

Part 17: Multiple Linear Regression

Text references: §3.2 in ISL, Chapters 12 and 13 in Ruppert.

In this part, we will extend from the simple linear regression model to **multiple regression**.

These models comprise the familiar linear models that take the form

$$Y = \beta_0 + \sum_{j=1}^p \beta_j x_j + \epsilon$$

where ϵ are assumed to be uncorrelated with mean zero and variance σ^2 . Assuming that the ϵ are normally distributed brings additional nice properties.

We will go through the fundamental steps of fitting such a model via an example.

Example: Factor Models

References to the “beta” for a stock derive from the slope of the regression line fit to a scatter plot where excess return for that equity is the response, and excess market return is the predictor. Models of this form are an important component of Capital Asset Pricing Model (CAPM).

The term **Factor Model** is used to indicate any linear model for excess return on an equity. The predictors in the model are the **factors**.

Quoting Ruppert, examples of factors include

1. returns on the market portfolio;
2. growth rate of the GDP;
3. interest rate on short term Treasury bills or changes in this rate;
4. inflation rate or changes in this rate;
5. interest rate spreads;
6. return on some portfolio of stocks;
7. the difference between the returns on two portfolios.

- 1 The CAPM only includes the factor for excess market return.
- 2 The **Fama-French Three Factor Model** adds two factors.
- 3 The first is **small minus big (SMB)**, the difference in returns between
- 4 portfolios of small and large stocks.
- 5 The second is **high minus low (HML)**, the difference in returns be-
- 6 tween portfolios of high and low book-to-market stocks.
- 7 See Ruppert for more details on these factors, and motivation on
- 8 their inclusion in the model.

1 Kenneth French maintains a web site with data on the three factors
2 `http://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html`

3 I wrote a simple R function that will allow for easy access to the daily
4 factor data.

5 This function, named `getFamaFrench()`, is located at
6 `http://www.stat.cmu.edu/~cschafer/MSCF/getFamaFrench.txt`

7 The syntax is much like the functions in `quantmod`:
8 `> ffhold = getFamaFrench(from="2016-1-1",`
9 `to="2016-6-30")`

10 (Note that data is typically unavailable for the most recent months.)

11 The first few rows of the data returned are as follows:¹

	Date	Mkt.RF	SMB	HML	RF
12					
13	23638 2016-01-04	-1.59	-0.83	0.53	0.000
14	23639 2016-01-05	0.12	-0.21	0.00	0.000
15	23640 2016-01-06	-1.35	-0.13	0.01	0.000
16	23641 2016-01-07	-2.44	-0.28	0.12	0.000
17	23642 2016-01-08	-1.11	-0.47	-0.04	0.000
18	23643 2016-01-11	-0.06	-0.65	0.35	0.000
19	...				

¹These factors are periodically revised, and hence the exact numbers may be different at the present time. This will affect the subsequent analysis, as well. The analysis shown here was performed on January 8, 2017.

1 Let's get the data for PNC for the same time period

2 > PNC = getSymbols("PNC", from="2016-1-1",

3 to="2016-6-30", auto.assign=F)

4 Now find the excess returns for PNC. Note that the Fama-French

5 data file includes the risk free rate, with name RF:

6 > ffhold\$PNCexret = 100*dailyReturn(PNC) - ffhold\$RF

1 The three factor model we fit is

$$2 \quad \text{PNC.Ex.Ret.} = \beta_0 + \beta_1 \text{Mkt.Ex.Ret.} + \beta_2 \text{SMB} + \beta_3 \text{HML} + \epsilon$$

3 The R command for fitting a linear regression is `lm()`:

```
4 > ff3modPNC = lm(PNCexret ~ Mkt.RF + SMB + HML,  
5 data=ffhold)
```

6 Note the use of the `data` argument to specify the data frame from
7 which the variables are taken.

8 The key information from the fit of the model is stored in the object
9 `ff3modPNC`. In particular, note that

```
10 > attributes(ff3modPNC)  
11 $names  
12 [1] "coefficients" "residuals" "effects" "rank"  
13 [5] "fitted.values" "assign" "qr" "df.residual"  
14 [9] "xlevels" "call" "terms" "model"
```

15 For example, you can obtain the residuals from

16 `as.numeric(ff3modPNC$residuals)`.²

²It is necessary to transform `ff3modPNC$residuals` using `as.numeric()` because it inherits the time series format from the original data.

1 Diagnostic Plots

2 Figure 1 shows three key diagnostic plots:

- 3 1. Plot of residuals versus fitted values.
- 4 2. Normal probability plot.
- 5 3. Plot of residuals versus time.

6 Below are the R commands to create this figure.

```
7 par(mfrow=c(2,2))
8
9 # Plot of residuals versus fitted values
10 plot(as.numeric(ff3modPNC$fit), as.numeric(ff3modPNC$resid),
11      pch=16, xlab="Fitted Values", ylab="Residuals",
12      cex.axis=1.3, cex.lab=1.3)
13
14 # Normal probability plot
15 qqnorm(as.numeric(ff3modPNC$resid), cex.axis=1.3, cex.lab=1.3,
16        pch=16, main="")
17 qqline(as.numeric(ff3modPNC$resid))
18
19 # Plot of residuals versus time
20 plot(ff3modPNC$resid, xlab="Time", ylab="Residuals", cex.axis=1.3,
21      cex.lab=1.3, pch=16, main="")
```

