# Chapter 1

# Simple Linear Regression Model

## 1.1 What is regression?

The word was *originally* used by an English statistician Galton. In his 1885 paper to Royal Anthropological Institute, he derived regression and used it interchangeably with the word reversion.

$$\text{regression} \equiv \text{regression to the mean } (e.g., \text{ Galton 1885}).$$

**Example 1.1.** *Regression towards the mean.*

- Sports stars tend to have a poorer season following a really good one.

- Split the class into two groups – top 50% and lower 50%. We can expect the lower group to do better and the top group to do worse. △

**Definition 1.1.** Regression Analysis *is a statistical methodology that analyzes the relation between two or more quantitative variables so that one variable can be predicted from other(s). That is to say, it is equivalent to find a function relation of the form*

$$Y = f(X).$$

## 1.2  Simple linear regression model

We will start with the simplest regression model (simple linear regression model),

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i.$$

Here *simple* means a single regressor variable $X$. The $Y$ variables are called *dependent* or *response* variables, and the $X$'s are called *independent variables*, *explanatory variables*, *regressors*, *covariates*, or *predictor variables*. The parameters $\beta_0$ and $\beta_1$ are called *regression coefficients*.

Our goal is to find $\beta_0$ and $\beta_1$. It is natural to find $\beta_0$ and $\beta_1$ which minimize the errors $\epsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$.

### 1.2.1  Assumptions and properties

Assumptions: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

- $X_i$ is a known constant.

- $\epsilon_i$ is a *random error* with:

  1. $E(\epsilon_i) = 0$ (the errors tend to compensate each other).

     If $E(\epsilon_i) = \alpha$, then $\beta_0^* = \alpha + \beta_0$ can absorb $\alpha$.

  2. $\text{Var}(\epsilon_i) = \sigma^2$ (for each observation the error has the same variance).

  3. $\text{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$ (the errors corresponding to difference observations are uncorrelated so they do not influence each other).

Properties:

- The response variable $Y_i$ is the sum of two components called constant term ($f(X_i) = \beta_0 + \beta_1 X_i$) and random term ($\epsilon_i$). The first term $\beta_0 + \beta_1 X_i$ repre-

sents the basic relation between the variables and the latter $\epsilon_i$ represents the disturbances (errors) which affect this relation.

- Even if $X_i$ is assumed to be known and fixed, $Y_i$ is a random variable for the presence of $\epsilon_i$.

- The expected value of $Y_i$ is equal to the constant term

$$E(Y_i) = E(\beta_0 + \beta_1 X_i + \epsilon_i) = \beta_0 + \beta_1 X_i,$$

  and so, in practice, there is a linear relation between $E(Y)$ and $X$ ;

- The variance of $Y_i$ is equal to that of the error term

$$\mathrm{Var}(Y_i) = \mathrm{Var}(\beta_0 + \beta_1 X_i + \epsilon_i) = \mathrm{Var}(\epsilon_i) = \sigma^2.$$

- $Y_i$ and $Y_j$ $(i \neq j)$ are uncorrelated,

$$\mathrm{Cov}(Y_i, Y_j) = \mathrm{Cov}(\epsilon_i, \epsilon_j) = 0,$$

  and so the measure of the response variable in correspondence of a subject does not affect the measure in correspondence of another subject.

## 1.3 Estimation of regression function

### 1.3.1 Parameter estimation by least squares methods

Let us denote

$$Q_1 = \sum_{i=1}^{n} |\epsilon_i| = \sum_{i=1}^{n} |Y_i - (\beta_0 + \beta_1 X_i)|$$
$$Q_2 = \sum_{i=1}^{n} \epsilon_i^2 = \sum_{i=1}^{n} \{Y_i - (\beta_0 + \beta_1 X_i)\}^2$$
$$Q_3 = \sum_{i=1}^{n} \epsilon_i = \sum_{i=1}^{n} \{Y_i - (\beta_0 + \beta_1 X_i)\}.$$

Minimizing $Q_1$ is called $L_1$ regression and minimizing $Q_2$ is $L_2$ or least-squares regression.

Can we use $Q_3$? *Answer*: no.

We will focus on a $L_2$ regression. Differentiate $Q_2$ with respect to $\beta_0$ and $\beta_1$ to get

$$\frac{\partial Q_2}{\partial \beta_0} = -2 \sum_{i=1}^{n} \{Y_i - (\beta_0 + \beta_1 X_i)\} = 0 \tag{1.1}$$

$$\frac{\partial Q_2}{\partial \beta_1} = -2 \sum_{i=1}^{n} X_i \{Y_i - (\beta_0 + \beta_1 X_i)\} = 0. \tag{1.2}$$

Let us denote $\hat{\beta}_0$ (or $b_0$) and $\hat{\beta}_1$ (or $b_1$) to be the solution of the above equations which are called normal equations. The solution $\hat{\beta}_0$ and $\hat{\beta}_1$ are called *point estimators* of $\beta_0$ and $\beta_1$.

The normal equations can be solved simultaneously:

$$\hat{\beta}_1 = b_1 = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2} = \frac{\sum X_i Y_i - n\bar{X}\bar{Y}}{\sum X_i^2 - n\bar{X}^2} = \frac{S_{xy}}{S_{xx}} \tag{1.3}$$

$$\hat{\beta}_0 = b_0 = \frac{1}{n}\{\sum_{i=1}^{n} Y_i - \hat{\beta}_1 \sum_{i=1}^{n} X_i\} = \bar{Y} - \hat{\beta}_1 \bar{X}, \tag{1.4}$$

where $\bar{X} = (1/n)\sum_{i=1}^{n} X_i$ and $\bar{Y} = (1/n)\sum_{i=1}^{n} Y_i$.

The $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased. That is

$$E(\hat{\beta}_0) = \beta_0 \quad \text{and} \quad E(\hat{\beta}_1) = \beta_1.$$

## 1.3.2 Normal error regression model and MLE

By minimizing $Q_2$, we found the parameters of the $L_2$ regression model under the assumption that $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$, where $\epsilon_i$ satisfy that

1. $E(\epsilon_i) = 0$,

2. $\mathrm{Var}(\epsilon_i) = \sigma^2$, and

3. $\mathrm{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$.

There are so many possible distributions of $\epsilon_i$ satisfying $E(\epsilon_i) = 0$, $\mathrm{Var}(\epsilon_i) = \sigma^2$, and $\mathrm{Cov}(\epsilon_i, \epsilon_j) = 0$ for $i \neq j$. One appealing choice is $\epsilon_i \sim N(0, \sigma^2)$. If the distribution of the error terms is specified, estimators of the parameters $\beta_0$, $\beta_1$ and $\sigma^2$ can be obtained by the method of *maximum likelihood.*

> **Fact.** *It is shown that* $\mathrm{Cov}(U, V) = 0$ *if the random variables* $U$ *and* $V$ *are independent. The converse, in general, is not true (see Example 1.2). However, if the random variables* $U$ *and* $V$ *are normal, then* $\mathrm{Cov}(U, V) = 0$ *guarantees that* $U$ *and* $V$ *are independent.*

**Example 1.2.** A probability mass function

$$f(u, v) = 1/4$$

for $(u, v) = (0, 1), (1, 0), (0, -1), (-1, 0)$. The marginal $f_1(1) = f_1(-1) = 1/4$, $f_1(0) = 1/2$ and $f_2(1) = f_2(-1) = 1/4$, $f_2(0) = 1/2$. Then we have $\mathrm{Cov}(U, V) = 0$ but $U$ and $V$ are not independent (i.e., $f(u, v) \neq f_1(u)f_2(v)$).                        △

