

## L8. Modeling Problems - Other Diagnostics

### Modeling Assumption Violations:

(1) Heteroscedasticity

(2) Non-Normality Residuals

(3) False assumption of linearity between response variable  $Y$  and any predictor.

#### 1. Heteroscedasticity.

Recall Model Assumption:  $\epsilon_i \stackrel{\text{i.i.d}}{\sim} N(0, \sigma^2), i=1, \dots, n$

$$\text{or } \vec{\epsilon} \sim MVN(\vec{0}, \sigma^2 I_n)$$

The OLS method requires all  $\epsilon_i$ s have a constant variance  $\sigma^2$ , s.t. the parameter estimate

$\hat{b} = (\vec{x}' \vec{x})^{-1} \vec{x}' \vec{y}$  is the "best" estimator  
in "BLUE".

- Heteroscedasticity occurs when this assumption is violated, which means the variance of ~~residuals~~<sup>errors</sup> are non-constant, more commonly, it refers to the spread of residuals changes systematically with predictors.
  - (a) Problem.
    - (i) The OLS estimate  $\hat{b}$  is still linear and unbiased, but not the "best" anymore: we can't say  $\hat{b}$  has the smallest variance among all unbiased linear estimators. There is another estimator with a smallest variance.
    - (ii)  $se(\hat{\beta}_k)$  for all OLS output are incorrect estimates of the standard deviation of  $\hat{\beta}_k$ . It will result in misleading t-test and C.I.
    - (iii) Predictions of  $y$  are still unbiased, but the prediction intervals are incorrect.

$$\begin{aligned}
 \text{Recall: } \text{Var}(\vec{b}) &= E(\vec{b} - \vec{\beta})(\vec{b} - \vec{\beta})^T \\
 &= E((\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{\varepsilon} \vec{\varepsilon}^T \vec{X} (\vec{X}^T \vec{X})^{-1}) \quad \text{lecture 4} \\
 &= (\vec{X}^T \vec{X})^{-1} \vec{X}^T \cdot E(\vec{\varepsilon} \vec{\varepsilon}^T) \cdot \vec{X} (\vec{X}^T \vec{X})^{-1}
 \end{aligned}$$

When applying OLS method:

$$\text{assume } E(\vec{\varepsilon} \vec{\varepsilon}^T) = \sigma^2 I_n = \begin{pmatrix} \sigma^2 & & \\ & \ddots & 0 \\ 0 & & \sigma^2 \end{pmatrix}$$

$$\text{so } \text{Var}(\vec{b})_{\text{OLS}} = \sigma^2 (\vec{X}^T \vec{X})^{-1}$$

and is estimated by  $\text{MSE}(\vec{X}^T \vec{X})^{-1}$

which is an unbiased estimator of  
 $\text{Var}(\vec{b})$  when the assumption of constant  
 variance holds.

When Heteroscedasticity exist:

$$\text{the real } \text{Var}(\vec{b}) = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & 0 \\ 0 & & \sigma_n^2 \end{pmatrix} \vec{X} (\vec{X}^T \vec{X})^{-1}$$

so  $\text{MSE}(\vec{X}^T \vec{X})^{-1}$  is now biased  
 estimate for  $\text{Var}(\vec{b})$

In this case, if we still apply OLS method:

(i)  $\text{Var}(\vec{b})_{\text{OLS}} \neq \text{Var}(\vec{b})$

Not "best" anymore

(ii) ~~the~~  $\overset{\wedge}{\text{Var}}(\vec{b})_{\text{OLS}} = \text{MSE} \cdot (\vec{x}^T \vec{x})^{-1}$  is

but unbiased anymore, so

$\text{se}(\beta_k)_{\text{OLS}}$  is incorrect.

(iii) Let's give some result about the prediction first

Textbook section 6.7

For a new observation with  $\vec{x}_h = \begin{pmatrix} 1 \\ x_{h1} \\ \vdots \\ x_{h,p-1} \end{pmatrix}$

The fitted value  $\hat{Y}_h = \vec{x}_h^T \vec{b}$

- unbiased :  $E(\hat{Y}_h) = \vec{x}_h^T \vec{\beta} = E(Y_h)$

- $\text{Var}(\hat{Y}_h) = \text{Var}(\vec{x}_h^T (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y})$

