

Introduction to OLS with R

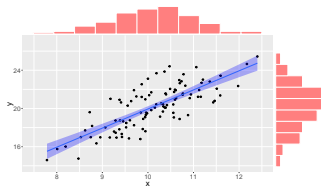
Erasmus+ visiting lecture at Faculty of Management & Economics RUB

Athanassios Stavrakoudis

<http://stavrakoudis.econ.uoi.gr>

Department of Economics, University of Ioannina, Greece

13 January 2016



Contents

- 1 Where to begin
- 2 Anscombe's quartet
- 3 Galton Example regression
- 4 Introduction of simple regression model
- 5 Coefficient of determination
- 6 Prediction

What is this about

- An introduction about basic linear models within **R**.
- **R** is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS.
- Download from : www.r-project.org
- **R** is most conveniently used with **RStudio**, an IDE for **R**.
- Download from : www.rstudio.com
- There are plenty of books, internet resources and MOOCs about **R**.
- **R** the *de-facto* statistical environment in academia. **R** is heavily used by the industry too.

Packages and libraries

There are several packages used in this lectures. If not present in your system you can install them:

```
1 if (!require("HH")) install.packages("HH")
2 if (!require("HistData")) install.packages("HistData")
3 if (!require("dplyr")) install.packages("dplyr")
4 if (!require("xtable")) install.packages("xtable")
5 if (!require("UsingR")) install.packages("UsingR")
6 if (!require("ggplot2")) install.packages("ggplot2")
7 if (!require("ggExtra")) install.packages("ggExtra")
8 if (!require("gridExtra")) install.packages("gridExtra")
```

or execute the following script : [install.R](#)

RStudio environment

The screenshot displays the RStudio interface with the following components:

- Source Editor:** Contains R code for loading data, plotting, and fitting a linear model.
- Console:** Shows the execution output of the code in the source editor.
- Environment Pane:** Lists the objects in the current environment: 'mtcars' (32 observations of 11 variables) and 'fit' (a list of 12).
- Plots Pane:** Displays a scatter plot of 'mtcars\$mpg' versus 'mtcars\$hp' with a blue regression line and a vertical reference line at approximately 130 hp.

```
1 rm(list=ls(all=T))
2 data(mtcars)
3
4 plot(mtcars$hp, mtcars$mpg)
5 fit <- lm(mpg ~ hp, data=mtcars)
6 abline(fit, col=4, lwd=2)
7
8 fit <- lm(mtcars$mpg ~ mtcars$hp)
9
10 attach(mtcars)
11 fit <- lm(mpg ~ hp)
12
```

```
> rm(list=ls(all=T))
> data(mtcars)
>
> plot(mtcars$hp, mtcars$mpg)
> fit <- lm(mpg ~ hp, data=mtcars)
> abline(fit, col=4, lwd=2)
>
> fit <- lm(mtcars$mpg ~ mtcars$hp)
>
> attach(mtcars)
> fit <- lm(mpg ~ hp)
> |
```

mtcars\$hp	mtcars\$mpg
53	24.4
54	26.7
57	30.4
61	27.8
68	26.7
75	24.4
76	21.4
78	20.1
85	21.5
97	19.2
103	17.8
104	16.4
105	15.2
106	14.7
108	15.8
110	14.7
113	14.7
115	15.2
118	14.7
121	15.2
125	15.2
130	15.2
135	15.2
140	15.2
145	15.2
150	15.2
155	15.2
160	15.2
165	15.2
170	15.2
175	15.2
180	15.2
185	15.2
190	15.2
195	15.2
200	15.2
205	15.2
210	15.2
215	15.2
220	15.2
225	15.2
230	15.2
235	15.2
240	15.2
245	15.2
250	15.2
255	15.2
260	15.2
265	15.2
270	15.2
275	15.2
280	15.2
285	15.2
290	15.2
295	15.2
300	15.2
305	15.2
310	15.2
315	15.2
320	15.2
325	15.2
330	15.2
335	15.2
340	15.2
345	15.2
350	15.2
355	15.2
360	15.2
365	15.2
370	15.2
375	15.2
380	15.2
385	15.2
390	15.2
395	15.2
400	15.2

