

Diagnostic and Remedial Measures

Most of the diagnostic procedures for multiple regression are same as those for the SLR.

* Scatter Plots:

Scatter plot of the response variable against each predictor variable can be used to determine the nature and the strength of bivariate relationships.

Scatterplot matrix can be used to see scatter plots together.

Correlation matrix can also be used to determine the strength of bivariate relationships.

Correlation matrix:

$$\begin{pmatrix} \text{corr}(Y, Y) & \text{corr}(Y, X_1) & \dots & \text{corr}(Y, X_{p-1}) \\ \text{corr}(X_1, Y) & \text{corr}(X_1, X_1) & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \text{corr}(X_{p-1}, X_{p-1}) \end{pmatrix}$$

Residual plots:

- * Residual vs fitted values
 - constant variance
 - outliers

* Residual vs time

- correlation between error terms.

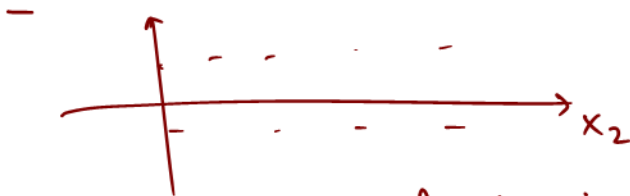
* Box-plots and Normal Probability plots

- Normality

* Residual vs predictor variables (each)

- adequacy of the regression function for that predictor.

* Residual vs predictors which are not in the model

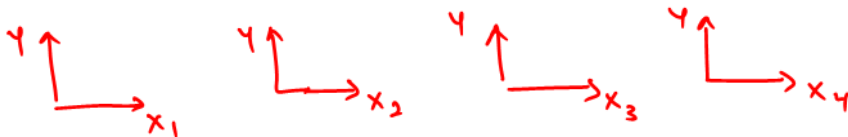


$$y = \beta_0 + \beta_1 x_1$$

x_2

- predictor should be in the model or not.

Chapter - 7: Extra Sum of Squares (ESS)



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

x_3

not SSE

ESS measures the marginal reduction in the error Sum of Squares (SSE) (ie marginal increase in regression Sum of Squares) when one or more predictors added to the model given that the other predictors are in the model.



Notation:

$$SSTO = SSR + SSE$$

* $SSR(X_1, X_2, X_3)$ - total variation explained by X_1, X_2, X_3 .

* $SSR(X_1 | X_2)$ - additional variation explained by X_1 when added to a model with X_2 .

* $SSR(X_1, X_4 | X_2, X_3)$ - additional variation explained by X_1 and X_4 when added to a model with X_2 and X_3 .

Note:

* ESS can also be viewed of SSE's.

* ESS represent the part of SSE that is explained by an added group of variables that was not previously explained by the rest.

$$\begin{aligned} SSR(X_1 | X_2) &= \textcircled{SSE(X_2)} - SSE(X_1, X_2) \\ &= SSR(X_1, X_2) - SSR(X_2). \end{aligned}$$

Decomposition of SSR into ESS

$$SSTO = SSR + SSE$$

$$SSTO = SSR(X_1) + \underbrace{SSE(X_1)}$$

$$= SSR(X_1) + SSR(X_2 | X_1) + \cancel{SSE(X_1, X_2)} \rightarrow \textcircled{1}$$

Further

$$SSTO = \cancel{SSR(X_1, X_2)} + SSE(X_1, X_2) \rightarrow \textcircled{2}$$

By ① und ②:

$$SSR(X_1, X_2) = SSR(X_1) + SSR(X_2|X_1)$$

Similarly,

$$SSR(X_1, X_2) = \overbrace{SSR(X_2) + SSR(X_1|X_2)}^{= SSR(X_2, X_1)}$$

$$SSR(X_1, X_2, X_3) = SSR(X_1) + SSR(X_2|X_1) + SSR(X_3|X_1, X_2)$$

Extended ANOVA table containing the decomposition of SSR

Source of variation	SS'	df	MS
Regression	- SSR(X ₁ , X ₂ , X ₃)	3	MSR(X ₁ , X ₂ , X ₃)
X ₁	SSR(X ₁)	1 1 1 } 1 1 1 }	MSR(X ₁)
X ₂ X ₁	SSR(X ₂ X ₁)		MSR(X ₂ X ₁)
X ₃ X ₁ , X ₂	SSR(X ₃ X ₁ , X ₂)		MSR(X ₃ X ₁ , X ₂)
Error	- SSE(X ₁ , X ₂ , X ₃)	n-4	MSE(X ₁ , X ₂ , X ₃)
Total	SSTO	n-1	

Note:

$$* MSR(X_2|X_1) = \frac{SSR(X_2|X_1)}{1}$$

$$* MSR(X_2, X_3|X_1) = \frac{SSR(X_2, X_3|X_1)}{2} = \frac{SSR(X_2|X_1) + SSR(X_3|X_1, X_2)}{2}$$

Test for Regression coefficients (using ESS)

Consider the model

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i$$

* Here we can use the general linear test to test the regression coefficients.

