# Math 362: Mathmatical Statistics II

Le Chen
le.chen@emory.edu

Emory University
Atlanta, GA

Last updated on April 2, 2020

2020 Spring

Chapter 11. Regression

# Plan

# Chapter 11. Regression

https://madhureshkumar.wordpress.com/

https://xkcd.com/

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

# Three ways to view the same thing

$$(x_1, y_1), \cdots, (x_n, y_n)$$

1. Purely data, no probability structure assumed.

$$(x_1, Y_1), \cdots, (x_n, Y_n)$$

2. A random sample of size $n$, where $Y_i$ follows a distribution depending on $x_i$ which is deterministic.

$$(X_1, Y_1), \cdots, (X_n, Y_n)$$

3. A random sample of size $n$, where $(X_i, Y_i)$ follow some joint distribution.

1.

# Plan

# Chapter 11. Regression

In linear regression, the observations (**red**) are assumed to be the result of random deviations (**green**) from an underlying relationship (**blue**) between a dependent variable (*y*) and an independent variable (*x*).

Goal: Find a blue line that minimizes
the sum of the square of the green lines

Thm. Given $n$ points $(x_1, y_1), \cdots, (x_n, y_n)$, the straight line $y = a + bx$ minimizing

$$L(a, b) = \sum_{i=1}^{n} [y_i - (a + bx_i)]^2$$

when

$$b = \frac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right) \left(\sum_{i=1}^{n} y_i\right)}{n \left(\sum_{i=1}^{n} x_i^2\right) - \left(\sum_{i=1}^{n} x_i\right)^2}$$

and

$$a = \frac{\sum_{i=1}^{n} y_i - b \sum_{i=1}^{n} x_i}{n} = \bar{y} - b\bar{x}.$$

Proof.

$$\begin{cases} \dfrac{\partial}{\partial a} L(a, b) = \displaystyle\sum_{i=1}^{n} (-2) [y_i - (a + bx_i)] = 0 \\ \dfrac{\partial}{\partial b} L(a, b) = \displaystyle\sum_{i=1}^{n} (-2x_i) [y_i - (a + bx_i)] = 0 \end{cases}$$ (Normal equations)

**Thm.** Given $n$ points $(x_1, y_1), \cdots, (x_n, y_n)$, the straight line $y = a + bx$ minimizing

$$L(a, b) = \sum_{i=1}^{n} [y_i - (a + bx_i)]^2$$

when

$$b = \frac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right) \left(\sum_{i=1}^{n} y_i\right)}{n \left(\sum_{i=1}^{n} x_i^2\right) - \left(\sum_{i=1}^{n} x_i\right)^2}$$

and

$$a = \frac{\sum_{i=1}^{n} y_i - b \sum_{i=1}^{n} x_i}{n} = \bar{y} - b\bar{x}.$$

**Proof.**

$$\begin{cases} \dfrac{\partial}{\partial a} L(a, b) = \sum_{i=1}^{n} (-2) [y_i - (a + bx_i)] = 0 \\ \dfrac{\partial}{\partial b} L(a, b) = \sum_{i=1}^{n} (-2x_i) [y_i - (a + bx_i)] = 0 \end{cases} \qquad \text{(Normal equations)}$$

$$\iff \begin{cases} \displaystyle\sum_{i=1}^{n} y_i - na - b\sum_{i=1}^{n} x_i = 0 & (1) \\ \displaystyle\sum_{i=1}^{n} x_i y_i - a\sum_{i=1}^{n} x_i - b\sum_{i=1}^{n} x_i^2 = 0 & (2) \end{cases}$$

$$(1) \implies a = \bar{y} - b\bar{x}$$

$$(1) \times \sum_{i=1}^{n} x_i - (2) \times n \implies b = \frac{n\sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)}{n\left(\sum_{i=1}^{n} x_i^2\right) - \left(\sum_{i=1}^{n} x_i\right)^2}$$

$\square$

# (Moore-Penrose) Pseudoinverse

1. Well determined system

$$Ax = b \implies x = A^{-1}y.$$

2. Overdetermined system

$$Ax = y$$
$$A^T Ax = A^T y$$
$$\underbrace{(A^T A)^{-1} A^T A}_{=I} x = (A^T A)^{-1} A^T y$$
$$x = \underbrace{(A^T A)^{-1} A^T}_{=:A^+} y$$

3. Under determined system

$$Ax = y \implies x = \underbrace{A^T (AA^T)^{-1}}_{=:A^+} y.$$

# (Moore-Penrose) Pseudoinverse

1. Well determined system

$$Ax = b \implies x = A^{-1}y.$$

2. Overdetermined system

$$Ax = y$$
$$A^T A x = A^T y$$
$$\underbrace{(A^T A)^{-1} A^T A}_{=I} x = (A^T A)^{-1} A^T y$$
$$x = \underbrace{(A^T A)^{-1} A^T}_{=:A^+} y$$