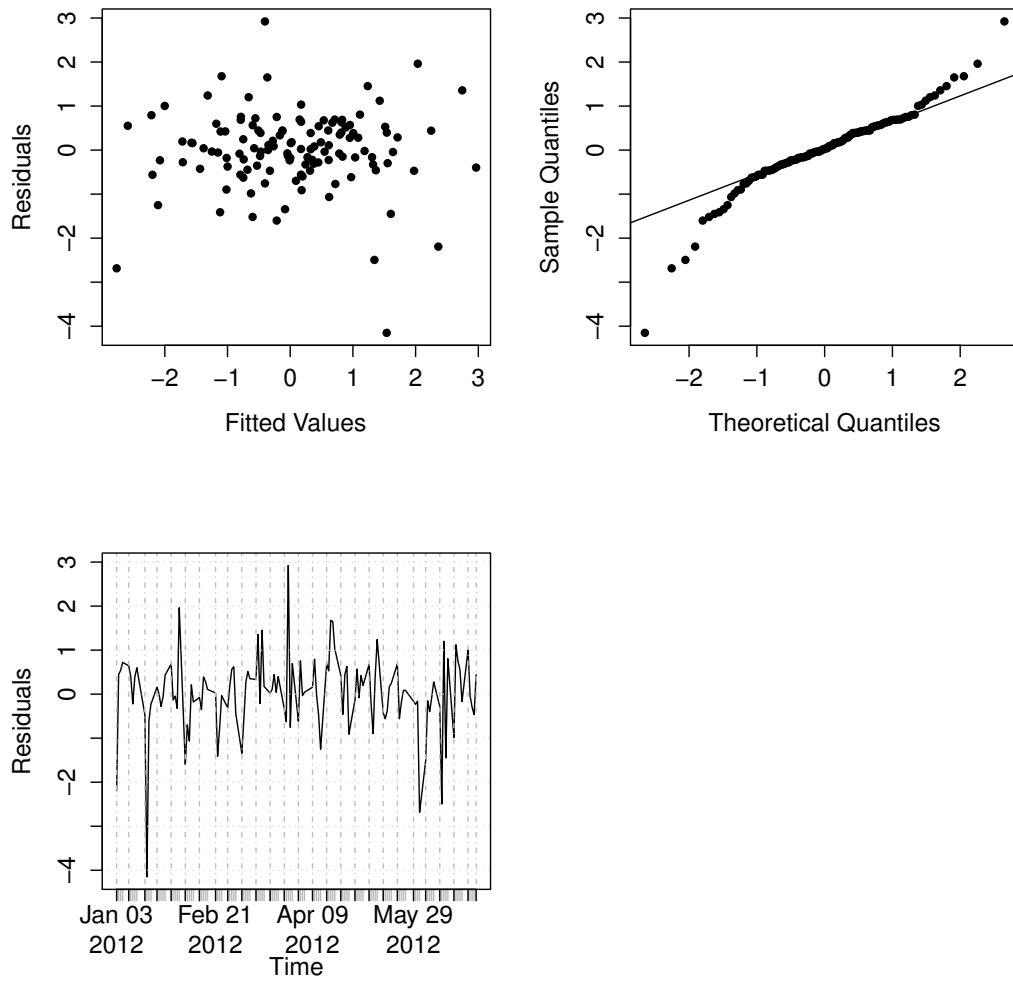


Figure 1: Diagnostic plots for the three factor model for PNC.

1 **Exercise:** Comment on the quality of the fit.

2 _____

3 _____

4 _____

5 _____

6 _____

- 1 Use the `summary()` function to see the parameter estimates and
2 their standard errors:³

```
3 Call:
4 lm(formula = PNCexret ~ Mkt.RF + SMB + HML, data = ffhold)
5
6 Residuals:
7      Min       1Q   Median       3Q      Max
8 -2.35913 -0.51226  0.09069  0.41264  3.08356
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) -0.16153    0.06945  -2.326   0.0217 *
13 Mkt.RF       1.20089    0.06952  17.275 < 2e-16 ***
14 SMB        -0.03428    0.13890  -0.247   0.8055
15 HML         0.78117    0.11937   6.544 1.52e-09 ***
16 ---
17 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
18
19 Residual standard error: 0.7755 on 121 degrees of freedom
20 Multiple R-squared:  0.7699, Adjusted R-squared:  0.7642
21 F-statistic: 135 on 3 and 121 DF,  p-value: < 2.2e-16
```

³Again, a reminder that your results could be slightly different since the Fama-French factors are periodically updated.

1 **Exercise:** There is strong evidence that β_3 , the coefficient for HML,
2 is larger than zero. How did I reach this conclusion, and what is the
3 interpretation?

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 _____

11 _____

12 _____

13 _____

14 _____

15 _____

16 _____

17 _____

18 _____

19 _____

Linear Regression in Matrix Notation

First, recall our simple linear regression model:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, \quad i = 1, 2, \dots, n$$

We can rewrite this model in matrix notation by letting

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

and then noting that

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

When moving to multiple regression, the main change we must make is to add columns to the matrix \mathbf{X} , commonly called the **design matrix**, and add entries to $\boldsymbol{\beta}$:

$$\mathbf{X} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,p} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,p} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

1 One can show that the least squares estimators (now a vector) are

2
$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

3 and can now define the vector of fitted values and residuals as

4
$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix} \quad \text{and} \quad \hat{\boldsymbol{\epsilon}} = \begin{bmatrix} \hat{\epsilon}_1 \\ \hat{\epsilon}_2 \\ \vdots \\ \hat{\epsilon}_n \end{bmatrix}$$

5 and it follows that

6
$$\hat{\mathbf{Y}} = \mathbf{X} \hat{\boldsymbol{\beta}} \quad \text{and} \quad \hat{\boldsymbol{\epsilon}} = \mathbf{Y} - \hat{\mathbf{Y}}.$$

1 The Hat Matrix

2 We note that we could write

$$3 \quad \widehat{\mathbf{Y}} = \mathbf{X}\widehat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} = \mathbf{H}\mathbf{Y}$$

4 where

$$5 \quad \mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

6 is called the **hat matrix** (because it puts a “hat” on \mathbf{Y}).

