

Residual analysis :

Two the error analysis for the following data

①	x	1	2	2.5	4	5	5.5	6	7	7.5
	y	2	4.5	3	6.5	6	10	9	11	5

calculate R^2 and $\sum e_i^2$ for the given data and also the same without considering the last value of x and y.

Sol: (with value)

x	y	$(x-\bar{x})$	$(x-\bar{x})^2$	$(y-\bar{y})$	$(y-\bar{y})^2$	$(x-\bar{x})(y-\bar{y})$
1	2	-3.5	12.25	-4.33	18.7489	15.155
2	4.5	-2.5	6.25	-1.83	3.3489	4.575
2.5	3	-2	4	-3.33	11.0889	6.66
4	6.5	-0.5	0.25	0.17	0.0289	-0.085
5	6	0.5	0.25	-0.33	0.1089	-0.165
5.5	10	1	1	3.67	13.4689	3.67
6	9	1.5	2.25	2.67	7.1289	4.005
7	11	2.5	6.25	4.67	21.8089	11.675
7.5	5	3	9	-1.33	1.7689	-3.99
<u>40.5</u>	<u>57</u>	<u>41.5</u>			<u>77.5001</u>	<u>41.5</u>

$$\bar{x} = \frac{\sum x_i}{n} = \frac{40.5}{9} = 4.5$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{57}{9} = 6.333$$

$$r = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x_i-\bar{x})^2} \sqrt{\sum (y_i-\bar{y})^2}} = \frac{41.5}{\sqrt{41.5} \sqrt{77.5001}} = \frac{41.5}{56.71202}$$

$$r^2 = 0.73176$$

$$r^2 = 0.53548$$

(without last value)

$$x = 33, y = 52$$

②

$$\bar{x} = \frac{\sum x_i}{n} = \frac{83}{8} = 10.375$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{52}{8} = 6.5$$

x	y	$(x - \bar{x})$	$(x - \bar{x})^2$	$(y - \bar{y})$	$(y - \bar{y})^2$	$(x - \bar{x})(y - \bar{y})$
1	2	-3.375	11.375	-4.5	20.25	14.8125
2	4.5	-2.375	5.556	-2	4	4.25
2.5	3	-1.875	3.516	-3.5	12.25	5.6875
4	6.8	-0.375	0.145	0	0	0
5	6	0.875	0.765	-0.5	0.25	-0.4375
5.5	10	1.875	3.516	3.5	12.25	4.8125
6	9	1.875	3.516	2.5	6.25	4.6875
7	11	2.875	8.265	4.5	20.25	12.9375
			31.3748		75.5	4.6

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{46}{\sqrt{31.3748} \sqrt{75.5}} = \frac{46}{48.67026}$$

$$r = 0.9451$$

$$r^2 = 0.89328$$

calculate $\sum e_i^2$:

$$x \quad y \quad x^2 \quad y^2 \quad xy$$

$$1 \quad 2 \quad 1 \quad 4 \quad 2$$

$$2 \quad 4.5 \quad 4 \quad 20.25 \quad 9$$

$$2.5 \quad 3 \quad 6.25 \quad 9 \quad 7.5$$

$$4 \quad 6.5 \quad 16 \quad 42.25 \quad 26$$

$$5 \quad 6 \quad 25 \quad 36 \quad 30$$

$$5.5 \quad 10 \quad 30.25 \quad 100 \quad 55$$

$$6 \quad 9 \quad 36 \quad 81 \quad 54$$

$$7 \quad 11 \quad 49 \quad 121 \quad 77$$

$$7.5 \quad 5 \quad 56.25 \quad 25 \quad 37.5$$

$$40.5 \quad 57 \quad 223.75 \quad 438.5 \quad 298$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n y_i x_i - (\sum y_i)(\sum_{i=1}^n x_i)/n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2/n}$$

$$= \frac{298 - (57)(40.5)/9}{223.75 - \frac{(40.5)^2}{9}}$$

$$= \frac{298 - 256.5}{223.75 - 182.25} = \frac{41.5}{41.5}$$

$$\boxed{\hat{\beta}_1 = 1}$$

$$x = 4.5, y = 6.33$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$= 6.33 - 1(4.5) = 6.33 - 4.5$$

$$\hat{\beta}_0 = 1.83$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

$$\hat{y}_1 = 1.83 + 1(1) = 2.83$$

$$\hat{y}_2 = 1.83 + 1(2) = 3.83$$

$$\hat{y}_3 = 1.83 + 1(2.5) = 4.33$$

$$\hat{y}_4 = 1.83 + 1(4) = 5.83$$

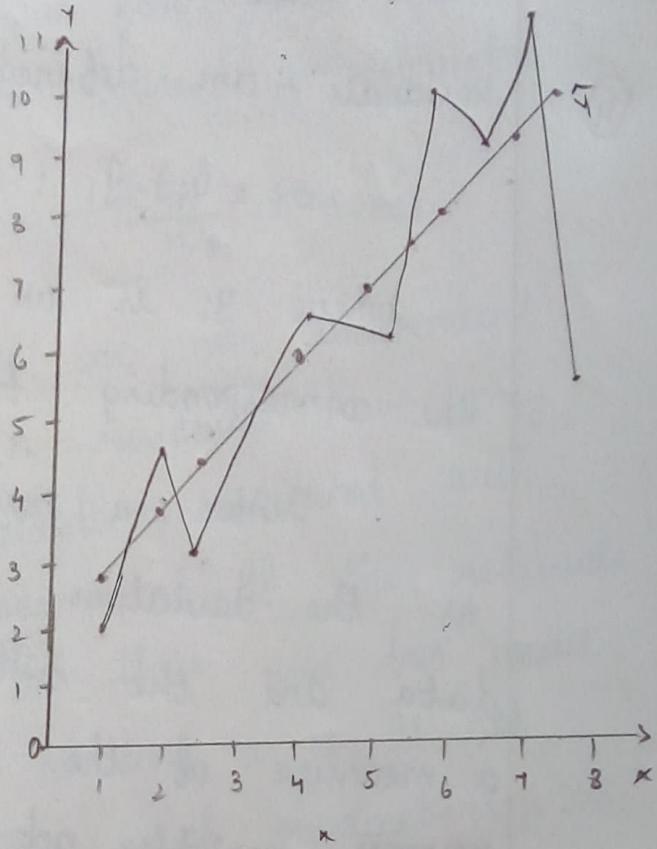
$$\hat{y}_5 = 1.83 + 1(5) = 6.83$$

$$\hat{y}_6 = 1.83 + 1(5.5) = 7.33$$

$$\hat{y}_7 = 1.83 + 1(6) = 7.83$$

$$\hat{y}_8 = 1.83 + 1(7) = 8.83$$

$$\hat{y}_9 = 1.83 + 1(7.5) = 9.33$$



Error:

$$e_1 = y_1 - \hat{y}_1 = 2 - 2.83 = -0.83 \quad e_1^2 = 0.6889$$

$$e_2 = y_2 - \hat{y}_2 = 4.5 - 3.83 = 0.67 \quad 0.4489$$

$$e_3 = y_3 - \hat{y}_3 = 3 - 4.33 = -1.33 \quad 1.7689$$

$$e_4 = y_4 - \hat{y}_4 = 6.5 - 5.83 = 0.67 \quad 0.4489$$

$$e_5 = y_5 - \hat{y}_5 = 6 - 6.83 = -0.83 \quad 0.6889$$

$$e_6 = y_6 - \hat{y}_6 = 10 - 7.33 = 2.67 \quad 7.1289$$

$$e_7 = y_7 - \hat{y}_7 = 9 - 7.83 = 1.17 \quad 1.3689$$

$$e_8 = y_8 - \hat{y}_8 = 11 - 8.83 = 2.17 \quad 4.7089$$

$$e_9 = y_9 - \hat{y}_9 = 5 - 9.33 = -4.33 \quad 18.7489$$

$$\overline{0.03}$$

$$\overline{56.0001}$$

Definition of Residual:

We have already known that the residuals are defined as,

$$e_i = y_i - \hat{y}_i ; i=1, 2, \dots, n \rightarrow ①$$

where y_i is an observation and \hat{y}_i is the corresponding fitted value.

Since, a residual may be viewed as the deviation between the observed data and the fitted data. It is also a measure of the variability in the response variable not explained by the regression model. It is also convenient to think of the residuals as the observed values of the assumptions on the errors should show up in the residuals.

Analysis of residuals is an effective way to discover several types of model inadequacies. Plotting residuals is a very effective way to investigate how well the regression model fits the data and also to check the assumptions of regression model.

The residuals have several important properties. They have zero mean and their approximate average variance is estimated

(5)

$$\text{by, } \frac{\sum (E_i - \bar{E})^2}{n-p} = \frac{\sum E_i^2}{n-p} = MS_{\text{Res}}$$

The residuals are not independent however, as the n residuals have only $n-p$ degrees of freedom associated with them. This non independence of the residuals has little effect on their use for model adequacy checking as long as n is not small relative to the no. of parameters p .

Properties of Residuals:

we can consider the properties we have seen in the topics of properties of least square fit.

Methods of scaling residuals:

There are four popular methods for scaling residuals.

These scaled residuals are helpful in finding observations that are outliers or extreme values, that is the observations that are separated in some fashion from the rest of the data. The methods are

⑥

(ii) standardised residuals.

(iii) studentized residuals.

(iv) Press residuals (prediction error sum of squares).

(v) R - student residuals.

(ii) **standardized residuals:**

since the approximate average variance of a residual is estimated by MS_{RES} , a logical scaling for the residuals would be the standardized residuals.

$$d_i = \frac{R_i}{\sqrt{MS_{RES}}}, i=1, 2, \dots, n$$

The standardized residuals have mean zero and approximately unit variance consequently a large standardized residuals ($d_i > 5$, say) potentially indicates an outlier.

(ii) **studentized Residuals:**

using MS_{RES} as the variance of the i th residual e_i is only an approximation. we can improve the residual scaling by dividing e_i by the exact standard deviation of the i th residuals.

Recall from Eq. that we may write the vector of residuals as

$$e = (I - H) y$$

where $H = X(X'X)^{-1}X'$ is the hat matrix.
The hat matrix has several useful properties.
It is symmetric ($H' = H$) and idempotent
($HH = H$). Similarly the matrix $I - H$ is
symmetric and idempotent. Substituting $y = X\beta + e$
into yields.

$$e = (I - H)(X\beta + e) = X\beta - HX\beta + (I - H)e$$

$$= X\beta - X(X'X)^{-1}X'X\beta + (I - H)e = (I - H)e$$

Thus, the residuals are the same
linear transformation of the observations
 y and the errors e . The covariance matrix
of the residual is

$$\text{var}(e) = \text{var}[(I - H)e] = (I - H) \text{var}(e) (I - H)' = \sigma^2(I - H)$$

since $\text{var}(e) = \sigma^2 I$ and $I - H$ is symmetric
and idempotent. The matrix $I - H$ is
generally not diagonal, so the residuals have
different variances and they are correlated.

The variance of the i^{th} residual is

$$\text{var}(e_i) = \sigma^2(1 - h_{ii})$$

where h_{ii} is the i^{th} diagonal element
of the hat matrix H . The covariance between
residuals e_i and e_j is

$$\text{cov}(e_i, e_j) = -\sigma^2 h_{ij}$$

where h_{ij} is the j^{th} element of
the hat matrix. Now since $0 \leq h_{ii} \leq 1$,

(8) using the residual mean square MS_{Res} to estimate the variance of the residuals actually overestimates $\text{var}(e_i)$. Furthermore, since h_{ii} is a measure of the location of the e_i depends on where the point x_i lies. Generally points near the centre of the x -space have larger variance than residuals at more remote locations.

Violation of model assumptions are more likely at remote points, and these violations may be hard to detect from inspection of the ordinary residuals e_i because their residuals will usually be smaller.

A logical procedure, then, is to examine the studentized residuals

$$r_i = \frac{e_i}{\sqrt{MS_{Res}(1-h_{ii})}} ; i=1, 2 \dots n$$

instead of e_i ($\text{or } d_i$). The studentized residuals have constant variance $\text{var}(r_i) = 1$ regardless of the location of x_i when the form of the model is correct. In many situations the variance of the residuals stabilizes, particularly for large data sets. In these case there may be little difference between the standardized and studentized residuals. Thus, standardized and studentized residuals often convey equivalent

information. However, since any point with a large residuals and a large $|h_{ii}|$ is potentially highly influential on the least-squares fit, examination of the studentized residuals is generally recommended.

Some of these points are very easy to see by examining the studentized residuals for a simple linear regression. It is easy to show that the studentized residuals are

$$r_i = \frac{e_i}{\sqrt{MS_{Res} \left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right) \right]}} \quad i = 1, 2, \dots, n$$

Notice that when the observation x_i is close to the midpoint of the x -data, $x_i - \bar{x}$ will be small, and the estimated standard deviation of e_i [The denominator of eqn 4.9] will be large. Conversely, when x_i is near the extremends of the range of the x -data, $x_i - \bar{x}$ will be large, and the estimated standard deviation of e_i will be small. Also, when the sample size n is really large, the so in big data sets, studentized residuals may not differ dramatically from standardized residuals.