Since $\epsilon_i \sim N(0, \sigma^2)$ and $\epsilon_i$ and $\epsilon_j$ $(i \neq j)$ are uncorrelated, $\epsilon_i$ are independent and identically distributed (iid) with $N(0, \sigma^2)$. The joint pdf of $(\epsilon_1, \epsilon_2, \ldots, \epsilon_n)$ is

$$f(\epsilon_1, \epsilon_2, \ldots, \epsilon_n) = \prod_{i=1}^{n} f(\epsilon_i).$$

Hence the likelihood function $L(\beta_0, \beta_1, \sigma^2)$ is given by

$$
\begin{aligned}
L(\beta_0, \beta_1, \sigma^2) &= \prod_{i=1}^{n} f(\epsilon_i) \\
&= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(Y_i - \beta_0 - \beta_1 X_i)^2\right] \\
&= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(Y_i - \beta_0 - \beta_1 X_i)^2\right].
\end{aligned}
$$

In general, it is easy to deal with the log likelihood $l(\cdot) = \ln L(\cdot)$ rather than $L(\cdot)$. The log likelihood is given by

$$l(\beta_0, \beta_1, \sigma^2) = \text{Constant} - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (Y_i - \beta_0 - \beta_1 X_i)^2.$$

Hence, if we assume $\epsilon_i \sim N(0, \sigma^2)$, then the MLE of the regression coefficients $\beta_0$ and $\beta_1$ are the same as the least-squares estimators. It should be made clear that if we can not make the normal assumption, the least squares estimator is not the MLE. Note that if we assume that $\epsilon_i$ comes from the double exponential distribution, the MLE of the $\beta_0$ and $\beta_1$ are the same as the $L_1$ estimators (*i.e..* minimizing $Q_1$).

### 1.3.3 Estimation of mean response

Given parameter estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ in the regression function:

$$\mu_Y = E(Y) = \beta_0 + \beta_1 X,$$

we can estimate $\mu_Y$ as $\hat{\mu}_Y = \hat{\beta}_0 + \hat{\beta}_1 X$. For convenience, we denote

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

instead of $\hat{\mu}_Y$ or $\widehat{E(Y)}$. We call a value of the response variable $(Y)$ a response value and $E(Y)$ the mean response. $\hat{Y}$ is a point estimator of the mean response when the level of the predictor variable is $X$. For the $i$th observed $X_i$, we call $\hat{Y}_i$ the fitted value for the $i$th case, where

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i, \qquad i = 1, \ldots, n.$$

### 1.3.4 Residuals

The $i$th residual is the difference between the observed value $Y_i$ and the corresponding fitted value $\hat{Y}_i$. This residual is denoted by $\hat{\epsilon}_i$ (or $e_i$) and is defined as

$$\hat{\epsilon}_i = e_i = Y_i - \hat{Y}_i.$$

Note that $\epsilon_i = Y_i - E(Y_i)$.

The estimated simple linear regression by the least-squares method has the following properties:

1. $\displaystyle\sum_{i=1}^{n} \hat{\epsilon}_i = 0.$

2. $\displaystyle\sum_{i=1}^{n} \hat{\epsilon}_i^2$ is a minimum.

3. $\displaystyle\sum_{i=1}^{n} Y_i = \sum_{i=1}^{n} \hat{Y}_i.$

4. $\displaystyle\sum_{i=1}^{n} X_i \hat{\epsilon}_i = 0.$

5. $\displaystyle\sum_{i=1}^{n} \hat{Y}_i \hat{\epsilon}_i = 0.$

6. $\displaystyle\sum_{i=1}^{n} Y_i \hat{\epsilon}_i = \sum_{i=1}^{n} \hat{\epsilon}_i^2.$

7. The regression line always passes through the point $(\bar{X}, \bar{Y})$.

## 1.4 Estimating $\sigma^2$

The variance $\sigma^2$ of the error term $\epsilon_i$ in regression needs to be estimated to obtain an indication of the variability of the probability distributions of $Y$.

Let us look at the simplest case, where we just have repeated observations $Y_1, \ldots, Y_n$. The variability in the $Y$'s is given by the usual sample variance:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (Y_i - \bar{Y})^2.$$

This is made unbiased by dividing by $n-1$ not $n$. In regression, we use

$$\text{MSE} = \frac{1}{n-2} \sum_{i=1}^{n} (\hat{\epsilon}_i - \bar{\hat{\epsilon}})^2 = \frac{1}{n-2} \sum_{i=1}^{n} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2.$$

Notice that $\bar{\hat{\epsilon}} = \frac{1}{n} \sum \hat{\epsilon}_i = 0$. Why $n-2$ is used instead of $n-1$ or $n$? To make MSE unbiased $i.e.$, $E(\text{MSE}) = \sigma^2$. The denominator term $n-2$ is degrees of freedom (the number of quantities that are free to vary).

We can also estimate $\sigma^2$ from the log likelihood $l(\beta_0, \beta_1, \sigma^2)$,

$$\frac{\partial l(\beta_0, \beta_1, \sigma^2)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (Y_i - \beta_0 - \beta_1 X_i)^2 = 0.$$

Solving for $\sigma^2$, we have

$$\hat{\sigma}_M^2 = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i)^2.$$

Thus, $\hat{\sigma}_M^2$ is biased.

## 1.5  Example I

A company produces refrigeration equipment and its replacement parts. In the past, one of the replacement parts has been produced periodically in different size lots. The company is interested in the optimum lot size. The data in the first column are different lot sizes and those in the second column are their corresponding work hours required to produce the lot.[1]

---

[1]Kutner, M. H. et al. Applied Linear Statistical Models. 5th edition. New York: McGraw-Hill, 2005.