Some introductory commands

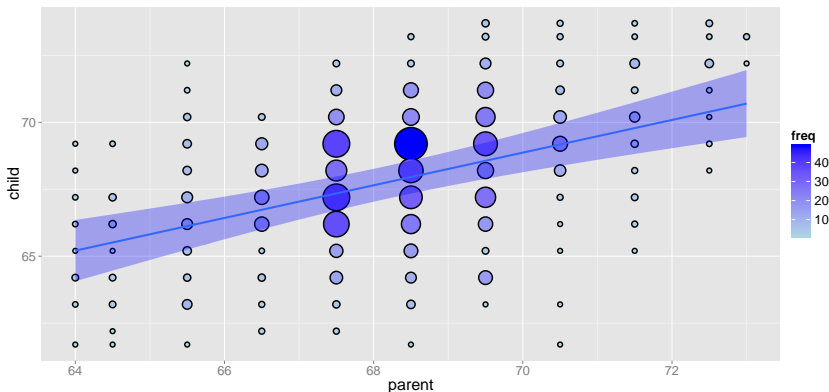
R has lots of methods to load data from local media or network sources. There are many embedded data sets useful for examples.

```
1 rm(list=ls(all=T))           # clear the enviroment
2 data(mtcars)                 # load a data frame
3
4 plot(mtcars$hp, mtcars$mpg)  # plot miles per gallon
5                               # vs horse power
6 fit <- lm(mpg ~ hp, data=mtcars) # fit within the data frame
7 abline(fit, col=4, lwd=2)    # add regression line
8
9 fit <- lm(mtcars$mpg ~ mtcars$hp) # use variable names
10
11 attach(mtcars)
12 fit <- lm(mpg ~ hp)         # use attached variables
```

Most of the lecture is about the **lm()** function.

Some history

Over a century ago, F. Galton observed that children's height correlates with their parents height. This is one of the first application of OLS and linear models to science.



F. Galton, Regression towards mediocrity in hereditary stature, *The Journal of the Anthropological Institute of Great Britain and Ireland* **1886**, 15:246–263.

doi:10.2307/2841582

Variables in OLS method

 $Y =$ $X\beta + \epsilon$

Response variable

Response variable (or dependent variable) must be **continuous**.

Predictor variable

Predictor variable(s) (or independent variable) can be:

- 1 continuous
- 2 discrete
- 3 categorical

Data in OLS

- Errors** check/insect the data possible errors
- Missing values** Maybe some data values are missed
- Patterns** not all data sets are suitable for OLS
- Outliers** unusual or extreme values

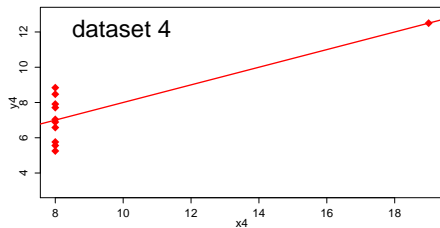
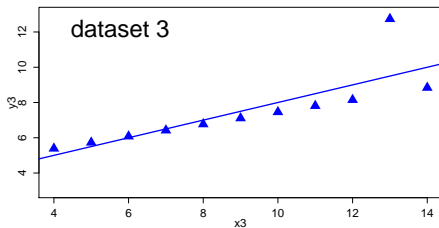
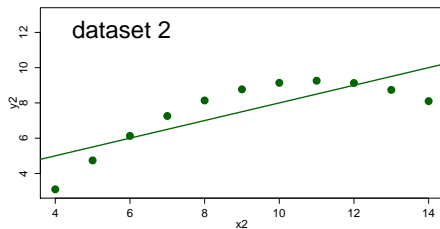
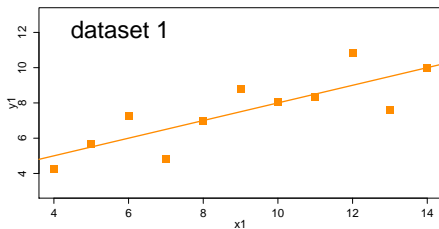
Plotting

Begin your analysis with plotting of data!

Contents

- 1 Where to begin
- 2 Anscombe's quartet**
- 3 Galton Example regression
- 4 Introduction of simple regression model
- 5 Coefficient of determination
- 6 Prediction

Anscombe's quartet: four different data sets



F.J. Anscombe, Graphs in Statistical Analysis, *The American Statistician*, 1973, 27:17–21

Same properties: mean and variance

```
1 data(anscombe)
2
3 > apply(anscombe, 2, mean)
4   x1   x2   x3   x4   y1   y2   y3   y4
5 9.000 9.000 9.000 9.000 7.501 7.501 7.500 7.501
6
7 > apply(anscombe, 2, var)
8   x1   x2   x3   x4   y1   y2   y3   y4
9 11.000 11.000 11.000 11.000 4.127 4.128 4.123 4.123
```

Attention!

All four data sets have the same mean and variance (to 2nd decimal point). This does not mean that come from the same distribution!

Same properties: correlation and OLS coefficients

```
1 > cor(x1, y1)
2 [1] 0.8164205
3 > cor(x2, y2)
4 [1] 0.8162365
5 > cor(x3, y3)
6 [1] 0.8162867
7 > cor(x4, y4)
8 [1] 0.8165214
```

```
1 > coef(lm.1)
2 (Intercept)          x1
3   3.0000909    0.5000909
4 > coef(lm.2)
5 (Intercept)          x2
6   3.000909    0.500000
7 > coef(lm.3)
8 (Intercept)          x3
9   3.0024545    0.4997273
10 > coef(lm.4)
11 (Intercept)          x4
12   3.0017273    0.4999091
```

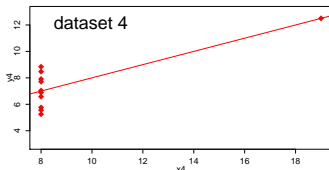
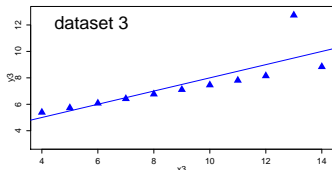
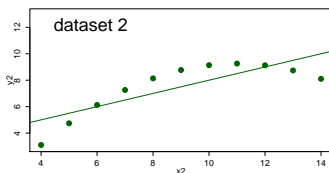
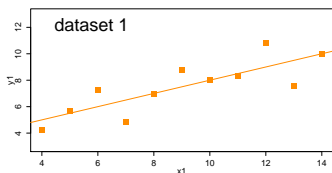
Attention!