Press Residuals :

(10) The standardized and studentized residuals are effective in detecting outliers. Another approach to making residuals useful in finding outliers is to examine the quantity that is computed from $y_i - \hat{y}_{(i)}$, where $\hat{y}_{(i)}$ is the fitted value of the i th response based on all observations except the i th one. The logic behind this is that if the i th observation y_i is really unusual, the regression model based on all observations may be overly influenced by this observation. This could produce a fitted value \hat{y}_i that is very similar to the observed value y_i , and consequently, the ordinary residuals e_i will be small. Therefore, it will be hard to detect the outlier. However, if the i th observation is deleted, then $\hat{y}_{(i)}$ cannot be influenced by that observation, so the resulting residuals should be likely to indicate the presence of the outlier.

If we delete the i th observation, fit the regression model to the remaining $n-1$ observations, and calculate the predicted value of y_i corresponding to the deleted observation, the corresponding

Prediction error is

$$e_{(i)} = y_i - \hat{y}_{(i)}$$

(ii)

This prediction error calculation is repeated for each observation $i=1, 2, \dots, n$. These prediction errors are usually called PRESS residuals. (because of their use in computing the prediction error sum of squares, discussed) some authors call the $e_{(i)}$ deleted residuals.

It would initially seem that calculating the PRESS residuals requires fitting n different regression. However, it is possible to calculate PRESS residuals from the results of a single least squares fit to all n observations we shown. How this is accomplished. It turns out that the i th PRESS residual is

$$e_i = \frac{e_i}{1-h_{ii}} ; i=1, 2, \dots, n$$

It is easy to see that the PRESS residuals is just the ordinary residual weighted according to the diagonal elements of the hat matrix h_{ii} . Residuals associated with points for which h_{ii} is large will have large PRESS residuals. These points will generally be high influence points. Generally a large difference between the ordinary residual and the

(12)

PRESS residuals will indicate a point where the model fits the data well, but a model built without that point predicts poorly.

Finally, the variance of the i^{th} PRESS residual is

$$\text{var}[e_{ii}] = \text{var}\left[\frac{e_i}{1-h_{ii}}\right] = \frac{1}{(1-h_{ii})^2} [\sigma^2(1-h_{ii})]$$

$$= \sigma^2 / (1-h_{ii})$$

so that a standardized PRESS residual is

$$\frac{e_{ii}}{\sqrt{\text{var}(e_{ii})}} = \frac{e / (1-h_{ii})}{\sqrt{\sigma^2(1-h_{ii})}} = \frac{e}{\sqrt{\sigma^2(1-h_{ii})}}$$

which, if we use MSRes to estimate σ^2 is just the studentized residuals discussed previously.

R-Student:

The studentized residual r_i discussed above is often considered an outlier diagnostic. It is customary to use Mres as our estimate of σ^2 in computing r_i . This is referred to as internal scaling of the residual because MSRes is an internally generated estimate of σ^2 obtained from fitting the model to all n observations. Another approach would be to use an estimate of σ^2 based on a

(13)

data set with the i^{th} observation removed. Denote the estimate of σ^2 so obtained by s_{res}^2 . We can show that

$$s_i^{(2)} = \frac{(n-p) \text{MSRes} - e_i^2 (1-h_{ii})}{n-p-1}$$

The estimate of σ^2 is used instead of MSRes to produce an externally studentized residual, usually called R-student given by

$$t_i = \frac{e_i}{\sqrt{s_{\text{res}}^2 (1-h_{ii})}} ; i = 1, 2, \dots, n$$

In many situations t_i will differ little from the studentized residual r_i . However, if the i^{th} observation is influential, then s_{res}^2 can differ significantly from MSRes , and thus the R-student statistic will be more sensitive to this point.

It turns out that under the usual regression assumption t_i will follow the t_{n-p-1} distribution. Appendix q establishes a formal hypothesis-testing procedure for outlier detection based on R-student. One could use a Bonjean-von Mises type approach and compare all n values of $|t_i|$ to $t_{(\alpha/2n)-n-p-1}$ to provide guidance regarding outliers. However it is our view that a formal approach is usually not necessary and that

(14) only relatively guide cutoff values need be considered. In general, a diagnostic view as opposed to a strict statistical hypothesis - testing view is best. Furthermore, detection of outliers often needs to be considered simultaneously with detection of influential observations.

Residual Plots:

Graphical analysis of residuals is a very effective way to investigate the adequacy of the fit of a regression model and to check the underlying assumption. Here we introduce and illustrate the basic residual plots. They are namely.

i) Normal probability plot.

ii) Plot of residuals against the fitted values;

iii) Plot of residuals against the regressor.

iv) Plot of residuals in time sequence.

These plots are difficultly generated by regression computer software packages. They should be examined routinely in all regression modeling problems. It is often a good idea to plot both the original residuals and one or more of the scaled residuals. We mainly plot studentized residuals because they have

constant variance.

(i) Normal probability plot:

Small departure from the normality assumption do not affect the model greatly, but gross non normality is potentially more serious as the t or F statistic and confidence and prediction intervals depend on the normality assumption. Further, more if the error is come from a distribution with heavier tails then the normal least square fit may be sensitive to a small subset of the data heavy tailed error distributions often generate outliers that pull the least square fit too much in their direction in this cases a very simple method of checking the normality assumption is to construct a normal probability plot of the residuals. This a graph designed to test the cumulative normal distribution that the cumulative normal distribution will plot as a straight line.