Minitab

```
1   MTB > #===================================
2   MTB > # (1) Reading Data
3   MTB > #===================================
4   MTB > # ----------------------------------
5   MTB > # Write the data into C1 and C2 variables directly.
6   MTB > # ----------------------------------
7   MTB > READ C1 C2 .
8   DATA >     80   399
9   DATA >     30   121
10   ..............................
11   ..............................
12   ..............................
13  DATA >     80   342
14  DATA >     70   323
15  DATA > END .
16  25 rows read.
17
18  MTB > # Write C1 only
19  MTB > SET C1 .
20  DATA >     80  30   50 90 70 60 120 80 100 50 40 70 90
21  DATA >     20 110 100 30 50 90 110 30   90 40 80 70
22  DATA > END .
23
24  MTB > # Write C2 only
25  MTB > SET C2 .
26  DATA >   399 121 221 376 361 224 546 352 353 157 160 252 389
27  DATA >   113 435 420 212 268 377 421 273 468 244 342 323
28  DATA > END .
29  MTB >
30  MTB > # Double check if they are read well
31  MTB > PRINT C1 C2.
32
33  Data Display
34  Row   C1    C2
35    1   80   399
36    2   30   121
37   ......................
38   ......................
39   25   70   323
40
41  MTB > READ C1 C2;
42  SUBC> file "S:\data\CH01TA01.txt" .
43  Entering data from file: S:\DATA\CH01TA01.TXT
44  25 rows read.
45
46  MTB > # Double check if they are read well
47  MTB > PRINT C1-C2 .
48
49  Data Display
50  Row   C1    C2
51    1   80   399
52    2   30   121
53   ......................
54   ......................
55   25   70   323
56
57  MTB > #===================================
58  MTB > # (2) Scatter plot
59  MTB > #===================================
60
61  MTB > GSTD .
62  * NOTE * The character graph commands are obsolete.
63  * NOTE * Standard Graphics are now enabled, and Professional Graphics are
64         * disabled. Use the GPRO command when you want to re-enable
```

```
 65           * Professional Graphics.
 66
 67   MTB > PLOT C2 C1 .
 68   Scatterplot
 69
 70            -
 71    C2      -                                                              *
 72            -
 73            -                                              *
 74        450+                                                    *
 75            -                                             *    *
 76            -                                  *      3
 77            -                          *      *      *
 78            -                          *      *
 79        300+
 80            -          *          *
 81            -               *              *
 82            -          *          *    *
 83            -
 84        150+               *      *
 85            -       *    *
 86            -
 87            ----+---------+---------+---------+---------+---------+--C1
 88                20        40        60        80       100       120
 89
 90   MTB >GPRO .
 91   * NOTE * Professional Graphics are now enabled , and Standard Graphics are
 92            * disabled. Use the GSTD command when you want to re-enable Standard
 93            * Graphics.
 94
 95   MTB > PLOT C2*C1 .   # NB: not PLOT C2 C1. (* is needed).
 96
 97   MTB > #===================================
 98   MTB > # (3) Parameter estimation
 99   MTB > #===================================
100   MTB > # Minitab provides estimation and ANOVA as well.
101   MTB > REGR C2 1 C1;
102   SUBC> FITS C3 .
103
104   Regression Analysis: C2 versus C1
105   The regression equation is
106   C2 = 62.4 + 3.57 C1
107
108   Predictor    Coef  SE Coef      T      P
109   Constant    62.37    26.18   2.38  0.026
110   C1         3.5702   0.3470  10.29  0.000
111
112   S = 48.8233   R-Sq = 82.2%   R-Sq(adj) = 81.4%
113
114
115   Analysis of Variance
116   Source          DF       SS       MS        F       P
117   Regression       1   252378   252378   105.88   0.000
118   Residual Error  23    54825     2384
119   Total           24   307203
120
121   Unusual Observations
122   Obs  C1      C2       Fit  SE Fit   Residual   St Resid
123    21  30  273.00   169.47   16.97     103.53       2.26R
124
125   R denotes an observation with a large standardized residual.
126
127   MTB > #----------------------------------------
128   MTB > # Scatter plot with the fitted line
129   MTB > #----------------------------------------
130   MTB > GPRO .
131   * NOTE * Professional Graphics are now enabled , and Standard Graphics are
132            * disabled. Use the GSTD command when you want to re-enable Standard
```

```
133       * Graphics.
134
135 MTB > PLOT C2*C1 ;
136 SUBC> Symbol ;
137 SUBC> Regress .
138
139 MTB > #==================================
140 MTB > # (4) Some Math
141 MTB > #==================================
142 MTB > LET C4 = C2 - C3  # No period(.) after LET command.
143 MTB > LET C5 = C4**2
144 MTB > PRINT C1-C5 .
145
146 Data Display
147 Row   C1   C2        C3        C4        C5
148   1   80  399   347.982    51.018    2602.8
149   2   30  121   169.472   -48.472    2349.5
150 ...........................................
151 ...........................................
152
153  25   70  323   312.280    10.720     114.9
154
155 MTB > LET C11 = SUM(C4)
156 MTB > LET C12 = SUM(C2)
157 MTB > LET C13 = SUM(C3)
158 MTB > LET C14 = C1*C4
159 MTB > LET C15 = C3*C4
160 MTB > LET C16 = C2*C4
161 MTB > LET C14 = SUM(C14)
162 MTB > LET C15 = SUM(C15)
163 MTB > LET C16 = SUM(C16)
164 MTB > PRINT C11-C16 .
165
166 Data Display
167 Row          C11   C12   C13         C14         C15       C16
168   1  -0.0000000  7807  7807  -0.0000000  -0.0000000  54825.5
169
170 MTB > ##- MSE
171 MTB > LET C21 = SUM(C5) / (25-2)
172 MTB > PRINT C21 .
173
174 Data Display
175 C21
176     2383.72
```