3. Under determined system

$$Ax = y \implies x = \underbrace{A^T (A A^T)^{-1}}_{=:A^+} y.$$

# (Moore-Penrose) Pseudoinverse

1. Well determined system

$$Ax = b \implies x = A^{-1}y.$$

2. Overdetermined system

$$Ax = y$$
$$A^T A x = A^T y$$
$$\underbrace{(A^T A)^{-1} A^T A}_{=I} x = (A^T A)^{-1} A^T y$$
$$x = \underbrace{(A^T A)^{-1} A^T}_{=:A^+} y$$

3. Under determined system

$$Ax = y \implies x = \underbrace{A^T (AA^T)^{-1}}_{=:A^+} y.$$

Proof. (Another proof based on pseudoinverse)

$$A = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2} , \qquad x = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1}, \qquad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{1 \times n}$$

$$A^T A = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix}$$

$$(A^T A)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix}$$

Proof. (Another proof based on pseudoinverse)

$$A = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2}, \qquad x = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1}, \qquad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{1 \times n}$$

$$A^T A = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^{n} x_i \\ \sum_{i=1}^{n} x_i & \sum_{i=1}^{n} x_i^2 \end{pmatrix}$$

$$(A^T A)^{-1} = \frac{1}{n \sum_{i=1}^{n} x_i^2 - \left( \sum_{i=1}^{n} x_i \right)^2} \begin{pmatrix} \sum_{i=1}^{n} x_i^2 & -\sum_{i=1}^{n} x_i \\ -\sum_{i=1}^{n} x_i & n \end{pmatrix}$$

Proof. (Another proof based on pseudoinverse)

$$A = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2}, \qquad x = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1}, \qquad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{1 \times n}$$

$$A^T A = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & \sum_{i=1}^{n} x_i \\ \sum_{i=1}^{n} x_i & \sum_{i=1}^{n} x_i^2 \end{pmatrix}$$

$$(A^T A)^{-1} = \frac{1}{n \sum_{i=1}^{n} x_i^2 - \left( \sum_{i=1}^{n} x_i \right)^2} \begin{pmatrix} \sum_{i=1}^{n} x_i^2 & -\sum_{i=1}^{n} x_i \\ -\sum_{i=1}^{n} x_i & n \end{pmatrix}$$

$$A^T y = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = x = (A^T A)^{-1} A^T y$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$= \begin{pmatrix} \dfrac{\left( \sum_{i=1}^n x_i^2 \right) \left( \sum_{i=1}^n y_i \right) - \left( \sum_{i=1}^n x_i \right) \left( \sum_{i=1}^n x_i y_i \right)}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \\[3ex] \dfrac{n \sum_{i=1}^n x_i y_i - \left( \sum_{i=1}^n x_i \right) \left( \sum_{i=1}^n y_i \right)}{n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2} \end{pmatrix}$$

$$A^T y = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$\begin{pmatrix} a \\ b \end{pmatrix} = x = (A^T A)^{-1} A^T y$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$= \begin{pmatrix} \dfrac{\left(\sum_{i=1}^n x_i^2\right)\left(\sum_{i=1}^n y_i\right) - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n x_i y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \\ \\ \dfrac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2} \end{pmatrix}$$

$$b = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i\right) \left(\sum_{i=1}^n y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}.$$

$$a = \frac{\left(\sum_{i=1}^n x_i^2\right) \left(\sum_{i=1}^n y_i\right) - \left(\sum_{i=1}^n x_i\right) \left(\sum_{i=1}^n x_i y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$= \frac{\left(\sum_{i=1}^n x_i^2\right) \left(\sum_{i=1}^n y_i\right) - \left(\sum_{i=1}^n x_i\right) \left[\left(\sum_{i=1}^n x_i y_i\right) - \frac{1}{n} \left(\sum_{i=1}^n x_i\right) \left(\sum_{i=1}^n y_i\right)\right]}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$- \frac{\frac{1}{n} \left(\sum_{i=1}^n x_i\right)^2 \left(\sum_{i=1}^n y_i\right)}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i\right)^2}$$

$$= \frac{1}{n} \sum_{i=1}^n y_i - b \frac{1}{n} \sum_{i=1}^n x_i = \bar{y} - b\bar{x}.$$

$\square$

# A probabilistic view ...

Def. The function $f(X)$ for which

$$\mathbb{E}\left[(Y - f(X))^2\right]$$

is minimized is called the **regression curve of $Y$ on $X$**.

Thm. Let $(X, Y)$ be two random variables such that $\mathrm{Var}(X)$ and $\mathrm{Var}(Y)$ both exist. Then the regression cure of $Y$ on $X$ is given (for all $x$) by

$$f(x) = \mathbb{E}\left[Y | X = x\right].$$

# A probabilistic view ...

Def. The function $f(X)$ for which

$$\mathbb{E}\left[(Y - f(X))^2\right]$$

is minimized is called the **regression curve of $Y$ on $X$**.