7 The diagonal entries of \mathbf{H} are referred to as the **leverages**.

8 **Comment:** Mathematically, \mathbf{H} is a **projection matrix**, as it “projects”
9 the vector of observed responses \mathbf{Y} onto the space of all vectors which
10 are linear combinations of the columns of \mathbf{X} . In fact, we note that \mathbf{H}
11 is symmetric and **idempotent**, meaning that $\mathbf{H}\mathbf{H} = \mathbf{H}$.

12 **Exercise:** Show that $\mathbf{I} - \mathbf{H}$ is also symmetric and idempotent.

13 _____

14 _____

15 _____

16 _____

17 _____

Additional Results

A range of important results can be proven:

1. The variance of $\hat{\beta}$ is $\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$.

2. $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta} = \mathbf{H}\mathbf{Y}$, where the (symmetric, idempotent) hat matrix \mathbf{H} is defined as above.

3. $\hat{\epsilon} = \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$.

4. If ϵ is assumed normal with mean zero and variance $\sigma^2 \mathbf{I}$,

(a) $\hat{\beta}$ is the maximum likelihood estimator.

(b) \mathbf{Y} is multivariate normal with mean $\mathbf{X}\beta$ and covariance $\sigma^2 \mathbf{I}$.

(c) $\hat{\beta}$ is multivariate normal with mean β and covariance $\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$.

(d) $\hat{\mathbf{Y}}$ is multivariate normal with mean $\mathbf{X}\beta$ and covariance $\sigma^2 \mathbf{H}$.

(e) $\hat{\epsilon}$ is multivariate normal with mean zero and covar. $\sigma^2(\mathbf{I} - \mathbf{H})$.

Degrees of Freedom

One key thing that does change when moving to multiple regression is the number of **degrees of freedom** in the residual vector.

Recall that we had seen previously that with simple linear regression, it held that

$$\sum_i \hat{\epsilon}_i = 0 \quad \text{and} \quad \sum_i \hat{\epsilon}_i X_i = 0.$$

This result is one way of seeing why the residuals only have $n - 2$ degrees of freedom in the case of simple linear regression: There are two constraints placed on the vector of residuals.

But, these facts could be written more compactly using matrix notation: $\mathbf{X}^T \hat{\boldsymbol{\epsilon}} = \mathbf{0}$, where $\mathbf{0}$ is a vector consisting entirely of zeros.

In fact, this result extends to multiple regression: The vector of residuals $\hat{\boldsymbol{\epsilon}}$ is orthogonal to **every column** of \mathbf{X} , i.e., $\mathbf{X}^T \hat{\boldsymbol{\epsilon}} = \mathbf{0}$.

1 **We now can state that there are $n - (p + 1)$ degrees of freedom in the**
2 **residuals**, and hence need to update some of our previous results.
3 (You will note that simple linear regression corresponds to the case
4 where $p = 1$.) In particular:

5 1. The unbiased estimator of σ^2 is

6
$$\hat{\sigma}^2 = \text{RSS} / (n - p - 1)$$

7 2. The statistic

8
$$\frac{\hat{\beta}_i - \beta_i}{\widehat{\text{SE}}(\hat{\beta}_i)}$$

9 has the t -distribution with $n - p - 1$ degrees of freedom. Hence,
10 the following hold:

11 (a) A $100(1 - \alpha)\%$ confidence interval for β_i is formed as

12
$$\hat{\beta}_i \pm t_{\alpha/2, n-p-1} \widehat{\text{SE}}(\hat{\beta}_i)$$

13 (b) Hypothesis tests concerning β_i should compare the test statis-
14 tic (the “t value”) with the t -distribution with $n - p - 1$ degrees
15 of freedom.

1

2 **Basic Variable Selection**

3 The basic, general strategy for avoiding overfitting with models is
4 to compare the values of AIC across different number/choices of co-
5 variates

6 Recall that we prefer the model which has the smallest value of AIC

7 Unfortunately, when the number of covariates is large, the number of
8 possible models is enormous. (There are 2^p possible models if there
9 are p covariates.)

- 1 We'll utilize the regression problem introduced earlier in which we
2 sought to estimate the forward-rate function based on the current
3 prices of zero-coupon bonds. The data are shown again as Figure 2.