R

```
1  > #==================================
2  > # (1) Reading Data
3  > #==================================
4  >
5  > # ---------------------------------
6  > # Write the data into C1 and C2 variables directly.
7  > # ---------------------------------
8  > x1 = c(80, 30, 50, 90, 70, 60,120, 80,100, 50,
9  +        40, 70, 90, 20,110,100, 30, 50, 90,110,
10 +        30, 90, 40, 80, 70)
11 > y1 = c(399,121,221,376,361,224,546,352,353,157,
12 +        160,252,389,113,435,420,212,268,377,421,
13 +        273,468,244,342,323)
14 >
15 > # ---------------------------------
```

```
16  > # From Hard disc
17  > # ----------------------------------
18  > # Note: use ANSI ascii text file. UTF-8 text is not supported.
19  >  mydata = read.table("S:/data/CH01TA01.txt")
20  >
21  > # The above is the same as:
22  > # setwd("S:/data")
23  > # mydata = read.table("CH01TA01.txt")
24  >
25  > # Double-check if they are read well
26  > x2 = mydata[,1]
27  > y2 = mydata[,2]
28  > cbind(x2,y2)
29        x2  y2
30   [1,]  80 399
31   [2,]  30 121
32   [3,]  50 221
33   [4,]  90 376
34   [5,]  70 361
35   [6,]  60 224
36   [7,] 120 546
37   [8,]  80 352
38   [9,] 100 353
39  [10,]  50 157
40  [11,]  40 160
41  [12,]  70 252
42  [13,]  90 389
43  [14,]  20 113
44  [15,] 110 435
45  [16,] 100 420
46  [17,]  30 212
47  [18,]  50 268
48  [19,]  90 377
49  [20,] 110 421
50  [21,]  30 273
```

```
51  [22,]   90 468
52  [23,]   40 244
53  [24,]   80 342
54  [25,]   70 323
55  >
56  > # ---------------------------------
57  > # From URL
58  > # ---------------------------------
59  > # See the URL: https://github.com/AppliedStat/LM
60  > # If your computer is connected to Internet, then the following should work:
61  > url = "https://raw.githubusercontent.com/AppliedStat/LM/master/CH01TA01.txt"
62  > mydata = read.table(url)
63  >
64  > # Check below
65  > # is.matrix(mydata)
66  > #   is.list(mydata)
67  > # as.matrix(mydata)
68  >
69  > # Double-check if they are read well
70  > x3 = mydata[,1]
71  > y3 = mydata[,2]
72  > cbind(x3,y3)
73         x3  y3
74   [1,]   80 399
75   [2,]   30 121
76   [3,]   50 221
77   [4,]   90 376
78   [5,]   70 361
79   [6,]   60 224
80   [7,]  120 546
81   [8,]   80 352
82   [9,]  100 353
83  [10,]   50 157
84  [11,]   40 160
85  [12,]   70 252
86  [13,]   90 389
87  [14,]   20 113
88  [15,]  110 435
89  [16,]  100 420
90  [17,]   30 212
91  [18,]   50 268
92  [19,]   90 377
93  [20,]  110 421
94  [21,]   30 273
95  [22,]   90 468
96  [23,]   40 244
97  [24,]   80 342
98  [25,]   70 323
99  >
100 > # ---------------------------------
101 > # For convenience,
102 > # ---------------------------------
103 > x = x1
104 > y = y1
105 >
106 > #=================================
107 > # (2) Scatter plot
108 > #=================================
109 > # ?Devices
110 > # postscript( "ex1.ps", width=4, height=4 )
111 > # pdf(file="ex1.pdf", width=4, height=4 )
112 > plot(x,y)
113 > dev.off()
114 null device
115           1
116 >
117 > # plot with text is optional in R
118 > # It may need to install txtplot package
```

```
119  > # install.packages("txtplot")
120  > library("txtplot")
121  > txtplot(x,y)
122       +-+---------+---------+---------+---------+--------+---+
123       |                                                *   |
124   500 +                                                    +
125       |                               *                    |
126       |                                    *     *         |
127   400 +                        *     *                     |
128       |                   *     *          *               |
129   300 +                        *                           +
130       |       *           *                                |
131       |          *     *     *     *                       |
132   200 +       *                                            +
133       |          *     *                                   |
134       | *     *                                            |
135   100 +-+---------+---------+---------+---------+--------+---+
136         20        40        60        80        100      120
137  >
138  > #===================================
139  > # (3) Parameter estimation
140  > #===================================
141  > lm( y ~ x)
142
143  Call:
144  lm(formula = y ~ x)
145
146  Coefficients:
147  (Intercept)              x
148        62.37           3.57
149
150  > output = lm (y ~ x)
151  > summary(output)
152
153  Call:
154  lm(formula = y ~ x)
155
156  Residuals:
157      Min      1Q  Median      3Q      Max
158  -83.876 -34.088  -5.982  38.826 103.528
159
160  Coefficients:
161              Estimate Std. Error t value Pr(>|t|)
162  (Intercept)   62.366     26.177   2.382   0.0259 *
163  x              3.570      0.347  10.290 4.45e-10 ***
164  ---
165  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1        1
166
167  Residual standard error: 48.82 on 23 degrees of freedom
168  Multiple R-squared:  0.8215,   Adjusted R-squared:  0.8138
169  F-statistic: 105.9 on 1 and 23 DF,  p-value: 4.449e-10
170
171  > anova(output)
172  Analysis of Variance Table
173
174  Response: y
175            Df Sum Sq Mean Sq F value    Pr(>F)
176  x          1 252378  252378  105.88 4.449e-10 ***
177  Residuals 23  54825    2384
178  ---
179  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1        1
180  >
181  > #===================================
182  > # (4) Some Math
183  > #===================================
184  > # To compare with Minitab codes
185  > # NB: Cx (x=1,2,..) are variables in Minitab
186  > C2 = y   # copy of y to C2
```

```
187  > C1 = x
188  >
189  > C3 = fitted(output)
190  >
191  > C4 = C2 - C3
192  > C5 = C4^2
193  >
194  > cbind(C1, C2, C3, C4, C5)
195       C1  C2       C3             C4             C5
196  1    80 399 347.9820   51.0179798 2.602834e+03
197  2    30 121 169.4719  -48.4719192 2.349527e+03
198  3    50 221 240.8760  -19.8759596 3.950538e+02
199  4    90 376 383.6840   -7.6840404 5.904448e+01
200  5    70 361 312.2800   48.7200000 2.373638e+03
201  6    60 224 276.5780  -52.5779798 2.764444e+03
202  7   120 546 490.7901   55.2098990 3.048133e+03
203  8    80 352 347.9820    4.0179798 1.614416e+01
204  9   100 353 419.3861  -66.3860606 4.407109e+03
205  10   50 157 240.8760  -83.8759596 7.035177e+03
206  11   40 160 205.1739  -45.1739394 2.040685e+03
207  12   70 252 312.2800  -60.2800000 3.633678e+03
208  13   90 389 383.6840    5.3159596 2.825943e+01
209  14   20 113 133.7699  -20.7698990 4.313887e+02
210  15  110 435 455.0881  -20.0880808 4.035310e+02
211  16  100 420 419.3861    0.6139394 3.769216e-01
212  17   30 212 169.4719   42.5280808 1.808638e+03
213  18   50 268 240.8760   27.1240404 7.357136e+02
214  19   90 377 383.6840   -6.6840404 4.467640e+01
215  20  110 421 455.0881  -34.0880808 1.161997e+03
216  21   30 273 169.4719  103.5280808 1.071806e+04
217  22   90 468 383.6840   84.3159596 7.109181e+03
218  23   40 244 205.1739   38.8260606 1.507463e+03
219  24   80 342 347.9820   -5.9820202 3.578457e+01
220  25   70 323 312.2800   10.7200000 1.149184e+02
221  >
222  > C11 = sum(C4)
223  > C12 = sum(C2)
224  > C13 = sum(C3)
225  > C14 = C1*C4
226  > C15 = C3*C4
227  > C16 = C2*C4
228  > C14 = sum(C14)
229  > C15 = sum(C15)
230  > C16 = sum(C16)
231  > C17 = sum(C5)
232  > cbind(C11, C12, C13, C14, C15, C16, C17)
233                C11  C12  C13          C14          C15      C16      C17
234  [1,] 2.273737e-13 7807 7807 1.921308e-11 7.958079e-11 54825.46 54825.46
235  >
236  > # MSE
237  > sum(C5) / (25-2)
238  [1] 2383.716
239  > sum ( (y-fitted(output))^2 ) / (25-2)
240  [1] 2383.716
```