Correlation and OLS coefficients are the same in all four data sets. This does not mean that imply that in all cases we !

Ancombe's quartet: plot your data

Do not use software as a black box

The four data sets have the same properties but obviously they are different: always plot your data as a first step of your analysis.



Question

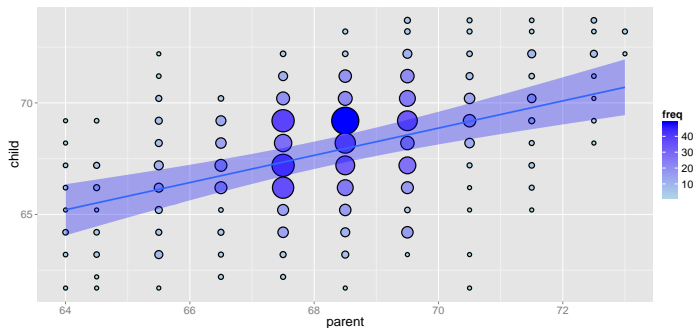
In which case can we reasonably accept the linear model hypothesis?

Contents

- 1 Where to begin
- 2 Anscombe's quartet
- 3 Galton Example regression**
- 4 Introduction of simple regression model
- 5 Coefficient of determination
- 6 Prediction

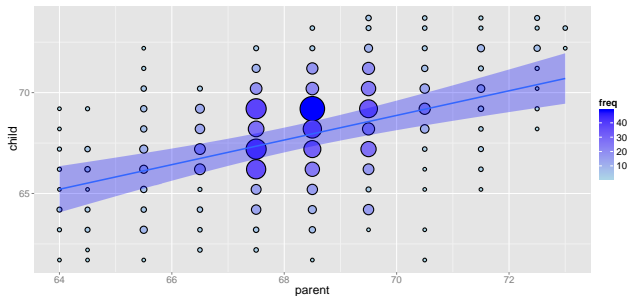
Galton regression

$$\text{childHeight} = \beta_0 + \beta_1 \text{parentHeight} + \epsilon$$



Galton, F. "Regression towards mediocrity in hereditary stature". *The Journal of the Anthropological Institute of Great Britain and Ireland* **1886**, 15:246–263. doi:10.2307/2841583.

Galton Data



- ① A study about parent and child height
- ② Units are in inches
- ③ Parent weight is computed as

$$parent = 0.5 \times (father + 1.08 \times mother)$$
- ④ The graph is overplotted, many points appear with the same coordinates
- ⑤ Size of bubble indicates frequency of appearance

Plot Galton data

```
1 library(HistData)
2 data(Galton)
3
4 plot(Galton$parent, Galton$child)
5 fit <- lm(child~parent, data=Galton)
6 abline(fit, col=4, lwd=3)
7 summary(fit)
8
9 library(ggplot2)
10 ggplot(Galton, aes(x=parent, y=child)) +
11   geom_point(shape=5) +
12   scale_colour_hue(l=50) +
13   stat_smooth(method=lm, se=T, level=0.95,
14             fill="blue", alpha=0.3, size=1.2) +
15   theme(text=element_text(size=20),
16         axis.title=element_text(size=20))
```

Child's height vs parent's height

```
1 > library(HistData)
2 > lm (child ~ parent, data=Galton)
3
4 Call:
5 lm(formula = child ~ parent, data = Galton)
6
7 Coefficients:
8 (Intercept)      parent
9      23.9415      0.6463
```

Interpretation

- 1 A child's height (in inches) is $\approx 23.9'' + 0.65$ times parent height.
- 2 A child's height is $\approx 23.9''$ if his/her parent height is $\approx 0''$.
- 3 If two parents height differ by $1''$ then their children height will differ by $\approx 0.65''$

A note about variables, lists and numbers

```

1 > coefficients(fit)           # get the coefficients
2 (Intercept)      parent
3 23.9415302      0.6462906
4 > fit$coefficients
5 (Intercept)      parent
6 23.9415302      0.6462906
7 > fit$coef
8 (Intercept)      parent
9 23.9415302      0.6462906
10 > fit$coef[1]              # the constant term
11 (Intercept)
12 23.94153
13 > fit$coef[2]              # the slope
14 parent
15 0.6462906
16 > fit$coef[[2]]            # the slope
17 [1] 0.6462906
18 > str(fit$coef[2])
19 Named num 0.646
20 - attr(*, "names")= chr "parent"
21 > str(fit$coef[[2]])
22 num 0.646

