

## L7. Modeling Problems

During Model Evaluation, we should try to identify modeling problems that could potentially hurt the model performance and inferences. The model problems can be summarized as

- Data Structural Problems :

- Multicollinearity
- Influential Points

- Model Assumption Violations :

- Heteroscedasticity
- Non-normality residuals
- $y$  and  $X$  are not linearly related.

but modeled that way :

False assumption of linearity  
between  $y$  and predictors.

For all the modeling problems, we will break down to

- a. Damage it might cause (less on math, more on symptoms)
- b. Detection
- c. Solution.

In LF, we will start with

### I. Multicollinearity

### II. Influential Points.

#### I. Multicollinearity.

Problem: two or more of the predictors are highly correlated,  
or in a broader description, one of the columns  
of  $\vec{X}^T \vec{X}$  is a linear combination of the others,  
i.e.  $\vec{X}^T \vec{X}$  is close to non-invertible.

(a) Damage it might cause:

Recall:  $\vec{b} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$      $\leftarrow$  becomes very "large"  
 $E(\vec{b}) = \vec{\beta}$     sensitive to the change in  $\vec{X}$ .

$\text{Var}(\vec{b}) = \sigma^2 (\vec{X}^T \vec{X})^{-1}$      $\leftarrow$  becomes very "large"

- The coefficient estimates:

can swing wildly based on which other independent variables are in the model. The coefficient  $\beta_1$  becomes very sensitive to small changes in the model.

**Symptom**: when you add and drop variables.

there are huge changes to the value of coefficients, sometimes even the sign of the fitted coefficients

**Ex**:

A perfect data not showing multicollinearity at all:

Model	$\hat{\beta}_1$	
$y \sim x_1$	-1	Adding $x_2$
$y \sim x_1 + x_2$	-1	or not, doesn't change the estimate

Model	$\hat{\beta}_2$	
$y \sim x_2$	-1.75	Adding $x_1$
$y \sim x_1 + x_2$	-1.75	changed $\hat{\beta}_2$ a lot.

A data with serious multicollinearity.

$$\begin{aligned} \hat{\beta}_1 &= 1.20 \\ y \sim x_1 &\Rightarrow 1.20 \\ y \sim x_1 + x_2 &\Rightarrow 1.039 \\ y \sim x_2 &\Rightarrow 34.44 \\ y \sim x_1 + x_2 &\Rightarrow 5.83 \end{aligned}$$

- Test result unreliable due to the inflation in  $se(\hat{\beta}_k)$

(1) In t test:  $t_{\text{stat}} = \frac{\hat{\beta}_k}{se(\hat{\beta}_k)} \downarrow$

the inflation in  $se(\hat{\beta}_k)$  caused smaller t statistics, harder to have a rejection / significant predictor when having a lot of predictors in the model.

symptom

: There are very few coefficients show significance from t test.

(2) In anova (typ=1), by changing the order of variables, the significance might change dramatically.

Ex

Data without Multicollinearity.

Model	$se(\hat{\beta}_1)$	$t_{\text{stat}}$ for $\hat{\beta}_1$	ANOVA (typ=1)
$y \sim x_1$	1.47	-0.68	$ss(x_1   \text{Null}) = 8$
$y \sim \cancel{x_2} + x_1$	1.41	-0.71	$ss(x_1   x_2) = 8$ .

## Data with serious Multicollinearity

Model	$se(\hat{\beta}_1)$	tstat	ANOVA (typ=1)
$y \sim x_1$	0.093	12.92	$ss(x_1   \text{Null}) = 505.5$
$y \sim x_2 + x_1$	6.193	5.39	$ss(x_1   x_2) = 88.43$

### Note about Multicollinearity:

(1) In reality, it always exists. When it's not too "serious" we don't need to worry about it too much.

(2) When the problem is serious, we can still use the model to do the prediction, since all estimation are unbiased. But the std. err of prediction would increase;

Not recommended; results might be not reliable.

It affects more in the model inference: testing. It becomes difficult for us to examine the relationship between each predictor and response variable independently.

(b) Detection:

To detect Multicollinearity in the data.

