

Dr.K.S.Bhanu

Department of Statistics

Institute of Science, Nagpur.

Semester III - Unit I

## **Outliers , leverage points and influential observations**

Contents of the topic: Divided into five parts

- A-** What are outliers and influential observations?
- B-** Reasons for the presence of outliers/Influential observations:
- C-** Detection of leverage and outliers (Diagnostic plots):
- D-** Deletion Diagnostics (for influential observations
- E-** Treatment – Remedial measures.

- 
- This topic deals with situations where one or a few observations are different-in some way-from the rest of the data.
  - We distinguish between two ways a few points may differ from the remaining points.
  - A data point might lie far from the general trend in the rest of the data: such a point is called an outlier.
  - Sometimes, a statistical analysis is very sensitive to a single (or a few) data point(s), in the sense that if the value of this point is changed even slightly, the outcome of the analysis alters greatly. Such points are called leverage points ..

[\[A\] Outliers & Influential observations: \(Discordant observations\):](#)

Discordant observations are those observation the presence of which causes deviations from the four assumptions of the regression model.

(a)Outliers:

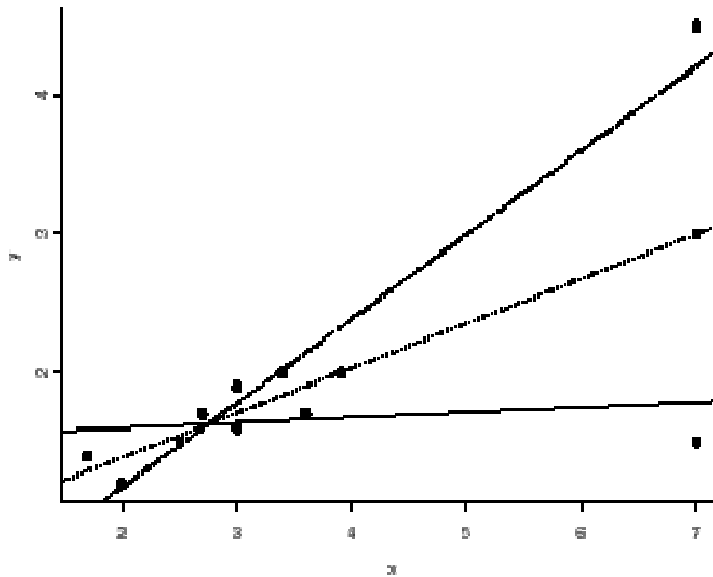
An outlier in a regression model is a data point that produces a large residual. In the space of the explanatory variables(X) if an observation is far removed (away) from the centre of the data it is referred to as an X-outlier. If the observation is such that the corresponding value of the response is far removed from the other values of the response variable that is called Y-outlier.

b) Leverage points:

**Leverage :**

In the regression setup, the values of the regressors (X) can be very high or very low. A single large or low value of X can pull the reg. line to one side .Such observations of X which are far away from the bulk of the data are called **leverage points**. These points can be treated as outliers in the space of the explanatory variables(x)

A leverage point is a point for which the observed value of this particular point has a great influence on the analysis. Suppose we have a cluster of data with  $x$ -values not too far apart; also, we have one observation corresponding to an  $x$ -value further away. The value of this isolated point is disproportionately influential on the least squares line: one might say that it works as a **lever**-if the value of this observation is changed, the least squares line changes considerably (as illustrated in the figure). In contrast, if the value of one of the points within the cluster is changed, the least squares line will not be affected to the same extent.