```

Sum of residuals

```
1 > mean(residuals(fit))
2 [1] -2.359884e-15
3 > mean(fit$resid)
4 [1] -2.359884e-15
5 > mean(residuals(fit))
6 [1] -2.359884e-15
7
8 > u <- fit$resid
9 > mean(u)
10 [1] -2.359884e-15
```

Notes

- 1 We verify that $\sum \hat{\epsilon}_i = 0$
- 2 There are several ways to obtain the result
- 3 Storing the residuals on a variable (for example **u**) helps if you want further calculations

Getting more results with summary()

```
1 > summary(fit)
2
3 Call:
4 lm(formula = child ~ parent, data = Galton)
5
6 Residuals:
7     Min       1Q   Median       3Q      Max
8 -7.8050 -1.3661  0.0487  1.6339  5.9264
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) 23.94153    2.81088   8.517  <2e-16 ***
13 parent      0.64629    0.04114  15.711  <2e-16 ***
14 ---
15 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
16
17 Residual standard error: 2.239 on 926 degrees of freedom
18 Multiple R-squared:  0.2105, Adjusted R-squared:  0.2096
19 F-statistic: 246.8 on 1 and 926 DF, p-value: < 2.2e-16
```

Coefficients, standard errors and t-statistics

```
1 > fit$coeff
2 (Intercept)      parent
3 23.9415302      0.6462906
4
5 > summary(fit)$coeff
6           Estimate Std. Error  t value    Pr(>|t|)
7 (Intercept) 23.9415302 2.81087834  8.517455 6.536845e-17
8 parent      0.6462906 0.04113588 15.711115 1.732509e-49
```

Summarizing the results

- 1 Calling **summary()** yields much more detailed results.
- 2 We do not need only the coefficient values, we also need their Std. Errors to interpret the results.
- 3 t-statistic and p-values helps us evaluate the significance of the estimation.

Tabulate the results of a regression

```

1 > library(xtable)
2 > xtable(fit)
3 \begin{table}[ht]
4 \centering
5 \begin{tabular}{rrrrr}
6 \hline
7 & Estimate & Std. Error & t value & Pr(>|t|) \\
8 \hline
9 (Intercept) & 23.9415 & 2.8109 & 8.52 & 0.0000 \\
10 parent & 0.6463 & 0.0411 & 15.71 & 0.0000 \\
11 \hline
12 \end{tabular}
13 \end{table}

```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	23.9415	2.8109	8.52	0.0000
parent	0.6463	0.0411	15.71	0.0000

Standard error of the coefficient

$$s_{\hat{\beta}_1} = \sqrt{\frac{\frac{1}{n-2} \sum \hat{\varepsilon}_i^2}{\sum (x_i - \bar{x})^2}}$$

```

1 > n <- length(Galton$parent)
2 > u2 <- fit$resid ** 2
3 > x2 <- (Galton$parent - mean(Galton$parent)) ** 2
4 >
5 > sqrt( sum(u2)/(n-2) / sum(x2) )
6 [1] 0.04113588
7 > sqrt(diag(vcov(fit)))
8 (Intercept)      parent
9 2.81087834 0.04113588

```

Standard error of the coefficient

General Model: $Y_i = \beta_0 + \beta_1 \cdot x_i + \epsilon_i$

Galton: $child_i = \beta_0 + \beta_1 \cdot parent_i + \epsilon_i$

$$H_0 : \beta_1 = 0 \quad \longleftrightarrow \quad H_1 : \beta_1 \neq 0$$

T-test statistic:

$$T_1 = \frac{\hat{\beta}_1 - 0}{\hat{\sigma}_{\beta_1}}$$

With normality of ϵ_i under the Null hypothesis:

$$T_1 \sim t_{n-2}$$

$$\text{p-value: } p = 2 \cdot (1 - F_{t_{n-2}}(|T_1|))$$

Obtain p-values analytically

Here is a good example of how to understand and interpret an analytical mathematical formula into R code.

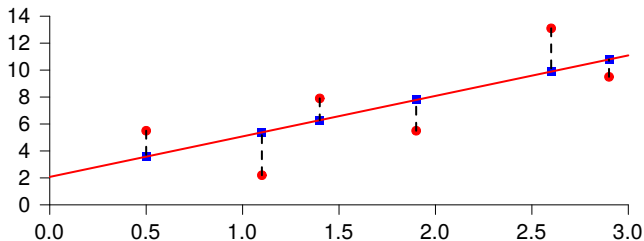
$$p = 2 \left(1 - F_{t_{n-2}}(|T_1|) \right)$$

```
1 > summary(fit)$coeff
2           Estimate Std. Error t value Pr(>|t|)
3 (Intercept)  23.9415    2.81088   8.517 6.537e-17
4 parent        0.6463    0.04114  15.711 1.733e-49
5
6 n      <- nrow(Galton)
7 sde    <- sqrt(diag(vcov(fit)))
8 tstat  <- fit$coeff / sde
9 pval   <- 2 * ( 1 - pt(tstat, df=n-2) )
```

