y, x_1, x_2, \dots, x_p where y is random realisation and subject to measurement error treated as dependent variable and observation on other p x -variables are actually known, easily available or known and highly correlated to y . If we assume that collected data are well explained by their second moments. Let us define the variance-covariance matrix as

$$\mathbf{C} = \begin{pmatrix} s_{yy} & s_{yx_1} & \cdots & s_{yx_p} \\ s_{x_1y} & s_{x_1x_1} & \cdots & s_{x_1x_p} \\ \cdots & \cdots & \cdots & \cdots \\ s_{x_py} & s_{x_px_1} & \cdots & s_{x_px_p} \end{pmatrix} \quad \text{Where } s_{uv} = n^{-1} \sum_{i=1}^n (u_i - \bar{u})(v_i - \bar{v})$$

and vector of unknown constants $\underline{w} = (a_0 \ a_1 \ \cdots \ a_p)'$

Then $\underline{w}' \mathbf{C} \underline{w} = n^{-1} \sum_1^n \{a_0(y - \bar{y}) + a_1(x_1 - \bar{x}_1) + \cdots + a_p(x_p - \bar{x}_p)\}^2 \geq 0$
 " = " holds iff $a_0(y - \bar{y}) + a_1(x_1 - \bar{x}_1) + \cdots + a_p(x_p - \bar{x}_p) = 0$ corresponding

to 0 eigenvalue, eigenvector is normal to this plane. Which may be assumed because of the dependence between \underline{y} and \mathbf{X} . we are to some extent justifies our linearity of the model.

As we mentioned that in any live data set observations are somehow correlated using information criteria namely Akaike information criterion (AIC) may be used for model selection as advocated by Barnett, et.al.(2010). Variance-Covariance matrix of the observation vector \underline{y} involves $n(n + 1)/2$ parameters. To reduce the number of parameters equality of the variances (σ^2) are assumed and four usually available alternative covariance structures for modeling different real life scenario. (i) Independent: all covariances are assumed

to 0. (ii) exchangeable: covariance between any two observations are same i.e. $Cov(u_i, u_j) = \sigma^2\rho$; Therefore, observations are exchanged (re-arranged) over time and/or different segments. (iii) Autoregressive: current action depends on the immediate past result. i.e. $u_t = \rho u_{t-1} + \epsilon_t$. (iv) Unstructured: all the covariances are different. Following is an extract of comparisons on the basis of AIC values. Cell values giving percentage successful selections, and bold are percentage of correct choices.

Table 1: Showing AIC values for different choices of model

True covariance	Selected covariance			
	Indep	Exch.	AR	Unst.
Independent	70	15	15	0
Exchangeable ($\rho = 0.2$)	0	97	2	1
Exchangeable ($\rho = 0.5$)	0	100	0	0
Autoregressive ($\rho = 0.3$)	0	3	97	0
Autoregressive ($\rho = 0.7$)	0	0	100	0
Unstructured	0	49	24	27

Hence, we proceed for the diagnostic results for exchangeable model.

For (2.2) let $\underline{\beta} = (\beta_1 \ \beta_2 \ \dots \ \beta_p)$ and \mathbf{X} may contain a column of unities and we minimise the sum of squares of errors i.e. $(\underline{y} - \mathbf{X} \underline{\beta})' \mathbf{V}^{-1} (\underline{y} - \mathbf{X} \underline{\beta})$, i.e. minimise $\begin{pmatrix} 1 & \underline{\beta}' \end{pmatrix} \begin{pmatrix} \underline{y}' \mathbf{V}^{-1} \underline{y} & -\underline{y}' \mathbf{V}^{-1} \mathbf{X} \\ -\mathbf{X}' \mathbf{V}^{-1} \underline{y} & \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \end{pmatrix} \begin{pmatrix} 1 \\ \underline{\beta} \end{pmatrix}$ or maximise $\underline{\beta}' \mathbf{X}' \mathbf{V}^{-1} \mathbf{X} \underline{\beta}$ Rao (1973). Estimator (2.6) satisfies this.