Thm. Let $(X, Y)$ be two random variables such that $\text{Var}(X)$ and $\text{Var}(Y)$ both exist. Then the regression cure of $Y$ on $X$ is given (for all $x$) by

$$f(x) = \mathbb{E}\left[Y | X = x\right].$$

Proof. Let $f(x) = \mathbb{E}\left[Y|X=x\right]$ and let $\phi(x)$ be a general function. Then

$$\mathbb{E}\left[(Y-\phi(X))^2\right] = \mathbb{E}\left[([Y-f(X)] + [f(x)-\phi(x)])^2\right]$$
$$= \mathbb{E}\left[(Y-f(X))^2\right] + \mathbb{E}\left[(f(x)-\phi(X))^2\right]$$
$$+ \mathbb{E}\left[(Y-f(X))(f(x)-\phi(X))\right].$$

Let $\psi(x)$ be either $f(x)$ or $\phi(x)$. We claim that

$$\mathbb{E}\left[(Y-f(X))\psi(X)\right] = 0.$$

Indeed,

$$\mathbb{E}[Y\psi(X)] = \iint_{\mathbb{R}^2} f_{X,Y}(x,y) y \psi(x) \mathrm{d}y \mathrm{d}x$$
$$= \int_{\mathbb{R}} \mathrm{d}x\, \psi(x) f_X(x) \underbrace{\int_{\mathbb{R}} \mathrm{d}y \frac{f_{X,Y}(x,y)}{f_X(x)} y}_{= \mathbb{E}[Y|X=x]}$$
$$= \mathbb{E}[f(X)\psi(X)].$$

Hence,

$$\mathbb{E}\left[(Y-\phi(X))^2\right] = \mathbb{E}\left[(Y-f(X))^2\right] + \mathbb{E}\left[(f(x)-\phi(X))^2\right]$$

which is minimized when $\phi(x) = f(x)$.

18

Proof. Let $f(x) = \mathbb{E}[Y|X = x]$ and let $\phi(x)$ be a general function. Then

$$\mathbb{E}\left[(Y - \phi(X))^2\right] = \mathbb{E}\left[([Y - f(X)] + [f(x) - \phi(x)])^2\right]$$
$$= \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right]$$
$$+ \mathbb{E}\left[(Y - f(X))(f(x) - \phi(X))\right].$$

Let $\psi(x)$ be either $f(x)$ or $\phi(x)$. We claim that

$$\mathbb{E}\left[(Y - f(X))\psi(X)\right] = 0.$$

Indeed,

$$\mathbb{E}[Y\psi(X)] = \iint_{\mathbb{R}^2} f_{X,Y}(x,y) y \psi(x) \mathrm{d}y \mathrm{d}x$$
$$= \int_{\mathbb{R}} \mathrm{d}x \psi(x) f_X(x) \underbrace{\int_{\mathbb{R}} \mathrm{d}y \frac{f_{X,Y}(x,y)}{f_X(x)} y}_{= \mathbb{E}[Y|X = x]}$$
$$= \mathbb{E}[f(X)\psi(X)].$$

Hence,

$$\mathbb{E}\left[(Y - \phi(X))^2\right] = \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right]$$

which is minimized when $\phi(x) = f(x)$.

18

Proof. Let $f(x) = \mathbb{E}[Y|X = x]$ and let $\phi(x)$ be a general function. Then

$$
\begin{aligned}
\mathbb{E}\left[(Y - \phi(X))^2\right] =&\, \mathbb{E}\left[([Y - f(X)] + [f(x) - \phi(x)])^2\right] \\
=&\, \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right] \\
&+ \mathbb{E}\left[(Y - f(X))(f(x) - \phi(X))\right].
\end{aligned}
$$

Let $\psi(x)$ be either $f(x)$ or $\phi(x)$. We claim that

$$
\mathbb{E}\left[(Y - f(X))\psi(X)\right] = 0.
$$

Indeed,

$$
\begin{aligned}
\mathbb{E}[Y\psi(X)] &= \iint_{\mathbb{R}^2} f_{X,Y}(x, y)y\psi(x)\mathrm{d}y\mathrm{d}x \\
&= \int_{\mathbb{R}} \mathrm{d}x\,\psi(x)f_X(x) \underbrace{\int_{\mathbb{R}} \mathrm{d}y \frac{f_{X,Y}(x, y)}{f_X(x)} y}_{= \mathbb{E}[Y|X = x]} \\
&= \mathbb{E}[f(X)\psi(X)].
\end{aligned}
$$

Hence,

$$
\mathbb{E}\left[(Y - \phi(X))^2\right] = \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right]
$$

which is minimized when $\phi(x) = f(x)$.