When assumption holds

$$= \sigma^2 \vec{x}_h^T (\vec{x}^T \vec{x})^{-1} \vec{x}^T$$

- Confidence Interval for the Mean Response.  $\mu_h$

$$\hat{Y}_h \pm t_{\alpha/2, df=n-p} \text{se}(\hat{Y}_h)$$

where  $\text{se}(\hat{Y}_h) = \sqrt{\text{MSE} \hat{X}_h^T (\hat{X}^T \hat{X})^{-1} \hat{X}^T}$

- Prediction Interval for the new observation  $\hat{Y}_h$

$$\hat{Y}_h \pm t_{\alpha/2, df=n-p} \sqrt{\text{se}(\hat{Y}_h)^2 + \text{MSE}}$$

When Heteroscedasticity exists.

$$\begin{aligned} \text{Var}(\hat{Y}_h) &= \text{Var}(\hat{X}_h^T (\hat{X}^T \hat{X})^{-1} \hat{X}^T \hat{y}) \\ &= \hat{X}_h^T (\hat{X}^T \hat{X})^{-1} \hat{X}^T \underbrace{\text{Var}(\hat{y})}_{\neq \sigma^2 I_n} \end{aligned}$$

So  $\text{MSE} \hat{X}_h^T (\hat{X}^T \hat{X})^{-1} \hat{X}^T$  is not unbiased estimate anymore.

Therefore prediction/confidence interval calculated this way would be incorrect.

(b) Detection.

(i) Residual v.s. Fitted Value Plot:

You would observe obvious change of bandwidth  
in the plot

(ii) Breusch-Pagan Test

Basic idea is the variance of error should not  
change given difference predictor values.

If the assumption of constant variance is  
violated, then the variance would be changed  
with predictor values.

1) Fit your model  $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{ip-1} + \epsilon_i$

2) Obtain the residuals  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$

3) Build an auxiliary regression model:

$$\epsilon_i^2 = \gamma_0 + \gamma_1 x_{i1} + \dots + \gamma_{p-1} x_{ip-1} + \zeta_i \quad (E)$$

4)  $H_0: \gamma_1 = \gamma_2 = \dots = \gamma_{p-1} = 0$  v.s  $H_a: \text{at least } \gamma_i \neq 0$

5) Test stat :  $\chi_s^2 = nR^2$

where  $R^2$  is the w.e. of determination of model (E)

Where  $H_0$  is true:  $\chi_s^2 \sim \chi^2(p)$

6) Test Result:

When Fail to reject  $H_0$ :  $\chi_s^2 < \chi_{d, df=p}^2$  or p-value >  $\alpha$

We conclude there's not significant heteroscedasticity.

when reject  $H_0$ :  $\chi_s^2 > \chi_{d, df=p}^2$  or p-value <  $\alpha$ .

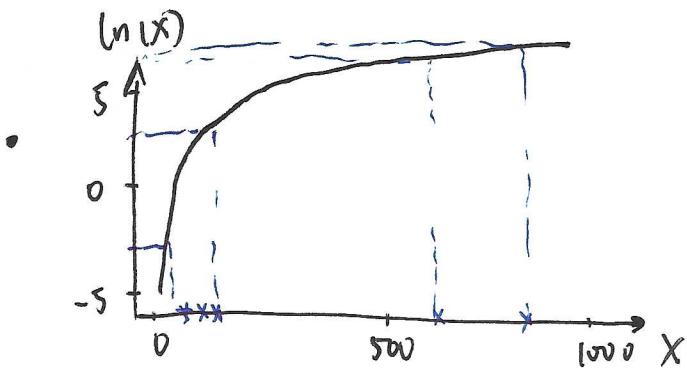
We conclude there's significant heteroscedasticity problem.

(iii) White Test \* is another common one.

