

Lecture 7

Multiple Linear Regression III

Reading: Chapter 13

STAT 8020 Statistical Methods II

September 10, 2020

Whitney Huang
Clemson University

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.1

Notes

Agenda

- 1 Multicollinearity
- 2 Variable Selection Criteria
- 3 Model Diagnostics

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.2

Notes

Multicollinearity

Multicollinearity is a phenomenon of high inter-correlations among the predictor variables

- Numerical issue \Rightarrow the matrix $X^T X$ is nearly singular
- Statistical issue
 - β^1 s are not well estimated
 - Spurious regression coefficient estimates
 - R^2 and predicted values are usually OK

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.3

Notes

Example

- Consider a two predictor model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$$

- We can show

$$\hat{\beta}_{1|2} = \frac{\hat{\beta}_1 - \sqrt{\frac{\hat{\sigma}_Y^2}{\hat{\sigma}_{X_1}^2}} r_{X_1, X_2} r_{Y, X_2}}{1 - r_{X_1, X_2}^2},$$

where $\hat{\beta}_{1|2}$ is the estimated partial regression coefficient for X_1 and $\hat{\beta}_1$ is the estimate for β_1 when fitting a simple linear regression model $Y \sim X_1$

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
74

Notes

An Simulated Example

Suppose the true relationship between response Y and predictors (X_1, X_2) is

$$Y = 4 + 0.8X_1 + 0.6X_2 + \varepsilon,$$

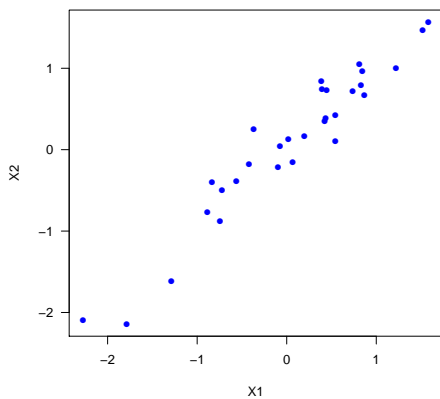
where $\varepsilon \sim N(0, 1)$ and X_1 and X_2 are positively correlated with $\rho = 0.95$. Let's fit the following models:

- Model 1: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$
- Model 2: $Y = \beta_0 + \beta_1 X_1 + \varepsilon_1$
- Model 3: $Y = \beta_0 + \beta_2 X_2 + \varepsilon_2$

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
75

Notes

Scatter Plot: X_1 vs. X_2



Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
76

Notes

Model 1 Fit

```
Call:
lm(formula = Y ~ X1 + X2)

Residuals:
    Min       1Q   Median       3Q      Max
-1.91369 -0.73658  0.05475  0.87080  1.55150

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.0710    0.1778   22.898 < 2e-16 ***
X1           2.2429    0.7187    3.121 0.00426 **
X2          -0.8339    0.7093   -1.176 0.24997
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9569 on 27 degrees of freedom
Multiple R-squared:  0.673,    Adjusted R-squared:  0.6488
F-statistic: 27.78 on 2 and 27 DF,  p-value: 2.798e-07
```

Multiple Linear Regression III



Multicollinearity

Variable Selection Criteria

Model Diagnostics

7.7

Notes

Model 2 Fit

```
Call:
lm(formula = Y ~ X1)

Residuals:
    Min       1Q   Median       3Q      Max
-2.09663 -0.67031 -0.07229  0.87881  1.49739

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.0347    0.1763   22.888 < 2e-16 ***
X1           1.4293    0.1955    7.311 5.84e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9634 on 28 degrees of freedom
Multiple R-squared:  0.6562,    Adjusted R-squared:  0.644
F-statistic: 53.45 on 1 and 28 DF,  p-value: 5.839e-08
```

Multiple Linear Regression III



Multicollinearity

Variable Selection Criteria

Model Diagnostics

7.8

Notes

Model 3 Fit

```
Call:
lm(formula = Y ~ X2)

Residuals:
    Min       1Q   Median       3Q      Max
-2.2584 -0.7398 -0.3568  0.8795  2.0826

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.9882    0.2014   19.80 < 2e-16 ***
X2           1.2973    0.2195    5.91 2.33e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.096 on 28 degrees of freedom
Multiple R-squared:  0.555,    Adjusted R-squared:  0.5391
F-statistic: 34.92 on 1 and 28 DF,  p-value: 2.335e-06
```

Multiple Linear Regression III



Multicollinearity

Variable Selection Criteria

Model Diagnostics

7.9

Notes

Variance Inflation Factor (VIF)

We can use the **variance inflation factor (VIF)**

$$VIF_i = \frac{1}{1 - R_i^2}$$

to quantify the severity of multicollinearity in MLR, where R_i^2 is the **coefficient of determination** when X_i is regressed on the remaining predictors

Multiple Linear Regression III
 CLEMSON UNIVERSITY
 Multicollinearity
 Variable Selection Criteria
 Model Diagnostics
 7.10

Notes

Variable Selection

Multiple Linear Regression Model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{p-1} X_{i,p-1} + \varepsilon_i, \quad \varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- What is the appropriate subset size?
- What is the best model for a fixed size?

