

# Outlier Detection

Recall LM  
LNp.7-1~7-7

## influential obs

- Minitab identifies two types of outliers denoted by R and X:
  - R: its standardized residual  $(y_i - \hat{y}_i)/se(\hat{y}_i)$  is large.
  - X: its X value gives large leverage (i.e., far away from majority of the X values).

dangerous  
to exclude  
in an  
automatic  
manner

- For the mortality data, the observation with  $T = 31.8$ ,  $M = 67.3$  (i.e., left most point in plot on LNp.2-2) is identified as both R and X.

report

- After removing this outlier and refitting the remaining data, the output is given on LNp.2-11. There is still an outlier identified as X but not R. This one (second left most point on LNp.2-2) should not be removed (why?)

Residual plots on LNp.2-12 show no systematic pattern.

obs. with large  
leverage exist  
in most data

Notes: Outliers are not discussed in the book, see standard regression texts.

Residual plots will be discussed in unit 3.

overall pattern  $\leftrightarrow$  cf. unusual observations

## Regression Results after Removing the Outlier

cf.

results in LNp.2-9

The regression equation is

$$M = -52.62 + 3.02 T$$

Predictor	Coef	SE Coef	T	P
Constant	-52.62	15.82	-3.33	0.005
T	$\hat{\beta}_1 \rightarrow 3.0152$	0.3466	8.70	0.000

$> 2.3577$  (LNp.2-9)

$> 6.76$  (LNp.2-9), d.f. = 1

$S = 5.93258$      $R\text{-Sq} = 85.3\%$      $R\text{-Sq}(\text{adj}) = 84.2\%$

$< 7.54466$  (LNp.2-9)

$> 76.5\%$  (LNp.2-9)

Analysis of Variance

Source	DF	SS	MS	F	P
Regression	1	2664.3	2664.3	75.70	0.000
Residual Error	13	457.5	35.2		
Total	14	3121.9			

$> 2599.5$  (LNp.2-9)

cf.  $< 3396.4$  (LNp.2-9)

Unusual Observations

Obs	T	M	Fit	SE Fit	Residual	St Resid
15	34.0	52.50	49.90	4.25	2.60	0.63 X

original 16th obs

X denotes an observation whose X value gives it large leverage.