* Test whether a single $\beta_k = 0$ or not ($k = 1, 2, 3$).

Steps:

1) Hypotheses:

$$H_0: \beta_3 = 0 \text{ (Reduced model is better)}$$

$$\text{vs } H_1: \beta_3 \neq 0 \text{ (Full model is better)}$$

2) Test Statistic (under H_1)

* Full model: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i$

$$SSE(F) = SSE(X_1, X_2, X_3)$$

$$df_F = n - 4$$

(under H_0)

* Reduced model: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \varepsilon_i$

$$SSE(R) = SSE(X_1, X_2)$$

$$df_R = n - 3$$

$$T = \frac{b_3 - 0}{S\{b_3\}} \sim t_{n-4}$$

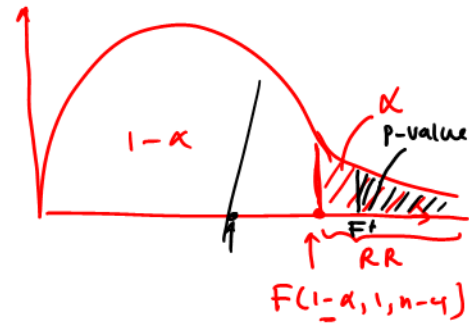
$$F = \frac{SSE(R) - SSE(F)}{df_R - df_F} \div \frac{SSE(F)}{df_F}$$

$$= \frac{SSE(X_1, X_2) - SSE(X_1, X_2, X_3)}{(n-3) - (n-4)} \div \frac{SSE(X_1, X_2, X_3)}{n-4}$$

$$= \frac{SSR(X_3 | X_1, X_2)}{1} \div \frac{SSE(X_1, X_2, X_3)}{n-4}$$

$$= \frac{MSR(X_3 | X_1, X_2)}{MSE(X_1, X_2, X_3)} \sim F_{1, n-4}$$

3) Calculate $F(1-\alpha, 1, n-4)$ (OR $p\text{-value} = P(F > F^*)$)
 observed value



4) If $F^* > F(1-\alpha, 1, n-4)$ } \Rightarrow Reject H_0
 (OR $p\text{-value} < \alpha$) } (i.e. X_3 should be in the model)

If $F^* \leq F(1-\alpha, 1, n-4)$ } \Rightarrow Fail to reject H_0
 (OR $p\text{-value} \geq \alpha$) } (i.e. X_3 can be dropped from the model).

Note:

* For the hypothesis $H_0: \beta_3 = 0$ vs $H_1: \beta_3 \neq 0$, the t -Statistic is
 $t = \frac{b_3}{S\{b_3\}}$ can also be used. It can be shown that $t^2 = F$.

$$* \quad p\text{-value} = \underbrace{P(F > F^*)}_{F\text{-test}} = \underbrace{P(T > |t^*|)}_{t\text{-test}}$$

* Test for several β_k s

Steps:

1) $H_0: \beta_2 = \beta_3 = 0$ vs $H_1: \text{not } H_0$

$$\begin{aligned} 2) \quad F &= \frac{SSE(R) - SSE(F)}{df_R - df_F} \div \frac{SSE(F)}{df_F} \\ &= \frac{SSE(X_1) - SSE(X_1, X_2, X_3)}{(n-2) - (n-4)} \div \frac{SSE(X_1, X_2, X_3)}{n-4} \\ &= \frac{SSR(X_2, X_3 / X_1)}{2} \div \frac{SSE(X_1, X_2, X_3)}{n-4} \end{aligned}$$

$$= \frac{\text{MSR}(X_2, X_3/X_1)}{\text{MSE}(X_1, X_2, X_3)} \sim F_{2, n-4}$$

$\frac{\text{SSR}(X_2/X_1) + \text{SSR}(X_3/X_1, X_2)}{2}$

3) calculate $F(1-\alpha, 2, n-4)$ (OR calculate p-value = $P(F > F^*)$)

4) If $F^* > F(1-\alpha, 2, n-4)$ } \Rightarrow Reject H_0
 (OR if p-value $\leq \alpha$) } (i.e. predictors X_2 and X_3 (at least one) should be in the model).