In the next few slides we will discuss some commonly used model selection criteria

Multiple Linear Regression III
 CLEMSON UNIVERSITY
 Multicollinearity
 Variable Selection Criteria
 Model Diagnostics
 7.11

Notes

Mallows' C_p Criterion

$$\begin{aligned} (\hat{Y}_i - \mu_i)^2 &= (\hat{Y}_i - E(\hat{Y}_i) + E(\hat{Y}_i) - \mu_i)^2 \\ &= \underbrace{(\hat{Y}_i - E(\hat{Y}_i))^2}_{\text{Variance}} + \underbrace{(E(\hat{Y}_i) - \mu_i)^2}_{\text{Bias}^2} \end{aligned}$$

where $\mu_i = E(Y_i | X_i = x_i)$

- Mean squared prediction error (MSPE):
 $\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2$
- C_p criterion measure:

$$\begin{aligned} \Gamma_p &= \frac{\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2}{\sigma^2} \\ &= \frac{\sum \text{Var}_{\text{pred}} + \sum \text{Bias}^2}{\text{Var}_{\text{error}}} \end{aligned}$$

Multiple Linear Regression III
 CLEMSON UNIVERSITY
 Multicollinearity
 Variable Selection Criteria
 Model Diagnostics
 7.12

Notes

C_p Criterion

- Do not know σ^2 nor numerator
- Use $MSE_{X_1, \dots, X_{p-1}} = MSE_F$ as the estimate for σ
- For numerator:
 - Can show $\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 = p\sigma^2$
 - Can also show $\sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2 = E(SSE_F) - (n - p)\sigma^2$

$$\Rightarrow C_p = \frac{SSE - (n-p)MSE_F + pMSE_F}{MSE_F}$$

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.13

Notes

C_p Criterion Cont'd

Recall

$$C_p = \frac{\sum_{i=1}^n \sigma_{\hat{Y}_i}^2 + \sum_{i=1}^n (E(\hat{Y}_i) - \mu_i)^2}{\sigma^2}$$

- When model is correct $E(C_p) \approx p$
- When plotting models against p
 - Biased models will fall above $C_p = p$
 - Unbiased models will fall around line $C_p = p$
 - By definition: C_p for full model equals p

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.14

Notes

Adjusted R^2 Criterion

Adjusted R^2 , denoted by R_{adj}^2 , attempts to take account of the phenomenon of the R^2 automatically and spuriously increasing when extra explanatory variables are added to the model.

$$R_{\text{adj}}^2 = 1 - \frac{SSE/(n - p - 1)}{SST/(n - 1)}$$

- Choose model which maximizes R_{adj}^2
- Same approach as choosing model with smallest MSE

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.15

Notes

Predicted Residual Sum of Squares *PRESS* Criterion

- For each observation i , predict Y_i using model generated from other $n - 1$ observations
- $PRESS = \sum_{i=1}^n (Y_i - \hat{Y}_{i(i)})^2$
- Want to select model with small *PRESS*

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.16

Notes

Other Approaches: Information criteria

- Akaike's information criterion (AIC)

$$n \log\left(\frac{SSE_k}{n}\right) + 2k$$

- Bayesian information criterion (BIC)

$$n \log\left(\frac{SSE_k}{n}\right) + k \log(n)$$

- Can be used to compare **non-nested** models

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.17

Notes

Automatic Search Procedures

- Forward Selection
- Backward Elimination
- Stepwise Search
- All Subset Selection

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.18

Notes

Model Assumptions

Model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{p-1} X_{i,p-1} + \varepsilon_i, \quad \varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

We make the following **assumptions**:

- Linearity:

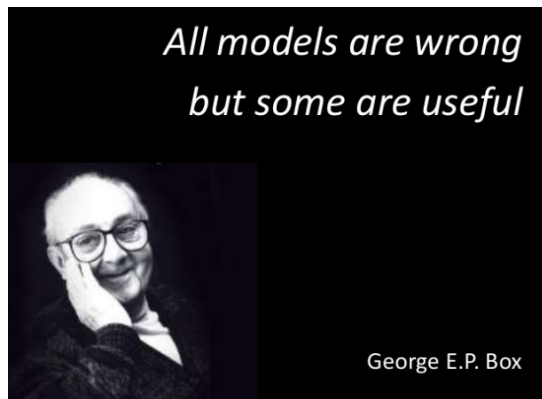
$$E(Y|X_1, X_2, \dots, X_{p-1}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1}$$

- Errors have constant variance, are independent, and normally distributed

$$\varepsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.19

Notes



Multiple Linear Regression III
CLEMSON UNIVERSITY
Multicollinearity
Variable Selection Criteria
Model Diagnostics
7.20

Notes

Notes
