

# **Regression Model Building - Diagnostics**

KNNL – Chapter 10

# Model Adequacy for Predictors

## Added Variable Plot (Partial Regression Plot)

- Graphical way to determine partial relation between response and a given predictor, after controlling for other predictors – shows form of relation between new  $X$  and  $Y$
- May not be helpful when other predictor(s) enter model with polynomial or interaction terms that are not included
- Algorithm (assume plot for  $X_3$ , given  $X_1, X_2$ ):
  - Fit regression of  $Y$  on  $X_1, X_2$ , obtain residuals =  $e_i(Y | X_1, X_2)$
  - Fit regression of  $X_3$  on  $X_1, X_2$ , obtain residuals =  $e_i(X_3 | X_1, X_2)$
  - Plot  $e_i(Y | X_1, X_2)$  (vertical axis) versus  $e_i(X_3 | X_1, X_2)$  (horizontal axis)
- Slope of the regression through the origin of  $e_i(Y | X_1, X_2)$  on  $e_i(X_3 | X_1, X_2)$  is the partial regression coefficient for  $X_3$

# Added Variable Plot

- If line is flat horizontal, then the additional variable is not helpful, controlling for the rest of the predictors.
- If the extra variable is important and has a linear relationship with  $Y$ , then we see a tilted straight-line scatterplot.
- If we see a curve, then we should add the extra variable with a transformation (like polynomial)

# Outlying Y Observations – Studentized Residuals

Model Errors (unobserved):

$$\varepsilon_i = Y_i - (\beta_0 + \beta_1 X_{i1} + \dots + \beta_{p-1} X_{i,p-1}) \quad E\{\varepsilon_i\} = 0 \quad \sigma^2\{\varepsilon_i\} = \sigma^2 \quad \sigma\{\varepsilon_i, \varepsilon_j\} = 0 \quad \forall i \neq j$$

Residuals (observed) where  $h_{ij} = (i, j)^{th}$  element of  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ :  
 $n \times n$

$$e_i = Y_i - \hat{Y}_i = Y_i - (b_0 + b_1 X_{i1} + \dots + b_{p-1} X_{i,p-1})$$

$$E\{e_i\} = 0 \quad \sigma^2\{e_i\} = \sigma^2(1 - h_{ii}) \quad \sigma\{e_i, e_j\} = -h_{ij}\sigma^2 \quad \forall i \neq j$$

$$s^2\{e_i\} = MSE(1 - h_{ii}) \quad s\{e_i, e_j\} = -h_{ij}MSE \quad \forall i \neq j$$

Semi-Studentized Residual (Residual divided by estimate of  $\sigma$ , trivial to compute):

$$e_i^* = \frac{e_i}{\sqrt{MSE}}$$

Studentized Residual (Residual divided by its standard error, messier to compute):

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

# Outlying $Y$ Observations – Studentized Deleted Residuals

Deleted Residual (Observed value minus fitted value when regression is fit on the other  $n-1$  cases):

$$d_i = Y_i - \hat{Y}_{i(i)} \quad \hat{Y}_{i(i)} = b_{0(i)} + b_{1(i)}X_{i1} + \dots + b_{p-1(i)}X_{i,p-1}$$

$b_{k(i)} \equiv$  regression coefficient of  $X_k$  when case  $i$  is deleted

Studentized Deleted Residual (makes use of having predicted  $i^{th}$  response from regression based on other  $n-1$  cases):

$$t_i = \frac{d_i}{s\{d_i\}} \sim t(n-p-1)$$

$$s^2\{d_i\} = s^2\{\text{pred}_i\} = MSE_{(i)} \left[ 1 + \mathbf{X}_i' (\mathbf{X}_{(i)}' \mathbf{X}_{(i)})^{-1} \mathbf{X}_i \right] \quad \mathbf{X}_i' = \begin{bmatrix} 1 & X_{i1} & \dots & X_{i,p-1} \end{bmatrix}$$

Note:  $SSE = (n-p)MSE = (n-p-1)MSE_{(i)} + \frac{e_i^2}{1-h_{ii}}$

$$\Rightarrow t_i = e_i \left[ \frac{n-p-1}{SSE(1-h_{ii}) - e_i^2} \right]^{1/2} \quad \text{Computed without re-fitting } n \text{ regressions}$$

Test for outliers (Bonferroni adjustment): Outlier if  $|t_i| \geq t \left( 1 - \left( \frac{\alpha}{2n} \right), n-p-1 \right)$

# Outlying X-Cases – Hat Matrix Leverage Values

$$\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' = \begin{bmatrix} h_{11} & h_{12} & \cdots & h_{1n} \\ h_{21} & h_{22} & \cdots & h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \cdots & h_{nn} \end{bmatrix} \quad h_{ij} = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_j \quad \mathbf{x}_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{i,p-1} \end{bmatrix}$$

Notes:  $0 \leq h_{ii} \leq 1$        $\sum_{i=1}^n h_{ii} = \text{trace}(\mathbf{H}) = \text{trace}(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}') = \text{trace}(\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}) = \text{trace}(\mathbf{I}_p) = p$

Cases with X-levels close to the “center” of the sampled X-levels will have small leverages. Cases with “extreme” levels have large leverages, and have the potential to “pull” the regression equation toward their observed Y-values. Large leverage values are  $> 2p/n$  (2 times larger than the mean)

