

10) Multiple Linear Regression

Vitor Kamada

January 2018

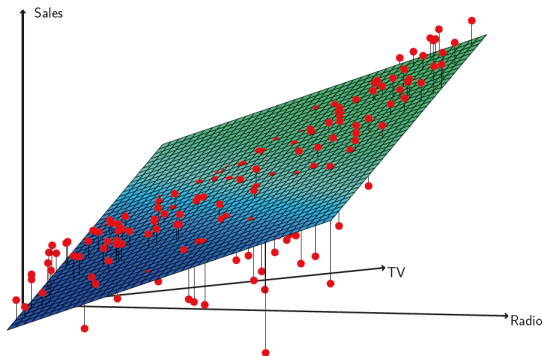
Tables, Graphics, and Figures from
An Introduction to Statistical Learning

James et al. (2017): Chapters: 3.2 and 3.6.3

Multiple Linear Regression

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$$



Simple vs Multiple Regression

	Coefficient	Std. error	t-statistic	p-value
Intercept	9.312	0.563	16.54	< 0.0001
radio	0.203	0.020	9.92	< 0.0001

	Coefficient	Std. error	t-statistic	p-value
Intercept	12.351	0.621	19.88	< 0.0001
newspaper	0.055	0.017	3.30	0.00115

	Coefficient	Std. error	t-statistic	p-value
Intercept	2.939	0.3119	9.42	< 0.0001
TV	0.046	0.0014	32.81	< 0.0001
radio	0.189	0.0086	21.89	< 0.0001
newspaper	-0.001	0.0059	-0.18	0.8599

Correlation Matrix

$$r = \frac{\text{cov}(x,y)}{s_x s_y}$$

	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5762
newspaper			1.0000	0.2283
sales				1.0000

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

$$H_a : \text{at least one } \beta_j \neq 0$$

$$F = \frac{\frac{TSS - RSS}{p}}{\frac{RSS}{n - p - 1}}$$

Quantity	Value
Residual standard error	1.69
R^2	0.897
F-statistic	570

Test if q Coefficients are Zero

$$H_0 : \beta_3 = \beta_4 = \beta_5 = 0$$

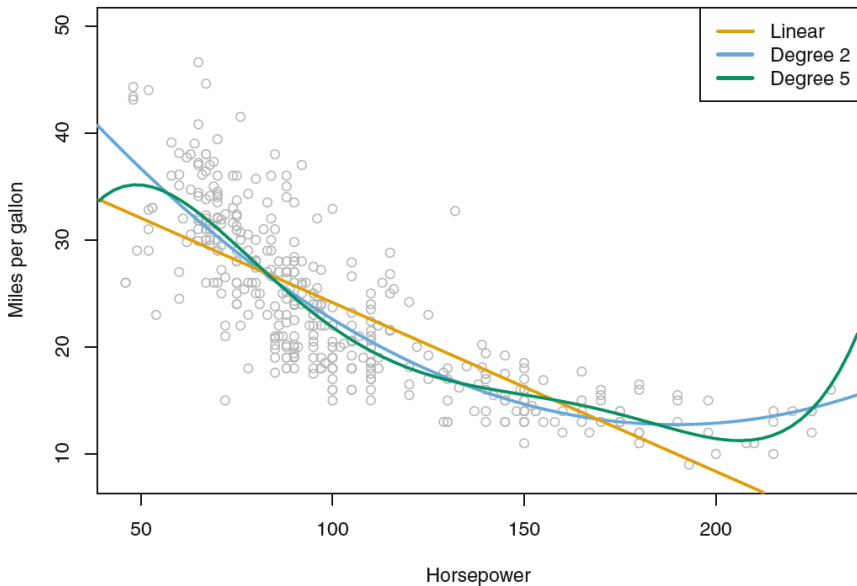
$$F = \frac{\frac{RSS_r - RSS_{ur}}{q}}{\frac{RSS_{ur}}{n-p-1}}$$

Interaction Effect

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$

	Coefficient	Std. error	t-statistic	p-value
Intercept	6.7502	0.248	27.23	< 0.0001
TV	0.0191	0.002	12.70	< 0.0001
radio	0.0289	0.009	3.24	0.0014
TV×radio	0.0011	0.000	20.73	< 0.0001

Auto Data Set



Non-linear Relationships

$$mpg = \beta_0 + \beta_1 hp + \beta_2 hp^2 + \epsilon$$

	Coefficient	Std. error	t-statistic	p-value
Intercept	56.9001	1.8004	31.6	< 0.0001
horsepower	-0.4662	0.0311	-15.0	< 0.0001
horsepower ²	0.0012	0.0001	10.1	< 0.0001

```
lm.fit=lm(medv~lstat +age,data=Boston)
```

Residuals:

Min	1Q	Median	3Q	Max
-15.981	-3.978	-1.283	1.968	23.158

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	33.22276	0.73085	45.458	< 2e-16	***
lstat	-1.03207	0.04819	-21.416	< 2e-16	***
age	0.03454	0.01223	2.826	0.00491	**

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.173 on 503 degrees of freedom