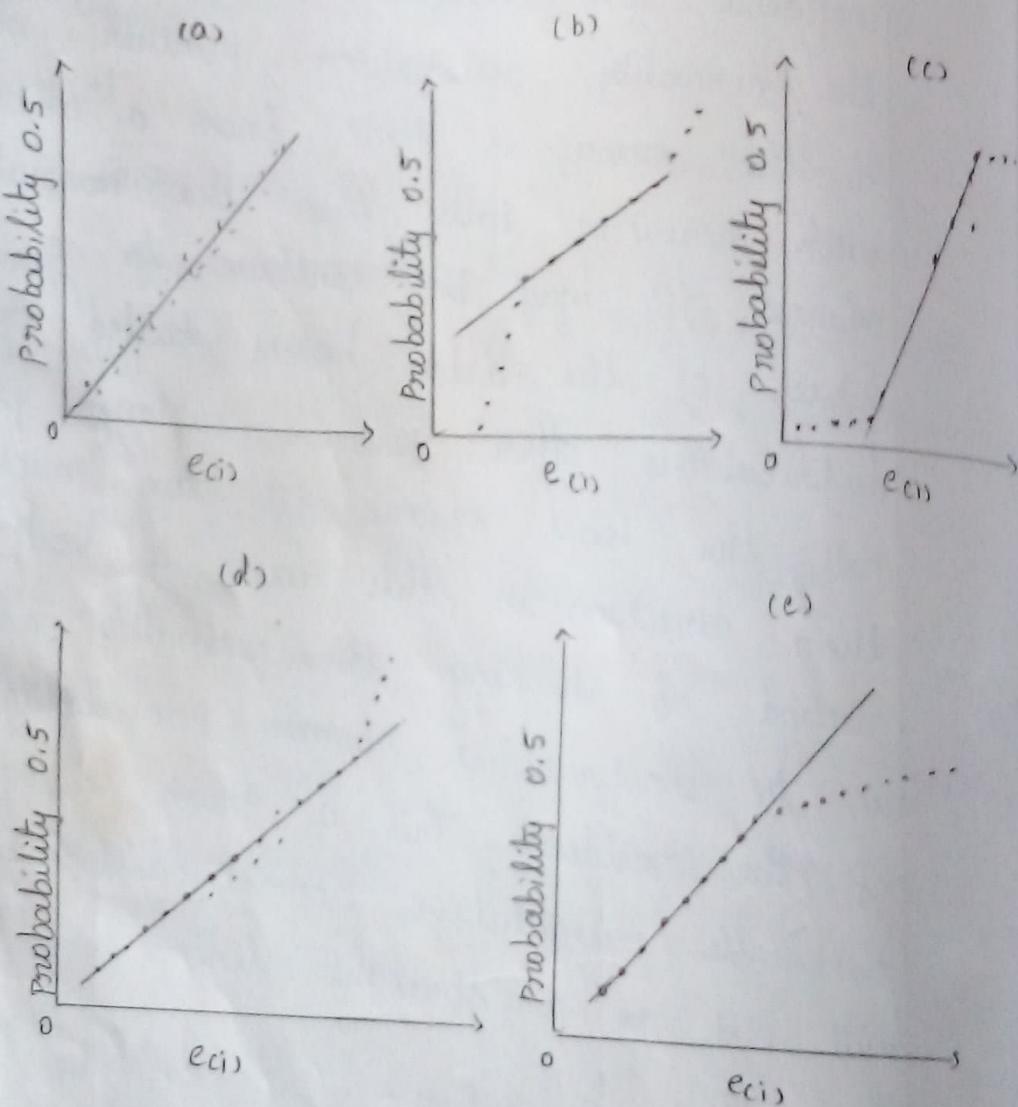
Let $e_{[1]} \leq e_{[2]} \leq \dots \leq e_{[n]}$ be the residuals ranked in the increasing order if we plot $e_{[i]}$ against the cumulative probability

$$P_i = \frac{(i - \frac{1}{2})}{n}, i = 1, 2, \dots, n$$

on the normal probability plot. The resulting points should lie approximately on a

(16)

straight line. The straight line is usually determined with emphasis on the central values rather than the extremes [The 0.33 and 0.67 cumulative probability points]. Substantial departures from a straight line indicates that the distribution is not normal.



Normal Probability plots

- Ideal
- heavy tailed distribution.
- Light - tailed distribution.
- positive skew.
- negative skew.

(17)

Fig (a) displays an idealised normal probability plot notice that the points lie approximately along a straight line. Fig (b-e) present other difficult problems. Fig (b) shows a sharp upward and downward curve at both extremes, indicating that the tails of these distinguish on two heavy for each to be consider normal. Fig (c) shows flattening at the extremes, which is a pattern difficult of samples from a distribution with thinner tails than the normal. Fig (d and e) exhibits patterns associated with positive and negative skew respectively.

Problem:

- 1) construct a normal probability plot for the example we have discussed in the earlier in residual analysis.