Proof. Let $f(x) = \mathbb{E}\left[Y | X = x\right]$ and let $\phi(x)$ be a general function. Then

$$\mathbb{E}\left[(Y - \phi(X))^2\right] = \mathbb{E}\left[([Y - f(X)] + [f(x) - \phi(x)])^2\right]$$
$$= \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right]$$
$$+ \mathbb{E}\left[(Y - f(X))(f(x) - \phi(X))\right].$$

Let $\psi(x)$ be either $f(x)$ or $\phi(x)$. We claim that

$$\mathbb{E}\left[(Y - f(X))\psi(X)\right] = 0.$$

Indeed,

$$\mathbb{E}[Y\psi(X)] = \iint_{\mathbb{R}^2} f_{X,Y}(x, y) y \psi(x) \mathrm{d}y \mathrm{d}x$$
$$= \int_{\mathbb{R}} \mathrm{d}x \psi(x) f_X(x) \underbrace{\int_{\mathbb{R}} \mathrm{d}y \frac{f_{X,Y}(x, y)}{f_X(x)} y}_{= \mathbb{E}[Y|X = x]}$$
$$= \mathbb{E}[f(X)\psi(X)].$$

Hence,

$$\mathbb{E}\left[(Y - \phi(X))^2\right] = \mathbb{E}\left[(Y - f(X))^2\right] + \mathbb{E}\left[(f(x) - \phi(X))^2\right]$$

which is minimized when $\phi(x) = f(x)$. $\qquad\square$

If one imposes that $f(x) = a + bx$, then

The following squared error:

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

is minimized at

$$b = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \quad \text{and} \quad a = \mathbb{E}[Y] - b\mathbb{E}[X]$$

with the mean squared error

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right] = \left(1 - \rho_{XY}^2\right)\sigma_Y^2.$$

If one imposes that $f(x) = a + bx$, then

Thm. The following squared error:

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

is minimized at

$$b = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \quad \text{and} \quad a = \mathbb{E}[Y] - b\mathbb{E}[X]$$

with the mean squared error

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right] = \left(1 - \rho_{XY}^2\right)\sigma_Y^2.$$

If one imposes that $f(x) = a + bx$, then

Thm. The following squared error:

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

is minimized at

$$b = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \quad \text{and} \quad a = \mathbb{E}[Y] - b\mathbb{E}[X]$$

with the mean squared error

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right] = \left(1 - \rho_{XY}^2\right)\sigma_Y^2.$$

If one imposes that $f(x) = a + bx$, then

Thm. The following squared error:

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

is minimized at

$$b = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \quad \text{and} \quad a = \mathbb{E}[Y] - b\mathbb{E}[X]$$

with the mean squared error

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right] = \left(1 - \rho_{XY}^2\right)\sigma_Y^2.$$

Proof.

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$
$$= \mathbb{E}\left[\left\{[Y - \mathbb{E}(Y)] - b[X - \mathbb{E}(X)] - [a - \mathbb{E}[Y] + b\mathbb{E}(X)]\right\}^2\right]$$

$$\|$$

$$\mathbb{E}\left[[Y - \mathbb{E}(Y)]^2\right] \qquad \qquad \text{Var}(Y)$$

$$+ b^2\mathbb{E}\left[[X - \mathbb{E}(X)]^2\right] \qquad \qquad + b^2\text{Var}(X)$$

$$+ \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 \qquad = \qquad + \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2$$

$$-2b\mathbb{E}\left[[Y - \mathbb{E}(Y)][X - \mathbb{E}(X)]\right] \qquad \qquad -2b\,\text{Cov}(X, Y)$$

$$-2\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[Y - \mathbb{E}(Y)\right] \qquad \qquad +0$$

$$+2b\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[X - \mathbb{E}(X)\right] \qquad \qquad +0$$

Proof.

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$
$$= \mathbb{E}\left[\left\{[Y - \mathbb{E}(Y)] - b[X - \mathbb{E}(X)] - [a - \mathbb{E}[Y] + b\mathbb{E}(X)]\right\}^2\right]$$

$$\|$$

$$\mathbb{E}\left[[Y - \mathbb{E}(Y)]^2\right] \qquad\qquad \text{Var}(Y)$$
$$+ b^2\mathbb{E}\left[[X - \mathbb{E}(X)]^2\right] \qquad\qquad + b^2\text{Var}(X)$$
$$+ \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 \qquad = \qquad + \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2$$
$$-2b\mathbb{E}\left[[Y - \mathbb{E}(Y)][X - \mathbb{E}(X)]\right] \qquad\qquad -2b\,\text{Cov}(X, Y)$$
$$-2\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[Y - \mathbb{E}(Y)\right] \qquad\qquad +0$$
$$+2b\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[X - \mathbb{E}(X)\right] \qquad\qquad +0$$

Proof.

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$
$$=\mathbb{E}\left[\left\{[Y - \mathbb{E}(Y)] - b[X - \mathbb{E}(X)] - [a - \mathbb{E}[Y] + b\mathbb{E}(X)]\right\}^2\right]$$

$$\|$$

$$\mathbb{E}\left[[Y - \mathbb{E}(Y)]^2\right] \qquad \text{Var}(Y)$$
$$+b^2\mathbb{E}\left[[X - \mathbb{E}(X)]^2\right] \qquad +b^2\text{Var}(X)$$
$$+\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 \qquad = \qquad +\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2$$
$$-2b\mathbb{E}\left[[Y - \mathbb{E}(Y)][X - \mathbb{E}(X)]\right] \qquad -2b\,\text{Cov}(X, Y)$$
$$-2\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[Y - \mathbb{E}(Y)\right] \qquad +0$$
$$+2b\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[X - \mathbb{E}(X)\right] \qquad +0$$

Proof.