## Residual Plots After Outlier Removal

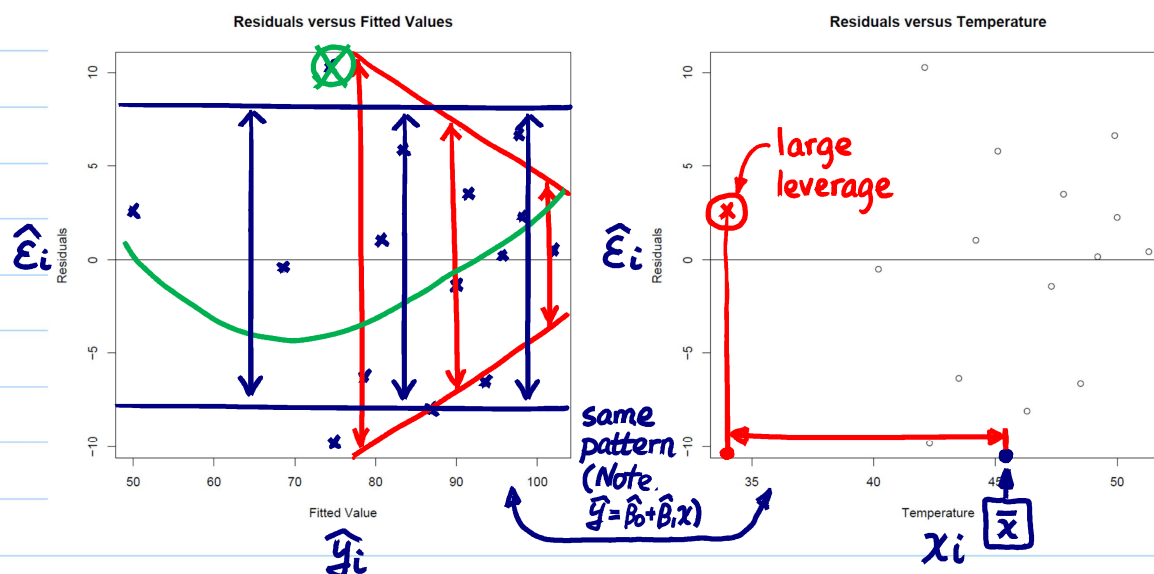


Figure 3: Residual Plots

**Comments :** No systematic pattern is discerned.

1. non-constant variance
2. curvature in the mean of residuals

objective of many regression analysis

**Prediction from the Breast Cancer Data** p. 2-13  
 ① for mean response  $E(y_x) = \mu_x$  ② for future obs.  $y_x$   
 ① interpolation ② extrapolation

- The fitted regression model is  $Y = -21.79 + 2.36X$ , where  $Y$  denotes the mortality rate and  $X$  denotes the temperature.  $\hat{\beta}_0$   $\hat{\beta}_1$
- The predicted mean of  $Y$  at  $X = x_0$  can be obtained from the above model. For example, prediction for the temperature of 49 is obtained by substituting  $x_0 = 49$ , which gives  $y_{x_0} = 93.85$ .  $\hat{y}_{x_0} \rightarrow y_{x_0} = \mu_{x_0} + \epsilon$   
 $\hat{\mu}_{x_0} \rightarrow \mu_{x_0} = \beta_0 + \beta_1 x_0$
- The standard error of  $\hat{\mu}_{x_0}$  is given by

$$\text{Var}(\hat{\mu}_{x_0}) = \text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0)$$
  
 (exercise, note.  $\text{cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$ )

$$s.e.(\hat{\mu}_{x_0}) = \hat{\sigma} \sqrt{\frac{1}{N} + \frac{1}{N} \frac{(\bar{x} - x_0)^2}{\sum_{i=1}^N (x_i - \bar{x})^2 / N}}$$

Note 1. It's a function of  $x_i$ 's only

Note 2. What happens if  $x_0$  is away from  $\bar{x}$ ?

- Here  $x_0 = 49$ ,  $1/N + (\bar{x} - x_0)^2 / \sum_{i=1}^N (x_i - \bar{x})^2 = 0.1041$ , and  $\hat{\sigma} = \sqrt{MSE} = 7.54$ . Consequently,  $s.e.(\hat{\mu}_{x_0}) = 2.432$ .

Note 3. It converges to zero when  $N \rightarrow \infty$

estimate  
 $\pm$  (critical value)  
 $\times$  s.e.(estimate)

## Confidence interval for mean and prediction interval for future observation

 $\mu_{x_0}$ 

$$y_{x_0} = \mu_{x_0} + \varepsilon$$

- A 95% confidence interval for the mean response  $\mu_{x_0} = \beta_0 + \beta_1 x_0$  at  $x = x_0$  is

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t_{N-2, 0.025} \times s.e.(\hat{\mu}_{x_0}).$$

Same

estimate  $\hat{\mu}_{x_0}$  critical value s.e.(estimate)

- Here the 95% confidence interval for the mean mortality corresponding to a temperature of 49 is [88.63, 99.07].

$$E[\hat{\beta}_0 + \hat{\beta}_1 x_0 - (\beta_0 + \beta_1 x_0 + \varepsilon)]^2$$

- A 95% prediction interval for an individual observation  $y_{x_0}$  corresponding to  $x = x_0$  is

Note It converges to  $\sigma^2$  when  $N \rightarrow \infty$

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t_{N-2, 0.025} \times \hat{\sigma} \sqrt{1 + \frac{1}{N} + \frac{(\bar{x} - x_0)^2}{\sum_{i=1}^N (x_i - \bar{x})^2}}$$

estimate  $\hat{y}_{x_0}$  critical value

s.e.(estimate)

$$E(\hat{y}_{x_0} - y_{x_0})^2 = \text{Var}(\hat{\mu}_{x_0}) + \text{Var}(\varepsilon) = \sigma^2$$

where 1 under the square root represents  $\sigma^2$ , variance of the new observation  $y_{x_0}$ .