ascending order : $P\left(\frac{i-1/2}{n}\right)$

$$e_1 = -4.333 \quad (1 - 1/2)/7 = 0.07142$$

$$e_2 = -1.333 \quad (2 - 1/2)/7 = 0.21428$$

$$e_3 = -0.833 \quad (3 - 1/2)/7 = 0.3571$$

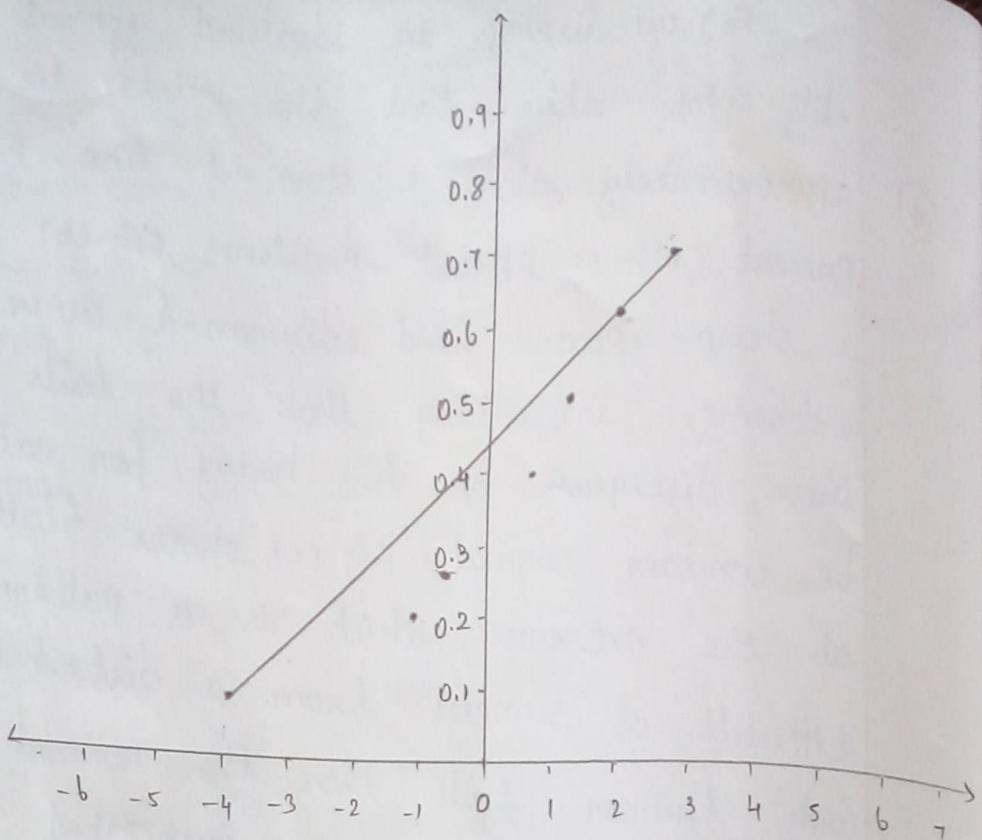
$$e_4 = 0.667 \quad (4 - 1/2)/7 = 0.5$$

$$e_5 = 1.167 \quad (5 - 1/2)/7 = 0.6428$$

$$e_6 = 2.167 \quad (6 - 1/2)/7 = 0.7858$$

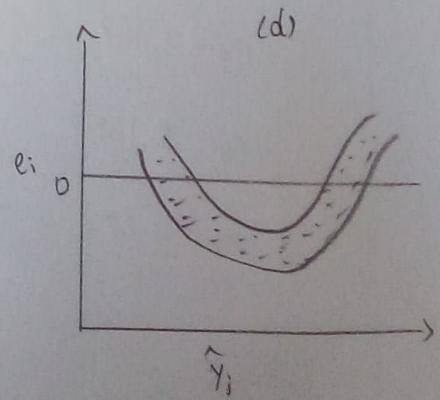
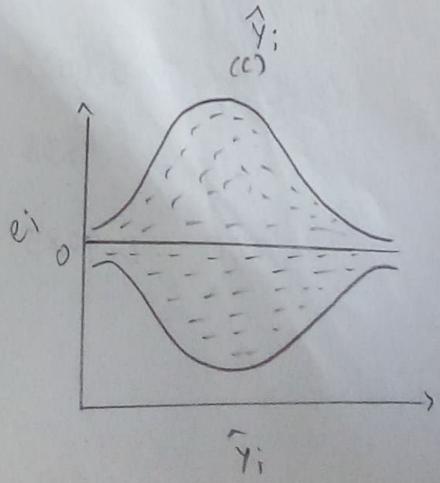
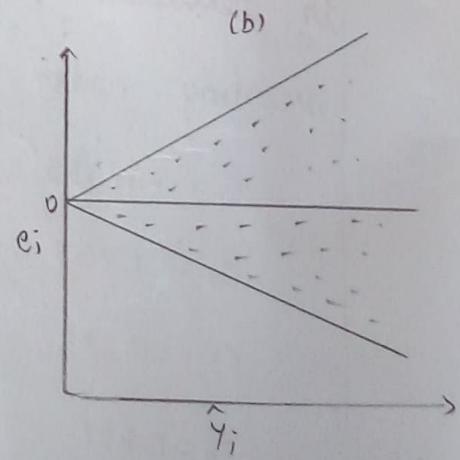
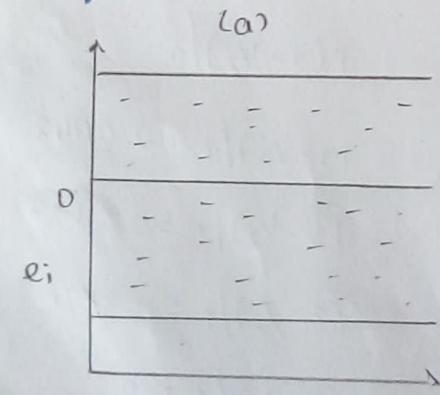
$$e_7 = 2.67 \quad (7 - 1/2)/7 = 0.9285$$

(18)



ii) plot of residuals against the fitted value \hat{y} :

A plot of residuals e_i versus the corresponding fitted values \hat{y}_i is usefull for detecting several common types of model inadequacies.



(19) Fig (a) which indicates that the residuals can be contained in a horizontal line (band), then there are no obvious model defects. plots of e_i versus \hat{y}_i against \hat{y}_i that resembles any of the patterns. In figures (b-d) are symptomatic of model deficiencies.

The patterns in figure (b) and (c) indicate that the variance of the errors is not constant. The outward opening funnel pattern in fig (b) implies that the variance is an increasing function of y [some time an inward opening funnel is also possible].

The double bow pattern in fig (c) often occurs when y is a proportion between 0 and 1. The variance of a binomial proportion near 0.5 is greater than near 0 or 1. The usual approach for dealing with inequality of variance is to apply a suitable transformation in either the regressor variable or the response variable or both. In practice, transformation on the response variable are generally employed to stabilize variance.

A curved plot such as in fig (d) indicates non-linearity. This could mean that other regressor variables are needed in the model. For example, a squared term may be necessary. Transformation on the

regression on the response variable may also be helpful in these cases.

A plot of the residuals against \hat{y}_i may also reveal one or more unusually large residuals. These points are of course, potential outliers. Large residuals that occur at the extreme \hat{y}_i values could also indicate that either the variance is not constant or the true relationship between x and y is not linear. These possibilities should be investigated before the points are considered as outliers.

(iii) plot of residuals against the regression (x_{ij}):

plotting the residuals against corresponding value of each regressor variable (x_{ij}), can also be helpful. For known the inadequacies. These plots often exhibits pattern such as shows in fig (a,b,c,d) in fitted value of \hat{y} . Except that the horizontal scale is x_{ij} for the j th regressor rather than \hat{y}_i . Once again an impression of a horizontal band containing the residuals is desirable. The funnel and double bow patterns in figures (b) and (c) indicate non constant variance. The wavy

(21) band in figure (d) on a non-linear pattern in general implies that the assumed relationship between the y and the regressor x_i is non-constant. Thus either higher order terms in x_i [such as x_i^2] or the transformation should be considered.

In the simple linear regression case it is not necessary to plot residuals versus both \hat{y}_i and the regressor variable x_i . The reason is that the fitted values \hat{y}_i are linear combination of the regressor value x_i , so that plots only would differ in the scale for the abscissa.

(iv) plot of residuals in time sequence:

If the time series in which the data were collected is known. It is a good idea to plot the residual against time order. Ideally, this plot will resemble in figures (a,b,c,d) that is a horizontal band will enclose all of the residuals, and the residuals will fluctuate in a more or less random fashion within this band. However, if this plot resembles the pattern in (b)-(d) this may indicate that the variances is changing with time, or that linear are quadratic terms in time should be added to the model. The

(2)

time sequence plots of residuals may indicate that the error at one time period are correlated with those at other time periods. The correlation between model errors at different time period is called auto correlation. A plot such as figure (a) indicates +ve auto correlation, while fig (b) is typical of negative auto correlation. The presence of auto correlation is a potential series violation of the basic regression assumption. More discussion about methods for detecting auto correlation and remedial measures are discussed.

Date
20.01.20

Diagnostics tests:

When we compute a sample average each observation in the sample has the same weight in determining the outcome. In the regression situation this is not the case. For example, we have considered that the location of observation is x can play unimportant role in determining the regression co-efficient. We have also focused attention (on) outliers or observation that have unusual y value.

On the other hand we sometimes find that a small subset of data exerts a dis-proportionate influence on the model

(23)

co-efficient and properties. In extreme case, the parameter estimates may depend more on the influential subset of points than on the majority of the data. Consequently we would like to find influential points based on the model if these influential points are indeed bad values then they should be eliminated from the samples. Here we

Here we present several diagnostic test for leverage and influence. These diagnostic are available in most multiple regression computer packages. It is important to use the diagnostic in conjunction with residuals analysis techniques.