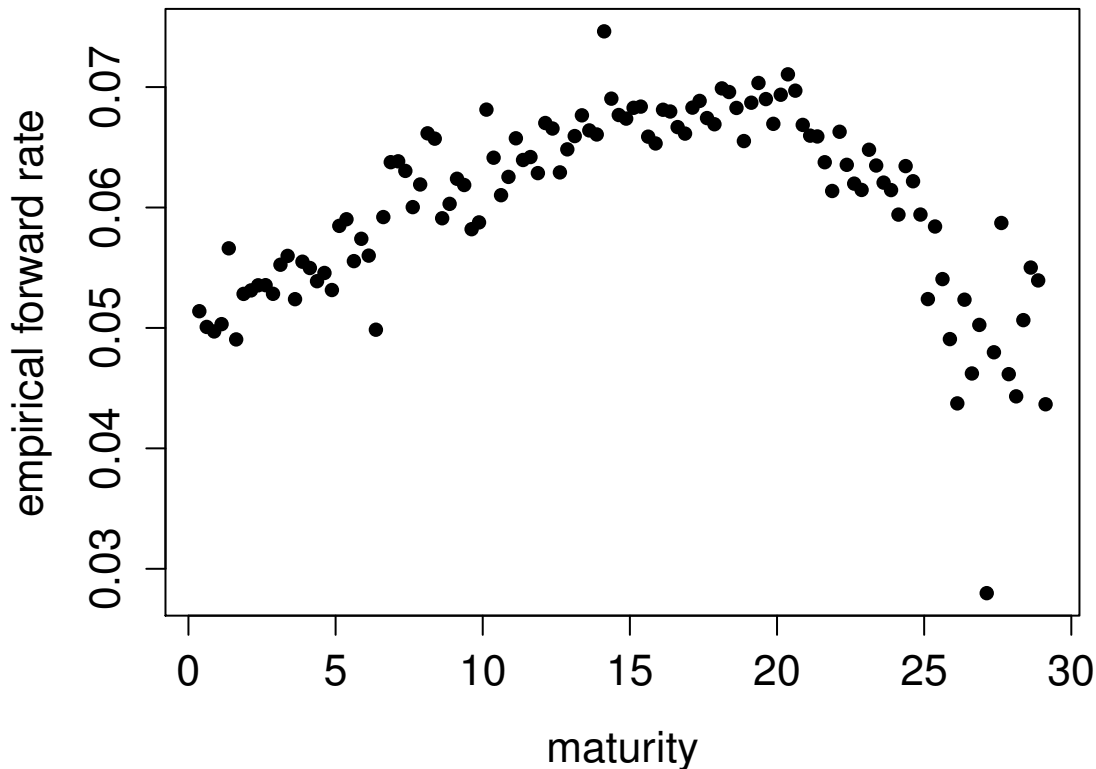


Figure 2: The response versus time for the estimation of the forward rate function.

- 4 You can read in this data set using the command

```
5 > forratedat =  
6   read.table("http://www.stat.cmu.edu/~cschafer/MSCF/forratedat.txt",  
7   header=T)  
8 > attach(forratedat)
```

1 We are considering polynomial models of the form

$$2 \quad Y_i = \sum_{j=0}^p \beta_j T_i^j + \epsilon_i$$

3 but our goal is to find the best subset of the p polynomial terms to
4 include in the model. (As is usually the case, we do not consider
5 removing the intercept β_0 .)

6 The “full” model (largest possible model) we’ll consider includes up
7 to the ninth power.

8 We will scale the predictor first, so that it has mean zero and standard
9 deviation one:

```
10 > maturity = scale(maturity)
```

11 Now, we can fit the full model:

```
12 > fullmod = lm(emp ~ maturity + I(maturity^2) +  
13     I(maturity^3) + I(maturity^4) + I(maturity^5) +  
14     I(maturity^6) + I(maturity^7) + I(maturity^8) +  
15     I(maturity^9))
```

16 Note the use of `I()` function to wrap the polynomial terms is re-
17 quired by R.

Exhaustive Search

In cases where the number of predictors is small, an exhaustive search over all possible models is possible.

One implementation of this in R is the function `bestglm()` which is part of the package `bestglm`.

Unfortunately, the syntax for this function requires that the predictors and response be in a single data frame, with the response as the last column.

Consider the following commands

```
> allpreds = cbind(maturity, maturity^2,  
  maturity^3, maturity^4, maturity^5,  
  maturity^6, maturity^7, maturity^8,  
  maturity^9)  
> Xyframe = data.frame(cbind(allpreds, emp))
```

Note that now `Xyframe` holds the predictors and the response.

```
> bestmod = bestglm(Xyframe, IC="AIC")
```

By specifying `IC = "AIC"` we are telling the function to use AIC to choose the model.

1 The R command

2 `> print(bestmod)`

3 will show the output below

4 AIC

5 BICq equivalent for q in (0.0443666836953486, 0.908095644854053)

6 Best Model:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.066820498	0.0004770183	140.079524	8.381181e-126
V1	0.004950250	0.0008628244	5.737262	8.591912e-08
V4	-0.008479843	0.0007636180	-11.104824	1.134645e-19
V5	-0.004005042	0.0008106249	-4.940684	2.801999e-06
V6	0.002275548	0.0002881148	7.898061	2.313625e-12
V7	0.001183985	0.0002729549	4.337659	3.206729e-05

14 We see that the optimal model takes the form

15
$$Y_i = \beta_0 + \beta_1 T_i + \beta_4 T_i^4 + \beta_5 T_i^5 + \beta_6 T_i^6 + \beta_7 T_i^7 + \epsilon_i$$

1 **Stepwise Regression**

2 An exhaustive search is not feasible when there is a large number of
3 predictors.

4 A classic strategy for dealing with the large number of possible mod-
5 els is to use a **stepwise variable selection** procedure.

6 With this method, one (typically) starts with the largest model under
7 consideration, and then takes a sequence of **steps**: At each step one
8 covariate is either **added** or **dropped**.

9 The decision for which step to take (i.e., which covariate to add or
10 drop) is based on AIC: You take the step which reduces AIC the most.