Python

```
1  ~/MyFiles/teaching/IE-34243> python3
2  Python 3.5.2 (default, Nov 12 2018, 13:43:14)
3  [GCC 5.4.0 20160609] on linux
4  Type "help", "copyright", "credits" or "license" for more information.
5  >>>
```

```
 6  >>> #================================
 7  ... # (1) Reading Data
 8  ... #================================
 9  ...
10  >>> # ----------------------------------
11  ... # Write the data into x1 and y1 variables directly.
12  ... # ----------------------------------
13  ... x1 = [ 80, 30, 50, 90, 70, 60,120, 80,100, 50,
14  ...          40, 70, 90, 20,110,100, 30, 50, 90,110,
15  ...          30, 90, 40, 80, 70 ]
16  >>> y1 = [399,121,221,376,361,224,546,352,353,157,
17  ...          160,252,389,113,435,420,212,268,377,421,
18  ...          273,468,244,342,323 ]
19  >>>
20  >>> # ----------------------------------
21  ... # Read From Hard disc and write into x2 and y2
22  ... # ----------------------------------
23  ... f = open("S:/data/CH01TA01.txt", "r")
24  >>> file2 = f.read().splitlines()
25  >>> print(file2)
26  ['   80   399', '   30   121', '   50   221', '   90   376', '   70   361', '   60   224',
    '  120   546', '   80   352', '  100   353', '   50   157', '   40   160', '   70
    252', '   90   389', '   20   113', '  110   435', '  100   420', '   30   212', '
    50   268', '   90   377', '  110   421', '   30   273', '   90   468', '   40   244',
    '   80   342', '   70   323']
27  >>> f.close()
28
29  >>> # Write the data in the file into x2 and y2 variables
30  ... x2 = []; y2 = []  # Make space for the data
31  >>>
32  >>> for line in file2[0:]:
33  ...      p = line.split()
34  ...      x2 = x2 + [float(p[0])]
35  ...      y2 = y2 + [float(p[1])]
36  ...
37  >>> # Double-check if they are read well
38  ... for i in range(len(x2)):
39  ...      print(x2[i],y2[i])
40  80.0 399.0
41  .....................
42  .....................
43  70.0 323.0
44  >>>
45  >>> # ----------------------------------
46  ... # From URL
47  ... # ----------------------------------
48  ... # See the URL: https://github.com/AppliedStat/LM
49  ... # If your computer is connected to Internet, then the following should work:
50  ... from urllib.request import urlopen
51  >>> link = "https://raw.githubusercontent.com/AppliedStat/LM/master/CH01TA01.txt"
52  >>> url  = urlopen(link)
53  >>> file3= url.readlines()
54  >>> print(file3)
55  [b'   80   399\n', b'   30   121\n', b'   50   221\n', b'   90   376\n', b'   70
    361\n', b'   60   224\n', b'  120   546\n', b'   80   352\n', b'  100   353\n', b'
    50   157\n', b'   40   160\n', b'   70   252\n', b'   90   389\n', b'   20   113\n',
    b'  110   435\n', b'  100   420\n', b'   30   212\n', b'   50   268\n', b'   90
    377\n', b'  110   421\n', b'   30   273\n', b'   90   468\n', b'   40   244\n', b'
    80   342\n', b'   70   323\n']
56  >>> url.close()
57  >>>
58  >>> # Write the data in the file into x3 and y3 variables
59  ... x3 = []   # Make space for the data
60  >>> y3 = []   # Make space for the data
61  >>> for line in file3[0:]:
62  ...      p = line.split()
63  ...      x3 = x3 + [float(p[0])]
64  ...      y3 = y3 + [float(p[1])]
```

```
65  ... #   y3 += [float(p[1])]       # Same as the above
66  ... #   y3.append(float(p[1]))
67  ...
68  >>> # Double-check if they are read well
69  ... print(x3)
70  [80.0, 30.0, 50.0, 90.0, 70.0, 60.0, 120.0, 80.0, 100.0, 50.0, 40.0, 70.0, 90.0,
        20.0, 110.0, 100.0, 30.0, 50.0, 90.0, 110.0, 30.0, 90.0, 40.0, 80.0, 70.0]
71  >>> print(y3)
72  [399.0, 121.0, 221.0, 376.0, 361.0, 224.0, 546.0, 352.0, 353.0, 157.0, 160.0, 252.0,
        389.0, 113.0, 435.0, 420.0, 212.0, 268.0, 377.0, 421.0, 273.0, 468.0, 244.0,
        342.0, 323.0]
73  >>> for i in range(len(x3)):
74  ...     print(x3[i],y3[i])
75
76  80.0 399.0
77  ....................
78  ....................
79  70.0 323.0
80
81  >>> # --------------------------------
82  ... # For convenience,
83  ... # --------------------------------
84  ... x = x1
85  >>> y = y1
86
87  >>> #================================
88  ... # (2) Scatter plot
89  ... #================================
90  ... import matplotlib.pyplot as plot # https://matplotlib.org/
91  >>> # Windows10: C:\> pip install matplotlib
92  ... # The below may be needed for upgrading pip.
93  ... # Windows10: C:\> python -m pip install  --upgrade pip  --user
94  ...
95  >>> plot.figure( figsize=(10,10) )
96  <matplotlib.figure.Figure object at 0x7f0de8411e48>
97  >>> plot.scatter(x, y, c="black" )
98  <matplotlib.collections.PathCollection object at 0x7f0dc958a6d8>
99  >>> # plot.savefig("ch01-example1.eps")
100 ... plot.show()
101
102 >>> #================================
103 ... # (3) Parameter estimation
104 ... #================================
105 ... # https://www.statsmodels.org
106 ... # Windows10: C:\> pip install -U statsmodels --user
107 ... from statsmodels.formula.api import ols
108 >>> from statsmodels.stats.anova import anova_lm
109 >>> import pandas  # https://pandas.pydata.org
110 >>>
111 >>> # This data structure is needed for ols and anova_lm
112 ... data = pandas.DataFrame({"x": x, "y": y})
113 >>>
114 >>> model = ols("y~x", data) # NB: ols("y~0+x", data)
115 >>> OUT = model.fit()
116 >>> OUT.summary()
117 <class 'statsmodels.iolib.summary.Summary'>
118 """
119                       OLS Regression Results
120 ==============================================================================
121 Dep. Variable:                     y   R-squared:                       0.822
122 Model:                           OLS   Adj. R-squared:                  0.814
123 Method:                Least Squares   F-statistic:                     105.9
124 Date:               Fri, 28 Jun 2019   Prob (F-statistic):           4.45e-10
125 Time:                       13:28:05   Log-Likelihood:                -131.64
126 No. Observations:                 25   AIC:                             267.3
127 Df Residuals:                     23   BIC:                             269.7
128 Df Model:                          1
129 Covariance Type:           nonrobust
```
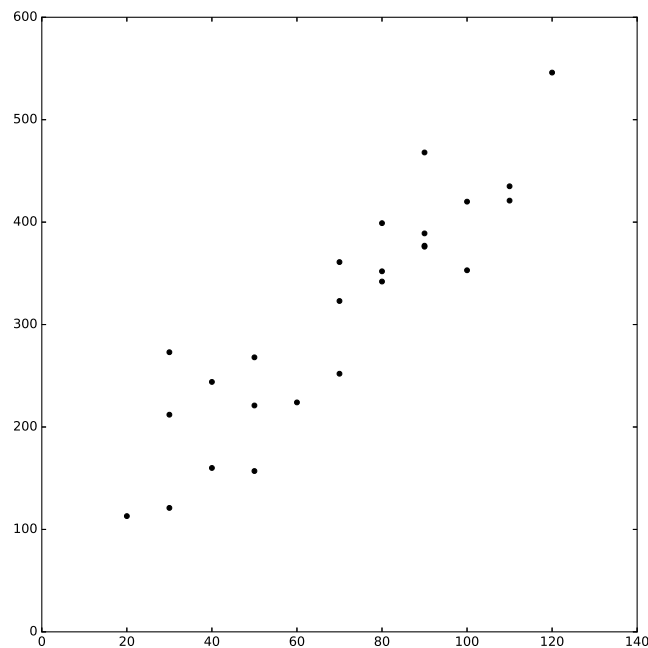
```
130   ==========================================================================
131                      coef      std err             t      P>|t|        [0.025      0.975]
132   --------------------------------------------------------------------------
133   Intercept       62.3659      26.177          2.382      0.026         8.214     116.518
134   x                3.5702       0.347         10.290      0.000         2.852       4.288
135   ==========================================================================
136   Omnibus:                              0.608   Durbin-Watson:                       1.432
137   Prob(Omnibus):                        0.738   Jarque-Bera (JB):                    0.684
138   Skew:                                 0.298   Prob(JB):                            0.710
139   Kurtosis:                             2.450   Cond. No.                            202.
140   ==========================================================================
141
142   Warnings:
143   [1] Standard Errors assume that the covariance matrix of the errors is correctly
          specified.
144   """
145
146   >>> # NB: Compare the following
147   ... # ols("y~  x", data).fit().summary()
148   ... # ols("y~0+x", data).fit().summary()
149   ... anova_lm(OUT)
150              df         sum_sq         mean_sq            F          PR(>F)
151   x         1.0   252377.580808   252377.580808   105.875709   4.448828e-10
152   Residual  23.0   54825.459192     2383.715617          NaN             NaN
153
154
155   >>> #================================
156   ... # (4) Some Math
157   ... #================================
158   ... # To compare with Minitab codes
159   ... # NB: Cx (x=1,2,..) are variables in Minitab
160   ... C2 = y
161   >>> C1 = x
162   >>>
163   >>> C3 = OUT.fittedvalues
164   >>>
165   >>> C4 = C2 - C3
166   >>> C5 = C4**2
167   >>>
168   >>> for i in range(len(C1)):   # ugly
169   ...     print(C1[i],C2[i],C3[i],C4[i],C5[i])
170   ...
171   80.0 399.0 347.98202020202024 51.01797979797976 2602.834262667071
172    ................................................
173    ................................................
174   70.0 323.0 312.28000000000003 10.71999999999997 114.91839999999937
175
176   >>> for i in range(len(C1)):   # better
177   ...     print("%5d %5d %10.5f %10.5f %10.5f" %(C1[i],C2[i],C3[i],C4[i],C5[i]))
178   ...
179      80    399    347.98202     51.01798 2602.83426
180      30    121    169.47192    -48.47192 2349.52695
181    ................................................
182    ................................................
183      70    323    312.28000     10.72000  114.91840
184   >>> C11 = sum(C4)
185   >>> C12 = sum(C2)
186   >>> C13 = sum(C3)
187   >>>
188   >>> import numpy as np  # We need np.multiply from numpy (https://www.numpy.org/)
189   >>> # Windows10: C:\> pip install -U numpy
190   ... C14 = np.multiply(C1,C4)
191   >>> C15 = np.multiply(C3,C4)
192   >>> C16 = np.multiply(C2,C4)
193   >>>
194   >>> C14 = sum(C14)
195   >>> C15 = sum(C15)
196   >>> C16 = sum(C16)
```

```
197  >>> C17 = sum(C5)
198  >>>
199  >>> print(C11, C12, C13, C14, C15, C16, C17)     # ugly
200  -1.1937117960769683e-12 7807.0 7807.000000000002 -8.151346264639869e-11
         -3.6288838600739837e-10 54825.459191918846 54825.4591919192
201
202  >>> print((" %11.7G"*7) %(C11, C12, C13, C14, C15, C16, C17))   # better
203   -1.193712E-12   7807   7807   -8.151346E-11 -3.628884E-10   54825.46   54825.46
```