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

$$=\mathbb{E}\left[\left\{[Y - \mathbb{E}(Y)] - b[X - \mathbb{E}(X)] - [a - \mathbb{E}[Y] + b\mathbb{E}(X)]\right\}^2\right]$$

$$||$$

| | |
|---|---|
| $\mathbb{E}\left[[Y - \mathbb{E}(Y)]^2\right]$ | $\text{Var}(Y)$ |
| $+b^2\mathbb{E}\left[[X - \mathbb{E}(X)]^2\right]$ | $+b^2\text{Var}(X)$ |
| $+\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2$ | $+\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2$ |
| $-2b\mathbb{E}\left[[Y - \mathbb{E}(Y)][X - \mathbb{E}(X)]\right]$ | $-2b\,\text{Cov}(X, Y)$ |
| $-2\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[Y - \mathbb{E}(Y)\right]$ | $+0$ |
| $+2b\left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]\mathbb{E}\left[X - \mathbb{E}(X)\right]$ | $+0$ |

(the middle column between the two above is joined by the "=" sign)

$$\Downarrow$$

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

$$||$$

$$\text{Var}(Y) + b^2\text{Var}(X) + \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 - 2b\,\text{Cov}(X, Y)$$

The best $a$, called $a^*$, should be such that

$$\left[a^* - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 = 0 \quad \iff \quad a^* = \mathbb{E}[Y] - b\mathbb{E}[X]$$

$$\Downarrow$$

$$\mathbb{E}\left[\{Y - (a + bX)\}^2\right]$$

$$||$$

$$\text{Var}(Y) + b^2\text{Var}(X) + \left[a - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 - 2b\,\text{Cov}(X, Y)$$

The best $a$, called $a^*$, should be such that

$$\left[a^* - \mathbb{E}[Y] + b\mathbb{E}(X)\right]^2 = 0 \quad \Longleftrightarrow \quad a^* = \mathbb{E}[Y] - b\mathbb{E}[X]$$

$$\Downarrow$$

$$\mathbb{E}\left[\{Y - (a^* + bX)\}^2\right]$$

$$||$$

$$\text{Var}(Y) + b^2\text{Var}(X) - 2b\,\text{Cov}(X, Y)$$

$$||$$

$$\sigma_Y^2 + b^2\sigma_X^2 - 2b\rho_{XY}\sigma_X\sigma_Y$$

$$||$$

$$\left(1 - \rho_{XY}^2\right)\sigma_Y^2 + \left(b\sigma_X - \rho_{XY}\sigma_Y\right)^2$$

The best $b$, called $b^*$, should be

$$\left(b^*\sigma_X - \rho_{XY}\sigma_Y\right)^2 = 0 \quad \iff \quad b^* = \rho_{XY}\frac{\sigma_Y}{\sigma_X}$$

$$\Downarrow$$

$$\mathbb{E}\left[\{Y - (a^* + bX)\}^2\right]$$

$$||$$

$$\text{Var}(Y) + b^2\text{Var}(X) - 2b\,\text{Cov}(X, Y)$$

$$||$$

$$\sigma_Y^2 + b^2\sigma_X^2 - 2b\rho_{XY}\sigma_X\sigma_Y$$

$$||$$

$$\left(1 - \rho_{XY}^2\right)\sigma_Y^2 + \left(b\sigma_X - \rho_{XY}\sigma_Y\right)^2$$

The best $b$, called $b^*$, should be

$$(b^*\sigma_X - \rho_{XY}\sigma_Y)^2 = 0 \quad \Longleftrightarrow \quad b^* = \rho_{XY}\frac{\sigma_Y}{\sigma_X}$$

$$\Downarrow$$
$$\mathbb{E}\left[\{Y - (a^* + b^*X)\}^2\right]$$
$$\|$$
$$\left(1 - \rho_{XY}^2\right)\sigma_Y^2$$

with

$$b^* = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \qquad \text{and} \qquad a^* = \mathbb{E}[Y] - b\mathbb{E}[X]$$

$\square$

$$\Downarrow$$

$$\mathbb{E}\left[\{Y - (a^* + b^*X)\}^2\right]$$

$$||$$

$$\left(1 - \rho_{XY}^2\right)\sigma_Y^2$$

with

$$b^* = \rho_{XY}\frac{\sigma_Y}{\sigma_X} = \frac{\sigma_{XY}}{\sigma_X^2} \qquad \text{and} \qquad a^* = \mathbb{E}[Y] - b\mathbb{E}[X]$$

$\square$

Remark In practice, we have data $(x_1, y_1), \cdots, (x_n, y_n)$ instead of the joint law of $(X, Y)$

$$\Downarrow$$

Replace

$$\mu_X, \mu_Y, \sigma_X^2, \sigma_Y^2, \rho_{XY}, \sigma_{XY}$$

by their maximum likelihood estimates

$$\bar{x}, \bar{y}, \hat{\sigma}_X^2, \hat{\sigma}_Y^2, r_{XY}, \hat{\sigma}_{XY}$$