IV. The Main Results

The dispersion matrix corresponding to an individual/unit of the model takes

$$\text{the form } \mathbf{V} = \sigma^2 \begin{pmatrix} 1 & \rho & \rho & \dots & \rho \\ \rho & 1 & \rho & \dots & \rho \\ \dots & \dots & \dots & \dots & \dots \\ \rho & \rho & \rho & \dots & 1 \end{pmatrix} = \sigma^2 \{ (1 - \rho) \mathbf{I} + \rho \mathbf{J} \} \text{ where } \sigma^2 = \sigma_\mu^2 +$$

σ_ν^2 and \mathbf{J} be the square matrix of 1's of order n again $\sigma^2 \{ n\rho \bar{\mathbf{J}}_n + (1 + (n-1)\rho) \mathbf{E}_n \}$ where $\bar{\mathbf{J}}_n = ((1/n))$ and $\mathbf{E}_n = \mathbf{I}_n - \bar{\mathbf{J}}_n$ are idempotent matrices and hence we got the result

Result 4.1 (Baltagi, 2005). For a matrix $\mathbf{V} = k\mathbf{A} + m\mathbf{B}$ where \mathbf{A} and \mathbf{B} are idempotent matrices and k, m are scalars such that $\mathbf{A} + \mathbf{B} = \mathbf{I}$ and $\mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A} = \mathbf{0}$, $\mathbf{V}^{-1} = \frac{1}{k}\mathbf{A} + \frac{1}{m}\mathbf{B}$ and in fact $\mathbf{V}^q = k^q\mathbf{A} + m^q\mathbf{B}$ for any number q .

Proof. The result can be obtained by direct verification. □

Let $z_j = (y_j - na\bar{y}) = \mathbf{a}'_j \underline{y}$ and $a = \frac{\rho}{1+(n-1)\rho}$, $\mathbf{a}'_j = (-a, -a, \dots, 1 - a, \dots, -a)'$ with j^{th} component is $1 - a$.

Result 4.2. When j^{th} is deleted from the data set, $DFBETA_j = \frac{(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{a}_j e_j}{\sigma^2(1-\rho)(1-a)-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j}$; $e_j = y_j - \tilde{x}'_j \hat{\beta}$; j^{th} residual.

and $DFFIT_j = \frac{\mathbf{c}'_j\mathbf{H}\mathbf{a}_j e_j}{\sigma^2(1-\rho)(1-a)-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j}$; \mathbf{c}_j is the j^{th} unit vector.

Proof.

$$\begin{aligned} \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} &= \frac{1}{\sigma^2(1-\rho)} \left\{ \mathbf{X}'\mathbf{X} - \frac{\rho}{1+(n-1)\rho} \mathbf{X}' \mathbf{1}\mathbf{1}' \mathbf{X} \right\} \\ &= \frac{1}{\sigma^2(1-\rho)} \left\{ \mathbf{X}'_j\mathbf{X}_j - \frac{\rho}{1+(n-1)\rho} \mathbf{X}'_j \mathbf{1}_j \mathbf{1}'_j \mathbf{X}_j \right\} \\ &\quad + \frac{1}{\sigma^2(1-\rho)} \left\{ \tilde{x}_j \tilde{x}'_j - \frac{\rho}{1+(n-1)\rho} (\mathbf{X}'_j \mathbf{1}_j \tilde{x}'_j + \tilde{x}_j \mathbf{1}'_j \mathbf{X}_j + \tilde{x}_j \tilde{x}'_j) \right\} \\ &= \frac{1}{\sigma^2(1-\rho)} \left\{ \mathbf{X}'_j\mathbf{X}_j - \frac{\rho}{1+(n-2)\rho} \mathbf{X}'_j \mathbf{1}_j \mathbf{1}'_j \mathbf{X}_j \right\} \\ &\quad + \frac{1}{\sigma^2(1-\rho)} \left\{ \frac{\rho^2}{(1+(n-2)\rho)(1+(n-1)\rho)} \mathbf{X}'_j \mathbf{1}_j \mathbf{1}'_j \mathbf{X}_j - \right. \\ &\quad \left. \frac{\rho}{1+(n-1)\rho} (\mathbf{X}'_j \mathbf{1}_j \tilde{x}'_j + \tilde{x}_j \mathbf{1}'_j \mathbf{X}_j) + \frac{1+(n-2)\rho}{1+(n-1)\rho} \tilde{x}_j \tilde{x}'_j \right\} \\ &= \mathbf{X}'_j\mathbf{V}_j^{-1}\mathbf{X}_j + \frac{1}{\sigma^2(1-\rho)} \left\{ \left(\frac{1+(n-1)\rho}{1+(n-2)\rho} \tilde{x}_j - \frac{n\rho}{1+(n-2)\rho} \bar{\mathbf{X}} \right) \left(\tilde{x}_j - \frac{n\rho}{1+(n-1)\rho} \bar{\mathbf{X}} \right)' \right\} \end{aligned}$$