If $F^* \leq F(1-\alpha, 2, n-4)$ } \Rightarrow Fail to reject H_0
 (OR If p-value $\geq \alpha$) } (i.e. predictors X_2 and X_3 can be dropped from the model).

Note:
 * Overall F-test is a special case of the previous test.

$$H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$$

$$H_1: \text{not } H_0.$$

$$F = \frac{\text{SSR}(X_1, X_2, \dots, X_{p-1})}{p-1} \div \frac{\text{SSE}(X_1, X_2, \dots, X_{p-1})}{n-p}$$

$$= \frac{\text{MSR}}{\text{MSE}} \leftarrow \text{Same as the test statistic in F-test.}$$

$\frac{\text{SSTO} - \text{SSE}}{\text{SSE}(R) - \text{SSE}(F)} \div \frac{\text{SSE}(F)}{\text{df}_F}$

* Other tests

1) consider the test

$$H_0: \beta_1 = \beta_2 \quad \text{vs} \quad H_1: \beta_1 \neq \beta_2$$

2) Full model: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \epsilon_i$

Reduced model: $Y_i = \beta_0 + \beta_c (X_{i1} + X_{i2}) + \beta_3 X_{i3} + \epsilon_i$, where $\beta_c = \beta_1 = \beta_2$.

* Here we can not use ESS. We have to fit Full and reduced models separately and calculate $SSE(F)$, $SSE(R)$, df_F and df_R .

2) Consider the test

$$H_0: \beta_1 = \beta_{10}, \beta_2 = \beta_{20} \quad \text{vs} \quad H_1: \text{not } H_0$$

\uparrow
constants

Reduced model:

$$Y_i = \beta_0 + \beta_{10} X_{i1} + \beta_2 X_{i2} + \beta_{30} X_{i3} + \epsilon_i$$

$$\Rightarrow \underbrace{Y_i - \beta_{10} X_{i1} - \beta_{30} X_{i3}}_{\text{new response}} = \beta_0 + \beta_2 X_{i2} + \epsilon_i$$

Coefficient of Partial Determination

$$\begin{bmatrix} R^2 \\ R^2, R_a^2 \end{bmatrix}$$

A coefficient of partial determination measures the marginal contribution of one X variable when all the others are already in the model.

Notation

$$R_{Y1/2}^2 = \frac{SSE(X_2) - SSE(X_1, X_2)}{SSE(X_2)} = \frac{SSR(X_1/X_2)}{SSE(X_2)}$$

- Percentage of the left over variation in Y (after regressing on X_2) that is explained by X_1 .

Similarly,

$$* R_{y2|1}^2 = \frac{SSR(X_2|X_1)}{SSE(X_1)}$$

$$* R_{y1|23}^2 = \frac{SSR(X_1|X_2, X_3)}{SSE(X_2, X_3)}$$

$$* R_{y4|123}^2 = \frac{SSR(X_4|X_1, X_2, X_3)}{SSE(X_1, X_2, X_3)}$$

$$* R_{y1}^2 = \frac{SSR(X_1)}{SSTO}$$

* Coefficient of Partial Correlation

The square root of coefficient of partial determination is called the coefficient of partial correlation.

$$* \underset{\uparrow}{r}_{y2|1} = \sqrt{R_{y2|1}^2} \quad \leftarrow \textcircled{\beta_2}$$

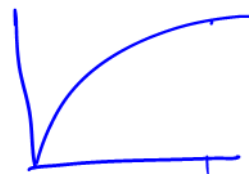
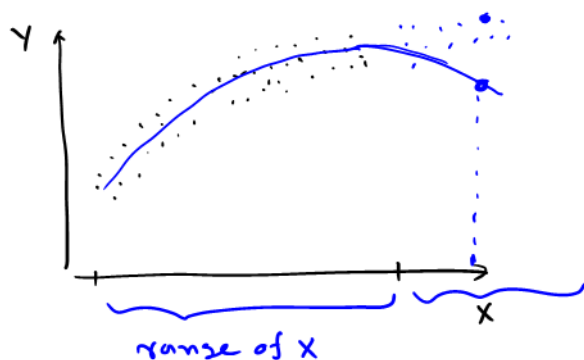
The sign of coefficient of partial correlation is same as the sign of corresponding regression coefficient.