(1) Plot - heatmap of pair-wise correlation  
between predictors (Naive!! but pretty)

Ex: Credit.csv.

Leave out the response variable "Balance"

some correlations that might be tricky:

"Income" v.s. "Limit" (0.79)

"Income" v.s. "Rating" (0.79)

Normally we say correlation  $\geq 0.8$  is strong,

those are on the boundaries.

"Limit" v.s. "Rating" (1!)

(2) Variance Inflation Factor (VIF) - most common.

VIF measures how much the variance is inflated in the coefficient estimates.

- $VIF_j$  for  $\hat{\beta}_j$  is

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

where  $R_j^2$  is the  $R^2$ -value obtained by fitting

the model without  $X_j$

i.e.  ~~$y \sim X_1 + X_2 + X_3$~~

$$VIF_j = \frac{1}{1 - R_j^2} \text{ where } R_j^2 \text{ is}$$

the  $R^2$  value from  $y \sim X_2 + X_3$

Recall  $R^2 = \frac{SSR}{SST} \leftarrow$  without  $X_j$ , how much variation in the model has lost.

- When  $VIF_j = 1$ : we don't lose any variation without  $X_j$ , AKA. It's not correlated with any other predictors. / the variance in  $X_j$  is not inflated at all.
- When  $1 < VIF_j \leq 4$ : light. } fine.
- When  $4 < VIF_j \leq 10$ : moderate } impacted by the multicollinearity.
- When  $VIF_j > 10$ :  $X_j$  is highly ~~correlated with other predictors~~.

Ex: Credit.CSV.

"Limit" "Rating" and "Education"

are highly inflated by the multicollinearity.

Recall from the heatmap; "Limit" and "Rating" are highly correlated.

### (C) Solutions.

Fast solutions:

- drop some of the highly correlated predictors based on practical understanding;  
i.e. I dropped "Limit" since I feel "Rating" is more important to stay.

- linearly combine some independent variables together when it makes sense: like when they are on the same scale. (dimension depth + length + width)  
or define a new one:  
use BMI ~~rat~~ instead of height and weight.

More complex solutions (~~Not~~ Not include here) :

- Use Regularization (Machine learning)
  - ↳ standardize the data
- Use other regression methods:  
Principal Components Analysis / Partial Least Squares.

