L26 DFFITS and Cook's D

- 1. Statistics for identifying outliers
 - (1) Residuals and deleted residuals

The *i*th residual $R_i = \hat{e}_i = y_i - \hat{y}_i \sim N\left(0, (1 - h_{ii})\sigma^2\right)$. The *i*th deleted residual $\text{PRESS}_i = \hat{e}_{i(i)} = y_i - \hat{y}_{i(i)} \stackrel{*}{=} \frac{\hat{e}_i}{1 - h_{ii}} \sim N\left(0, \frac{\sigma^2}{1 - h_{ii}}\right)$.

(2) Studentized residuals and sudentized deleted residuals

$$\begin{split} \widehat{e}_{i} &= y_{i} - \widehat{y}_{i} \sim N\left(0, (1 - h_{ii})\sigma^{2}\right) \Longrightarrow \frac{\widehat{e}_{i}}{\sqrt{(1 - h_{ii})\sigma^{2}}} \sim N(0, 1^{2}).\\ \text{The ith studentized residual STUDENT}_{i} &= t_{i} = \frac{\widehat{e}_{i}}{\sqrt{(1 - h_{ii})MSE}} \sim t(n - p).\\ \widehat{e}_{i(i)} &= y_{i} - \widehat{y}_{i(i)} = \frac{\widehat{e}_{i}}{1 - h_{ii}} \sim N\left(0, \frac{\sigma^{2}}{1 - h_{ii}}\right) \Longrightarrow \frac{\widehat{e}_{i(i)}}{\sqrt{\sigma^{2}/(1 - h_{ii})}} = \frac{\widehat{e}_{i}}{\sqrt{(1 - h_{ii})\sigma^{2}}} \sim N(0, 1^{2}).\\ \text{The ith studentized deleted residual RSTUDENT} &= t_{(i)} = \frac{\widehat{e}_{i}}{\sqrt{(1 - h_{ii})MSE}} \sim t(n - 1 - p).\\ \text{Comment: Standardizing } \widehat{e}_{i} \text{ and } \widehat{e}_{i(i)} \text{ lead to the same } \frac{\widehat{e}_{i}}{\sqrt{(1 - h_{ii})\sigma^{2}}} \sim N(0, 1^{2}).\\ \text{In studentizing } \widehat{e}_{i}, \sigma^{2} \text{ is replaced by MSE that resulted in} \end{split}$$

STUDENT_i =
$$\frac{\widehat{e}_i}{\sqrt{(1-h_{ii})MSE}} \sim t(n-p).$$

In studentizing $\hat{e}_{i(i)}, \sigma^2$ is replaced by $MSE_{(i)}$ that resulted in

RSTUDEHT_i =
$$\frac{\widehat{e}_i}{\sqrt{(1-h_{ii})MSE_{(i)}}} \sim t(n-1-p).$$

2. DFFITS

(1) Difference of fitted values The *i*th difference of fitted values $\hat{y}_i - \hat{y}_{i(i)} = \hat{y}_i - y_i + y_i - \hat{y}_{i(i)} = -\hat{e}_i + \hat{e}_{i(i)} = -\hat{e}_i + \frac{1}{1 - h_{ii}}\hat{e}_i$

$$= \frac{h_{ii}}{1-h_{ii}}\widehat{e}_i \sim \frac{h_{ii}}{1-h_{ii}}N(0, (1-h_{ii})\sigma^2) = N\left(0, \frac{h_{ii}^2}{1-h_{ii}}\sigma^2\right)$$

So $\widehat{y}_i - \widehat{y}_{i(i)} = \begin{cases} \frac{h_{ii}}{1-h_{ii}}\widehat{e}_i \\ h_{ii}\widehat{e}_{i(i)} \end{cases} = \frac{h_{ii}}{1-h_{ii}}\widehat{e}_i \sim N\left(0, \frac{h_{ii}^2}{1-h_{ii}}\sigma^2\right)$

(2) Two studentized difference of fitted values

Standardizing $\hat{y}_i - \hat{y}_{i(i)}$, $\frac{\hat{y}_i - \hat{y}_{i(i)}}{\sqrt{h_{ii}^2 \sigma^2 / (1 - h_{ii})}} = \frac{\hat{e}_i}{\sqrt{(1 - h_{ii})\sigma^2}} \sim N(0, 1^2).$ Two studentized $\hat{y}_i - \hat{y}_{i(i)}$, $\frac{\hat{e}_i}{\sqrt{(1 - h_{ii})MSE}} \sim t(n - p)$ and $\frac{\hat{e}_i}{\sqrt{(1 - h_{ii})MSE_{(i)}}} \sim t(n - 1 - p)$ are studentized residual and studentized deleted residual.