- The 95% prediction interval for the predicted mortality of an individual corresponding to the temperature of 49 is [76.85, 110.85]. becomes wider

many concept/idea/formula in MLR are similar to those in SLR

## Multiple Linear Regression : Air Pollution Data

more than one predictor / effect

<http://lib.stat.cmu.edu/DASL/Stories/AirPollutionandMortality.html>

- Data collected by General Motors.

\* true model

$$Y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$$

- Response is age-adjusted mortality.

\* fitted model

$$Y = X_1 \beta_1 + \varepsilon'$$

- Predictors :

Why?

- Variables measuring demographic characteristics.

- Variables measuring climatic characteristics.

- Variables recording pollution potential of 3 air pollutants.

their roles are similar to block factors

- Objective : To determine whether air pollution is significantly related to mortality.

relationship

Q: Why should we put demographic & climatic variables in the model when we are only interested in ~

## Predictors

climate

1. **JanTemp** : Mean January temperature (degrees Fahrenheit)
2. **JulyTemp** : Mean July temperature (degrees Fahrenheit)
3. **RelHum** : Relative Humidity
4. **Rain** : Annual rainfall (inches)

demographic

5. **Education** : Median education
6. **PopDensity** : Population density
7. **%NonWhite** : Percentage of non whites
8. **%WC** : Percentage of white collar workers
9. **pop** : Population
10. **pop/house** : Population per household
11. **income** : Median income

pollutant

12. **HCPot** : HC pollution potential
13. **NOxPot** : Nitrous Oxide pollution potential
14. **SO2Pot** : Sulphur Dioxide pollution potential

## Initial Data Analysis ← Getting Started →

p. 2-17

- data cleaning
- need transformation?
- scatter plots
- histogram
- descriptive statistics
- ...

- There are 60 data points.

• Pollution variables are highly skewed, log transformation makes them nearly symmetric. The variables HCPot, NOxPot and SO2Pot are replaced by  $\log(\text{HCPot})$ ,  $\log(\text{NOxPot})$  and  $\log(\text{SO2Pot})$ .

- Observation 21 (Fort Worth, TX) has two missing values, so this data point will be discarded from the analysis. → 59 data points.

Why? Note,  
response is  
not skewed  
in the case.





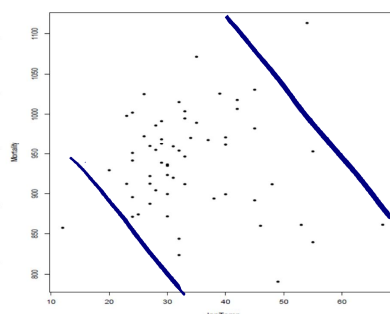
# Scatter Plots

"rain" ④  
significant  
(LNp.2-24)

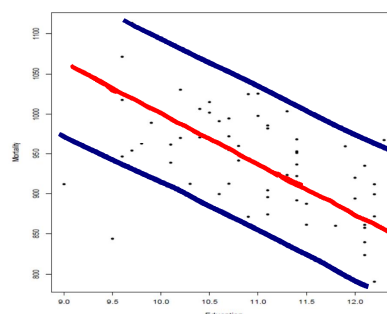
cf., why?  
check ★  
in LNp.2-15

Figure 4: Scatter Plots of mortality against selected predictors

(a) JanTemp



(b) Education

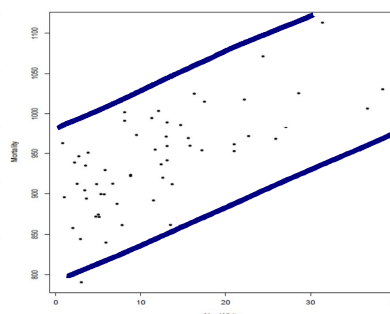


not significant  
(LNp.2-24)

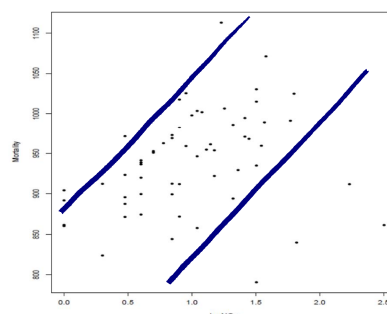
simple  
regression:

$$y = \beta_0 + \beta_1 x_5 + \varepsilon$$

(c) NonWhite



(d) Log(NOxPot)