Note: (1) Don't drop intercept easily. It's the constant  
~~~~~  
that adjust the scale of your response.

(2) If dropping one or two variables don't help, & you don't have to drop more.  
It's a balance between inflation and deflation!  
(battle between overfitting and underfitting again)

## II. Influential Points.

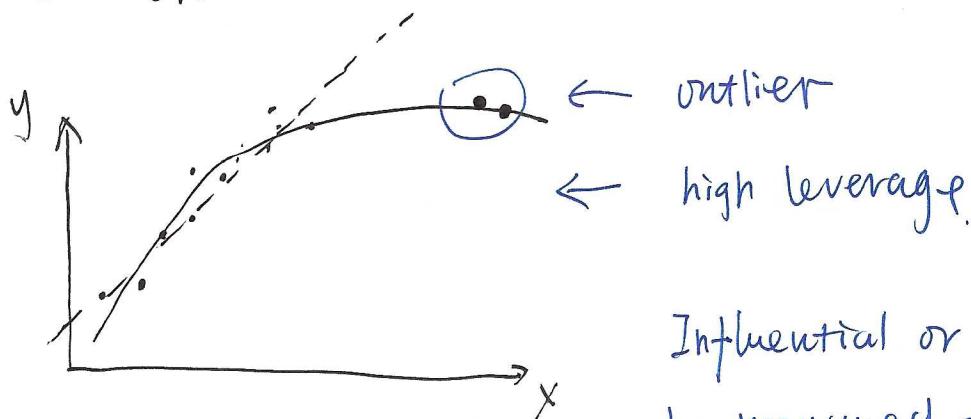
Data points can potentially be influential to your model

in different ways:

(i) Outlier: an observation has a response value  $y_i$  that is very far from the fitted value  $\hat{y}_i$ .

(ii) High leverage point: an observation has a particular unusual combination of predictor values is said to have high leverage, that could influence  $\hat{y}_i$ .

(iii) Observations with that's an outlier with high leverage.



Influential or not needed to be measured for the impact on your model.

(a) Problem: The influential points significantly affect the model estimation: with or without the observations lead to significantly different models.

## (b) Detection

(i) Just identifying outliers: obs with  $y_i$  "far away" from  $\hat{y}_i$

- Ordinary residuals? restricted by the units of measurement  
i.e. on weight,  $e_i = 20$  (lb)  
on length,  $e_i = 1$  (inch)  
which one is "bigger"?

- Standardize the residuals:

— semi-studentized residuals:  $e_i^* = \frac{e_i}{\sqrt{MSE}}$

"rule of thumb":  $|e_i^*| > 3$  can be identified as outlier.

(ii) Just identify high leverage: obs with  $X$  values that impact the  $\hat{y}_i$  a lot.

Recall:

$$\vec{y} = \vec{x}\vec{\beta} + \vec{\epsilon}$$

$$\vec{\beta} = (\vec{x}^\top \vec{x})^{-1} \vec{x}^\top \vec{y}$$

$$\hat{y} = \vec{x}\vec{\beta} = \underbrace{\vec{x}(\vec{x}^\top \vec{x})^{-1} \vec{x}^\top \vec{y}}_{= H\vec{y}} = H\vec{y}$$

$\hat{A}$ : "hat" matrix calculated by  $\vec{x}$

and measure when  $\vec{y}$  changes, how much it changes  $\vec{\hat{y}}$ .

look closer:

$$\begin{pmatrix} \vec{y}_1 \\ \vdots \\ \vec{y}_n \end{pmatrix} = \begin{pmatrix} h_{11} & h_{12} & \cdots & h_{1n} \\ h_{21} & h_{22} & \cdots & h_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ h_{n1} & h_{n2} & \cdots & h_{nn} \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$
$$\vec{\hat{y}} = \hat{A} * \vec{y}$$

Therefore:  $\vec{y}_i = h_{i1}y_1 + h_{i2}y_2 + \cdots + h_{in}y_n$

