

Diagnostic checks of regression models

The model we arrive at when performing regression analysis depends strongly on the quality of the data.

Problems

Multicollinearity: Two or more columns in the X -matrix are almost linear dependent.

Checked by the correlation between regression variables,

Estimated correlation between x_i and x_j is given by

$$r_{ij} = \frac{\sum_{k=1}^m (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^m (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^m (x_{jk} - \bar{x}_j)^2}}$$

Multicollinearity often occurs when the variation interval of two or more variables are approximately equal.

A good measure of multicollinearity is the variance inflation factor.

$$VIF = \frac{1}{1-R_j^2}$$

where R_j^2 is the multiple determination coefficient when x_j is regressed on the other regression variables.

Important $VIF < 10$. If it occurs, consider.

1. Remove column

2. Collect more data

3. Respecify the model. Instead of x_1, x_2, x_3 , $\frac{x_1 + x_2 + x_3}{3}$ or x_1, x_2, x_3 may be used.

4. Center variables $(x_i - \bar{x})$ polynomial models, or standardize.

5. Use PCA regression or PLS regression, eventually Ridge regression.

Influential observations.

leverage points. $\hat{y} = H y \Rightarrow \hat{y}_i = \sum_{j=1}^m h_{ij} y_j = h_{ii} y_i + \sum_{i \neq j} h_{ij} y_j$

If h_{ii} is large, y_i could have a large influence on \hat{y}_i .

rank(H) = $k+1 = \sum_{i=1}^m h_{ii}$. Average $h_{ii} = \frac{k+1}{m}$. Observations

where $h_{ii} > \frac{2(k+1)}{m}$ are said to be leverage points.

Cook's distance D_i

Cook's distance

$$(\hat{y} - \hat{y}_{(i)})^T (\hat{y} - \hat{y}_{(i)}) = (x_b - x_{b(i)})^T (x_b - x_{b(i)}) = (\beta - \beta_{(i)})^T X^T (b - b_{(i)})$$

The (i) notation means that the i -th set of observations $(y_i, x_{i1}, \dots, x_{ik})$ is taken out.

Cook's distance D_i is defined as

$$D_i = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T X^T (\hat{\beta} - \hat{\beta}_{(i)})}{\hat{\sigma}^2 (k+1) \hat{\sigma}^2}, \quad i = 1, 2, \dots, n$$

with substituted estimates $\hat{h}_{ii} d_i^2$ where $d_i = \frac{e_i}{\sqrt{1-h_{ii}}}, \quad \hat{\sigma} = \hat{\sigma}_{LS}$

$$\frac{\hat{h}_{ii} d_i^2}{(1-h_{ii}) k+1}$$

We consider observation $(y_i, x_{i1}, \dots, x_{ik})$ to be influential if $D_i > 1$

Study of residuals

If the model is correct, we expect the residuals to be uncorrelated with mean zero and constant variance.

$$\hat{\epsilon} = \underline{y} - \hat{\underline{y}} = (\underline{I} - \underline{H})\underline{y} = (\underline{I} - \underline{H})\underline{\epsilon}$$

Therefore $E[\hat{\epsilon}] = (\underline{I} - \underline{H})E[\underline{\epsilon}] = 0$

and $\text{Cov}[\hat{\epsilon}] = E[(\underline{I} - \underline{H})\underline{\epsilon}\underline{\epsilon}^T(\underline{I} - \underline{H})] = \sigma^2(\underline{I} - \underline{H})$

which shows that $\hat{\epsilon}_1, \hat{\epsilon}_2, \dots, \hat{\epsilon}_m$ are not uncorrelated.

$$\text{Var}[\hat{\epsilon}_i] = \sigma^2(1-h_{ii})$$

Hence the variances of $\hat{\epsilon}_1, \dots, \hat{\epsilon}_m$ need not to be equal.

Since $\text{Var}\left[\frac{\hat{\epsilon}_i}{\sigma\sqrt{1-h_{ii}}}\right] = 1$, we often use $h_i = \frac{e_i}{\sigma\sqrt{1-h_{ii}}}$ in the study of the residuals. These are called standardized residuals. Some books call them studentized or internal studentized.

The residuals $\delta_i = y_i - \hat{y}_{i,-i}$ ($i = 1, 2, \dots, n$) where $\hat{y}_{i,-i}$ is the estimator for the i -th fitted value when the observation $(y_i, x_{i1}, \dots, x_{ik})$ is taken out, ~~are~~ are called the Press residuals (Predictive sum of squares). These can be shown to equal $\frac{\hat{\epsilon}_i}{1-h_{ii}}$ and

$$\text{Var}(\delta_i) = \frac{\text{Var}[\hat{\epsilon}_i]}{(1-h_{ii})^2} = \frac{\sigma^2(1-h_{ii})}{(1-h_{ii})^2} = \frac{\sigma^2}{1-h_{ii}}$$

$$\text{Hence } \frac{s_i}{s_{d(s_i)}} = \frac{\hat{\epsilon}_i}{\sqrt{(1-h_{ii})}} = \frac{\hat{\epsilon}_i}{\sqrt{1-h_{ii}}}$$

Assume we estimate σ^2 by $s_{(i)}^2$ which is the ~~estimate~~
estimate for the variance when $(y_i, x_{i1}, \dots, x_{ik})$ is
taken out. Then $s_i = y_i - \hat{y}_{i(i)}$ and $s_{(i)}^2$ are
independent ($\hat{\beta}$ and σ^2 are independent)

and $\frac{\hat{\epsilon}_i}{s_{(i)}\sqrt{1-h_{ii}}} \sim t_{n-k-2}$ and can be used to

investigate if the expected value of a residual is
different from zero, i.e. if we have an outlier.
The residuals $\frac{\hat{\epsilon}_i}{s_{(i)}\sqrt{1-h_{ii}}}$ are called R-studentized or

sometimes ~~internal~~ studentized.

Check of residuals provide information about

1. Outliers

2. Heterogeneity of variance

3. Mispecification of models (or if other variables
should be included).

4. Normal distribution

The most common residual plots are plots of residuals
against fitted values to check for heterogeneity in variance,
against other variables to see if they should be included,
and normal-pwt to check for normally distributed data and
outliers.