# 1.6   Example II

Over the past decades, the data useful to assess the impact of climate change have been collected from satellite. Especially since 1979, satellites have measured the area or extended area of sea ice in the Arctic Ocean on a daily basis. It should be noted that the extended areas (the fifth column in the data set) include the areas near the pole not imaged by the sensor. It is assumed to be entirely ice covered with at least 15% concentration. On the other hand, the areas (the sixth column) *exclude* the area not imaged by the sensor.

The averages of areas (or extended areas) of sea ice on a daily basis are calculated for each month. We consider the average *extended* area of sea ice in September. It should be noted that September is the month when the ice stops melting each summer and reaches its minimum extent. Witt[2] analyzed a time series of September Arctic sea ice extent from 1979 until 2015. The origianl data[3] can be obtained at:

www.amstat.org/publications/jse/v21n1/witt.pdf

http://dx.doi.org/10.7265/N5736NV7

Minitab

```
1   MTB > #===================================
2   MTB > # (1) Reading Data
3   MTB > #===================================
4   MTB > SET C1 .
5   DATA> 1979 1980 1981 1982 1983 1984 1985 1986 1987 1988
6   DATA> 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998
7   DATA> 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008
8   DATA> 2009 2010 2011 2012 2013 2014 2015
9   DATA> END .
10  MTB > SET C2 .
11  DATA> 7.22 7.86 7.25 7.45 7.54 7.11 6.93 7.55 7.51 7.53
12  DATA> 7.08 6.27 6.59 7.59 6.54 7.24 6.18 7.91 6.78 6.62
13  DATA> 6.29 6.36 6.78 5.98 6.18 6.08 5.59 5.95 4.32 4.73
14  DATA> 5.39 4.93 4.63 3.63 5.35 5.29 4.68
15  DATA> END .
16  MTB > PRINT C1 C2.
17
18  Data Display
19  Row    C1     C2
20    1  1979  7.22
21  ...................
22
23   37  2015  4.68
24
25  MTB > #===================================
26  MTB > # (2) Scatter plot
27  MTB > #===================================
28  MTB > GSTD .
29  * NOTE * The character graph commands are obsolete.
30  * NOTE * Standard Graphics are now enabled, and Professional Graphics are
31         * disabled. Use the GPRO command when you want to re-enable
32         * Professional Graphics.
33
```

[2]Witt, Gary Using Data from Climate Science to Teach Introductory Statistics. Journal of Statistics Education, 21 2013 ⟨URL: www.amstat.org/publications/jse/v21n1/witt.pdf⟩.