## Fitting the Multiple Regression Equation

estimate  $\beta$

unknown parameter

$$y_x = \sum_{j=0}^k \beta_j \times g_j(x_1, \dots, x_m) + \varepsilon$$

known function ( $g_0 = 1$ , intercept)

functional form

- Underlying Model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2).$$

$E(y_x)$

matrix form

$$Y = X\beta + \varepsilon$$

$$\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}$$

- Coefficients are estimated by minimizing

$$S(\beta) \equiv \sum_{i=1}^N \left( y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}) \right)^2 = (y - X\beta)'(y - X\beta).$$

$\partial S / \partial \beta = 0$

$$\sum_i \beta_i g_i = X\beta$$

- Least Squares estimates:

$$\hat{\beta} = (X'X)^{-1}X'y$$

$\Rightarrow E(\hat{\beta}) = \beta$  if model correct

- Variance-Covariance matrix of  $\hat{\beta}$ :  $\Sigma_{\hat{\beta}} = \sigma^2 (X'X)^{-1}$

depend only on  $x_i$ 's

$(X'X)_{jl}$ : inner product of the  $j$ th &  $l$ th columns in  $X$ .

a row: one group of observations

a column: one effect

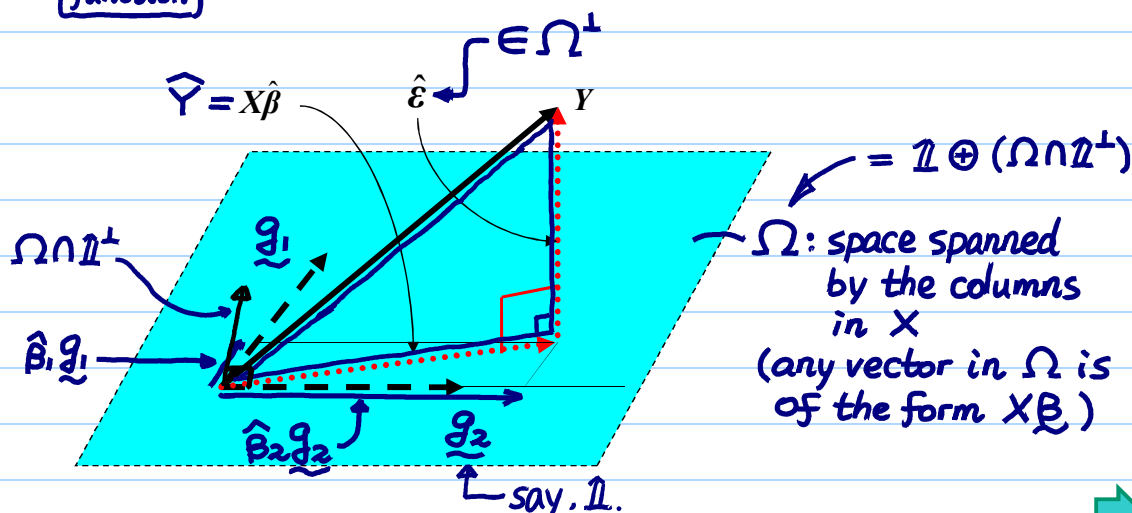
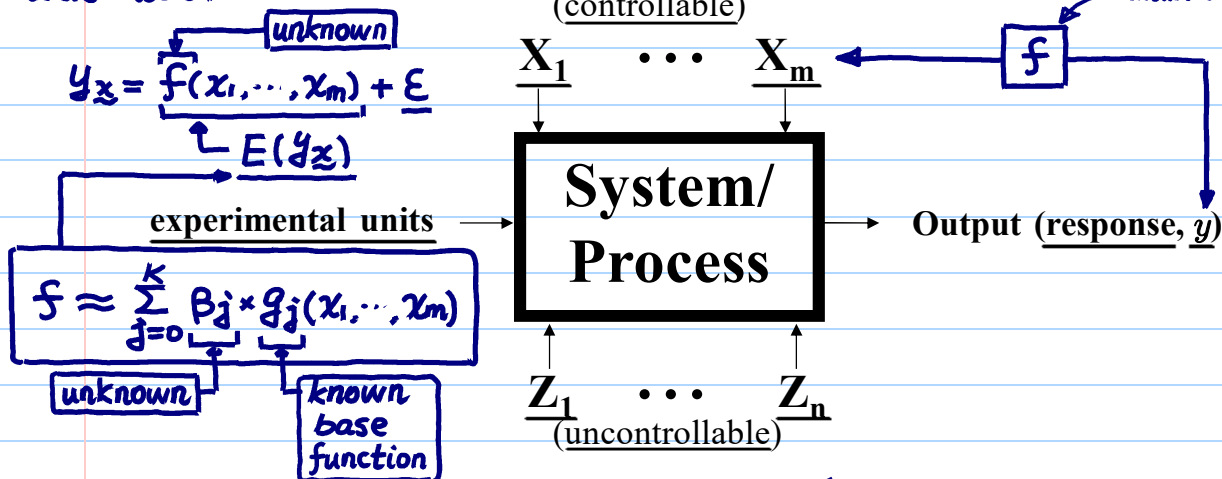
	$y$	$x_1$	$\dots$	$x_k$
$y_1$	$y_1$	$x_{11}$	$\dots$	$x_{k1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$y_n$	$y_n$	$x_{n1}$	$\dots$	$x_{nk}$