Using the identity $1 = 1 + \frac{\rho}{1+(n-2)\rho} - \frac{\rho}{1+(n-2)\rho}$. Hence we get

$$\mathbf{X}'_j\mathbf{V}_j^{-1}\mathbf{X}_j = \mathbf{X}'\mathbf{V}^{-1}\mathbf{X} - \frac{1}{\sigma^2(1-\rho)(1-a)} \mathbf{w}_j \mathbf{w}'_j$$

Similarly $\mathbf{X}'_j\mathbf{V}_j^{-1} \tilde{y}_j = \mathbf{X}'\mathbf{V}^{-1} \tilde{y} - \frac{1}{\sigma^2(1-\rho)(1-a)} \mathbf{w}_j z_j$

Where $\mathbf{w}_j = \tilde{x}_j - na\bar{\mathbf{X}} = \mathbf{a}'_j\mathbf{X}$,

Therefore, $DFBETA_j = \frac{(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{a}_j e_j}{\sigma^2(1-\rho)(1-a)-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j}$; $e_j = y_j - \tilde{x}'_j \hat{\beta}$; j^{th} residual.

and $DFFIT_j = \frac{\mathbf{c}'_j\mathbf{H}\mathbf{a}_j e_j}{\sigma^2(1-\rho)(1-a)-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j}$; \mathbf{c}_j is the j^{th} unit vector.

If we consider the deletion of j^{th} case that may arise in looking for an outlier. The corresponding results of $DFBETA$, $DFFIT$. It is more appropriate and statistically justified to study standardised version of the above measures. variances of the above measures given in following theorem.

Result 4.3.

$$Var(DFBETA_j) = \frac{\text{When the } j^{th} \text{ observation is deleted}}{(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{a}_j\mathbf{a}'_j\mathbf{X}(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}} \frac{1}{(1-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j)}$$

$$Var(DFFIT_j) = \frac{(\mathbf{a}'_j\mathbf{H}\mathbf{a}_j)^2}{(1-\mathbf{a}'_j\mathbf{H}\mathbf{a}_j)}$$

$$\begin{aligned} \text{Var}(DFBETA_j) &= \text{Var}(\hat{\beta} - \hat{\beta}_j) = \text{Var}(\hat{\beta}) + \text{Var}(\hat{\beta}_j) - 2\text{Cov}(\hat{\beta}, \hat{\beta}_j) \\ &= \sigma^2\{(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} + (\mathbf{X}'_j\mathbf{V}_j^{-1}\mathbf{X}_j)^{-1} - 2(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}^{-1}\text{Cov}(y, y'_j)\mathbf{V}_j^{-1}\mathbf{X}_j(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\} \end{aligned}$$

Where $\text{Cov}(y, y'_j) = \mathbf{V}^*$ obtained from \mathbf{V} by deleting j^{th} column. $\mathbf{X}'\mathbf{V}^{-1}\mathbf{V}^* = \mathbf{X}'\mathbf{E}_j = \mathbf{X}'_j$. Hence the result for $\text{Var}(DFBETA_j)$, and similarly for $\text{Var}(DFFIT_j)$.

V. NUMERICAL RESULT

Let us consider the data set given in the following table whose first two columns are (x_1 & x_2) are independent variable and third column y the dependent variable. For different values of ρ standardised values of $DFFIT$ are calculated and shown in last three columns. It is observed that corresponding to seventh observation $DFFIT$ values are high become more prominent as the model approaches correct specification as with values of ρ .