Contents

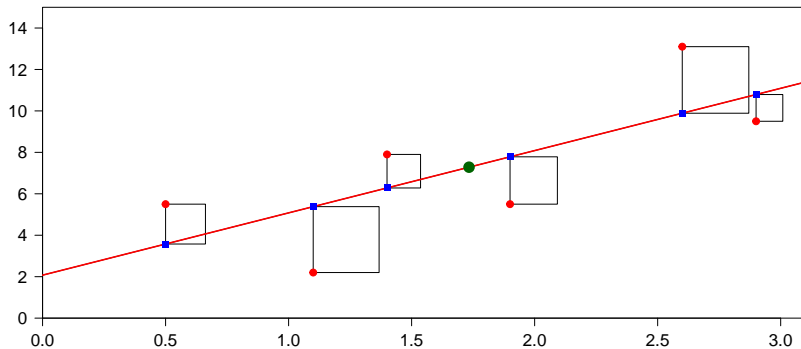
- ① Where to begin
- ② Anscombe's quartet
- ③ Galton Example regression
- ④ Introduction of simple regression model**
- ⑤ Coefficient of determination
- ⑥ Prediction

Basic Terminology



- ① **red points:** $\{x_i, y_i\}$, observable variables (price, GDP, temperature, years in work, etc), out of line.
- ② **blue points:** $\{x_i, \hat{y}_i\}$, fitted (estimated) data, on the line..
- ③ **red line:** regression line (fitted from data)
- ④ **dashed lines:** residuals, $\{\hat{y}_i - y_i\}$ (estimation of errors)
- ⑤ **y:** outcome or dependent variable
- ⑥ **x:** explanatory variable or independent variable or regressor

Squared residuals



- 1 **squares**: represent $(\hat{y}_i - y_i)^2$
- 2 **Regression line** is the one that minimizes:

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

- 3 **green point**: (\bar{x}, \bar{y}) , regression line passes through the mean values.

OLS, some reminders

$$\widehat{Y}_i = \beta_0 + \beta_1 \widehat{X}_i + \epsilon_i$$

$$E(\epsilon_i) = 0 \quad \forall i \quad \text{mean value of residuals}$$

$$V(\epsilon_i) = \sigma^2 \quad \forall i \quad \text{homoscedasticity}$$

$$\text{Cov}(\epsilon_i, \epsilon_j) = 0 \quad \forall i \neq j \quad \text{non correlated residuals}$$

$$\widehat{\beta}_1 = \text{cor}(Y, x) \frac{\sigma_Y}{\sigma_x} \quad \text{slope of the regression line}$$

$$\widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{x} \quad \text{constant term}$$

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2} \quad \text{if } \beta_0 = 0$$

Simple linear regression model, estimated in R

Regression equation

Regression equation $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$

```
lm(y ~ x)
```

Fit the model and estimation of the coefficients

```
Formula <- y ~ x  
lm(Formula)
```

Alternative way to call **lm** function

```
fit <- lm(y ~ x)  
summary(fit)
```

Store the model to a variable for further usage

```
lm(y ~ x, data=DF)  
lm(DF$y ~ DF$x)
```

Variables from a data frame (DF)

```
lm(y ~ x - 1)  
lm(y ~ 0 + x)
```

Exclude constant term, two equivalent methods

Estimation in R: quick reference

```
fit <- lm(y ~ x)
```

Store the results of the estimation to a variable

```
summary(fit)
```

Examine the results

```
1 plot(x, y)
```

```
2 abline(fit)
```

Quick scatter plot of the data with the regression line

```
coefficients(fit)
```

List the coefficients

```
residuals(fit)
```

Extract the residuals ($\hat{\epsilon}_i$)

```
fitted.values(fit)
```

Extract the fitted values (\hat{y}_i)

```
confint(fit)
```

Confidence interval

```
vcov(fit)
```

Variance/covariance matrix

Contents

- ① Where to begin
- ② Anscombe's quartet
- ③ Galton Example regression
- ④ Introduction of simple regression model
- ⑤ Coefficient of determination**
- ⑥ Prediction

Coefficient of determination R^2

The coefficient of determination (**Multiple R-squared** in **R**) indicates the proportion of the variance in Y , which is explained by the model.

$$R^2 = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = \frac{ESS}{TSS} = \frac{S_{xY}^2}{S_x^2 S_Y^2} = r_{xY}^2$$