- $h_{ii}$  quantifies the influence that observed  $y_i$  has on its predicted value  $\vec{y}_i$ .
- $h_{ii}$  completely decided by  $\vec{x}$
- If  $h_{ii}$  is large,  $y_i$  plays a big role in estimating  $\vec{y}_i$

$h_{ii}$  is defined as the leverage of the  $i^{th}$  observation.

- Properties of leverage (NO math proof required)

$$h_{ii} = \vec{H}_{ii} = (\vec{x}(\vec{x}^T \vec{x})^{-1} \vec{x}^T)_{ii}$$

(i)  $0 \leq h_{ii} \leq 1$

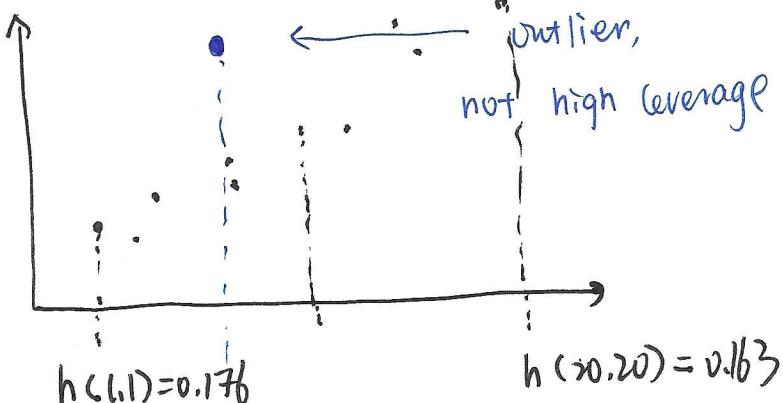
(ii)  $\sum h_{ii} = P$

(iii)  $\bar{h} = \frac{P}{n}$

(iv)  $h_{ii}$  measure the "distance" between  ~~$\vec{x}_i$  and  $\vec{\bar{x}}$~~

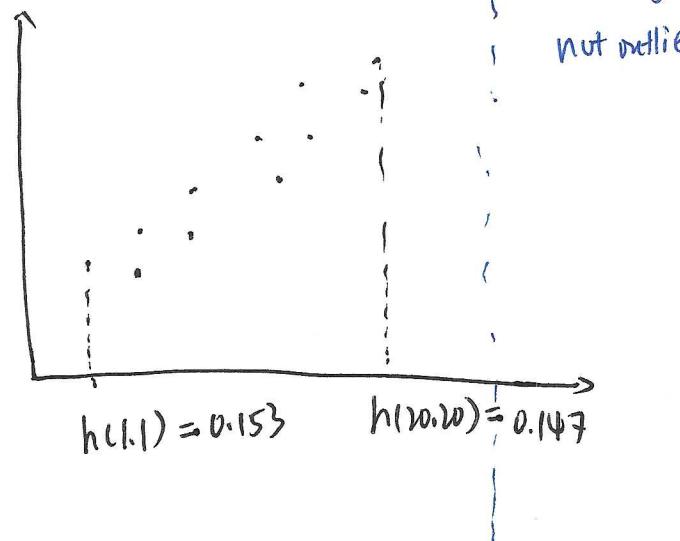
$\vec{x}_i$  and  $\vec{\bar{x}}$ .

high leverage,  
not outlier



$h(11,11) = 0.048$

not far from sample mean of  $X$



- Rule of Thumb: Obs with  $h_{ii} > \frac{2P}{n}$

flag as high leverage points.

(iii) Just being an outlier or high leverage points doesn't necessarily mean it's disturbing the model. We care more about how the observation influence the regression model — Influence Points

- Externaly Studentized Residuals

Recall semi-studentized residuals :  $e_i^* = \frac{e_i}{\sqrt{MSE}}$

"semi" because MSE is not the estimate of  $\text{Var}(e_i)$ , i.e.,  $\text{Var}(e_i) \neq \sigma^2$

Recall  $\vec{e} = \vec{y} - \hat{\vec{y}} = (I_n - \hat{H}) \vec{y}$ .

then  $\text{Var}(\vec{e}) = \text{Var}((I_n - \hat{H}) \vec{y}) = \sigma^2 (I_n - \hat{H}) \text{Var}(\vec{y}) (I_n - \hat{H})$

$$= \sigma^2 (I_n - \hat{H})$$

- so we can estimate  $\text{Var}(\vec{e})$  by  $MSE(I_n - \hat{H})$

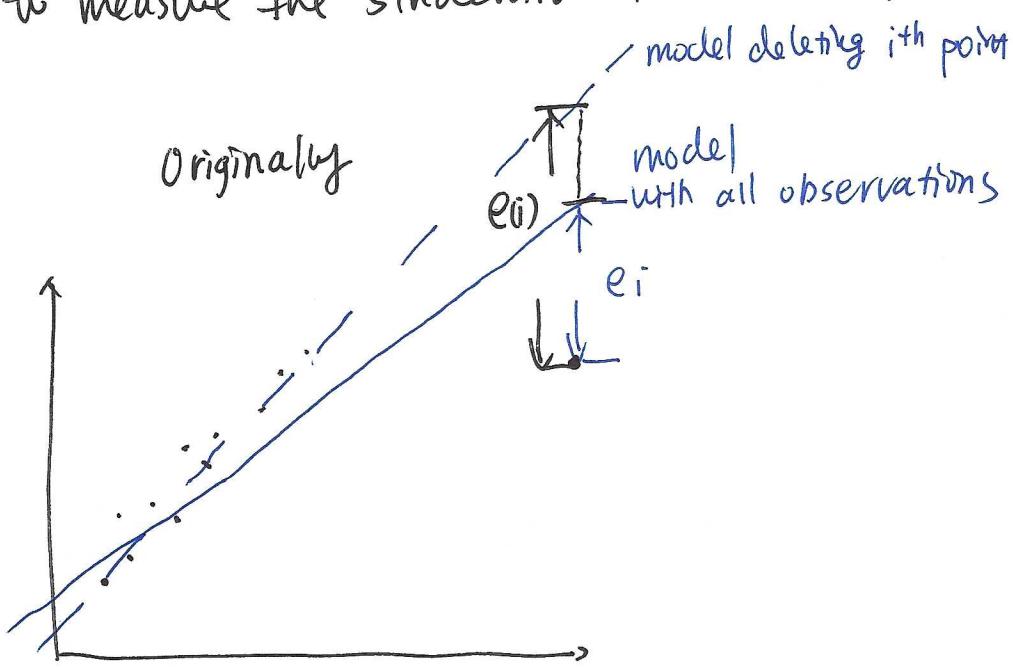
- $\boxed{\hat{\text{Var}}(e_i) = MSE(I - h_{ii})}$

- The "real" Studentized Residual =  $\frac{e_i}{\sqrt{MSE(I - h_{ii})}}$   
(Intervall)

Now we are trying to measure how much an obs impacts the regression line. When the influence is substantial, the line is "pulled" towards the observation. So we can measure this by deleting one obs at a time and fitting the model on the remaining ( $n-1$ ) observations.

Deleting one obs at a time and fitting the model on the remaining ( $n-1$ ) observations, then use  $e(i) = \hat{y}_i - \underbrace{\hat{y}_{(i)}}_{\substack{\uparrow \\ \text{from the model} \\ \text{with } i^{\text{th}} \text{ obs. deleted}}}$