- The leverage of the  $i$ th observation is defined as the  $i$ th diagonal element of the hat matrix.  $H = X(X'X)^{-1}X'$  (i.e.  $h_{ii}$ )
- The  $i$ th leverage is a no. between 0 and 1 and it indicates the amount of leverage, or influence, the  $i$ th observation has on the least-squares line.
- The larger the value of  $h_{ii}$ , the more influence the observation has on the least squares line.
- A rule of thumb says that an observation is a leverage point if it has a hat-diagonal  $h_{ii}$  is greater than  $2p/n$  where  $p=k+1$ .
- Note that the hat-matrix  $H$  only depends on the design matrix ( $X$ ) and not on the response variables  $Y_i$ . That is, the observed value of the response variable is irrelevant with regard to whether or not a point  $(x_i, y_i)$  is a leverage point.
- The leverage point may be remote in the  $x$ -space yet may lie almost on the regression line passing through the rest of the sample points.
- Note that leverage points do not necessarily constitute a problem. If the observation  $y_i$  corresponding to a leverage point lies close to the general trend in the data, the point is called a good leverage point, and there is no reason to do anything about the data point.
- (Such points do not affect the estimates of the regression coefficients but may affect the value of  $R^2$ )

- However, if  $y_i$  differs from the main trend-in particular, if  $y_i$  corresponds to an outlier-the point is called a bad leverage point, and should be removed from the dataset.

### (c) Influential observations:

Certain observations can considerably influence the analysis of the data due to their presence. They need not necessarily be the outliers or leverage points. These observations which individually or together with several other observations significantly influence the reg. analysis are called influential observations. An observation is influential if removing it would significantly change the position of the reg. line.

### [B] Reasons for the presence of outliers/Influential observations:

1. Gross errors in measuring either the response or explanatory variables.
2. Need for additional variables in the model.
3. Need of an alternate model like a non linear model or a mixture model.

### [C] Detection of leverage and outliers (Diagnostic plots):

1) using leverage (for X outliers):

In the case of simple linear regression we can plot the data point (X,Y) of the fitted model & identify the outlying observation. If the regression model is of the form

The hat matrix plays an important role in identifying influential observations.  $h_{ii}$  is taken as a measure of the distance of the  $i$ th observation from the center of the x-space. Large hat diagonals reveal observations that are potentially influential.

Any observation for which the hat diagonal exceeds  $2p/n$  ( where  $p = k+1$  is considered to be a leverage point.

Note that not all leverage points will be influential.

Observations with large hat diagonals( for which  $h_{ii}$  is large) and large residuals are likely to be influential.

## 2) Residual Analysis (to identify Y-outliers):

Y-outliers can be identified by looking for an unusual pattern in the residual. The  $i$ th residual is given by,

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

Plotting the residuals is a very effective way of investigating how well the regression model fits the data. If the model is correct, then the residuals  $\varepsilon_i$  should behave as zero mean and uncorrelated random variables with constant variance.

The residual  $\varepsilon_i$  measures the degree of misfit of the  $i$ th observation. Thus, residuals contain information which can be used to identify the outliers.

### 1. Standardized Residual:

$$d_i = \frac{\varepsilon_i}{\sqrt{MSRes}}$$

$$\text{where, } MSRes = \frac{SSRes}{n-p} = \frac{\sum \varepsilon_i^2}{n-p}; \quad p = k + 1$$

$$SSRes = \sum (Y_i - \hat{y}_i)^2$$

If  $|d_i| > 3$  ;  $i$ th data point is an outlier.

### 2. Studentised Residual:

#### a) Internally studentised residual:

$$r_i = \frac{\varepsilon_i}{\sqrt{MSE(1 - h_{ii})}}$$

A very large value of  $r_i$  implies that  $i$ th observation can be an outlier.

A plot of the above two residuals against the fitted value  $\hat{y}_i$  (hat) can give an indication about the presence of outliers in the model.