### (c) Solution

#### (i) log-transformation on y

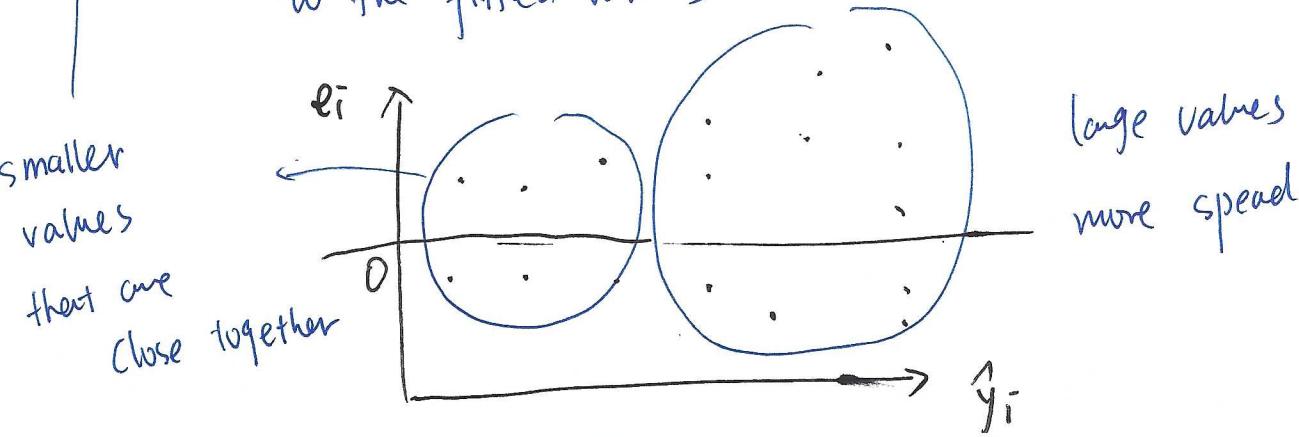
- Natural logarithm,  $\ln$ , transformation is the most common transformation in MLR



From the plot, the effects of taking natural log transformation are:

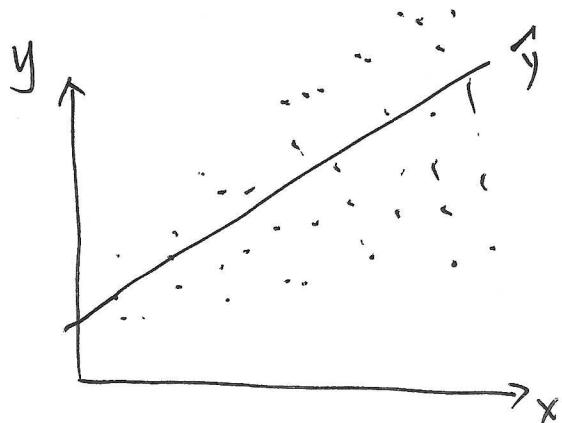
- small values that are close together are spread further out
- large values that are spread out are brought close together

Therefore, it works when the residuals changes proportionally to the fitted values



$$\text{since } e_i = y_i - \hat{y}_i$$

the heteroscedasticity problem lands on "y" in the data.



distribution of y around the fitted line  
also changes from clusters of  
smaller values to the spread  
of large values.

- In this case, a natural-log transformation on y

could help "even-out" the points:

$$y_{\text{new}} = \ln(y)$$

smf.ols("y\_{\text{new}} \sim x\_1 + \dots + x\_{p-1}")

$$\hat{y}_{\text{new}} = \ln(\hat{y}) \rightarrow \hat{y} = e^{\hat{y}_{\text{new}}}$$

Since WLS only works for the case of  $e_i$  is proportional with  $\hat{y}_i$ , or "Funnel" shape, other solutions can be

(ii) Weighted least square regression.

Recall in OLS: assumes  $\text{Var}(\vec{\varepsilon}) = \begin{pmatrix} \sigma^2 & & \\ & \ddots & 0 \\ 0 & & \sigma^2 \end{pmatrix}$

$$\text{find } \vec{\beta}_{\text{OLS}} = \arg \min_{\vec{\beta}} \text{SS}(\vec{\beta})_{\text{OLS}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$$

$$\text{where } \text{SS}(\vec{\beta})_{\text{OLS}} = \sum_{i=1}^n e_i^2$$

$$\text{and } \text{Var}(\vec{\beta}_{\text{OLS}}) = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \cdot \underbrace{E(\vec{\varepsilon} \vec{\varepsilon}^T)}_{= \text{Var}(\vec{\varepsilon})} \vec{X} (\vec{X}^T \vec{X})^{-1}$$

$$= \sigma^2 (\vec{X}^T \vec{X})^{-1}$$

and  $\beta$  estimated by  $\text{MSE}(\vec{X}^T \vec{X})^{-1}$

$$\text{so } \text{se}(\vec{\beta}_{\text{OLS}})_{k+1} = [\text{MSE}(\vec{X}^T \vec{X})^{-1}]_{k+1, k+1}$$

When heteroscedasticity exists.

Since  $\ln y$  only works for the case of  $e_i$  is proportional with  $y_i$ .

or "Funnel" shape, other solutions can be

(iii) Weighted least-square regression. — count different variances in the estimation.

