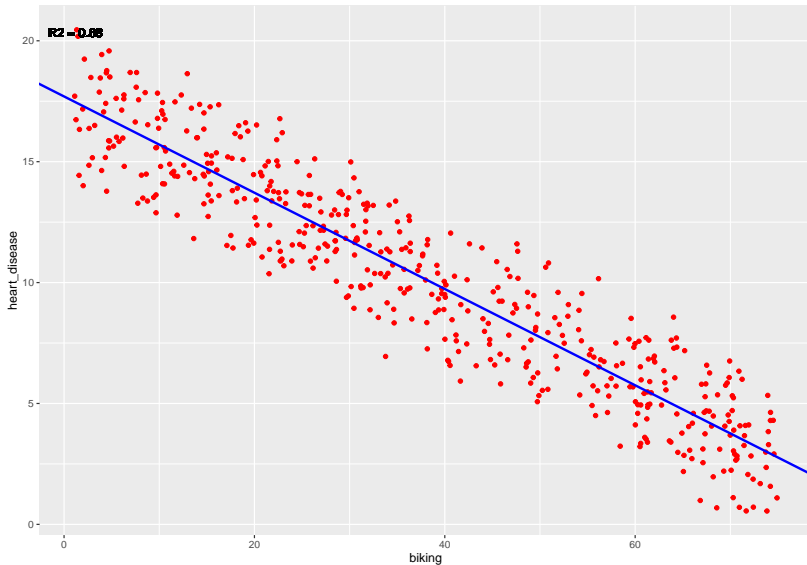


Multiple Linear Regression

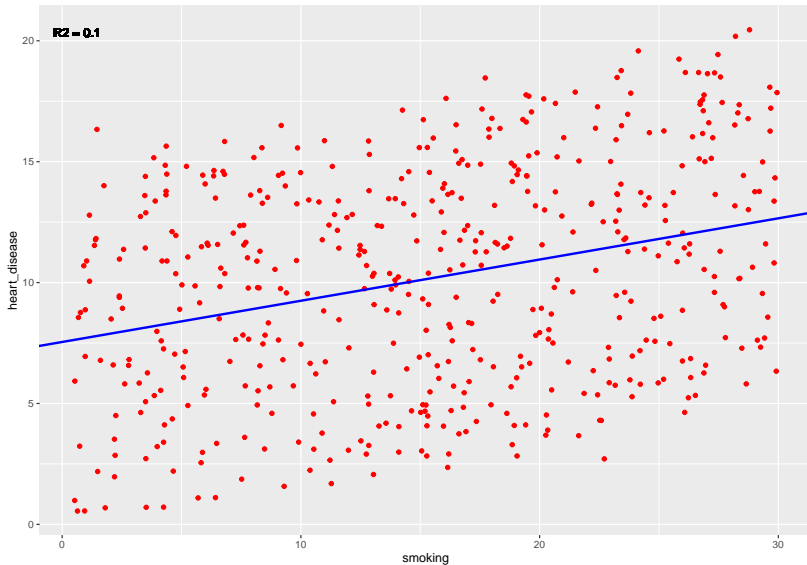
Example Data

biking	smoking	heart_disease
30.801246	10.896608	11.769423
65.129215	2.219563	2.854081
1.959664	17.588331	17.177803
44.800196	2.802559	6.816647
69.428454	15.974505	4.062223
54.403626	29.333175	9.550046
49.056162	9.060846	7.624507
4.784604	12.835021	15.854654
65.730788	11.991297	3.067462
35.257449	23.277683	12.098484
51.825567	14.435118	6.430248
52.936197	25.074869	8.608272
48.767479	11.023271	6.722524
26.166801	6.645749	10.597807
10.553075	5.990506	14.079478

Regress heart_disease individually



Regress heart_disease individually



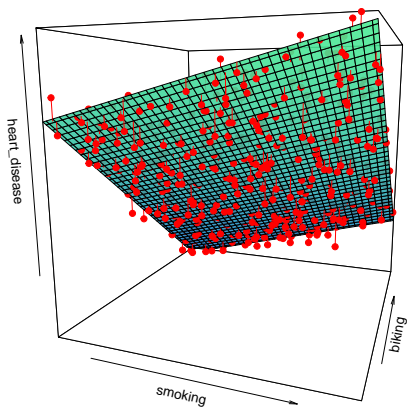
► Is there a better way? better model?