$R^2 \rightarrow 1$ indicates a good fit

$R^2 \rightarrow 0$ indicates a poor fit

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

Total sum of squares

$$ESS = \sum (y_i - \hat{y})^2$$

Error sum of squares

$$RSS = \sum (\hat{y}_i - \bar{y})^2$$

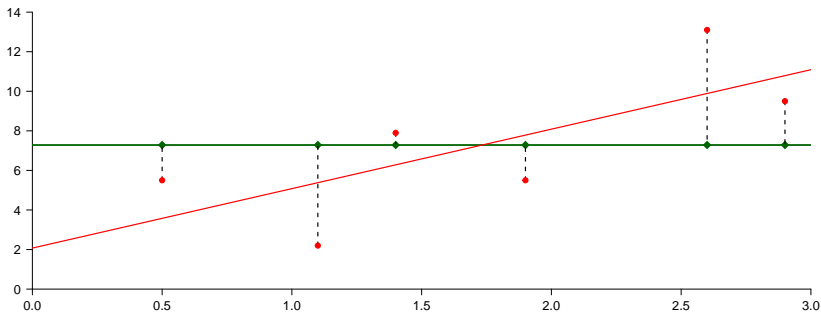
regression sum of squares

$$TSS = RSS + ESS$$

$$R^2 = RSS/TSS$$

Fraction of unexplained variance (FUV)

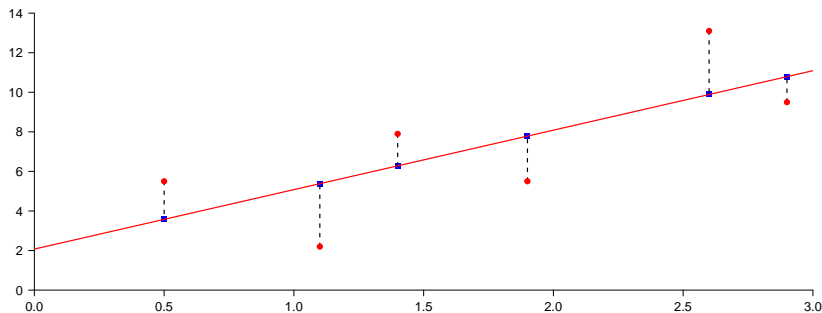
Total sum of squares (TSS)



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

- 1 **red points:** y_i , observations
- 2 **green points:** \bar{y} , mean value
- 3 **dashed lines:** $y_i - \bar{y}$

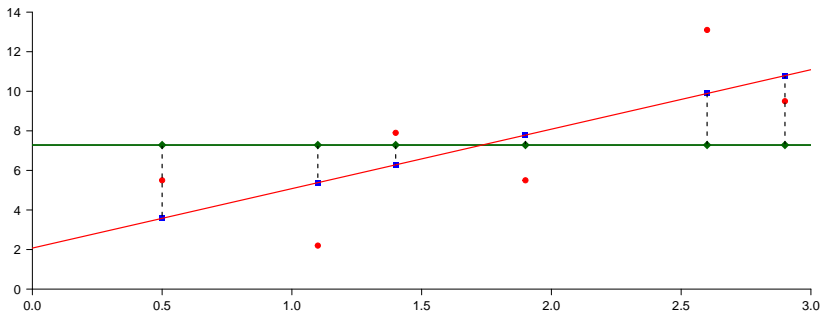
Error sum of squares (ESS)



$$\begin{aligned} ESS &= \sum (y_i - \hat{y})^2 = \sum \epsilon_i^2 \\ &= \sum \left(y_i - (\beta_0 + \beta_1 x_i) \right)^2 \end{aligned}$$

- ① **red points:** y_i , observations
- ② **blue points:** \hat{y}_i , fitted values
- ③ **dashed lines:** $y_i - \hat{y}_i$

Regression sum of squares (RSS)



$$RSS = \sum (\hat{y} - \bar{y})^2$$

- ① **red points:** y_i , observations
- ② **blue points:** \hat{y}_i , fitted values
- ③ **green points:** \bar{y} , mean value of y
- ④ **dashed lines:** $\hat{y} - \bar{y}$

R^2 value

```
1 > f <- lm (y ~ x)
2 > summary(f)
3
4 Call:
5 lm(formula = y ~ x)
6
7 Coefficients:
8             Estimate Std. Error t value Pr(>|t|)
9 (Intercept)    2.074      2.730   0.760   0.490
10 x              3.006      1.419   2.118   0.102
11
12 Residual standard error: 2.899 on 4 degrees of freedom
13 Multiple R-squared:  0.5285, Adjusted R-squared:  0.4107
14 F-statistic: 4.484 on 1 and 4 DF,  p-value: 0.1016
```

$$R^2 = 0.5285$$

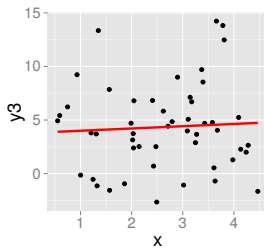
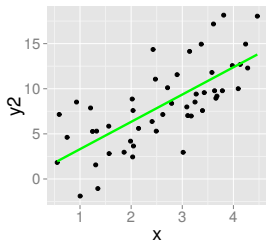
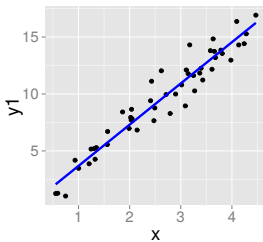
Compute TSS, ESS, RSS and R^2 values

```

1 > tss <- sum( ( y - mean(y) )^2 )
2 > rss <- sum( ( f$fit - mean(y) )^2 )
3 > ess <- sum( ( f$fit - y )^2 )
4 > r2 <- 1 - ess/tss
5 > r2
6 [1] 0.5285435
7 > r2 <- rss/tss
8 > r2
9 [1] 0.5285435
10 > tss
11 [1] 71.32833
12 > rss+ess
13 [1] 71.32833