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y} \quad \Rightarrow \quad \hat{Y}_i = \sum_{j=1}^n h_{ij}Y_j = \sum_{j=1}^{i-1} h_{ij}Y_j + h_{ii}Y_{ii} + \sum_{j=i+1}^n h_{ij}Y_j$$

Leverage values for new observations:  $h_{\text{new,new}} = \mathbf{X}'_{\text{new}} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_{\text{new}}$

New cases with leverage values larger than those in original dataset are extrapolations

# Identifying Influential Cases – Fitted Values

Influential Cases in Terms of Their Own Fitted Values - DFFITS:

$$DFFITS_i = \frac{\hat{Y}_i - \hat{Y}_{i(i)}}{\sqrt{MSE_{(i)} h_{ii}}} \equiv \# \text{ of standard errors own fitted value is shifted when case included vs excluded}$$

Computational Formula (avoids fitting all deleted models):

$$DFFITS_i = e_i \left[ \frac{n-p-1}{SSE(1-h_{ii}) - e_i^2} \right]^{1/2} \left( \frac{h_{ii}}{1-h_{ii}} \right)^{1/2} = t_i \left( \frac{h_{ii}}{1-h_{ii}} \right)^{1/2}$$

Problem cases are  $>1$  for small to medium sized datasets,  $> 2\sqrt{\frac{p}{n}}$  for larger ones

Influential Cases in Terms of All Fitted Values - Cook's Distance:

$$D_i = \frac{\sum_{j=1}^n \left( \hat{Y}_j - \hat{Y}_{j(i)} \right)^2}{pMSE} = \frac{\left( \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)} \right)' \left( \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_{(i)} \right)}{pMSE} = \frac{e_i^2}{pMSE} \left[ \frac{h_{ii}}{(1-h_{ii})^2} \right]$$

Problem cases are  $> F(0.50; p, n-p)$

# Influence on the Regression Coefficients

A measure of the influence of the  $i^{\text{th}}$  case on each regression coefficient  $b_k$  is the difference between the estimated regression coefficient  $b_k$  based on all  $n$  cases and the regression coefficient  $b_{k(i)}$  obtained when the  $i^{\text{th}}$  case is omitted:

$$(DFBETAS)_{k(i)} = \frac{b_k - b_{k(i)}}{\sqrt{MSE_{(i)} c_{kk}}}$$

where

$$(X'X)^{-1} = \begin{bmatrix} c_{00} & c_{01} & \cdots & c_{0,p-1} \\ c_{10} & c_{11} & \cdots & c_{1,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ c_{p-1,0} & c_{p-1,1} & \cdots & c_{p-1,p-1} \end{bmatrix}$$

Problem cases when:

$$(DFBETAS)_{k(i)} > \frac{2}{\sqrt{n}}$$

Note: Cook's Distance can also be computed as aggregate measure on the the influence on all regression coefficients and is algebraically equivalent to the overall influence on all fitted values.



# Multicollinearity - Variance Inflation Factors

- Problems when predictor variables are correlated among themselves
  - Regression Coefficients of predictors change, depending on what other predictors are included
  - Extra Sums of Squares of predictors change, depending on what other predictors are included
  - Standard Errors of Regression Coefficients increase when predictors are highly correlated
  - Individual Regression Coefficients are not significant, although the overall model is
  - Width of Confidence Intervals for Regression Coefficients increases when predictors are highly correlated
  - Point Estimates of Regression Coefficients are wrong sign (+/-)

# Multicollinearity – Informal Diagnostics

- Large changes in the estimated regression coefficients when a predictor variable is added or deleted.
- Nonsignificant results in individual t-tests on the regression coefficients.
- Estimated regression coefficients with an algebraic sign that is the opposite of that expected.
- Large coefficients of simple correlation between pairs of predictor variables in the correlation matrix.
- Strong linear associations between pairs of predictors in the scatterplot matrix.

# Variance Inflation Factor

Original Units for  $X_1, \dots, X_{p-1}, Y$ :  $\sigma^2 \{\mathbf{b}\} = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$

Correlation Transformed Values:  $X_{ik}^* = \frac{1}{\sqrt{n-1}} \left( \frac{X_{ik} - \bar{X}_k}{s_k} \right) \quad Y_i^* = \frac{1}{\sqrt{n-1}} \left( \frac{Y_i - \bar{Y}}{s_Y} \right)$

$$\sigma^2 \{\mathbf{b}^*\} = (\sigma^*)^2 \mathbf{r}_{xx}^{-1} \quad \sigma^2 \{b_k^*\} = (\sigma^*)^2 (VIF)_k$$

$$\text{where: } (VIF)_k = \frac{1}{1 - R_k^2}$$

with  $R_k^2 \equiv$  Coefficient of Determination for regression of  $X_k$  on the other  $p-2$  predictors

$$R_k^2 = 0 \Rightarrow (VIF)_k = 1 \quad 0 < R_k^2 < 1 \Rightarrow (VIF)_k > 1 \quad R_k^2 = 1 \Rightarrow (VIF)_k = \infty$$

Multicollinearity is considered problematic wrt least squares estimates if:

$$\max \left( (VIF)_1, \dots, (VIF)_{p-1} \right) > 10 \quad \text{or if} \quad (\overline{VIF}) = \frac{\sum_{k=1}^{p-1} (VIF)_k}{p-1} \text{ is much larger than } 1$$