1. $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$, $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \frac{n \sum_{i=1}^{n} x_i^2 - \left( \sum_{i=1}^{n} x_i \right)^2}{n^2}$

   $\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \frac{n \sum_{i=1}^{n} y_i^2 - \left( \sum_{i=1}^{n} y_i \right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^{n} x_i y_i - \bar{x}\bar{y}$

   $= \frac{n \sum_{i=1}^{n} x_i y_i - \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right)}{n^2}$

4. $r_{XY} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$b = r_{XY} \frac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}$$

1. $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$, $\bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \frac{n \sum_{i=1}^{n} x_i^2 - \left(\sum_{i=1}^{n} x_i\right)^2}{n^2}$

$\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \frac{n \sum_{i=1}^{n} y_i^2 - \left(\sum_{i=1}^{n} y_i\right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^{n} x_i y_i - \bar{x}\bar{y}$

$= \frac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)}{n^2}$

4. $r_{XY} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$b = r_{XY} \frac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}$$

1. $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \frac{n \sum_{i=1}^{n} x_i^2 - \left( \sum_{i=1}^{n} x_i \right)^2}{n^2}$

   $\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \frac{n \sum_{i=1}^{n} y_i^2 - \left( \sum_{i=1}^{n} y_i \right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^{n} x_i y_i - \bar{x} \bar{y}$

   $= \frac{n \sum_{i=1}^{n} x_i y_i - \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} y_i \right)}{n^2}$

4. $r_{XY} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$b = r_{XY} \frac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}$$

1. $\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i, \bar{y} = \frac{1}{n}\sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2 = \frac{1}{n}\sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \frac{n\sum_{i=1}^{n} x_i^2 - \left(\sum_{i=1}^{n} x_i\right)^2}{n^2}$

   $\hat{\sigma}_Y^2 = \frac{1}{n}\sum_{i=1}^{n}(y_i - \bar{y})^2 = \frac{1}{n}\sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \frac{n\sum_{i=1}^{n} y_i^2 - \left(\sum_{i=1}^{n} y_i\right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n}\sum_{i=1}^{n} x_i y_i - \bar{x}\bar{y}$

   $= \frac{n\sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right)\left(\sum_{i=1}^{n} y_i\right)}{n^2}$

4. $r_{XY} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$b = r_{XY}\frac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}$$

1. $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \frac{n \sum_{i=1}^{n} x_i^2 - \left(\sum_{i=1}^{n} x_i\right)^2}{n^2}$

   $\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \frac{n \sum_{i=1}^{n} y_i^2 - \left(\sum_{i=1}^{n} y_i\right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}) (y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^{n} x_i y_i - \bar{x}\bar{y}$

   $= \frac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right) \left(\sum_{i=1}^{n} y_i\right)}{n^2}$

4. $r_{XY} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$b = r_{XY} \frac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \frac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}$$

1. $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^{n} y_i$

2. $\hat{\sigma}_X^2 = \dfrac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 = \dfrac{1}{n} \sum_{i=1}^{n} x_i^2 - \bar{x}^2 = \dfrac{n \sum_{i=1}^{n} x_i^2 - \left(\sum_{i=1}^{n} x_i\right)^2}{n^2}$

   $\hat{\sigma}_Y^2 = \dfrac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \dfrac{1}{n} \sum_{i=1}^{n} y_i^2 - \bar{y}^2 = \dfrac{n \sum_{i=1}^{n} y_i^2 - \left(\sum_{i=1}^{n} y_i\right)^2}{n^2}$

3. $\hat{\sigma}_{XY} = \dfrac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y}) = \dfrac{1}{n} \sum_{i=1}^{n} x_i y_i - \bar{x}\bar{y}$

   $\qquad = \dfrac{n \sum_{i=1}^{n} x_i y_i - \left(\sum_{i=1}^{n} x_i\right) \left(\sum_{i=1}^{n} y_i\right)}{n^2}$

4. $r_{XY} = \dfrac{\hat{\sigma}_{XY}}{\hat{\sigma}_X \hat{\sigma}_Y}$

$$\Downarrow$$

$$\boxed{b = r_{XY} \dfrac{\hat{\sigma}_Y}{\hat{\sigma}_X} = \dfrac{\hat{\sigma}_{XY}}{\hat{\sigma}_X^2}, \qquad a = \bar{y} - b\bar{x}}$$

Maximum likelihood estimates

$$\hat{\sigma}_X^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

$$\hat{\sigma}_Y^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2$$

$$\hat{\sigma}_{XY} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

Sample (co)variances

$$s_X^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^{n} (y_i - \bar{y})^2$$

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

E.g. 1 Producing air conditioners. $x =$ rough weight of a rod. $y =$ finished weight. Find the best linear approximation of $xy$-relationship. Predict the weight when $x = 2.71$