[3]Fetterer, F. et al. Sea Ice Index (updated daily), Version 2. Boulder, Colorado USA: NSIDC: National Snow and Ice Data Center, 2016 ⟨URL: http://dx.doi.org/10.7265/N5736NV7⟩.

```
34  MTB > PLOT C2 C1 .
35  Scatterplot
36
37            -
38            -         *                              *
39      7.5+            * *    * **       *
40            -      *   *    *       *         *
41   C2       -              *                *      *
42            -                  *   *       *
43            -              *         *      **     *
44      6.0+                               *   *   *
45            -                                *
46            -                                   *       **
47            -
48            -                                  *  *        *
49      4.5+                                          *
50            -                              *
51            -
52            -                                   *
53            -
54            ------+---------+---------+---------+---------+---------+C1
55             1981.0     1988.0     1995.0     2002.0     2009.0     2016.0
56
57  MTB > GPRO .
58  * NOTE * Professional Graphics are now enabled, and Standard Graphics are
59         * disabled. Use the GSTD command when you want to re-enable Standard
60         * Graphics.
61
62  MTB > PLOT C2*C1 .
63
64  MTB > #================================
65  MTB > # (3) Parameter estimation
66  MTB > #================================
67  MTB > REGR C2 1 C1;
68  SUBC> FITS C3 .
69
70  Regression Analysis: C2 versus C1
71  The regression equation is
72  C2 = 181 - 0.0873 C1
73
74  Predictor        Coef    SE Coef        T       P
75  Constant       180.73      17.40    10.39   0.000
76  C1           -0.087321   0.008711   -10.02   0.000
77
78  S = 0.565773   R-Sq = 74.2%   R-Sq(adj) = 73.4%
79
80  Analysis of Variance
81  Source            DF        SS        MS        F       P
82  Regression         1    32.162    32.162   100.48   0.000
83  Residual Error    35    11.203     0.320
84  Total             36    43.366
85
86  Unusual Observations
87  Obs    C1      C2      Fit   SE Fit   Residual   St Resid
88   18   1996   7.9100   6.4362   0.0934     1.4738      2.64R
89   29   2007   4.3200   5.4757   0.1274    -1.1557     -2.10R
90   34   2012   3.6300   5.0391   0.1604    -1.4091     -2.60R
91
92  R denotes an observation with a large standardized residual.
93
94  MTB > #----------------------------------------
95  MTB > # Scatter plot with the fitted line
96  MTB > #----------------------------------------
97  MTB > GPRO .
98  * NOTE * Professional Graphics are now enabled, and Standard Graphics are
99         * disabled. Use the GSTD command when you want to re-enable Standard
100        * Graphics.
101
```

```
102  MTB > PLOT C2*C1 ;
103  SUBC> Symbol ;
104  SUBC> Regress .
105
106  MTB > #====================================
107  MTB > # (4) To compare with Figure 2 of Witt (2013).
108  MTB > #====================================
109  MTB > DELETE 35:37 C1.
110  MTB > DELETE 35:37 C2.
111  MTB > REGR C2 1 C1;
112  SUBC> FITS C3 .
113
114  Regression Analysis: C2 versus C1
115  The regression equation is
116  C2 = 189 - 0.0913 C1
117
118  Predictor       Coef   SE Coef       T       P
119  Constant      188.60     20.25    9.31   0.000
120  C1           -0.09128   0.01015   -8.99   0.000
121
122  S = 0.580615   R-Sq = 71.7%   R-Sq(adj) = 70.8%
123
124  Analysis of Variance
125  Source          DF      SS       MS       F       P
126  Regression       1   27.265   27.265   80.88   0.000
127  Residual Error  32   10.788    0.337
128  Total           33   38.053
129
130  Unusual Observations
131  Obs    C1      C2      Fit   SE Fit   Residual  St Resid
132   18  1996   7.9100  6.4129  0.0997    1.4971      2.62R
133   34  2012   3.6300  4.9525  0.1948   -1.3225     -2.42R
134
135  R denotes an observation with a large standardized residual.
136
137  MTB > #----------------------------------------
138  MTB > # Scatter plot with the fitted line
139  MTB > #----------------------------------------
140  MTB > GPRO .
141  * NOTE * Professional Graphics are now enabled, and Standard Graphics are
142         * disabled. Use the GSTD command when you want to re-enable Standard
143         * Graphics.
144
145  MTB > PLOT C2*C1 ;
146  SUBC> Symbol ;
147  SUBC> Regress .
```
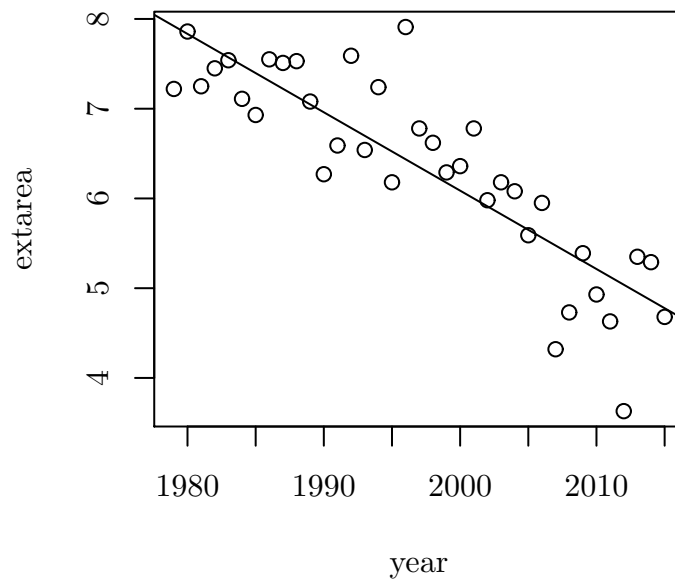
R

```
1   > #===================================
2   > # (1) Reading Data
3   > #===================================
4   > year = c(1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988,
5   +          1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998,
6   +          1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008,
7   +          2009, 2010, 2011, 2012, 2013, 2014, 2015 )
8   > area = c(7.22, 7.86, 7.25, 7.45, 7.54, 7.11, 6.93, 7.55, 7.51, 7.53,
9   +          7.08, 6.27, 6.59, 7.59, 6.54, 7.24, 6.18, 7.91, 6.78, 6.62,
10  +          6.29, 6.36, 6.78, 5.98, 6.18, 6.08, 5.59, 5.95, 4.32, 4.73,
11  +          5.39, 4.93, 4.63, 3.63, 5.35, 5.29, 4.68)
12  >
13  >
14  > #===================================
```

```
15  > # (2) Scatter plot
16  > #==================================
17  > plot( year, area )
18  >
19  >
20  > #==================================
21  > # (3) Parameter estimation
22  > #==================================
23  > OUT = lm (area ~ year)
24  > summary(OUT)
25
26  Call:
27  lm(formula = area ~ year)
28
29  Residuals:
30       Min        1Q    Median        3Q       Max
31  -1.40910  -0.34356   0.03251   0.35840   1.47376
32
33  Coefficients:
34                Estimate Std. Error t value Pr(>|t|)
35  (Intercept) 180.728966   17.396966    10.39 3.10e-12 ***
36  year         -0.087321    0.008711   -10.02 7.97e-12 ***
37  ---
38  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1         1
39
40  Residual standard error: 0.5658 on 35 degrees of freedom
41  Multiple R-squared:  0.7417,   Adjusted R-squared:  0.7343
42  F-statistic: 100.5 on 1 and 35 DF,  p-value: 7.968e-12
43
44  > # NB: Compare the following
45  > # lm (area ~ year)
46  > # lm (area ~ 0 + year)
47  >
48  > anova(OUT)
49  Analysis of Variance Table
```

```
50
51  Response: area
52            Df Sum Sq Mean Sq F value    Pr(>F)
53  year       1 32.162  32.162  100.48 7.968e-12 ***
54  Residuals 35 11.203   0.320
55  ---
56  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1        1
57  >
58  > #----------------------------------------
59  > # Scatter plot with the fitted line
60  > #----------------------------------------
61  > plot( year, area )
62  > abline(OUT)
63  >
64  >
65  > #======================================
66  > # (4) To compare with Figure 2 of Witt (2013).
67  > #======================================
68  > yr = year[1:34]
69  > ar = area[1:34]
70  >
71  > OUT2 = lm (ar ~ yr)
72  > summary(OUT2)
73
74  Call:
75  lm(formula = ar ~ yr)
76
77  Residuals:
78       Min       1Q   Median       3Q      Max
79  -1.32245 -0.37971  0.01339  0.38897  1.49711
80
81  Coefficients:
82              Estimate Std. Error t value Pr(>|t|)
83  (Intercept) 188.60240   20.25374   9.312 1.26e-10 ***
84  yr           -0.09128    0.01015  -8.993 2.85e-10 ***
85  ---
86  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1        1
87
88  Residual standard error: 0.5806 on 32 degrees of freedom
89  Multiple R-squared:  0.7165,   Adjusted R-squared:  0.7076
90  F-statistic: 80.88 on 1 and 32 DF,  p-value: 2.845e-10
91
92  > anova(OUT2)
93  Analysis of Variance Table
94
95  Response: ar
96            Df Sum Sq Mean Sq F value    Pr(>F)
97  yr         1 27.265 27.2650  80.878 2.845e-10 ***
98  Residuals 32 10.788  0.3371
99  ---
100 Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1        1
101 >
102 > #----------------------------------------
103 > # Scatter plot with the fitted line
104 > #----------------------------------------
105 > plot(yr, ar, ylim=c(0,9) )
106 > abline( OUT2 )
```