to measure the studentized residual.



The Externally Studentized Residual is defined as

$$\text{stud}(e_i) = \frac{e_i}{s(e_{(i)})}$$

Text book p395~396

$$s(e_{(i)}) = \sqrt{\frac{MSE_{(i)}}{1-h_{ii}}}$$

$$e_{(i)} = \frac{e_i}{1-h_{ii}}$$

Note: index  $i$  indicates the model using all data

( $i$ ) indicates the model using data deleting obs.  $i$ .

Textbook p396.

$$\text{stud}_i(e_i) = \frac{e_i}{\sqrt{\frac{MSE(i)}{1-h_{ii}}}}$$

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

(model with data  
deleting obs i)

or

$$= \frac{e_i}{\sqrt{\frac{SSE(1-h_{ii}) - e_i^2}{n-p-1}}}$$

(model with full data)

- Rule of using  $\text{stud}_i(e_i)$  to identify outliers with high influence:

$$\text{stud}_i(e_i) \sim t_{\underbrace{n-1-p}_{\text{sample size.}}}$$

When  $|\text{stud}_i(e_i)| \geq t_{0.5}, df = n-1-p$ , we can detect the  $i^{th}$  obs. as an outlier with high influence.

Another might be more common one to detect influential points:

Cook's distance

It again consider refitting the data without observation;  
and it's a combination of  $e_i$  and  $h_{ii}$ :

The COOK's distance  $D_i$  for  $i^{th}$  obs:

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{PMSE} = \underbrace{\frac{e_i^2}{PMSE}} \cdot \underbrace{\frac{h_{ii}}{(1-h_{ii})^2}}$$

- $D_i$  measures directly how much ALL the fitted values changes after deleting  $i^{th}$  observation.
- Rule of thumbs: flag the observation when  $D_i > 1$   
(practical) Correction:  $D_i > 4/n$   
the exact comparison can be decided by a F distribution, not discussed here.

## Summary :

### • Multicollinearity

#### Problem

- unstable estimation of  $\vec{b}$
- unreliable t-test : harder to have rejections / significant predictors
- unreliable anova test.

#### Detection

- Heatmap plot for correlations in  $X$
- $VIF > 10$

#### Solution (so far)

- Delete one of the highly correlated variables
- Combine some of the variables together if applies.

### • Influential Points

#### Problem

- Observations that are outliers or with high leverage might influence the fitted model

#### Detection

- $|stds(e_i)| > t_{\alpha/2, df=n-1-p}$
- Cook's distance  $D_i > 1$   
Correction:  $D_i > 4/n$

#### Solution

- Report the observations
- Report the model result with and without the influential points.