Recall OLS assumes  $\text{Var}(\vec{\epsilon}) = \begin{pmatrix} \sigma^2 & & \\ & \ddots & 0 \\ 0 & & \sigma^2 \end{pmatrix}$

while reality in heteroscedasticity is

$$\text{Var}(\vec{\epsilon}) = \begin{pmatrix} \sigma_1^2 & & 0 & \\ & \sigma_2^2 & & 0 \\ 0 & & \ddots & \\ & & & \sigma_n^2 \end{pmatrix} = \Sigma$$

$$\text{Var}(\vec{b}) = (\vec{X}^\top \vec{X})^{-1} \vec{X}^\top \cdot \underbrace{E(\vec{\epsilon}\vec{\epsilon}^\top)}_{\text{aka } \text{Var}(\vec{\epsilon})} \cdot \vec{X} (\vec{X}^\top \vec{X})^{-1}$$

$$\text{aka } \text{Var}(\vec{\epsilon}) = \Sigma$$

We need to estimate  $\Sigma$  by  $\hat{\Sigma} = \begin{pmatrix} e_1^2 & & 0 & \\ & e_2^2 & & 0 \\ 0 & & \ddots & \\ & & & e_n^2 \end{pmatrix}$

then estimate  $\text{Var}(\vec{b})$  by  $(\vec{X}^\top \vec{X})^{-1} \vec{X}^\top \hat{\Sigma} \vec{X} (\vec{X}^\top \vec{X})^{-1}$

$$\text{Var}(\vec{\epsilon}) = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{pmatrix} = \Sigma$$

- Weighted-least-square estimates take account of the different variances in the way that
  - An observation with small error variation should weigh more in the model.
  - An observation with bigger error variation should weigh less in the model.
  - It includes the  $\Sigma$  in the estimation in the way that

Define the Weight Matrix

$$W = \begin{pmatrix} w_1 & & & \\ & w_2 & & \\ & & \ddots & \\ & 0 & & w_n \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & \\ & & \ddots & \\ & 0 & & \frac{1}{\sigma_n^2} \end{pmatrix}$$

- Defined the weighted sum of squares:

$$SS_{WLS}(\vec{\beta}) = \sum w_i e_i^2$$

and  $\vec{b}_{WLS} = \arg\min SS_{WLS} = (\vec{X}^\top \hat{W} \vec{X})^{-1} \vec{X}^\top \hat{W} \vec{Y}$

Where  $\hat{W} = \boxed{?}$

Practically how to decide the weights can be tricky:

If we assume  $e_1 \sim N(0, \sigma_1^2), \dots, e_n \sim N(0, \sigma_n^2)$

we only have one observed sample for  $e_i$  from each  $\epsilon_i$ .

Here we just give an simple example, it's NOT required to run WLS yourself:

When  $e_i$  v.s.  $X_k$  has a funnel shape,

you can run model = " $e_i^2 \sim X_k$ "

then use fitted values of the model as the estimated variance of  $\text{Var}(e_i)$

$$\text{then } w_i = 1/e_i^2$$

(iii) Generalized Least squares: again deal with the estimate of  $\Sigma$  in a broader choices of cases.

logistic model is a case of GLM.

(iv) Use more Robust Regression: which is not sensitive to assumptions.

(2) Sometimes, after correcting Non-Normality, Non-linearity and model selection, heteroscedasticity can be greatly improved as well.

## Summary:

- Heteroscedasticity : The assumption of constant variance of  $\epsilon_i$  is violated.

### Problem

- In OLS result  
 $se(\hat{\beta}_k)$  is incorrect, resulting in misleading t-test and anova test results.
- $Var(\hat{Y}_h)$  is incorrect, resulting in misleading C.I. and prediction intervals.
- The estimation of coefficients and fitted values are still unbiased, therefore reliable.

### Detection

- Residual v.s. Fitted plot  
the bandwidth changes
- Breusch-Pagan Test or White Test :  
p-value  $\leq 0.05$  indicated serious heteroscedasticity problem.

### Solution:

- (i) If  $\epsilon_i$  v.s.  $\hat{y}_i$  is funnel shape, perform natural log transformation on  $y$ .
- (ii) Perform Weighted-least-square regression.

\* Practically hard to decide a good weight

(iii) Use Robust Regression methods in ML:

Huber, RANSAC, Theil Sen