Python

```
1  >>> #======================================
2  ... # (1) Reading Data
```

```
3   ... #=================================
4   ... x = [1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988,
5   ...       1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998,
6   ...       1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008,
7   ...       2009, 2010, 2011, 2012, 2013, 2014, 2015]
8   >>> y = [7.22, 7.86, 7.25, 7.45, 7.54, 7.11, 6.93, 7.55, 7.51, 7.53,
9   ...       7.08, 6.27, 6.59, 7.59, 6.54, 7.24, 6.18, 7.91, 6.78, 6.62,
10  ...       6.29, 6.36, 6.78, 5.98, 6.18, 6.08, 5.59, 5.95, 4.32, 4.73,
11  ...       5.39, 4.93, 4.63, 3.63, 5.35, 5.29, 4.68]
12  >>>
13  >>> #=================================
14  ... # (2) Scatter plot
15  ... #=================================
16  ... import matplotlib.pyplot as plot # https://matplotlib.org/
17  >>> import numpy as np           # https://www.numpy.org
18  >>> year = np.array(x).reshape((-1,1))
19  >>> area = np.array(y).reshape((-1,1))
20  >>>
21  >>> plot.figure( figsize=(10,10) )
22  <matplotlib.figure.Figure object at 0x7fec71fa7860>
23  >>> plot.scatter(year, area, c="black" )
24  <matplotlib.collections.PathCollection object at 0x7fec51acf518>
25  >>> plot.show()
26  >>>
27  >>> #=================================
28  ... # (3) Parameter estimation
29  ... #=================================
30  ... # https://www.statsmodels.org
31  ... # Windows10: C:\> pip install -U statsmodels --user
32  ... from statsmodels.formula.api import ols
33  >>> from statsmodels.stats.anova import anova_lm
34  >>> import pandas   # https://pandas.pydata.org
35  >>>
36  >>> # This data structure is needed for  ols and anova_lm
37  ... data = pandas.DataFrame({"year": x, "area": y})
38  >>> model = ols("area~year", data=data) # NB: ols("y~0+x", data)
39  >>> OUT = model.fit()
40  >>> print(OUT.summary())
41                          OLS Regression Results
42  ==============================================================================
43  Dep. Variable:                   area   R-squared:                       0.742
44  Model:                            OLS   Adj. R-squared:                  0.734
45  Method:                 Least Squares   F-statistic:                     100.5
46  Date:                Wed, 26 Jun 2019   Prob (F-statistic):           7.97e-12
47  Time:                        19:27:27   Log-Likelihood:                -30.399
48  No. Observations:                  37   AIC:                             64.80
49  Df Residuals:                      35   BIC:                             68.02
50  Df Model:                           1
51  Covariance Type:            nonrobust
52  ==============================================================================
53                  coef    std err          t      P>|t|      [0.025      0.975]
54  ------------------------------------------------------------------------------
55  Intercept    180.7290     17.397     10.389      0.000     145.411     216.047
56  year          -0.0873      0.009    -10.024      0.000      -0.105      -0.070
57  ==============================================================================
58  Omnibus:                        1.794   Durbin-Watson:                   1.739
59  Prob(Omnibus):                  0.408   Jarque-Bera (JB):                0.811
60  Skew:                          -0.133   Prob(JB):                        0.667
61  Kurtosis:                       3.675   Cond. No.                     3.74e+05
62  ==============================================================================
63
64  Warnings:
65  [1] Standard Errors assume that the covariance matrix of the errors is correctly
        specified.
66  [2] The condition number is large, 3.74e+05. This might indicate that there are
67  strong multicollinearity or other numerical problems.
68  >>>
69  >>> # NB: Compare the following
```
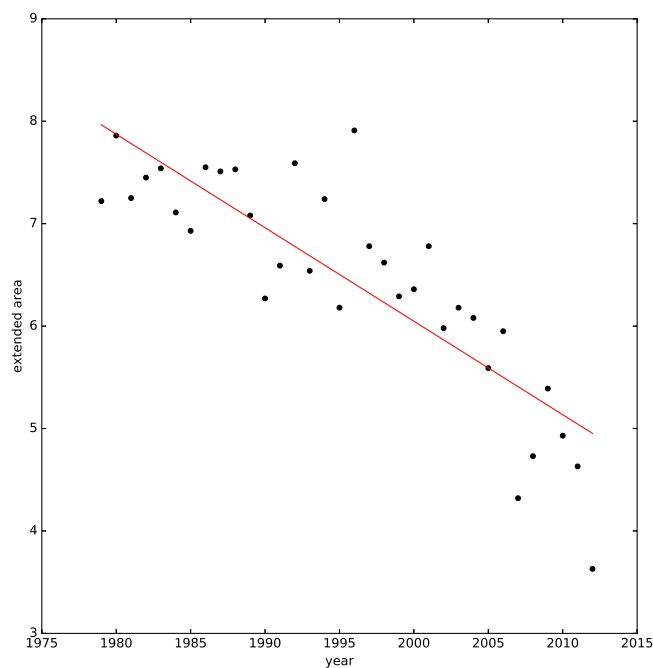
```
70  ... # ols("area~year", data).fit().summary()
71  ... # ols("area~0+year", data).fit().summary()
72  ...
73  >>> print(anova_lm(OUT))
74             df      sum_sq     mean_sq           F          PR(>F)
75  year       1.0   32.162073   32.162073   100.475221   7.967937e-12
76  Residual  35.0   11.203484    0.320100          NaN            NaN
77  >>>
78  >>>
79  >>> #----------------------------------------
80  ... # Scatter plot with the fitted line
81  ... #----------------------------------------
82  ... plot.figure( figsize=(10,10) )
83  <matplotlib.figure.Figure object at 0x7fec2e7a1be0>
84  >>> plot.scatter(year, area, c="black" )
85  <matplotlib.collections.PathCollection object at 0x7fec2e744eb8>
86  >>> plot.plot(year, OUT.fittedvalues, c="red", linewidth=1 )
87  [<matplotlib.lines.Line2D object at 0x7fec2e798cf8>]
88  >>> plot.xlabel("Year")
89  <matplotlib.text.Text object at 0x7fec2e73ff60>
90  >>> plot.ylabel("Extended area")
91  <matplotlib.text.Text object at 0x7fec2e74c668>
92  >>> plot.show()
93  >>>
94  >>> #================================
95  ... # (4) To compare with the paper
96  ... #================================
97  ... # NB: Shorten the data by slicing.
98  ... #    (it starts with zero and ends with 34).
99  ... data2 = pandas.DataFrame({"yr": x[0:34], "ar": y[0:34]})
100 >>> model2 = ols("ar~yr", data=data2) # NB: ols("y~0+x", data)
101 >>> OUT2 = model2.fit()
102 >>> print(OUT2.summary())
103                         OLS Regression Results
104 ==============================================================================
105 Dep. Variable:                     ar   R-squared:                       0.717
106 Model:                            OLS   Adj. R-squared:                  0.708
107 Method:                 Least Squares   F-statistic:                     80.88
108 Date:                Wed, 26 Jun 2019   Prob (F-statistic):           2.85e-10
109 Time:                        19:27:36   Log-Likelihood:                -28.729
110 No. Observations:                  34   AIC:                             61.46
111 Df Residuals:                      32   BIC:                             64.51
112 Df Model:                           1
113 Covariance Type:            nonrobust
114 ==============================================================================
115                  coef    std err          t      P>|t|      [0.025      0.975]
116 ------------------------------------------------------------------------------
117 Intercept    188.6024     20.254      9.312      0.000     147.347     229.858
118 yr            -0.0913      0.010     -8.993      0.000      -0.112      -0.071
119 ==============================================================================
120 Omnibus:                        0.986   Durbin-Watson:                   1.477
121 Prob(Omnibus):                  0.611   Jarque-Bera (JB):                0.228
122 Skew:                           0.032   Prob(JB):                        0.892
123 Kurtosis:                       3.396   Cond. No.                     4.06e+05
124 ==============================================================================
125
126 Warnings:
127 [1] Standard Errors assume that the covariance matrix of the errors is correctly
        specified.
128 [2] The condition number is large, 4.06e+05. This might indicate that there are
129 strong multicollinearity or other numerical problems.
130 >>> print(anova_lm(OUT2))
131            df      sum_sq     mean_sq           F          PR(>F)
132 yr         1.0   27.264989   27.264989   80.877733   2.845158e-10
133 Residual  32.0   10.787637    0.337114         NaN            NaN
134 >>>
135 >>> #----------------------------------------
136 ... # Scatter plot with the fitted line
```

```
137 ... #------------------------------------------
138 ... yr = year[0:34] # NB: slicing (takes 34 observations)
139 >>> ar = area[0:34] # NB: it starts with zero and ends with 34.
140 >>> plot.figure( figsize=(10,10) )
141 <matplotlib.figure.Figure object at 0x7fec2e2c5dd8>
142 >>> plot.scatter(yr, ar, c="black" )
143 <matplotlib.collections.PathCollection object at 0x7fec2de5b780>
144 >>> plot.plot(yr, OUT2.fittedvalues, c="red", linewidth=1 )
145 [<matplotlib.lines.Line2D object at 0x7fec2de61358>]
146 >>> plot.xlabel("Year")
147 <matplotlib.text.Text object at 0x7fec2e2e0198>
148 >>> plot.ylabel("Extended area")
149 <matplotlib.text.Text object at 0x7fec2e2e7860>
150 >>> plot.show()
```

## 1.7  Spurious Correlation

There are examples that attempt to demonstrate how no cause-and-effect (causation) is necessarily implied by the regression model.

For example, an observational study of elementary school children aged 6 – 11 finds a high positive correlation between shoe size $X$ and score $Y$ on a test of reading comprehension. Note that such data are observational data since the explanatory variable, shoe size, is not controlled.

This relation does not imply, however, that an increase in shoe size causes reading ability. There are hidden factors (adequate explanatory variables) such as age of child and amount of education which affect both the shoe size $(X)$ and reading ability $(Y)$.

Even when a strong statistical relationship reflects causal conditions, the causal conditions may act in the *opposite* direction, from $Y$ to $X$.

# 1.8 Miscellaneous