$$g_{ij} = g_j(x_{i1}, \dots, x_{ik})$$

model matrix  $X$

$Y$	$1$	$g_1$	$g_2$	$\dots$	$g_p$
$y_1$	1	$g_{11}$	$g_{12}$	$\dots$	$g_{1p}$
$y_2$	1	$g_{21}$	$g_{22}$	$\dots$	$g_{2p}$
$\dots$					
$y_n$	1	$g_{n1}$	$g_{n2}$	$\dots$	$g_{np}$

true model



## Analysis of Variance

overall F-test

$$(\omega) H_0: y = \beta_0 + \varepsilon$$

$$(\Omega) H_A: y = \beta_0 + \sum_{j=1}^k \beta_j g_j + \varepsilon$$

- The total variation in  $y$ , i.e., corrected total sum of squares,

$$CTSS = \sum_{i=1}^N (y_i - \bar{y})^2 = \mathbf{y}^T \mathbf{y} - N\bar{y}^2, \text{ can be decomposed into two parts}$$

(Analysis of Variance (ANOVA))

規律

隨機

RSS under  $\Omega$ RSS under  $\omega$ 

$$CTSS = \text{RegrSS} + \text{RSS}$$

where  $\text{RSS} = \text{Residual sum of squares} = \sum (y_i - \hat{y}_i)^2 = (\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})$  $\text{RegrSS} = \text{Regression sum of squares} = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2$ mean of  $\hat{y}_i$ 's =  $\bar{y}$ 

ANOVA Table

$$y_i \text{'s} \rightarrow y, \hat{y}_i \text{'s} \rightarrow X\hat{\beta}, \bar{y} \rightarrow \bar{y}$$

check the graph in Lnp.2-4

Source of variation

Source	Degrees of Freedom	Sum of Squares	Mean Squares
regression	$k$ $\left[ \beta_1, \dots, \beta_k \right]$	$\hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2$	$(\hat{\beta}^T \mathbf{X}^T \mathbf{X} \hat{\beta} - N\bar{y}^2) / k$
residual	$N - (k + 1)$	$(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta})$	$(\mathbf{y} - \mathbf{X}\hat{\beta})^T (\mathbf{y} - \mathbf{X}\hat{\beta}) / (N - k - 1)$
total	$N - 1$	$\mathbf{y}^T \mathbf{y} - N\bar{y}^2$	

F-statistic

## Explanatory Power of the Model

規律

隨機

- $R^2 = \frac{RegrSS}{CTSS} = 1 - \frac{RSS}{CTSS}$  measures of the proportion of variation in  $y$  explained by the fitted model.  $R$  is called the multiple correlation coefficient.

