

Exercises

Section 10.2 Fitting the Simple Linear Regression Model

10.4 The time between eruptions ...

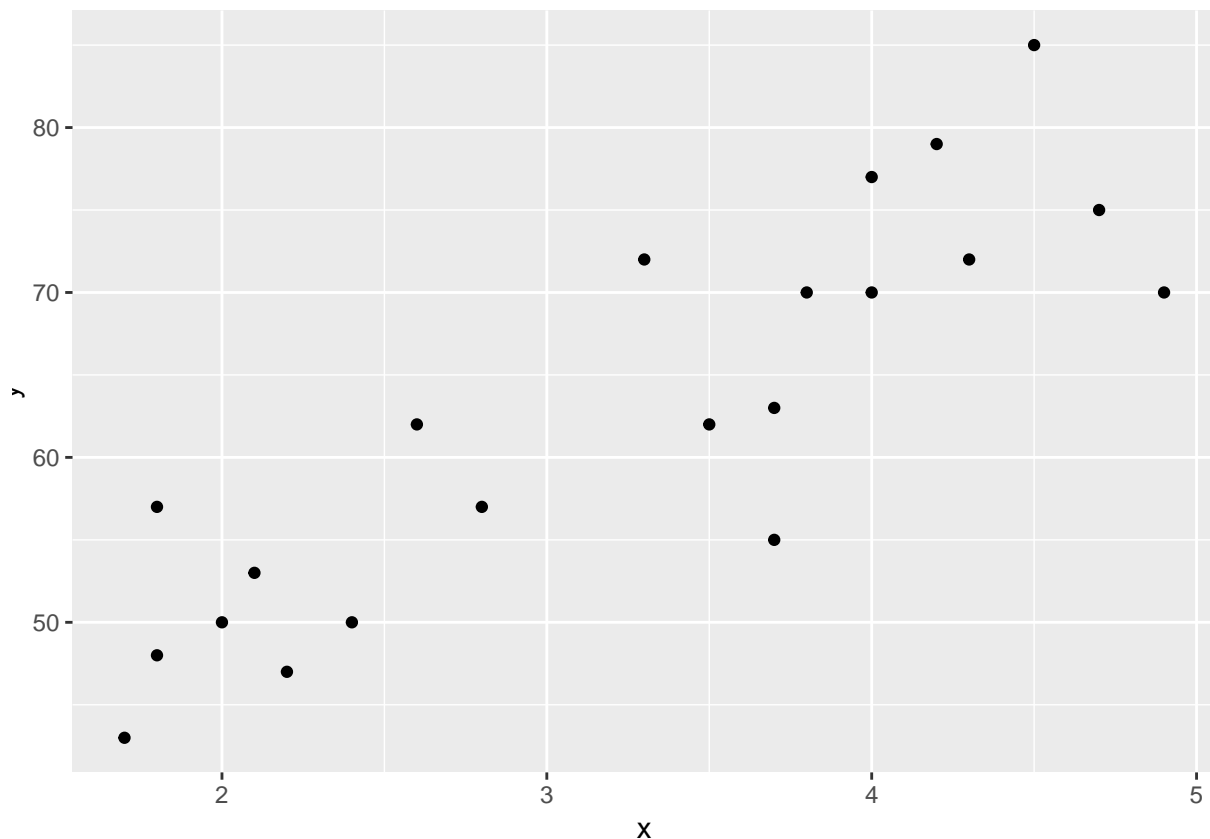
```
library(ggplot2)
```

```
oldfaithful <- read.csv(file='data/1987oldfaithful.csv',header=TRUE,sep=',')  
head(oldfaithful)
```

```
##   Observation.Number Duration.of.Eruption Time.Between.Eruptions  
## 1                   1                2.0                50  
## 2                   2                1.8                57  
## 3                   3                3.7                55  
## 4                   4                2.2                47  
## 5                   5                2.1                53  
## 6                   6                2.4                50
```

Assume the time between eruptions is linearly dependent on the duration of the last eruption.

```
x <- oldfaithful$'Duration.of.Eruption'  
y <- oldfaithful$'Time.Between.Eruptions'  
ggplot(data=oldfaithful, aes(x,y)) + geom_point()
```



Let

- $n = 21$
- $(x_i)_{i \in [n]}$ be the values of the independent variable.
- $(Y_i)_{i \in [n]}$ be random variables for which $(y_i)_{i \in [n]}$ are the observations.

Assume for each $i \in [n]$ that Y_i is linearly dependent on x_i . In particular that there exists two real values β_0 and β_1 and a random error $\epsilon_i \sim N(0, \sigma^2)$ such that

$$Y_i = \beta_0 + \beta_1 \cdot x_i + \epsilon_i$$

We can estimate β_0 and β_1 by the least squares method:

$$\min_{(\beta_0, \beta_1) \in \mathbb{R}^2} Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 \cdot x_i))^2$$

From calculus we know a function Q has a critical point at $(\hat{\beta}_0, \hat{\beta}_1)$ if $\nabla Q(\hat{\beta}_0, \hat{\beta}_1) = 0$.

$$\frac{\partial Q}{\partial \beta_0} = -2 \cdot \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 \cdot x_i))$$

$$\frac{\partial Q}{\partial \beta_1} = -2 \cdot \sum_{i=1}^n x_i \cdot (y_i - (\beta_0 + \beta_1 \cdot x_i))$$

Therefore

$$\begin{aligned} n \cdot \hat{\beta}_0 + \left(\sum_{i=1}^n x_i \right) \cdot \hat{\beta}_1 &= \left(\sum_{i=1}^n y_i \right) \\ \left(\sum_{i=1}^n x_i \right) \cdot \hat{\beta}_0 + \left(\sum_{i=1}^n x_i^2 \right) \cdot \hat{\beta}_1 &= \left(\sum_{i=1}^n x_i \cdot y_i \right) \end{aligned}$$

By defining $S_{xy} = \sum_{i=1}^n (x_i - \bar{x}) \cdot (y_i - \bar{y})$ one can show

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \cdot \bar{x}, \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

```
S <- function(x,y) {
  sum( (x-mean(x))*(y-mean(y)) )
}
```

```
betahat_1 <- S(x,y)/S(x,x); betahat_1 # slope
```

```
## [1] 9.790069
```

```
betahat_0 <- mean(y) - betahat_1*mean(x); betahat_0 # y-intercept
```

```
## [1] 31.01311
```