11 Once AIC no longer changes, the process stops, and you have your
12 final model.

13 This algorithm is far from perfect, but is useful in the appropriate
14 situations.

15 This is built into R using the function `step()`

1 Now we will use the `step()` function on the example above.

2 The syntax is simple:

3 `> finalmod = step(fullmod, direction="both")`

4 Recall that `fullmod` holds the output of the full model fit.

5 The use of `direction = "both"` tells R to consider not only re-
6 moving predictors, but also adding predictors to the model. (Once
7 a predictor is removed, it could be returned to the model in a later
8 step.)

1 The output of the first step of algorithm appears below.

```
2 Start:  AIC=-1296.97
3 emp ~ maturity + I(maturity^2) + I(maturity^3) + I(maturity^4) +
4       I(maturity^5) + I(maturity^6) + I(maturity^7) + I(maturity^8) +
5       I(maturity^9)
6
7           Df  Sum of Sq      RSS      AIC
8 - I(maturity^6)  1 3.0370e-07 0.0013612 -1298.9
9 - I(maturity^9)  1 3.3910e-07 0.0013613 -1298.9
10 - I(maturity^8)  1 9.2100e-07 0.0013619 -1298.9
11 - I(maturity^3)  1 9.7950e-07 0.0013619 -1298.9
12 - I(maturity^7)  1 1.7992e-06 0.0013627 -1298.8
13 - I(maturity^2)  1 2.1071e-06 0.0013630 -1298.8
14 - I(maturity^5)  1 3.4433e-06 0.0013644 -1298.7
15 - I(maturity^4)  1 4.7036e-06 0.0013656 -1298.6
16 - maturity      1 2.2945e-05 0.0013839 -1297.0
17 <none>           0.0013609 -1297.0
```

18 Note how the table lists each of the covariates currently in the model,
19 and gives what AIC would be if that covariate were removed.

20 We also see the line labelled <none>, indicating that the AIC would
21 stay at -1297.0 if none were removed.

22 The covariates are ordered by increasing values of AIC. (Unfortu-
23 natley, due to rounding it's not possible to see some distinctions in
24 this output.)

1 So, the term $I(\text{maturity}^6)$ is removed at the first step.

2 The second step is shown as:

3 Step: AIC=-1298.94

4 `emp ~ maturity + I(maturity^2) + I(maturity^3) + I(maturity^4) +`
5 `I(maturity^5) + I(maturity^7) + I(maturity^8) + I(maturity^9)`

```
6
7           Df Sum of Sq      RSS      AIC
8 - I(maturity^9)  1 3.3900e-07 0.0013616 -1300.9
9 - I(maturity^3)  1 9.8000e-07 0.0013622 -1300.9
10 - I(maturity^7)  1 1.7990e-06 0.0013630 -1300.8
11 - I(maturity^5)  1 3.4430e-06 0.0013647 -1300.7
12 - I(maturity^2)  1 1.9184e-05 0.0013804 -1299.3
13 - maturity      1 2.2945e-05 0.0013842 -1299.0
14 <none>                                0.0013612 -1298.9
15 + I(maturity^6)  1 3.0400e-07 0.0013609 -1297.0
16 - I(maturity^4)  1 8.0262e-05 0.0014415 -1294.3
17 - I(maturity^8)  1 1.2814e-04 0.0014894 -1290.5
```

18 The decision of this step is to remove $I(\text{maturity}^9)$. Note that
19 also under consideration was to add back in the term that was re-
20 moved.

1 This process continues until a final model is reached:

```
2 Step:  AIC=-1303.2
3 emp ~ maturity + I(maturity^4) + I(maturity^5) + I(maturity^7) +
4       I(maturity^8)
5
6           Df Sum of Sq      RSS      AIC
7 <none>                 0.0013818 -1303.2
8 + I(maturity^2)    1 0.00001918 0.0013626 -1302.8
9 + I(maturity^6)    1 0.00001738 0.0013644 -1302.7
10 + I(maturity^3)   1 0.00000105 0.0013808 -1301.3
11 + I(maturity^9)   1 0.00000041 0.0013814 -1301.2
12 - I(maturity^7)   1 0.00023342 0.0016152 -1287.1
13 - I(maturity^5)   1 0.00030284 0.0016847 -1282.2
14 - maturity       1 0.00040836 0.0017902 -1275.2
15 - I(maturity^8)   1 0.00075672 0.0021385 -1254.5
16 - I(maturity^4)   1 0.00212230 0.0035041 -1197.3
```

17 So, we have arrived (after only four steps) at a model that includes
18 powers 1, 4, 5, 7, and 8.

19 **Exercise:** How does this model compare to the one we stated as being
20 optimal using the exhaustive search? What happened?

21 _____

22 _____

23 _____

24 _____

25

Cross Validation

We've discussed the use of AIC for the purpose of **model selection**, i.e. the process of determining which covariates should be included in the model.

An important alternative is **leave-one-out cross validation**.

In this approach, one imagines refitting the model n times, each time excluding one of the n observations.

At iteration i , observation i is excluded. The fitted model is then used to predict the response for this observation. Call this $\hat{y}_{(-i)}$.

Finally, we calculate the quantity

$$\text{PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{(-i)})^2$$

where **PRESS** stands for **Prediction Error Sum of Squares**.

1 The motivation for this is as follows: Once observation i is removed
2 from the fit, we have removed the **influence** that this observation has
3 on the parameter estimates.