```

Comparison with R2



```
1 lm1 <- lm(y1 ~ x)
2 lm2 <- lm(y2 ~ x)
3 lm3 <- lm(y3 ~ x)
4 r2.1 <- summary(lm1)$r.sq
5 r2.2 <- summary(lm2)$r.sq
6 r2.3 <- summary(lm3)$r.sq
7 > cbind(r2.1, r2.2, r2.3)
8           r2.1      r2.2      r2.3
9 [1,] 0.92725 0.51492 0.0031825
```

R^2 values

Bigger R^2 values **indicate** better fit of the data.

Caution: it is not a strong criterion.

Adjusted R^2

$$\bar{R}^2 = 1 - \frac{n-1}{n-k} (1 - R^2)$$

Extract R^2 and \bar{R}^2

```

1 summary(f)$r.sq
2 summary(f)$adj.r.sq
3
4 fs      <- summary(f)
5 r.sq    <- fs$r.squared
6 r.sq.adj <- fs$adj.r.sq
7 > cbind(r.sq, r.sq.adj)
8         r.sq  r.sq.adj
9 [1,] 0.5285435 0.5893206

```

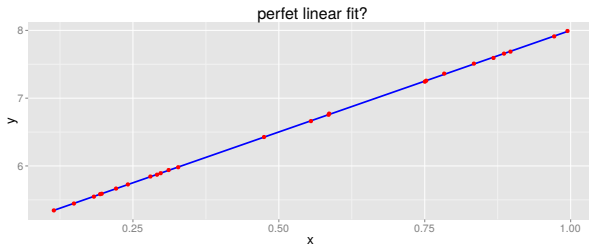
Analytically computing \bar{R}^2

```

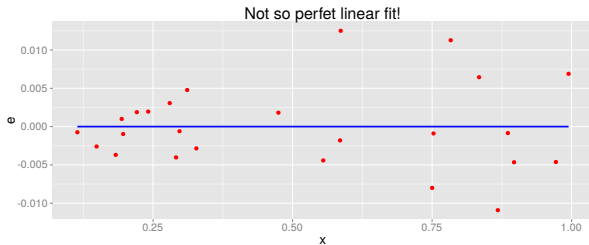
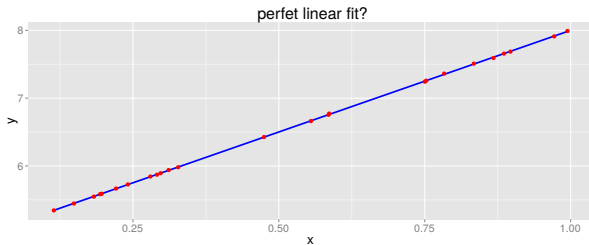
1 n <- length(x)
2 k <- f$rank
3 r.sq.adj2 <- (n-1)/(n-k)*(1-r.sq)
4 > cbind(r.sq.adj, r.sq.adj2)
5         r.sq.adj  r.sq.adj2
6 [1,] 0.5893206 0.5893206

```

The heteroscedasticity problem



The heteroscedasticity problem



Comment about correlations in errors

```
1 > n <- 25
2 > x <- c(runif(n, 0, 1))
3 > y <- 5 + 3*x + rnorm(n, sd=x*0.01)
4 > f <- lm( y~ x )
5 > summary(f)
6
7 Call:
8 lm(formula = y ~ x)
9
10             Estimate Std. Error t value Pr(>|t|)
11 (Intercept)  4.998397   0.001348   3708   <2e-16 ***
12 x            3.004048   0.002698   1113   <2e-16 ***
13 ---
14 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15
16 Residual standard error: 0.003483 on 23 degrees of freedom
17 Multiple R-squared: 1, Adjusted R-squared: 1
18 F-statistic: 1.24e+06 on 1 and 23 DF, p-value: < 2.2e-16
```

Contents

- ① Where to begin
- ② Anscombe's quartet
- ③ Galton Example regression
- ④ Introduction of simple regression model
- ⑤ Coefficient of determination
- ⑥ Prediction

Predicting in **lm** models

```
1 x <- c(0.5, 1.1, 1.4, 1.9, 2.6, 2.9)
2 y <- c(5.5, 2.2, 7.9, 5.5, 13.1, 9.5)
3 f <- lm(y ~ x)
4 b0 <- f$coef[[1]]
5 b1 <- f$coef[[2]]
6 xp <- 2
7 yp <- b0 + b1*xp
8 > yp
9 [1] 8.084824
```

Point prediction

Suppose we have a point $X_p = 2$. We want to know what is the predicted value by the model. This is:

$$Y_p = \hat{\beta}_0 + \hat{\beta}_1 \times X_p$$

The above code makes the needed computation.