Table 2: Calculated DFFITS values

x_1	x_2	y	$\rho = 0$	$\rho = 0.3$	$\rho = 0.5$
35.5	36.5	30.6	-0.6893	-0.9005	-0.9188
35.3	37.9	30.4	-0.0228	-0.0322	-0.0368
36.4	38.6	37.6	0.1393	0.1811	0.1953
37	38.4	39.8	-0.1344	-0.1675	-0.1761
37.7	38.2	40.4	-0.5861	-0.7141	-0.7425
37.6	38.4	42.2	-0.3140	-0.3764	-0.3903
37.8	37.2	103	3.5744	4.2480	4.4129
36.7	39.1	39.2	0.2855	0.3379	0.3530
38.3	38.4	44.1	-0.5948	-0.7064	-0.7440
37.9	38.8	39.6	-0.5279	-0.6315	-0.6729
39.5	38.7	48.7	-0.8531	-1.0334	-1.1181
38.5	38.4	41.9	-0.8877	-1.0925	-1.2061
40.3	40.2	52.7	-0.4186	-0.5263	-0.5966
40.5	41.6	53.8	0.1256	0.1621	0.1904
39.4	42.8	42.3	0.4258	0.5673	0.6988
40.1	40.7	53.6	-0.0215	-0.0297	-0.0392
40.7	41.5	53.3	-0.0760	-0.1102	-0.1619
39.3	41	50	0.3221	0.4940	0.8830
41.2	42.7	55	0.2538	0.4180	1.3201

Further Study

As there is a concept of association between the observations specified by ρ it would be better to have an estimator of ρ for unclustered observations.

REFERENCES

- [1]. C.M. Beach, J.G. MacKinnon, A maximum likelihood procedure for regression with auto correlated errors, *Econometrica*, 46(1978) 51-58.
- [2]. D.A. Belsley, E. Kuh, R.E. Welsch, *Regression Diagnostics*, John Wiley, New York(1980) .
- [3]. A.G. Barnett, N. Koper, A.J. Dobson, F. Schmiegelow, M. Maseau, Using Information Criteria to Select the Correct Variance-Covariance Structure for longitudinal data in ecology, *Methods in Ecology and Evolution*, 1 (1)(2010) 15-24.
- [4]. S. Chatterjee, A.S. Hadi, *Sensitivity Analysis in Linear Regression*, John Wiley, New York(1988) .
- [5]. D. Cochran, G.H. Orcutt, Applications of Least Square Regressions to Relationships containing Auto correlated Error Terms, *J. Amer. Statist. Soc.*, 44(1949) 32-61.
- [6]. R.D. Cook, Assessment of Local Analysis, *J. Roy. Statist. Soc. Ser. B*, 48(1986) 133-169.
- [7]. G.G. Judge, R.C. Hill, W.E. Griffiths, H. Lutkepohl, T.C. Lee, *Introduction to the Theory and Practice of Econometrics*, John Wiley, New York(1988) .
- [8]. S.W. Kim, R. Huggins, Diagnostics for Auto correlated Regression Models, *Austral. & New Zealand J. Statist.*, 40 (1)(1998) 65-71.
- [9]. J.W. Neter, W. Wasserman, M.H. Kutner, *Applied Linear Statistical Models*, Homewood, IL: Irwin(1985) . S.N. Guria 13
- [10]. S.J. Prais, C.B. Winsten *Trend Estimators and Serial Correlation*, Cowles Commission Discussion Paper No. 383, Chicago,(1954) .
- [11]. M.L. Putterman, Leverage and Influence in auto correlated regression models, *Appl. Statist.*, (1988) .
- [12]. C.R. Rao, *Linear Statistical Inference and its Application*, John Wiley, New York(1973) .
- [13]. R. Schall, T.T. Dunne, Diagnostics for Regression-ARMA Time Series, *Directions in Robust Statistics & Diagnostics*, ed. W. Stahel & S. Weis-berg, part 2(1991) 205-221.
- [14]. Sen, A., Srivastava, M., *Regression Analysis: Theory, Methods, and Applications*, Springer-Verlag, (1990) .
- [15]. S. Sen Roy, S.N. Guria, Estimation of Regression Parameters in the Presence of Outliers in the Response, *Statistics*, 43,6(2009) 531-539.
- [16]. C.L. Tsai, X. Wu, Assessing local influence in linear regression models with 1st-order autoregressive or heteroscedastic error structure, *Statist. & Probab. Letters*, 14(1992) 247-252.