Comment: Standardizing $y_i - \hat{y}_i$, $y_i - \hat{y}_{i(i)}$ and $\hat{y}_i - \hat{y}_{i(i)}$ lead to the same $\frac{\hat{e}_i}{\sqrt{(1-h_{ii})\sigma^2}} \sim N(0, 1^2)$ with two studentized forms: Studentized residual and Studentized deleted residual.

(3) DFFITS

$$\begin{array}{ll} \text{Define} & \text{DFFITS}_{i} = \frac{\widehat{y}_{i} - \widehat{y}_{i(i)}}{\sqrt{h_{ii} \, MSE_{(i)}}}. \\ \text{Then} & \text{DFFITS}_{i} = \sqrt{\frac{h_{ii}}{MSE_{(i)}}} \frac{\widehat{e}_{i}}{1 - h_{ii}} = \sqrt{\frac{h_{ii}}{MSE_{(i)}}} \, \widehat{e}_{i(i)} = \sqrt{\frac{h_{ii}}{1 - h_{ii}}} t_{(i)}. \\ \text{So} & \text{DFFITS}_{i} \sim \sqrt{\frac{h_{ii}}{1 - h_{ii}}} t(n - 1 - p) \end{array}$$

By the rule of thumb, if $|\text{DFFITS}_i| > 2\sqrt{\frac{p}{n}}$, then the *i*th observation is an outlier.

Ex1: Express DFFITS_i through
$$\hat{e}_i$$
, $\hat{e}_{i(i)}$ and $t_{(i)}$.
DFFITS_i = $\sqrt{\frac{h_{ii}}{MSE_{(i)}}} \frac{\hat{e}_i}{1-h_{ii}} = \sqrt{\frac{h_{ii}}{MSE_{(i)}}} \hat{e}_{i(i)} = \sqrt{\frac{h_{ii}}{1-h_{ii}}} t_{(i)}$

- 3. Cook's D
 - (1) Squared norm of difference of fitted vectors $\|\widehat{y} - \widehat{y}_{(i)}\|^2 = \sum_{j=1}^n (\widehat{y}_j - \widehat{y}_{j(i)})^2.$
 - (2) Cook's D

Cook's D
Define
$$D_i = \frac{\|\widehat{y} - \widehat{y}_{(i)}\|^2}{pMSE}$$

(3) Appication

Large D_i indicates that the *i*th observation is an outlier. By the rule of thumb, the cut-off value for D_i is $\frac{4}{n}$.

 \mathbf{Ex} : SAS

```
proc reg;
   model y=x1 x2;
   output out=b P=P R=R H=H STUDENT=STUDENT
          PRESS=PRESS RSTUDENT=RSTUDENT
          DFFITS=DFFITS COOKD=COOKD;
   run;
proc print;
   run;
```

L27 Multicollinearity

- 1. Multicollinearity
 - (1) An assumption

Based on sample $y \sim N(X\beta, \sigma^2 I_n)$ from a regression model, β can be estimated by its MLE and LSE $\hat{\beta} = (X'X)^{-1}X'y \sim N(\beta, \sigma^2(X'X)^{-1})$. To get this estimator there is an underlying assumption: X'X is non-singular. This assumption has many different but equivalent forms. For example

- (i) The columns of $X \in \mathbb{R}^p$ are linearly independent;
- (ii) No column of X is a linear combination of others;
- (iii) |X'X| > 0.
- (2) Multicollinearity

There is multicollinearity in $X \iff C$ The columns of X are almost linearly dependent $\iff C$ One column of X is almost a LC of others $\iff |X'X| > 0$ is almost 0.

Let $X'X = Q\Lambda Q'$ be the eigenvalue decomposition of X'X where the columns of Q are eigenvectors of X'X, Q is an orthogonal matrix $(Q' = Q^{-1})$; $\Lambda = \operatorname{diag}(\lambda_1, ..., \lambda_p), \lambda_i > 0$, i = 1, ..., p, are eigenvalues of X'X. Because $|X'X| = \lambda_1 \cdots \lambda_p$, the multicollinearity is reflected as the values of $\lambda_i > 0$ close to 0.

(3) Consequence of multicollinearity

If there is multicollinearity in X, then the total variance in $\widehat{\beta} \sim N(0, \sigma^2(X'X)^{-1})$,

$$\sum_{i=1}^{p} \operatorname{var}(\widehat{\beta}_{i}) = \operatorname{tr}\left(\operatorname{Cov}(\widehat{\beta})\right) = \operatorname{tr}[\sigma^{2}(X'X)^{-1}] = \sigma^{2}\operatorname{tr}[(Q\Lambda Q')^{-1}]$$
$$= \sigma^{2}\operatorname{tr}(Q\Lambda^{-1}Q') = \sigma^{2}\operatorname{tr}(\Lambda^{-1}) = \sigma^{2}\left(\frac{1}{\lambda_{1}} + \dots + \frac{1}{\lambda_{p}}\right)$$

will be vary large. So even though $\hat{\beta}$ is an unbiased estimator, but it is not stable.