4 Now, we can think of observation i as being a “new” observation,
5 and we can ask: “How well does the model, containing only these
6 covariates, predict the response for this case?”

7 The quantity $(y_i - \hat{y}_{(-i)})^2$ quantifies the amount of error in this pre-
8 diction.

9 **Exercise:** Which is larger: PRESS or RSS?

10 _____

11 _____

12 _____

13 _____

14 _____

15 _____

16 _____

17 _____

18 _____

19 _____

1 Fortunately, it is not actually necessary to fit n models in order calcu-
2 late PRESS.

3 The remarkable result is as follows:

5 Let \mathbf{Y} be a vector consisting of the n responses, and let $\widehat{\mathbf{Y}}$ be a
6 vector consisting of the n fitted values (from the full model).

7 If it is the case that $\widehat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$, then

8
$$y_i - \hat{y}_{(-i)} = \frac{\hat{\epsilon}_i}{1 - h_{ii}}$$

9 where h_{ii} is the i^{th} diagonal element of \mathbf{H} .

11 In the case of linear regression, it is the case that

12
$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

13 where \mathbf{X} whose columns hold the values of the predictors. See Rup-
14 pert for the details.

15 This matrix \mathbf{H} is previously mentioned **hat matrix**.

16 The diagonal elements h_{ii} are called the **leverages**.

1 The leverages can be found in R using

```
2 > levs = hatvalues(fullmod)
```

3 Then we can find PRESS using

```
4 > PRESS = sum((fullmod$resid/(1-levs))^2)
```

5 We find PRESS to be 0.001974 for this model.

6 Of course, our goal is to find the model with lowest PRESS. In R we
7 can do this using the function `bestglm()` with `IC="LOOCV"`:

```
8 > bestmod2 = bestglm(Xyframe, IC="LOOCV")
```

9 LOOCV

10 BICq equivalent for q in (0.0443666836953498, 0.908095644854043)

11 Best Model:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.066820498	0.0004770183	140.079524	8.381181e-126
V1	0.004950250	0.0008628244	5.737262	8.591912e-08
V4	-0.008479843	0.0007636180	-11.104824	1.134645e-19
V5	-0.004005042	0.0008106249	-4.940684	2.801999e-06
V6	0.002275548	0.0002881148	7.898061	2.313625e-12
V7	0.001183985	0.0002729549	4.337659	3.206729e-05

19 We arrive at the same model as when we used AIC.

1 **AIC versus PRESS**

2 AIC and PRESS will often give similar, if not the same, result for the
3 model choice.

4 There are a couple main reasons to prefer PRESS to AIC: First, the
5 quantity being estimated is by PRESS is a very natural measure of
6 prediction performance. Second, the theory behind AIC is based on
7 the assumption that the distributional assumption for the response is
8 correct.

9 A main reason to choose AIC over PRESS is that AIC is more stable.
10 The estimate that PRESS provides of the expected prediction error is
11 subject to large variance. Hence, PRESS may not perform well for
12 small to moderate sample sizes.

13 Although PRESS is simple to calculate in the case of linear regression,
14 this will not always be the case. In general, AIC will be easier to
15 calculate than PRESS.

16 As is often the case, our past experience will often guide our choice
17 of criterion.

Comment: AIC Calculations in R

One must use great care when comparing values of AIC reported by different functions in R. There are additive constants that are sometimes arbitrarily discarded. Of course, when using the same R function for comparing two models, there will be no problem.

Here is an important example. Suppose that the response values Y_1, Y_2, \dots, Y_n are modeled such that they are independent, normally-distributed random variables with mean $m(\boldsymbol{\beta}, \mathbf{x}_i)$ and variance σ^2 . Here, \mathbf{x}_i is the vector of predictors.

Exercise: Show that the value of the log-likelihood is

$$\ell(\boldsymbol{\beta}, \sigma^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \sum_{i=1}^n \left[\frac{(y_i - m(\boldsymbol{\beta}, \mathbf{x}_i))^2}{2\sigma^2} \right].$$

in this case.

1 Then, one can show that when you evaluate the log-likelihood at its
2 peak,

$$3 \quad -2\ell(\hat{\beta}, \hat{\sigma}^2) = n \log(2\pi) + n \log(\text{RSS}/n) + n.$$

4 The R function `AIC()` returns this full expression, plus two times the
5 number of **total** number of parameters (those in β , and σ^2).

6 The R function `extractAIC()` returns only

$$7 \quad n \log(\text{RSS}/n)$$

8 plus two times the number of parameters in β .

9 Hence, the values returned by `AIC()` and `extractAIC()` differ by

$$10 \quad n \log(2\pi) + n + 2.$$

11 This is constant across all models being considered, and hence is
12 of no practical importance. Continuing our example from above,
13 (where $n = 125$),

```
14 > AIC(ff3modPNC)
15 [1] 297.1191
16 > extractAIC(ff3modPNC)
17 [1] 4.00000 -59.61554
```

18 The function `step()` uses `extractAIC()` when calculating AIC.

Analysis of Variance

The term **analysis of variance (ANOVA)** is somewhat old-fashioned, but refers broadly to tools and tests which attempt to decompose the **total variability in the response into different components**.

First, we will define the **total sum of squares (SSTO or SST)** as

$$\text{SSTO} = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Exercise: Interpret the quantity SSTO in the context of regression. Think about the simplest possible model for predicting the response.