↑  $cor(\underline{y}, \underline{\hat{y}}) \leftarrow$  check the graph in LN p. 2-20

- Adjusted  $R^2$  :

$$\underline{R_a^2} = 1 - \frac{\frac{RSS}{N-(k+1)}}{\frac{CTSS}{N-1}} = 1 - \left( \frac{N-1}{N-k-1} \right) \frac{RSS}{CTSS} \leq R^2$$

↑  $\geq 1$

∴  $RSS$  always decreases

- When an additional predictor is included in the regression model,  $R^2$  always increases. This is not a desirable property for model selection. However,  $R_a^2$  may decrease if the included variable is not an informative predictor.

Usually  $R_a^2$  is a better measure for comparing different model fits.

Q: When can  $R^2$  reach 1?

Ans.  $\hat{y}_i = y_i$ , only when data have no replicates & # of parameters  $\beta$  = # of obs

total df available  
N  
↑  
sample size

df available for studying  $\beta$   
df available for estimating  $\sigma^2$

Want to have a final fitted model such that

1. # of parameters : small
2.  $RSS$  : small

## Testing significance of coefficients : t-Statistic

Examine whether  $\beta_j = 0$  when other effects  $\beta_i$ 's ( $g_i$ 's),  $i \neq j$ , are still in the model

- To test the null hypothesis  $H_0 : \beta_j = 0$  against the alternative hypothesis  $H_A : \beta_j \neq 0$  under the full model, use the test statistic

Note: collinearity

$$t_j = \frac{\hat{\beta}_j - 0}{s.d.(\hat{\beta}_j)}$$

↑

null model

$$(\omega) H_0 : y = \beta_0 + \dots + \beta_{j-1} g_{j-1} + \beta_{j+1} g_{j+1} + \dots + \epsilon$$

$$(\Omega) H_A : y = \beta_0 + \sum_{i=1}^k \beta_i g_i + \epsilon$$

alternative model

- The higher the value of  $|t_j|$ , the more significant is the coefficient.
- In practice, if  $p$ -value is less than  $\alpha = 0.05$  or  $0.01$ ,  $H_0$  is rejected.
- Confidence Interval :  $100(1 - \alpha)\%$  confidence interval for  $\beta_j$  is given by

null dist., under  $H_0$ ,  
 $t_j \sim t_{N-(k+1)}$

$$\hat{\beta}_j \pm t_{N-(k+1), \frac{\alpha}{2}} \times s.d.(\hat{\beta}_j)$$

estimate      critical value      s.d.(estimate)

$$H_0 : \beta_j = \beta_j^*$$

$$\left| \frac{\hat{\beta}_j - \beta_j^*}{s.d.(\hat{\beta}_j)} \right| \leq t_{N-(k+1), \frac{\alpha}{2}}$$

acceptance region

where  $t_{N-k-1, \frac{\alpha}{2}}$  is the upper  $\alpha/2$  point of the  $t$  distribution with  $N - k - 1$  degrees of freedom.

df<sub>RSS</sub>

If the confidence interval for  $\beta_j$  does not contain 0, then  $H_0$  is rejected.



check graphs  
in LNp.2-18

## Analysis of Air Pollution Data

Predictor	$\hat{\beta}_j$ 's	Coef	SE Coef	$s.d(\hat{\beta}_j)$ 's	T	t-test	P
⑧ $\beta_0$ Constant		1332.7	291.7		4.57	0.000 ✓	
① JanTemp		-2.3052	0.8795		-2.62	0.012 ✓	
② JulyTemp		-1.657	2.051		-0.81	0.424	
③ RelHum		0.407	1.070		0.38	0.706	
④ Rain		1.4436	0.5847		2.47	0.018 ✓	
⑤ Educatio		-9.458	9.080		-1.04	0.303	
⑥ PopDensi		0.004509	0.004311		1.05	0.301	
⑦ %NonWhit		5.194	1.005		5.17	0.000 ✓	
⑧ %WC		-1.852	1.210		-1.53	0.133	
⑨ pop		0.00000109	0.00000401		0.27	0.788	
⑩ pop/hous		-45.95	39.78		-1.16	0.254	
⑪ income		-0.000549	0.001309		-0.42	0.677	
⑫ logHC		-53.47	35.39		-1.51	0.138	
⑬ logNOx		80.22	32.66		2.46	0.018 ✓	
⑭ logSO2		-6.91	16.72		-0.41	0.681	

Q: Can we remove all insignificant effects simultaneously to simplify the model?  
Ans. No, in general, but OK if orthogonality exists.

how to interpret it?