- 2. Detecting multicollinearity
 - (1) Sample correlation coefficient method

proc corr; var x1 x2 x3; run;

calculates and presents sample correlation coefficient matrix $\begin{pmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{pmatrix}$. Because $|r_{ij}| \leq 1$ and $|r_{ij}| = 1 \iff x_i = ax_j + b$, the value of $|r_{ij}|$ close to 1 indicates that x_i and x_j are almost correlated consequently there is multicollinearity in X.

(2) Variance inflation factor method

Columns of X are almost linearly dependent if and only if one column, say x_i , is almost a linear combination of others. Thus regression model $x_i = \beta_0 + \sum_{j \neq i} \beta_j x_j + \epsilon$ will produce the coefficient of determination R_i^2 will value close to 1. Consequently $\text{VIF}_i = \frac{1}{1-R_i^2}$ called the variance inflation factor will also have large value.

```
proc reg;
model y=x1 x2 x3/vif;
run;
```

calculates and presents VIF_i , i = 1, 2, 3. If one $\text{VIF}_i > 2.5$, then we conclude the problem of multicollinearity in X.

3. Ridge regression: A remedy

(1) Ridge regression

With eigenvalue decomposition $X'X = Q\Lambda Q'$, the LSE and MLE of β is

$$\widehat{\beta} = (X'X)^{-1}X'y = (Q\Lambda Q')^{-1}X'y.$$

Small values of $\lambda_i > 0$ in $\Lambda = \text{diag}(\lambda_1, ..., \lambda_p)$ caused multicollinearity. An naive and intuitive idea is to add $K = \text{diag}(k_1, ..., k_p)$ where $k_i > 0$ for all i to Λ to have a new estimator $\widehat{\beta}(K) = [Q(\Lambda + K)Q']^{-1}X'y$.

The diagonal elements of Λ form the ridge of Λ , the proposed approach elevates the ridge and hence is called a ridge regression.

(2) Compare performance of estimators

Let $\hat{\eta}$ be an estimator for $\eta \in \mathbb{R}^k$. Then $E \| \hat{\eta} - \eta \|^2$ is the mean squared error. When comparing two estimators, the one with smaller mean squared error is the better one.

$$\begin{split} E\|\widehat{\eta}-\eta\|^2 &= E[(\widehat{\eta}-\eta)'(\widehat{\eta}-\eta)] \\ &= E\left\{[\widehat{\eta}-E(\widehat{\eta})+E(\widehat{\eta})-\eta]'[\widehat{\eta}-E(\widehat{\eta})+E(\widehat{\eta})-\eta]'\right\} \\ &= E\left\{[\widehat{\eta}-E(\widehat{\eta})]'[\widehat{\eta}-E(\widehat{\eta})]\right\} + \|E(\widehat{\eta})-\beta\|^2 \\ &= \sum_{i=1}^k \operatorname{var}(\widehat{\eta}_i) + \|E(\widehat{\eta})-\beta\|^2 = \text{Toal variance in } \widehat{\eta} + (\text{Bias})^2. \end{split}$$

(3) Performance of $\widehat{\beta}(K)$

(i) Ridge estimator is a biased estimator

$$\begin{split} E[\widehat{\beta}(K)] &= E\left\{ [Q(\Lambda + K)Q']^{-1}X'y \right\} \\ &\neq E[(Q\Lambda Q')^{-1}X'y] = E[(X'X)^{-1}X'y] = (X'X)^{-1}X'X\beta = \beta. \end{split}$$

(ii) Ridge estimator reduces the total variance

$$\begin{split} \sum_{i} \operatorname{var}[\widehat{\beta}_{i}(K)] &= \operatorname{tr}[\operatorname{Cov}(\widehat{\beta}(K))] \\ &= \operatorname{tr}\{\sigma^{2}[Q(\Lambda + K)^{-1}Q']^{-1}X'X[Q(\Lambda + K)^{-1}Q']\} \\ &= \sigma^{2}\operatorname{tr}[(\Lambda + K)^{-1}\Lambda(\Lambda + K)^{-1}] = \sigma^{2}\sum_{i}\frac{\lambda_{i}}{(\lambda_{i} + k_{i})^{2}} \\ &< \sigma^{2}\sum_{i}\frac{1}{\lambda_{i}} = \sum_{i}\operatorname{var}(\widehat{\beta}_{i}) \end{split}$$

(iii) Optimization

It has been shown that there are k_i , k = 1, ..., p, such the mean squared error of ridge estimator is less than that of MLE of β . However these k_i depend on parameters and must be estimated.