1 We have already defined the **residual sum of squares (RSS)** or **sum**
2 **of squared errors (SSE)** as

3
$$\text{SSE} = \sum_{i=1}^n (Y_i - \widehat{Y}_i)^2$$

4 We will also define the **regression sum of squares (SSR)** as

5
$$\text{SSR} = \sum_{i=1}^n (\widehat{Y}_i - \bar{Y})^2$$

6 **Exercise:** Show that $\text{SSTO} = \text{SSR} + \text{SSE}$.

7 _____

8 _____

9 _____

10 _____

11 _____

12 _____

13 _____

14 _____

15 _____

16 _____

17 _____

18 _____

19 _____

20

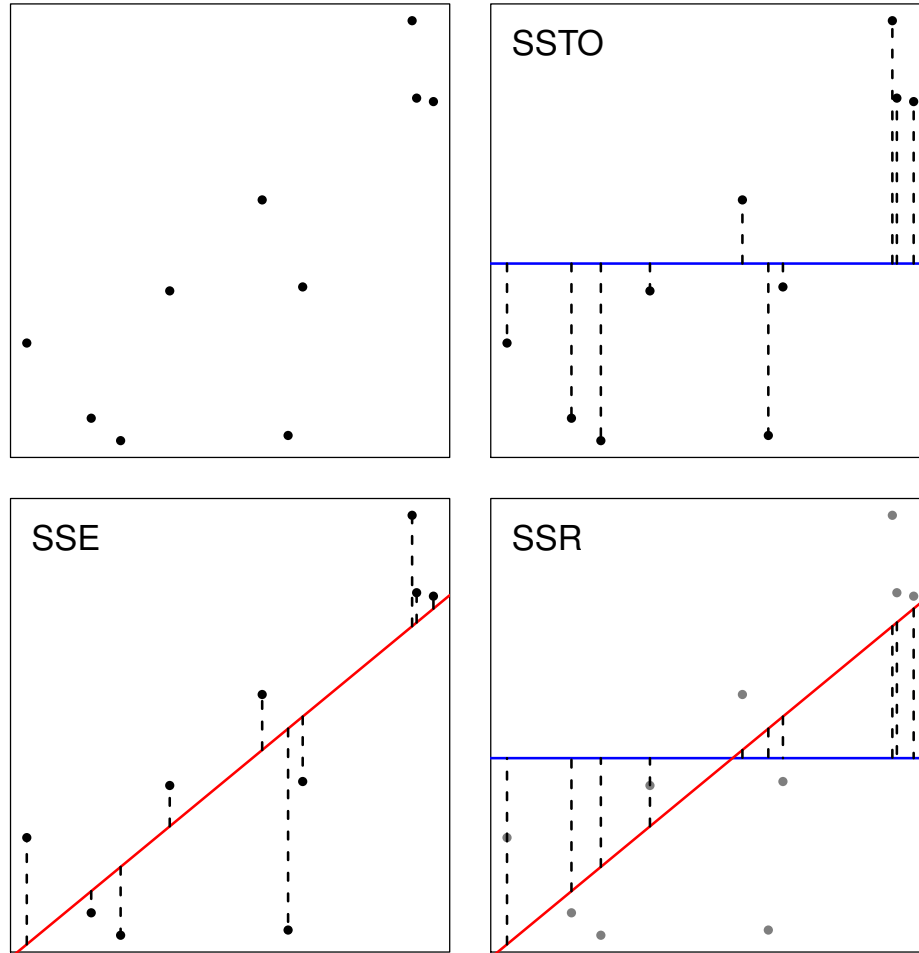


Figure 3: This figure shows the decomposition of the total sum of squares (SSTO) into the sum of squared errors (SSE) and the regression sum of squares (SSR). The top left plot shows the original scatter plot. The top right plot depicts how SSTO is calculated from the deviations of the points around the mean of the response values. The bottom left plot shows the residuals, these are squared and summed to find SSE. Finally, the bottom right plot shows how SSR is calculated from the deviations from the regression line to the mean line.

The ANOVA Table

Each of the three “sums of squares” defined previously has associated with it a certain number of **degrees of freedom**.

In particular, we say that “SSTO has $n - 1$ degrees of freedom” because the list of n deviations

$$Y_1 - \bar{Y}, Y_2 - \bar{Y}, \dots, Y_n - \bar{Y}$$

is subject to one restriction, namely that they must sum to zero.

We say that “SSE has $n - p - 1$ degrees of freedom” because the n residuals

$$Y_1 - \hat{Y}_1, Y_2 - \hat{Y}_2, \dots, Y_n - \hat{Y}_n$$

is subject to $p + 1$ restrictions, i.e. they are orthogonal to each of the columns of \mathbf{X} .

Finally, the remaining p degrees of freedom are attributed to SSR because there are p free parameters in the regression model.

1 This information is classically summarized in the **analysis of vari-**
2 **ance table**, as shown below.

3

Source	SS	DF	MS	F
Regression	SSR	p	SSR/p	MSR/MSE
4 Error	SSE	$n - p - 1$	$SSE/(n - p - 1)$	
Total	SSTO	$n - 1$		

5 The column **MS** provides the **mean squares** for each of regression
6 and error. You will note that the **mean squared error (MSE)** is simply
7 our unbiased estimator of σ^2 , i.e., $E(\text{MSE}) = \sigma^2$.

Global F Test for Regression

Assume that the irreducible errors are i.i.d. normal with mean zero and variance σ^2 , as is standard. Then we can prove the following:

Important Theoretical Result: If it is the case that

$$\beta_1 = \beta_2 = \cdots = \beta_p = 0$$

then SSR/σ^2 has the chi-squared distribution with p degrees of freedom. Further, SSE/σ^2 has the chi-squared distribution with $n - p - 1$ degrees of freedom. Finally, SSE and SSR are independent.