Multiple Regression Model

- ▶ $\text{heart_disease} = \beta_0 + \beta_1 \cdot \text{biking} + \beta_2 \cdot \text{smoking} + \epsilon$
- ▶ $\epsilon \sim N(0, \sigma^2)$

Graphing the solution

RSS: 211.74, R2 = 0.98



► $\text{heart_disease} = 14.98 + -0.2 \cdot \text{biking} + 0.18 \cdot \text{smoking}$

Model Definition

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

► Model Assumptions

- (A1) The response variable y is a random variable and the predictor x_1, x_2, \dots, x_n is non-random
- (A2) $\epsilon \sim N(0, \sigma^2)$

Parameters Estimation

Data Presentation

Observation	Response Variable	Predictors			
	y	x_1	x_2	\cdots	x_p
1	y_1	x_{11}	x_{12}	\cdots	x_{1p}
2	y_2	x_{21}	x_{22}	\cdots	x_{2p}
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
n	y_n	x_{n1}	x_{n2}	\cdots	x_{np}

Matrix Equation of MLR

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Least Squares Estimators

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

Example

An automobile insurance company wants to use gender ($x_1 = 0$, if female and $x_1 = 1$, if male) and traffic penalty point (x_2) to predict the number of claims (y). The observed values of these variables for a sample of six motorists are given by:

Motorist	1	2	3	4	5	6
x_1	0	0	0	1	1	1
x_2	0	1	2	0	1	2
y	1	0	2	1	3	5

You are using the following model:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i, \quad i = 1, 2, \dots, 6$$

Example (Continue)

You have determine

$$(X'X)^{-1} = \frac{1}{12} \begin{bmatrix} 7 & -4 & -3 \\ -4 & 8 & 0 \\ -3 & 0 & 3 \end{bmatrix}$$

Write the best fitted linear equation.

Goodness of Fit

Coefficient of Determination

► Similarly to the case of SLR, we have

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

TSSRSSReg SS

► And

$$R^2 = \frac{RegSS}{TSS} = 1 - \frac{RSS}{TSS}$$

as good as the baseline model
 $\hat{y}_i = \bar{y}$
 $R^2 = 0$
perfect fit

F-test

- Full Model: *($p \geq 2$) uses all the predictors*

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

- Baseline Model or i.i.d model: *uses no predictors*

$$y = \beta_0 + \epsilon$$

- The baseline model is equivalent to

$$\beta_1 = \beta_2 = \dots = \beta_p = 0$$

H_0

- We would like to test for the joint significant of all predictors, or if the full model is a significant improvement over the baseline model, or

$$\underline{H_0 : \underbrace{\beta_1 = \beta_2 = \cdots = \beta_p = 0}_{\text{i.i.d. model}}} \text{ vs. } \underline{H_a : \underbrace{\text{at least one } \beta_j \text{ is non-zero.}}_{\text{MLR model}}}$$

► Test Statistics

number of
predictors
↓

$$F = \frac{(\text{TSS} - \text{RSS}_1)/p}{\text{RSS} / (n - p - 1)} = \frac{\text{Reg SS} / p}{\text{RSS} / (n - p - 1)},$$

ANOVA Table

- The results of MLR are usually summarized in the ANOVA table

Source	Sum of Squares	df	Mean Square	F
Regression	Reg SS	p	Reg SS/ p	$\frac{\text{Reg SS}/p}{\text{RSS}/[n - (p + 1)]}$
Error	RSS	$n - (p + 1)$	$s^2 = \text{RSS}/[n - (p + 1)]$	
Total	TSS	$n - 1$		

Example

An actuary uses multiple regression model with three predictors and 20 observations and has the following results. $p = 3$

$$n = 20$$

	Sum of Squares
Regression	150
Total	200

$$RSS = 150$$

$$TSS = 200$$

$$RSS = 200 - 150 = 50$$

$$F = \frac{RSS / p}{RSS / (n - p - 1)} = \frac{150 / 3}{50 / (20 - 3 - 1)} = \frac{50}{50} \cdot 16 = 16$$

He wants to test the following hypothesis

$$H_0 : \beta_1 = \beta_2 = \beta_3 = 0$$

H_1 : At least one of β_1 , β_2 , and β_3 is zero

Calculate the F-statistics of the test.