E.g. 1 Producing air conditioners. $x$ = rough weight of a rod. $y$ = finished weight. Find the best linear approximation of $xy$-relationship. Predict the weight when $x = 2.71$
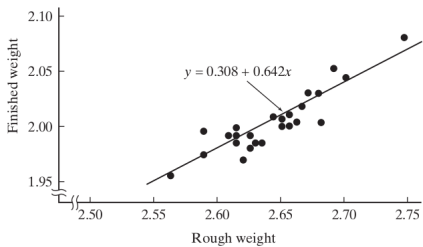
**Table 11.2.1**

| Rod Number | Rough Weight, $x$ | Finished Weight, $y$ | Rod Number | Rough Weight, $x$ | Finished Weight, $y$ |
|---|---|---|---|---|---|
| 1 | 2.745 | 2.080 | 14 | 2.635 | 1.990 |
| 2 | 2.700 | 2.045 | 15 | 2.630 | 1.990 |
| 3 | 2.690 | 2.050 | 16 | 2.625 | 1.995 |
| 4 | 2.680 | 2.005 | 17 | 2.625 | 1.985 |
| 5 | 2.675 | 2.035 | 18 | 2.620 | 1.970 |
| 6 | 2.670 | 2.035 | 19 | 2.615 | 1.985 |
| 7 | 2.665 | 2.020 | 20 | 2.615 | 1.990 |
| 8 | 2.660 | 2.005 | 21 | 2.615 | 1.995 |
| 9 | 2.655 | 2.010 | 22 | 2.610 | 1.990 |
| 10 | 2.655 | 2.000 | 23 | 2.590 | 1.975 |
| 11 | 2.650 | 2.000 | 24 | 2.590 | 1.995 |
| 12 | 2.650 | 2.005 | 25 | 2.565 | 1.955 |
| 13 | 2.645 | 2.015 | | | |

...

...  □

**Sol.** ...



...

Def. Let *a* and *b* be the least squares coefficients with the sample
$(x_1, y_1), \cdots, (x_n, y_n)$.

$\hat{y} = a + bx$: **predicted value** of y

$y_i - \hat{y}_i = y_i - (a + bx_i)$: *i***th residual**

Rem. Use the resigual plots to assessing the model.

Def. Let *a* and *b* be the least squares coefficients with the sample $(x_1, y_1), \cdots, (x_n, y_n)$.

$\hat{y} = a + bx$: **predicted value** of y

$y_i - \hat{y}_i = y_i - (a + bx_i)$: *i*th **residual**
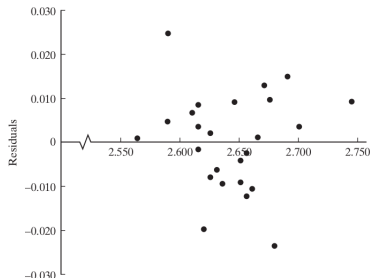
Rem. Use the resigual plots to assessing the model.

Def. Let *a* and *b* be the least squares coefficients with the sample $(x_1, y_1), \cdots, (x_n, y_n)$.

$\hat{y} = a + bx$: **predicted value** of y

$y_i - \hat{y}_i = y_i - (a + bx_i)$: *i***th residual**

Rem. Use the resigual plots to assessing the model.

Def. Let *a* and *b* be the least squares coefficients with the sample
$(x_1, y_1), \cdots, (x_n, y_n)$.

$\hat{y} = a + bx$: **predicted value** of y

$y_i - \hat{y}_i = y_i - (a + bx_i)$: *i***th residual**

Rem. Use the resigual plots to assessing the model.

Here are the residues and their plots:

E.g. 1' Here are the residues and their plots:

**Table 11.2.2**

| $x_i$ | $y_i$ | $\hat{y}_i$ | $y_i - \hat{y}_i$ |
|-------|-------|-------------|-------------------|
| 2.745 | 2.080 | 2.070 | 0.010 |
| 2.700 | 2.045 | 2.041 | 0.004 |
| 2.690 | 2.050 | 2.035 | 0.015 |
| 2.680 | 2.005 | 2.029 | −0.024 |
| 2.675 | 2.035 | 2.025 | 0.010 |
| 2.670 | 2.035 | 2.022 | 0.013 |
| 2.665 | 2.020 | 2.019 | 0.001 |
| 2.660 | 2.005 | 2.016 | −0.011 |
| 2.655 | 2.010 | 2.013 | −0.003 |
| 2.655 | 2.000 | 2.013 | −0.013 |
| 2.650 | 2.000 | 2.009 | −0.009 |
| 2.650 | 2.005 | 2.009 | −0.004 |
| 2.645 | 2.015 | 2.006 | 0.009 |
| 2.635 | 1.990 | 2.000 | −0.010 |
| 2.630 | 1.990 | 1.996 | −0.006 |
| 2.625 | 1.995 | 1.993 | 0.002 |
| 2.625 | 1.985 | 1.993 | −0.008 |
| 2.620 | 1.970 | 1.990 | −0.020 |
| 2.615 | 1.985 | 1.987 | −0.002 |
| 2.615 | 1.990 | 1.987 | 0.003 |
| 2.615 | 1.995 | 1.987 | 0.008 |
| 2.610 | 1.990 | 1.984 | 0.006 |
| 2.590 | 1.975 | 1.971 | 0.004 |
| 2.590 | 1.995 | 1.971 | 0.024 |
| 2.565 | 1.955 | 1.955 | 0.000 |

E.g. 2 Predict the Social Security expenditures.