Sometimes we find that a regression co-efficient may have a sign that does not make engineering (or) scientific sense, a regressor known to be important may be statistically insignificant. This situation may be the result of one (or) few influential observation residuals finding the observations we can use any of the diagnostic test.

~~PRESS~~ PRESS STATISTIC :

(24) we defined the PRESS residuals as $e_{(i)} = y_i - \hat{y}_{(i)}$, where $\hat{y}_{(i)}$ is the predicted value of the i^{th} observed response based on a model fit to the remaining $n-1$ sample points. We noted that large PRESS residuals are potentially useful in identifying observations where the model does not fit the data well, or observations for which the model is likely to provide poor future predictions.

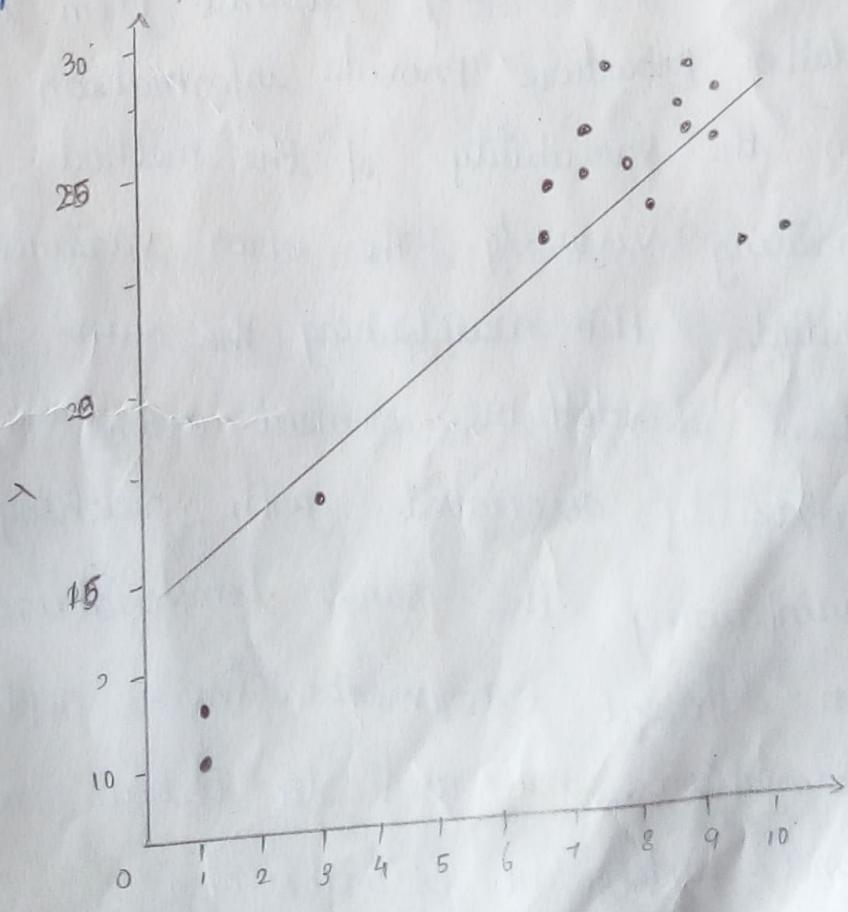
Allen [1971, 1974] has suggested using the prediction error sum of squares (or the PRESS statistic), defined as the sum of the squared PRESS residuals, as a measure of model quality. The PRESS statistic is

$$\begin{aligned} \text{PRESS} &= \sum_{i=1}^n [y_i - \hat{y}_{(i)}]^2 \\ &= \sum_{i=1}^n \left(\frac{e_i}{1-h_{ii}} \right)^2 \end{aligned}$$

PRESS is generally regarded as a measure of how well a regression model will perform in predicting new data. A model with a small value of PRESS is desired. □

A formal test of lack of fit:

(25) we will now present a formal statistical test for the lack of fit of a regression model. The Produce ~~assumptions~~ that the normality independence, and constant variance requirements are met and that only the first-order or straight-line character of the relationship is in doubt. For example, consider the data in figure. There is some indication that the straight-line fit is not very satisfactory, and it would be helpful to have a test procedure to determine if there is systematic curvature present.



(2)

The lack of fit requires that we have replicate observations on the response y for at least one level of x .

We emphasize that there should be true replications, not just duplicate readings or measurements of y .

For example,

Suppose that y is Product viscosity and x is temperature. True replication consists of running n_i separate experiments at $x=x_i$ and observing viscosity, not just running a single experiment at x_i and measuring viscosity n_i times.

The readings obtained from the latter procedure provide information only on the variability of the method of measuring viscosity. The error variance σ^2 includes this maintaining the same temperature level in diff measurement error and the variability associated with reaching and maintaining the same temperature level in different experiments. These replicated observations are used to obtain a model-independent estimate of σ^2 .

(21) Suppose we have n_i observations on the response at the i^{th} level of the regressor $x_i : i = 1, 2, \dots, m$. Let y_{ij} denote the j^{th} observation on the response at x_i , $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, n$. There are $n = \sum_{i=1}^m n_i$ total observations. The test procedure involves partitioning the residual sum of squares in two components.

$$SS_{\text{Res}} = SS_{\text{PE}} + SS_{\text{LOF}}$$

where SS_{PE} is the sum of squares due to pure error and SS_{LOF} is the sum of squares due to lack of fit.

To develop this partitioning of SS_{Res} , note that the $(i, j)^{\text{th}}$ residual is

$$y_{ij} - \hat{y}_i = (y_{ij} - \bar{y}_i) + (\bar{y}_i - \hat{y}_i) \rightarrow ①$$

where \bar{y}_i is the average of the n_i observation at x_i . Squaring both sides of Eq ① and summing over i and j yields

$$\sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \hat{y}_i)^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2 \rightarrow ②$$

since the cross-product term equals zero.

The left hand side of equation ② is usual residual sum of squares.

The two components on the right hand side of the equ (3) measure pure error and lack of fit.

We see that the pure sum of squares

$$SSPE = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 \rightarrow (4)$$

we consider the equation (4) is obtained by computing the corrected sum of square of the repeat observation at each level of x and then pooling over the m levels of x . Since there are n_{i-1} d.f for pure error at each level of x then the total no. of d.f associated with pure error sum of square is

$$\sum_{i=1}^m (n_{i-1}) = n-m \rightarrow (5)$$

The sum of square for lack of fit

$$\sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2 \rightarrow (6)$$

is a weighted sum of squared deviation between the mean response \bar{y}_i at each x level and the corresponding the fitted value. If the fitted values \hat{y}_i are close to the corresponding average response \bar{y} then there is a strong indication that the regression function is linear. There are $m-2$ d.f associated with SS_{LOF} , since there are m levels of x and two d.f

are lost because two parameters estimated must be estimated to obtain the \bar{Y}_i .