Multiple R-squared: 0.5513, Adjusted R-squared: 0.5495

F-statistic: 309 on 2 and 503 DF, p-value: < 2.2e-16

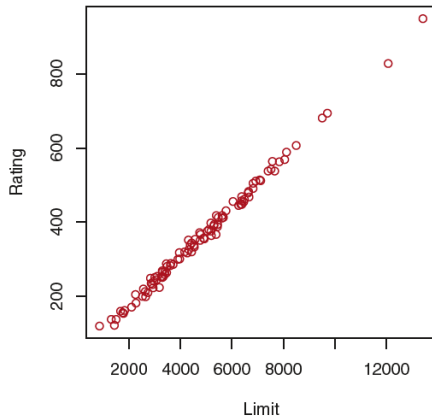
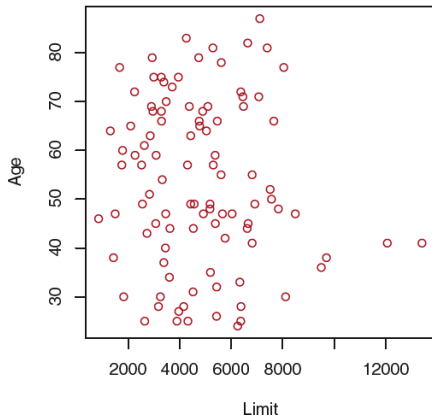
lm.fit=lm(medv~.,data=Boston)

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	3.646e+01	5.103e+00	7.144	3.28e-12	***
crim	-1.080e-01	3.286e-02	-3.287	0.001087	**
zn	4.642e-02	1.373e-02	3.382	0.000778	***
indus	2.056e-02	6.150e-02	0.334	0.738288	
chas	2.687e+00	8.616e-01	3.118	0.001925	**
nox	-1.777e+01	3.820e+00	-4.651	4.25e-06	***
rm	3.810e+00	4.179e-01	9.116	< 2e-16	***
age	6.922e-04	1.321e-02	0.052	0.958229	
dis	-1.476e+00	1.995e-01	-7.398	6.01e-13	***
rad	3.060e-01	6.635e-02	4.613	5.07e-06	***
tax	-1.233e-02	3.760e-03	-3.280	0.001112	**
ptratio	-9.527e-01	1.308e-01	-7.283	1.31e-12	***
black	9.312e-03	2.686e-03	3.467	0.000573	***
lstat	-5.248e-01	5.072e-02	-10.347	< 2e-16	***

$$Y = \text{balance}$$

		Coefficient	Std. error	t-statistic	p-value
Model 1	Intercept	-173.411	43.828	-3.957	< 0.0001
	age	-2.292	0.672	-3.407	0.0007
	limit	0.173	0.005	34.496	< 0.0001
Model 2	Intercept	-377.537	45.254	-8.343	< 0.0001
	rating	2.202	0.952	2.312	0.0213
	limit	0.025	0.064	0.384	0.7012

Scatterplots: Credit Data Set



Variance Inflation Factor (VIF)

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R_{X_j|X_{-j}}^2}$$

```
library(car); vif(lm.fit)
```

```
      crim      zn      indus      chas      nox      rm
1.792192 2.298758 3.991596 1.073995 4.393720 1.933744
      age      dis      rad      tax      ptratio      black
3.100826 3.955945 7.484496 9.008554 1.799084 1.348521
      lstat
2.941491
```

> 5 or 10 indicates a problematic collinearity

summary(lm(medv~lstat*age,data=Boston))

Residuals:

Min	1Q	Median	3Q	Max
-15.806	-4.045	-1.333	2.085	27.552

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	36.0885359	1.4698355	24.553	< 2e-16	***
lstat	-1.3921168	0.1674555	-8.313	8.78e-16	***
age	-0.0007209	0.0198792	-0.036	0.9711	
lstat:age	0.0041560	0.0018518	2.244	0.0252	*

Signif. codes:

0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.149 on 502 degrees of freedom

Multiple R-squared: 0.5557, Adjusted R-squared: 0.5531

F-statistic: 209.3 on 3 and 502 DF, p-value: < 2.2e-16


```
lm.fit2=lm(medv~lstat+I(lstat^2))
```

Residuals:

Min	1Q	Median	3Q	Max
-15.2834	-3.8313	-0.5295	2.3095	25.4148

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	42.862007	0.872084	49.15	<2e-16	***
lstat	-2.332821	0.123803	-18.84	<2e-16	***
I(lstat^2)	0.043547	0.003745	11.63	<2e-16	***

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

Residual standard error: 5.524 on 503 degrees of freedom
Multiple R-squared: 0.6407, Adjusted R-squared: 0.6393
F-statistic: 448.5 on 2 and 503 DF, p-value: < 2.2e-16

```
lm.fit=lm(medv~lstat)
```

```
anova(lm.fit,lm.fit2)
```

```
Model 1: medv ~ lstat
Model 2: medv ~ lstat + I(lstat^2)
  Res.Df  RSS Df Sum of Sq    F      Pr(>F)
1     504 19472
2     503 15347   1    4125.1 135.2 < 2.2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```