From R^2 to F-test

► The R^2 and the F - statistics have the following relation

$$F = \frac{\text{RegSS}/p}{\text{RSS}/(n-p-1)} = \frac{R^2/p}{(1-R^2)/(n-p-1)}$$

and

$$R^2 \rightarrow \frac{1}{R^2} \rightarrow \frac{1}{F} \rightarrow F \uparrow$$

p value \downarrow

$$R^2 = \frac{Fp}{Fp + n - p - 1}$$

$$\Rightarrow \frac{1}{R^2} = \frac{Fp + n - p - 1}{Fp}$$

$$\Rightarrow \left[\frac{1}{R^2} = 1 + \frac{n-p-1}{p} \cdot \left(\frac{1}{F} \right) \right]$$

Example

Sarah performs a regression of the return on a mutual fund (y) on four predictors plus an intercept. She uses monthly returns over 105 months. Her software calculates the $R^2 = .8$ but then it quits working before it calculates the value of F . Calculate the F-statistics for Sarah.

$$p = 4$$

$$n = 105$$

$$R^2 = .8$$

$$F = \frac{R^2 / p}{(1 - R^2) / (n - p - 1)} = \frac{.8 / 4}{.2 / 105 - 4 - 1} = \boxed{160}$$

Generalized F-test

► Full Model:

(model 2) uses all p predictors

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

► Reduced Model:

(model 2)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{p-q} x_{p-q} + \epsilon$$

uses only some of the predictors.

($p-q$) not use the other q predictors.

$$TSS = \sum (y_i - \bar{y})^2$$

	Reduced model		Full model
RSS	<u>RSS</u> ₀	≥	<u>RSS</u> ₁
Reg SS	(Reg SS) ₀	≤	(Reg SS) ₁
TSS	TSS	=	TSS

$$\begin{array}{lcl} \text{Reduced model:} & RSS_0 + (Reg SS)_0 & = TSS \\ \text{Full model} & RSS_1 + (Reg SS)_1 & = TSS \end{array} \quad ||$$

Model	RSS	$RegSS$
Reduced	RSS_0	$RegSS_0$
Full	RSS_1	$RegSS_1$

$$RSS_0 + RegSS_0 = TSS$$

$$RSS_1 + RegSS_1 = TSS$$

- ▶ $H_0 : \beta_{p-q+1} = \beta_{p-q+2} = \dots = \beta_{p-q} = 0$ or Reduced model is adequate
- ▶ Test Statistics

$$F = \frac{\text{Extra SS}/q}{\text{RSS}_1/(n - p - 1)} = \frac{(\text{RSS}_0 - \text{RSS}_1)/q}{\text{RSS}_1/(n - p - 1)}.$$

Example

$n=26$

- Model 1: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$ (Reduced)
- Model 2: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon$ (Full)

The results of the regression are as follows:

Model Number	Residual Sum of Squares	Regression Sum of Squares
1	$RSS_0 = 13.47$	22.75
2	$RSS_1 = 10.53$	25.70

The null hypothesis is $H_0 : \beta_3 = \beta_4 = 0$ with the alternative hypothesis that the two betas are not equal to zero.

Calculate the statistic used to test H_0 .

$$p = 4, q = 2$$

$$\Rightarrow F = \frac{(RSS_0 - RSS_1) / q}{RSS_1 / (n - p - 1)} = \frac{(13.47 - 10.53) / 2}{10.53 / (26 - 4 - 1)} = \boxed{2.09}$$

Example

You wish to find a model to predict insurance sales, using 27 observations and 8 variables x_1, x_2, \dots, x_8 . The analysis of variance (ANOVA) tables are below. Model A contains all 8 variables and Model B contains x_1 and x_2 only.

Calculate the F-statistics for testing

$$H_0 : \beta_3 = \beta_4 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = 0$$

$$F = \frac{(RSS_0 - RSS_1) / d}{RSS_1 / (n - p - 1)}$$

$$\left. \begin{array}{l} p = 8 \\ q = 6 \\ n = 27 \\ RSS_0 = 126,471 \\ RSS_1 = 76,893 \end{array} \right\}$$

Model A

Source	SS	df	MS
Regression	115,175	8	14,397
Error	76,893	18	4,272
Total	192,068	26	

Model B

Source	SS	df	MS
Regression	65,597	2	32,798
Error	126,471	24	5,270
Total	192,068	26	

t-test

- ▶ Similar to SLR, the t-test can be used to test for the magnitude of the coefficients. $\beta_0, \beta_1, \beta_2 \dots$
- ▶ Coefficients with larger p-values are less significant in the presence of other predictors and may be considered to be dropped. }

Example

	Coefficient	Standard Error	Stat	p-value	
Intercept	-	-	-2.24	0.0303	$\beta_0 = 0$
x_1	513,280.76	233,143.23	2.20	0.0330	$\beta_1 = 0$
x_2	280,148.46	483,001.55	0.58	0.5649	$\beta_2 = 0$
x_3	38.64	6.42	6.01	0.0000	$\beta_3 = 0$

At a 1% significance level, which of the following hypothesis that it's fail to reject: $\beta_1 = 0$, $\beta_2 = 0$ and $\beta_3 = 0$?

If $p\text{-value} < .01 \Rightarrow \text{Reject } H_0$

$> .01 \Rightarrow \text{fail to reject } H_0$