Chapter-8: Regression models for Quantitative and Qualitative Predictors:

* Polynomial Regression Models

Polynomial Regression models are used to model curvilinear relationships,

- 1) when the true curvilinear response function is a polynomial.
- 2) when the response function is unknown, but a polynomial model is a good approximation to the true function.



Polynomial regression may provide good fit for the data at hand but may turn in unexpected directions when extrapolated beyond the range of the data.

* Second order model with one predictor variable.

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i$$

where $x_i = X_i - \bar{X}$ - i^{th} centered predictor.

$$(X'X)^{-1} \quad X = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{pmatrix}$$

Note: Since X and X^2 are highly correlated, here we use

$x_i = X_i - \bar{X}$, centered predictor. This avoids computational difficulties.

(i.e. calculating $(X'X)^{-1}$ is easier).

Note:

$$y_i = \beta_0 + \beta_1 x_i + \beta_{11} x_i^2 + \varepsilon_i$$

* β_{11} — coefficient of x_i^2 (i.e. $\underbrace{x_i \cdot x_i}_{\beta_{11}}$)

* β_{12} — coefficient of $\underbrace{x_1 x_2}$

$\beta_{21} \times$

R Codes for Chapter-7 Extra Sums of Squares in R

We will work again with the data from Problem 6.9, “Grocery Retailer.” You can obtain the ANOVA table using the function “`anova(model)`”. Here you get sum of squares for each predictor variable in the model:

`> anova(lrm)`
Analysis of Variance Table

$SSR = SSR(X_1) + SSR(X_2|X_1) + SSR(X_3|X_1, X_2)$

$H_0: \beta_1 = 0 \quad \text{vs} \quad H_1: \beta_1 \neq 0$

Response: Retailer

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
X_1 Cases	1	136366	136366	6.6417	0.01309 *
X_2 Costs	1	5726	5726	0.2789	0.59987
X_3 Holiday	1	2034514	2034514	99.0905	2.941e-13 ***
Residuals	48	985530	20532		

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Handwritten annotations:
 $SSR(X_1)$ points to 136366
 $SSR(X_2|X_1)$ points to 5726
 $SSR(X_3|X_1, X_2)$ points to 2034514
 SSE points to 985530

The ANOVA table given by R provides the extra sum of squares for each predictor variable, *given that* the previous predictors are already in the model. Thus the Sum of Squares given for “Cases” is $SSR(X_1) = 136366$, while the Sum of Squares given for “Costs” is $SSR(X_2 | X_1) = 5726$, and the Sum of Squares given for “Holiday” is $SSR(X_3 | X_1, X_2) = 2034514$. This corresponds to Table 7.3 on p.261 of the text.

Now you have $SSR(X_1)$, $SSR(X_2 | X_1)$, and $SSR(X_3 | X_1, X_2)$ and their corresponding degrees of freedom and mean squares. If you sum them together you get $SSR(X_1, X_2, X_3)$, which has 3 degrees of freedom. Divide $SSR(X_1, X_2, X_3)$ by 3 to get $MSR(X_1, X_2, X_3)$. To get $SSE(X_1, X_2, X_3)$, its degrees of freedom, and $MSE(X_1, X_2, X_3)$, use the line beginning with “Residuals.” To calculate and store these in R, use the commands

```
> SSR = sum( anova(lrm)[1:3,2] )
> SSR
[1] 2176606
> MSR = SSR / 3
> MSR
[1] 725535.4
> SSE = anova(lrm)[4,2]
> SSE
[1] 985529.7
> MSE = anova(lrm)[4,3]
> MSE
[1] 20531.87
```

You can obtain alternate decompositions of the regression sum of squares into **extra sum of squares** by running new linear models with the predictors entered in a different order. For an example, if we want $SSR(X_3)$, $SSR(X_1|X_3)$ and $SSR(X_2|X_1, X_3)$, we could try:

```

> Model2 <- lm( Retailer ~ X3 Holiday + X1 Cases + X2 Costs)
> anova(Model2)
Analysis of Variance Table