#### b) Externally Studentised residual (R-student statistics):

$$t_i = \frac{e_i}{\hat{\sigma}_i^2(1 - h_{ii})}$$

Where,  $\hat{\sigma}_i^2$  is the estimate of  $\sigma^2$  when the  $i$ th observation is left out.

$$3. R^2_{\text{adjusted}} = 1 - \frac{\frac{SS_{\text{Res}}}{n-p}}{\frac{SST}{n-1}}$$

$$4. \text{Var}(\hat{\beta}_j) = \sigma^2 C_{jj} \quad S.E(\hat{\beta}_j) = \sqrt{\sigma^2 C_{jj}}$$

[D] DELETION DIAGNOSTICS (for influential observations): The influential observations can be identified using deletion diagnostics whereby the regression analysis is carried out both with and without the observation and the effect of deletion is investigated. The deletion of such an observation may change

- 1) The estimates of the regression coefficient  $\hat{\beta}$ .
  - 2) The fit of that observation or other observations
  - 3) The precision of the estimate i.e.  $V(\hat{\beta})$ .
- These can be identified using deletion diagnostics.

#### 1. DIAGNOSTICS FOR PARAMETER ESTIMATES:

They measure the effect of the influential observation on the estimates of the regression coefficient  $\hat{\beta}$ . The influence of deleting the  $i$ th observation on the  $j$ th regression coefficient  $\hat{\beta}_j$  can be detected using the following measure:

$$\text{DFBETAS} = S_{ij}(\beta) = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{\hat{\sigma}_{(i)}^2 C_{jj}}}$$

Where  $C = (C_{ij}) = (X'X)^{-1}$

And  $\hat{\beta}_{j(i)}$  and  $\hat{\sigma}_{(i)}^2$  are the estimates of  $\hat{\beta}_j$  and  $\hat{\sigma}^2$  after deleting the  $i$ th observation. A large value of  $S_{ij}(\beta)$  indicates that the  $i$ th observation has considerable influence on the  $j$ th regression coefficient.

The suggested cutoff point for  $S_{ij}(\beta)$  is  $2/\sqrt{n}$ .

i.e if  $|S_{ij}(\beta)| > 2/\sqrt{n}$  it implies that the  $i$ th observation could be an influential observation.

## 2. DIAGNOSTICS FOR FIT:

This diagnostic method is used to investigate the effect of deleting the  $i$ th observation on the predicted or fitted value. The change in the fit of the  $i$ th observation when it is deleted is obtained using.

$$DFFITSi = Si(FIT) = \frac{\hat{y}_i - \hat{y}_{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2 h_{ii}}}$$

Where,  $\hat{y}_{(i)}$  is the fitted value of  $y_i$  obtained by deleting the  $i$ th observation. It can be shown that

$$Si(\text{fit}) = t_i \sqrt{\frac{h_{ii}}{1 - h_{ii}}}$$

Where  $t_i$ =R-student statistic

$h_{ii}$ =leverage of the  $i$ th observation.

If the data point is an outlier, then  $t_i$  will be large in magnitude while if the data point has high leverage, then  $h_{ii}$  will be close to 1. In either case,  $Si(\text{fit})$  will be large.

Thus  $Si(\text{fit})$  depends both on leverage as well as prediction error.

The suggested cutoff point for  $Si(\text{fit})$  is  $2\sqrt{p/n}$ .

i.e. if  $Si(\text{fit}) > 2\sqrt{p/n}$  it implies that the  $i$ th observation could be an influential observation.

## 3. COOK'S DISTANCE:

To find whether a point is influential or not, it is desirable to consider the location of the point in both the X-space as well as the Y-space. The cook's distance is based upon the squared distance between the least square estimates  $\hat{\beta}$

based on all n points and the estimates obtained by deleting the ith point say  $\hat{\beta}_{(i)}$ . This distance measure is given by(in vector form)

$$Di = \frac{(\hat{\beta}_{(i)} - \hat{\beta})' X' X (\hat{\beta}_{(i)} - \hat{\beta})}{(p)(MSRes)} \dots \dots (1) \quad ; i = 1, 2, \dots, n$$

Where,  $MSRes = \frac{SSRes}{n-p} = \frac{\sum e_i^2}{n-p}$

Point with large values of Di will have considerable influence on the least square estimate  $\beta$ . Points for which  $Di \geq 1$  are considered to be influential.