Prediction in R with multiple values

Using a vector

```
1 xp <- c(1, 2, 3)
2 yp <- b0 + b1*xp
3 yp
4 [1] 5.079233 8.084824 11.090415
```

Using the `predict()` function

```
1 xp <- data.frame(x=c(1, 2, 3))
2 > predict.lm(f, xp)
3           1           2           3
4 5.079233 8.084824 11.090415
```

The second option has to be a **data frame**.

Prediction in R with **predict**

```
1 > predict(f, interval="prediction")
2         fit         lwr         upr
3 1  3.576438 -6.3849494 13.53782
4 2  5.379792 -3.6665849 14.42617
5 3  6.281470 -2.5124817 15.07542
6 4  7.784265 -0.9357985 16.50433
7 5  9.888179  0.5462288 19.23013
8 6 10.789856  0.9539719 20.62574
9
10 > newX <- data.frame(x=c(1, 2, 3))
11 > predict(f, newX)
12         1         2         3
13 5.079233 8.084824 11.090415
14
15 > predict(f, newX, interval="prediction")
16         fit         lwr         upr
17 1  5.079233 -4.0836929 14.24216
18 2  8.084824 -0.6737388 16.84339
19 3 11.090415  1.0642844 21.11655
```

Prediction interval

```
1 fit <- lm (y ~ x)
2 predict(fit, newdata, interval, level)
```

fit An lm object

newdata data for prediction (must be in data frame format)

interval can be **"confidence"** for confidence interval or **"prediction"** for prediction interval

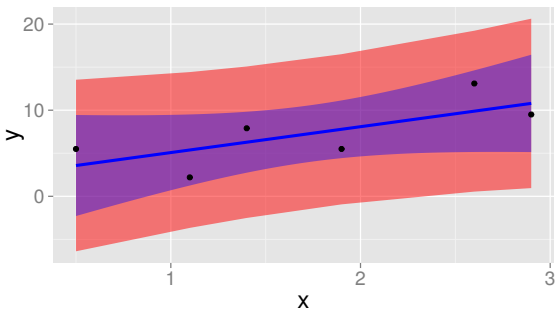
level tolerance/confidence level

Under the assumption of normally distributed residuals the prediction interval can be calculated as:

$$\hat{Y}_P = \hat{Y} \pm \hat{\sigma} \sqrt{1 + \frac{1}{n} \left(1 + \frac{(x_P - \bar{x})^2}{\tilde{S}_x^2} \right)} \cdot t_{n-2; 1-\alpha/2}$$

with $\hat{\sigma} = \frac{1}{n-2} \cdot \sum \hat{\epsilon}_i^2$ and $\tilde{S}^2 = \frac{n-1}{n} \cdot S^2$

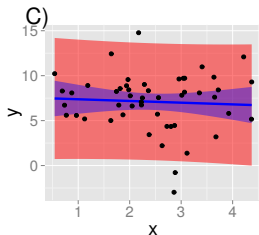
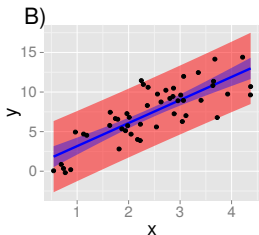
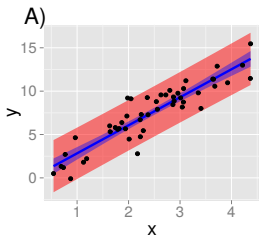
Confidence interval vs Prediction interval



Notes and explanations

- 1 Confidence interval tells us how good is the fitting.
- 2 Prediction interval tells where new data points are expected.
- 3 Confidence : about existed data
Prediction : about non existed data
- 4 Prediction interval area is always wider than confidence interval area.

Confidence interval vs Prediction interval






Discussion

Examine the three plots above and comment on the following:

- 1 Variance of y and ϵ
- 2 R^2
- 3 Std. Error of the coefficient (slope, $\widehat{\beta}_1$)
- 4 P-value of the coefficient (slope, $\widehat{\beta}_1$)
- 5 Band width of confidence interval
- 6 Band width of prediction interval

Reading list

-  J. Verzani g, Using R for Introductory Statistics , 2rd ed., 2014, Chapman & Hall/CRC
-  J. Fox & S. Weisberg, An R Companion to Applied Regression, 3rd ed., 2016, Sage Publications
-  C. Kleiber & A. Zeileis, Applied Econometrics with R, 2008, Springer

Online and electronic resources

- 1 <http://www.ats.ucla.edu/stat/r/>
- 2 <http://www.stat.wisc.edu/~larget/stat302/>
- 3 <https://www.coursera.org/learn/regression-models/>

Comments and Questions

Thank you very much
for your attention!