```

Response: Retailer

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
X₃ Holiday	1	2077646	2077646	101.1913	2.086e-13 ***
X₁ Cases	1	92285	92285	4.4947	0.0392 *
X₂ Costs	1	6675	6675	0.3251	0.5712
Residuals	48	985530	20532		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

GENERAL LINEAR TEST

$$H_0: \beta_2 = 0 \text{ vs } H_1: \beta_2 \neq 0$$

If we are considering dropping Costs (X_2) from the **lrm** model, we run a reduced model which uses only the other two predictors **Cases** and **Holiday**:

```

> Reduced <- lm( Retailer ~ Holiday + Cases) #fitting the reduced model
> Reduced

```

Call:

```
lm(formula = Retailer ~ Holiday + Cases)
```

Coefficients:

(Intercept)	Holiday	Cases
4.058e+03	6.196e+02	7.704e-04

Then to perform the F test, just type

```
> anova(Reduced, Retailer)
```

To get the ANOVA comparison:

```

> anova(Reduced, lrm)
Analysis of Variance Table

```

Model 1: Retailer ~ Holiday + Cases

Model 2: Retailer ~ Cases + Costs + Holiday

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	49	992204				
2	48	985530	1	6674.6	0.3251	0.5712

$$\alpha = 0.05$$

Test Statistic

Note that the first argument to the **anova()** function must be the **reduced model**, and the second argument must be the full model (the one with all the original predictors).

General Linear Test for the other reduced models:

Now suppose we want to test $H_0: \beta_2 = 0, \beta_3 = 600$ against its alternative. In this case the reduced model, corresponding to H_0 , is $Y_i = \beta_0 + \beta_1 X_{i1} + 600 X_{i3} + \varepsilon_i$, which may be rewritten as $Y_i - 600 X_{i3} = \beta_0 + \beta_1 X_{i1} + \varepsilon_i$. To obtain the reduced model in R, use the formulation:

```

> RetailerN = Retailer - 600 * Holiday ← new response
> Reduced2 <- lm(RetailerN ~ Cases)
> Reduced2

```

```
Call:
lm(formula = RetailerN ~ Cases)
```

```
Coefficients:
(Intercept)      Cases
  4.059e+03    7.756e-04
```

```
> anova(Reduced2)
Analysis of Variance Table
```

```
Response: RetailerN
      Df Sum Sq Mean Sq F value Pr(>F)
Cases    1  93738    93738   4.714 0.03469 *
Residuals 50 994244    19885
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

However, you will get an error message if you attempt to use the `anova()` function to compare this model with the full model, because the two models do not have the same response variable. Instead, you will need to obtain the *SSE* for this reduced model, along with its degrees of freedom, from its ANOVA table, and the *SSE* from the full model, along with its degrees of freedom, from the ANOVA table for the full model, then calculate F^* using equation (2.9) in the textbook.

```
> SSE_R=anova(Reduced2)[2,2]
> #SSE_R
> DF_R=anova(Reduced2)[2,1]
> #DF_R
>
> SSE_F = anova(lrm)[4,2]
> #SSE_F
> DF_F=anova(lrm)[4,1]
> #DF_F
>
> #Test atstatistics
> F=((SSE_R-SSE_F)/(DF_R-DF_F))/(SSE_F/DF_F)
> F
[1] 0.2122226
```

```
> #P-value
> Pvalue=1-pf(F,DF_R-DF_F,DF_F)
> Pvalue
[1] 0.8095395
```

Coefficients of Partial Determination

To obtain the coefficients of partial determination, you will need to use formulae like those in section 7.4. You may also need to run several different models, with the predictors in various different orders, in order to obtain values for the needed forms of *SSE* and the extra sums of squares.

To calculate $R^2_{Y1|23}$, first the following model with X_2 and X_3 should be fit to calculate $SSE(X_2, X_3)$

```
> M1=lm(Retailer~ Costs+Holiday)
> SSEX_2X_3=anova(M1)[3,2]
> SSEX_2X_3
[1] 1081237
```

A model with X_1 , X_2 and X_3 should be fit to calculate $SSE(X_1, X_2, X_3)$

```
> M2=lm(Retailer~ Cases+Costs+Holiday)
> SSEX_1X_2X_3=anova(M2)[4,2]
> SSEX_1X_2X_3
[1] 985529.7
```

Then $R^2_{Y1|23}$

```
> RSQ_Y1_23=(SSEX_2X_3-SSEX_1X_2X_3)/SSEX_2X_3
> RSQ_Y1_23
[1] 0.08851609
```