Since  $\hat{Y} = X\hat{\beta}$

$$\hat{Y}_{(i)} = X\hat{\beta}_{(i)}$$

Expression (1) can be rewritten as

$$Di = \frac{(\hat{y}_{(i)} - \hat{y})' (\hat{y}_{(i)} - \hat{y})}{(p)(MSRes)} \dots \dots \dots (2)$$

The numerator of Di is the squared Euclidean distance that the vector of fitted values moves when the ith observation is deleted.

The Di statistic can also be written as

$$Di = \frac{r_i^2}{p} \left( \frac{h_{ii}}{1 - h_{ii}} \right) = \frac{1}{p}$$

*(square of the studentised residual)*  $\left( \frac{h_{ii}}{1 - h_{ii}} \right)$



Thus  $D_i$  combines the residual magnitude of the  $i$ th observation given by  $r_i$  and the leverage of the  $i$ th point  $h_{ii}$  in the  $X$ -space to find the influence of the  $i$ th observation.

The magnitude of  $D_i$  is compared with  $F_{p,n-p}$ .

i.e if  $D_i \geq F_{p,n-p}$  it implies that the  $i$ th point is highly influential.

#### 4.PRECISION DIAGNOSTICS:

These measures study the influence of deleting the  $i$ th observation on the precision of an estimator through its variance and are called precision diagnostics.

We define the generalized variance of  $\beta$  as

$$GV(\hat{\beta}) = |\text{var}(\hat{\beta})| = |\sigma^2 (X'X)^{-1}|$$

To express the role of the  $i$ th observation on the precision of estimation, we define

$$\text{COVRATIO}(i) = \frac{|\hat{\sigma}_{(i)}^2 (X'_{(i)} X_{(i)})^{-1}|}{|\text{MSRES}(X'X)^{-1}|}$$

Where,  $\text{MSRES} = \hat{\sigma}^2 = \text{U. E. of } \sigma^2$

$$\begin{aligned} &= \frac{|\hat{\sigma}_{(i)}^2 (X'_{(i)} X_{(i)})^{-1}|}{|\hat{\sigma}^2 (X'X)^{-1}|} \\ &= \left( \frac{\hat{\sigma}_{(i)}^2}{\hat{\sigma}^2} \right)^p \left( \frac{1}{1-h_{ii}} \right) \end{aligned}$$

If  $\text{COVRATIO}(i) > 1$ , the  $i$ th observation improves the precision of estimation, while if  $\text{COVRATIO}(i) < 1$ , inclusion of the observation degrades the precision.

CUTOFF: The suggested cutoff point for  $\text{COVRATIO}$  is  $1 + 3 p/n$ .

I.e. if  $\text{COVRATIO} > 1 + 3 P/N$ . or if  $\text{COVRATIO} < 1 - 3 P/N$ ., then the  $i$ th point should be considered influential.

NOTE: A high leverage point will make the COVRATIO(i) large. This is because a high leverage point will improve the precision unless the point is an outlier in the Y-space.

## EJ TREATMENT OF OUTLIERS / INFLUENTIAL OBSERVATIONS.

1. The presence of outliers / influential observations should be carefully investigated to see if a reason for their unusual behavior can be found. Care should be taken to see if the outlier is important to the model before it is discarded. The effect of outliers on the regression model may be checked by dropping these points and refitting the regression model. If the estimates of the parameters of the model do not change much, we can conclude that the points are not influential observations. Hence these points can be safely deleted. Deleting these points, the precision of the estimates can be improved and the width of the C.I of  $\beta$  can be decreased.
2. If the influential observation is due to an error in measurement or if the sample point is really not needed for the model, then it can be safely discarded. However, if the analysis reveals that an influential point is a valid observation, then it should be retained.

These decisions are taken using robust estimation procedures

---