It becomes insignificant when other effects are in the model (cf graph in LNp.2-18)  
Possible reason: collinearity

S = 34.58      R-Sq = 76.7%      R-Sq(adj) = 69.3%

Analysis of Variance

$$R^2 = 1 - \frac{1}{1 + \frac{k}{N-(k+1)} F}$$

Source	DF	SS	MS	F	P
Regression	14	173383	12384	10.36	0.000
Residual Error	44	52610	1196		
Total	58	225993			

59 obs.

- formulation of hypothesis testing from the view of comparing models

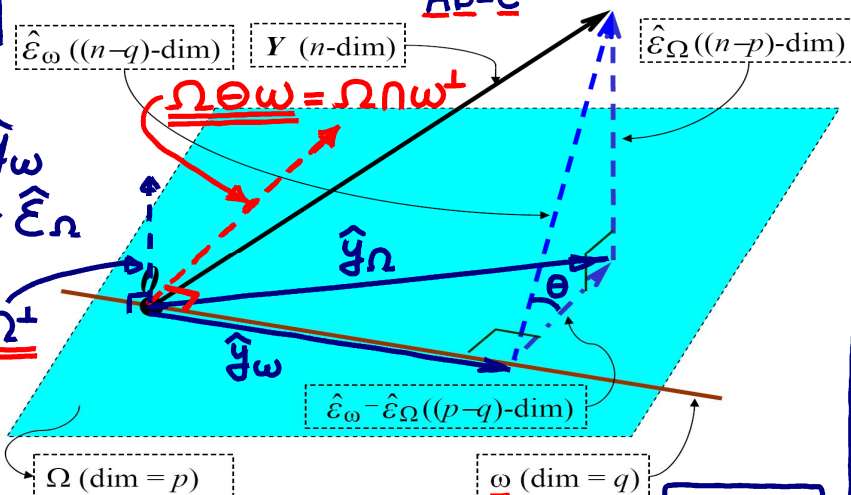
➤ a model space  $\equiv$  the space spanned by columns of some  $X$  effects

➤ consider a large model space,  $\Omega$ , and a smaller model space,  $\omega$ , where  $\omega \subset \Omega$ , i.e.,  $\omega$  represents a subset/a subspace of  $\Omega$ . Suppose dimension (# of parameters) of  $\Omega$  is  $p$  and  $\dim(\omega) = q$ , where  $p > q$ .  $\rightarrow df_{\Omega} = n - p, df_{\omega} = n - q$

➤ to answer "which of the model spaces is more adequate" in statistical language

$\Rightarrow$  perform the test  $H_0: \omega$  v.s.  $H_A: \Omega \setminus \omega$  check examples in LNp.2-21&23

$H_0: \beta_j = \beta_j^* = 0$



•  $\theta$  large  $\Rightarrow \hat{e}_{\omega} \approx \hat{e}_{\Omega}$   
 $\Rightarrow \hat{y}_{\Omega} \approx \hat{y}_{\omega}$   
 $\Rightarrow$  prefer  $\omega$

•  $\theta$  small  
 $\Rightarrow \hat{e}_{\omega}$  quite different from  $\hat{e}_{\Omega}$   
 $\Rightarrow \hat{y}_{\Omega}$  quite different from  $\hat{y}_{\omega}$   
 $\Rightarrow$  prefer  $\Omega$

$$F = \frac{\frac{\|\hat{e}_{\omega} - \hat{e}_{\Omega}\|^2}{(p-q)}}{\frac{\|\hat{e}_{\Omega}\|^2}{(n-p)}} \sim F_{p-q, n-p} \text{ (under } \omega \text{)}$$

$\cot^2(\theta)$

$df_{\omega} - df_{\Omega} \rightarrow df_{\Omega}$

a special case: t-test