(29) computationally we usually obtain SS_{LOF} by subtracting SS_{PE} from SS_{RES} .

The test statistic for lack of fit is

$$F_0 = \frac{(SS_{LOF} / m-2)}{(SS_{PE} / n-m)}$$

$$F_0 = \frac{SS_{LOF} / (m-2)}{SS_{PE} / (n-m)}$$

i.e $MS_{LOF} / MS_{PE} \rightarrow \textcircled{3}$

variance stabilizing transformation:

The assumption of constant variance is a basic requirement of regression analysis. A common reason for a violation of this assumption is for the response variable y to follow a probability distribution in which the variance is functionally related to the mean. For example, if y is a Poisson r.v. in a simple linear regression model then the variance y is equal to the mean. Since the mean of y is related to the regression variable x then the variance of y will be proportional to x . Variance stabilizing transformation are often useful in this case. Thus if the distribution of y is

Poisson we could regress against x since the variance of the square root of a Poisson r.v. variable is independent of mean.

Several commonly used variance stabilizing used are summarised in the following table.

The strength of a transformation depends upon the amount of curvature that indicates

useful & variance stabilizing transformation

Relationship of σ^2 to expectation of y

Transformation

1. $\sigma^2 \propto$ constant.

$y' = y$ (no transformation)

2. $\sigma^2 \propto E(y)$

$y' = \sqrt{y}$ (square root, Poisson data).

3. $\sigma^2 \propto E(y)[\underline{E(y)}]$

$y' = \sin^{-1}(\sqrt{y})$ (arc sin binomial proportions).

$0 \leq y_i \leq 10$

4. $\sigma^2 \propto (E(y))^2$

$y' = \ln(y)$ (log)

5. $\sigma^2 \propto (E(y))^3$

$y' = y^{1/2} = \sqrt{y}$

(Reciprocal square root).

6. $\sigma^2 \propto (E(y))^4$

$y' = y^{-1} = 1/y$ (Reciprocal)

Sometimes we can use prior experience or theoretical consideration to guide

to appropriate transformation. However in many case we have no a priory reason to suspect the error is not constant. Out first indication of the problem is from inspection of the scatter diagram or residual analysis. In these case the appropriate transformation must be empirically.

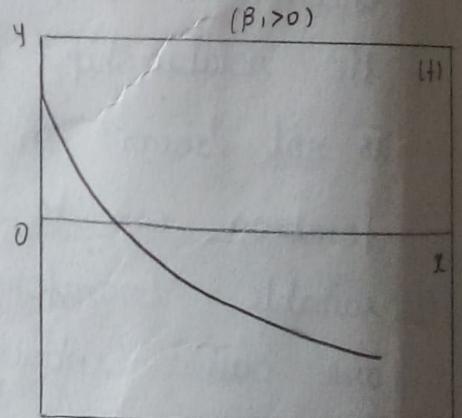
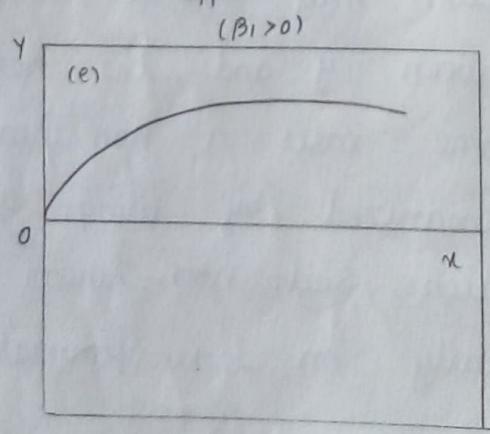
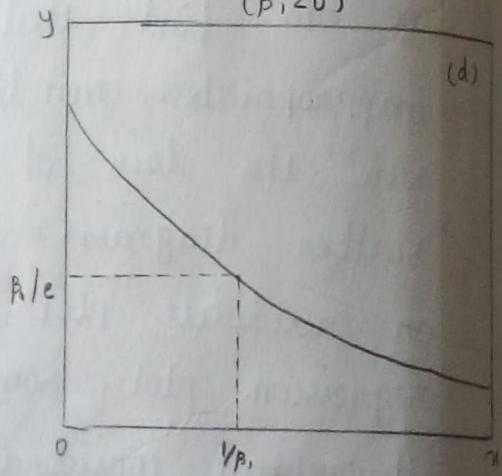
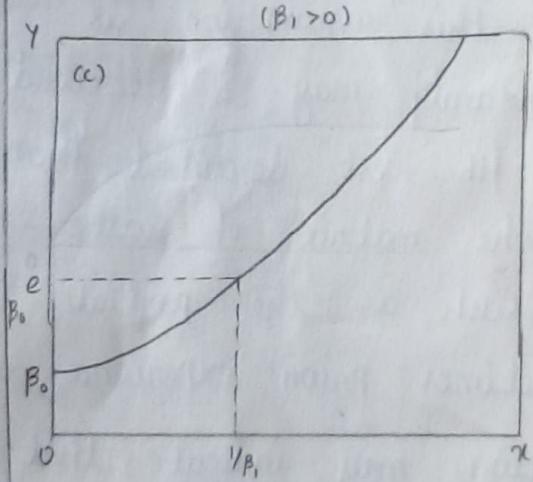
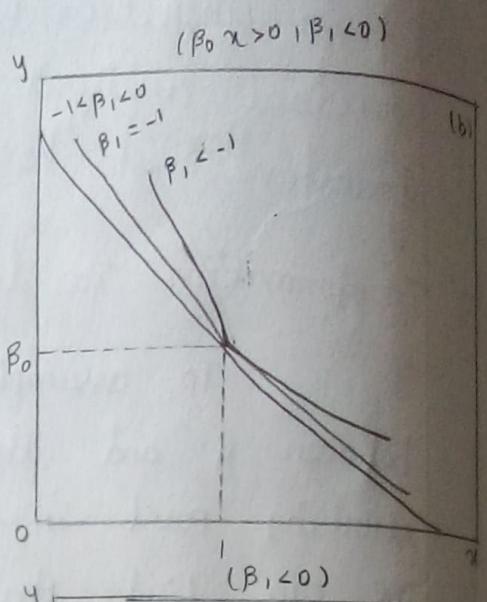
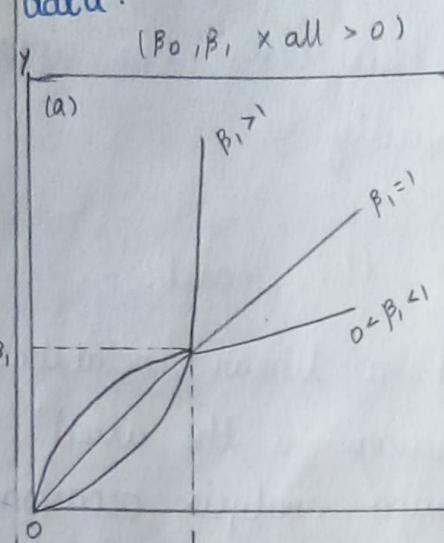
Transformation To linearize the model:

To assumption of a linear relationship between y and the regressors is the usual starting point in regression analysis occasionally we can find that this assumption is inappropriate. Non linearity may be detected via the lack-of-fit test described from scatter diagrams, the matrix of scatter-plots or residuals plot such as the partial regression plot. Sometimes prior experience or theoretical consideration may indicate that the relationship between y and the regressors is not linear. In some cases a non linear function can be linearized by using a suitable transformation. Such non-linear models are called intrinsically or transformably linear.

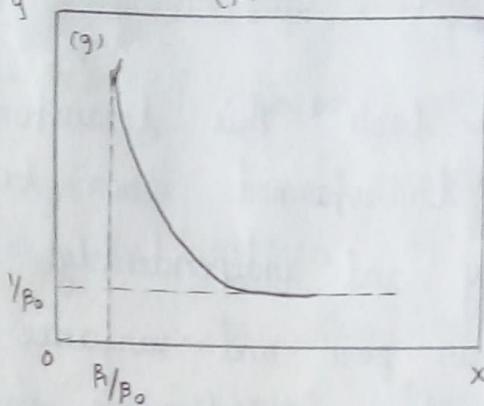
Several linearizable functions are shown in figure. The corresponding non-linear functions, transformations, and resulting linear forms are shown in table. which the scatter

diagram of y against x indicates curvature
 we may be able to match the observed of the
 plot to one of the curves in fig and use the
 linear form of the function to represent the
 data.

(32)

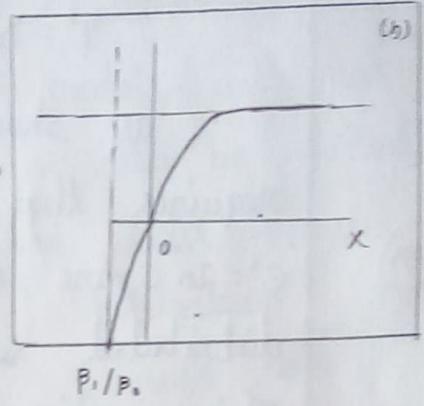


33 (B₁ > 0)



positive curvature

(B₁ > 0)



Negative curvature

Linearizable function (From Daniel and Wood [1980], used with permission of the publishers)

Linearizable function and corresponding Linear form

Figure	Linearizable function	Transformation	Linear form
a,b	$y = P_0 x^{P_1}$	$y' = \log y, x' = \log x$	$y' = \log P_0 + P_1 x'$
c,d	$y = P_0 e^{P_1 x}$	$y' = \ln y, x' = x$	$y' = \ln P_0 + P_1 x'$
e,f	$y = P_0 + P_1 \log x$	$x' = \log x$	$y' = P_0 + P_1 x'$
g,h	$y = \frac{x}{P_0 x - P_1}$	$y' = \frac{1}{y}, x' = 1/x$	$y' = P_0 - P_1 x'$

To illustrate a non-linear model that is intrinsically linear, consider the exponential function.

$$y = P_0 e^{P_1 x} \epsilon$$

This function is intrinsically linear since it can be transformed to a straight line by a logarithmic transformation.

$$\ln y = \ln P_0 + P_1 x + \ln \epsilon$$

$y' = \beta_0 + \beta_1 x + \epsilon'$

(34)

as shown in table. This transformation requires that the transformed error terms $\epsilon' = \ln \epsilon$ are normally and independently distributed with mean zero and variance σ^2 . This implies that the multiplicative error ϵ in the original model is log normally distributed. We should look at the residuals from the transformed model to see if the assumptions are valid. Generally if x and y are in the proper matrix, the usual least squares assumptions are more likely to be satisfied, although it is not unusual to discover at this stage that a non-linear model is preferable.

Various types of reciprocal transformations are also useful. For example, the model.

$$y = \beta_0 + \beta_1 (\frac{1}{x}) + \epsilon$$

can be linearized by using the reciprocal transformation $x' = \frac{1}{x}$. The resulting linearized model is

$$y = \beta_0 + \beta_1 x' + \epsilon$$

Other models that can be linearized by reciprocal transformations are

$$\frac{1}{y} = \beta_0 + \beta_1 x + \epsilon$$

$$y = \frac{x}{\beta_0 x - \beta_1 + \epsilon}$$

Generalized and weighted least squares:

(35) Linear regression models with non constant error variance can also be fitted by the method of weighted least squares

In this method of estimation the deviation between the observed and the expected values of y_i is multiplied by a weight w_i , chosen inversely proportional to the variance of y_i . For the case of simple linear regression, the weighted least squares function is

$$S(\beta_0, \beta_1) = \sum_{i=1}^n w_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \rightarrow ①$$

The resulting least squares normal equation are

$$\hat{\beta}_0 \sum w_i + \hat{\beta}_1 \sum w_i x_i = \sum w_i y_i \quad \rightarrow ②$$

$$\hat{\beta}_0 \sum w_i x_i + \hat{\beta}_1 \sum w_i x_i^2 = \sum w_i x_i y_i$$

Solving the equation ② will produce weighted least square estimates are $\hat{\beta}_0, \hat{\beta}_1$.

Thus we given development of weighted least square for the multiple regression model, i.e., generalized weighted least squares. Here we begin by considering a slightly more general situation concerning the structure of the model errors.

When the error ϵ are uncorrelated but have unequal variances so that the co-variance matrix Σ is

$$\sigma^2 \Sigma = \sigma^2 \begin{bmatrix} 1/w_1 & 0 & 0 & \cdots & 0 \\ 0 & 1/w_2 & 0 & \cdots & 0 \\ \vdots & & & & \\ 0 & 0 & 0 & \cdots & 1/w_n \end{bmatrix}$$

The estimation is usually called weighted least squares. Let $W = \Sigma^{-1}$. Since Σ is a diagonal matrix, W is also diagonal with diagonal elements w_1, w_2, \dots, w_n .

The weighted least squares normal equations are

$$(X' W X) \hat{\beta} = X' W Y$$

This is multiple regression of the weighted least squares normal equation therefore

$$\hat{\beta} = (X' W X)^{-1} X' W Y$$

is the weighted least squared estimator.