Exercise: What does this result tell us about the distribution of

$$F = \frac{MSR}{MSE}$$

when $\beta_1 = \beta_2 = \cdots = \beta_p = 0$?

1 **Exercise:** Explain how the results on the previous page form the basis
2 for a useful hypothesis test.

3 _____

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 _____

11 _____

12 _____

13 _____

14 _____

15 _____

16 _____

17 _____

18 _____

19 _____

20 _____

21

Testing a Subset of the β 's

The Global F Test described above takes an “all-or-nothing” approach to testing β parameters. Suppose instead one wants to test the null hypothesis

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$

versus

$$H_1: \text{at least one of } \beta_1, \beta_2, \dots, \beta_k \neq 0.$$

where $k < p$.⁴

Exercise: Explain why this test would be particularly useful when working with factors as predictors.

⁴The ordering of the predictors is arbitrary for the purposes of this test, so we can, without affecting anything, just assume that we want to test the first k β coefficients.

1 To see how this test is constructed, first refer to the regression model
2 with $\beta_1 = \beta_2 = \cdots = \beta_k = 0$ as the **reduced model**. It is “reduced” in
3 the sense that it has fewer parameters than the **full model**.

4 Imagine fitting this reduced model, and saving the SSE from this fit
5 as SSE_{red} . Then, under the normality assumption for the irreducible
6 errors, and when the null hypothesis

7
$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0.$$

8 is true, the statistic

9
$$F = \frac{(\text{SSE}_{\text{red}} - \text{SSE}_{\text{full}})/k}{\text{MSE}_{\text{full}}}$$

10 has the F -distribution with k and $n - p - 1$ degrees of freedom.

11 The test is conducted as above: If the value of F is too large, reject H_0
12 in favor of the alternative hypothesis that at least one of the β_i under
13 consideration is not equal to zero.

1 Conducting the Test in R

2 R has a built-in function `anova()` which, among other things, will
3 report the results of this test. The syntax is simple:

```
4 proslmfull = lm(log(PSA)~log(CV)+log(PW)+Age+sqrt(BPH)+  
5     sqrt(CP), data=prosdats)  
6  
7 proslmred = lm(log(PSA)~log(CV)+log(PW)+Age, data=prosdats)  
8  
9 > anova(proslmred, proslmfull)  
10 Analysis of Variance Table  
11  
12 Model 1: log(PSA) ~ log(CV) + log(PW) + Age  
13 Model 2: log(PSA) ~ log(CV) + log(PW) + Age + sqrt(BPH) + sqrt(CP)  
14  
15      Res.Df    RSS Df Sum of Sq      F Pr(>F)  
16 1         93 52.600  
17 2         91 50.729  2    1.8707 1.6779 0.1925
```

18 **Exercise:** Describe the conclusion of the test shown above.

Example: Predicting Insurance Claims Severity

This challenge comes from Kaggle, with details that can be found at <https://www.kaggle.com/c/allstate-claims-severity/>

How severe is an insurance claim?

When you've been devastated by a serious car accident, your focus is on the things that matter the most: family, friends, and other loved ones. Pushing paper with your insurance agent is the last place you want your time or mental energy spent. This is why [Allstate](#), a personal insurer in the United States, is continually seeking fresh ideas to improve their claims service for the over 16 million households they protect.



Allstate is currently developing automated methods of predicting the cost, and hence severity, of claims. In this recruitment challenge, Kagglers are invited to show off their creativity and flex their technical chops by creating an algorithm which accurately predicts claims severity. Aspiring competitors will demonstrate insight into better ways to predict claims severity for the chance to be part of Allstate's efforts to ensure a worry-free customer experience.

Figure 4: Screen shot from Kaggle.

1 The full training set can be downloaded from their site. We will work
2 with a reduced version, consisting of 5000 observations, each with 76
3 predictor values. Sixty-two of these predictors are categorical, while
4 the remaining 14 can be treated as continuous.

5 This data can be read in using

```
6 > insdat = read.table(  
7   "http://www.stat.cmu.edu/~cschafer/MSCF/Allstatesub.txt",  
8   sep=" ", header=T)
```

9 By looking at `names(insdat)`, one can see that the predictors have
10 been given generic names to protect the intellectual property of All-
11 state. The final column, named `loss`, is the response.

12 Our objective at this point is to learn some more advanced R syntax
13 for dealing with cases when there are a large number of predictors.

1 Consider the R code below.

```
2 fulllm = lm(loss ~ ., data=insdat)
3
4 predlist = attributes(terms(formula(fulllm)))$term.labels
5
6 pvalue = numeric(length(predlist))
7
8 for(i in 1:length(predlist))
9 {
10     holdout = drop1(fulllm, predlist[i], test="F")
11     pvalue[i] = holdout$`Pr(>F)`[2]
12 }
13
14 cbind(predlist, signif(pvalue))
```

15 **Exercise:** Go through the above code line-by-line to explore and un-
16 derstand the R syntax.

17 _____

18 _____

19 _____

20 _____

21 _____

22 _____

23 _____

24 _____

Part 17: Multiple Linear Regression

1 _____

2 _____

3 _____

4 _____

5 _____

6 _____

7 _____

8 _____

9 _____

10 _____

11 _____

12 _____

13 _____

14 _____

15 _____

16 _____

17 _____

18 _____

19 _____

20 _____