Does the the least squares line $y = -38.0 + 12.9x$ a good model to predict the cost in 2010 would be \$543, i.e., the case $x = 45$?

Sol.
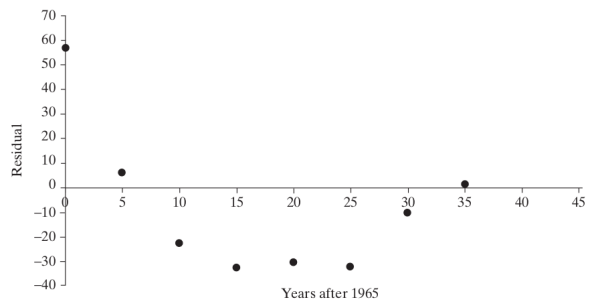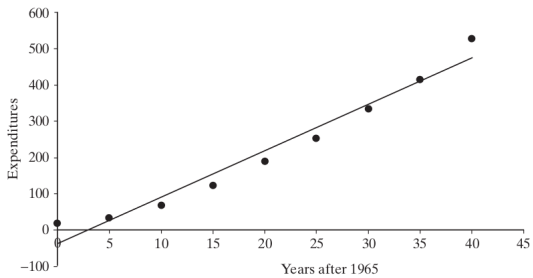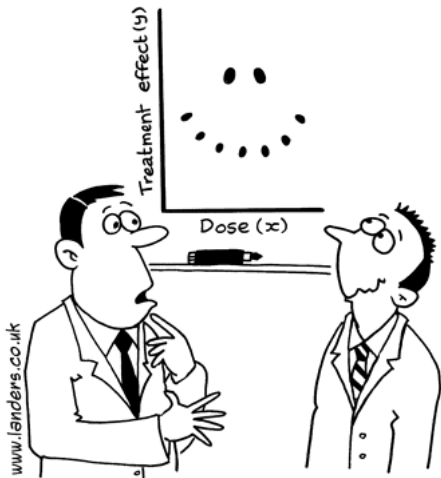
E.g. 2 Predict the Social Security expenditures.

**Table 11.2.3**

| Year | Years after 1965, $x$ | Social Security Expenditures ($ billions), $y$ |
|------|------|------|
| 1965 | 0 | 19.2 |
| 1970 | 5 | 33.1 |
| 1975 | 10 | 69.2 |
| 1980 | 15 | 123.6 |
| 1985 | 20 | 190.6 |
| 1990 | 25 | 253.1 |
| 1995 | 30 | 339.8 |
| 2000 | 35 | 415.1 |
| 2005 | 40 | 529.9 |

*Source*: www.socialsecurity.gov/history/trustfunds.html.

Does the the least squares line $y = -38.0 + 12.9x$ a good model to predict the cost in 2010 would be \$543, i.e., the case $x = 45$?
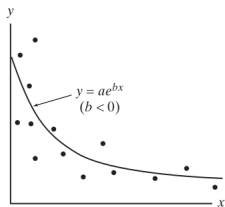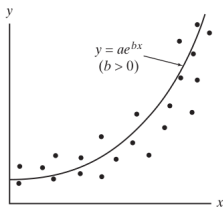
Sol.

"It's a non-linear pattern with outliers.....but for some reason I'm very happy with the data."

# Exponential Regression



$$y = ae^{bx} \quad \Longleftrightarrow \quad \ln y = \ln a + bx$$

$$b = \frac{n \sum_{i=1}^{n} x_i \ln y_i - \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} \ln y_i \right)}{n \left( \sum_{i=1}^{n} x_i^2 \right) - \left( \sum_{i=1}^{n} x_i \right)^2} \qquad \ln a = \frac{\sum_{i=1}^{n} \ln y_i - b \sum_{i=1}^{n} x_i}{n}$$

**E.g.** Moore's law:

Gordon Moore predicted in 1965 that the number of transistors per chip would double every 18 months.

Based on the real data, check:
1) Whether is the chip capacity doubling at a fixed rate?
2) Find out the rate.

Moore's law:

Gordon Moore predicted in 1965 that the number of transistors per chip would double every 18 months.

Based on the real data, check:
1) Whether is the chip capacity doubling at a fixed rate?
2) Find out the rate.

Moore's law:

Gordon Moore predicted in 1965 that the number of transistors per chip would double every 18 months.

Based on the real data, check:

1) Whether is the chip capacity doubling at a fixed rate?

2) Find out the rate.

Moore's law:

Gordon Moore predicted in 1965 that the number of transistors per chip would double every 18 months.

Based on the real data, check:
1) Whether is the chip capacity doubling at a fixed rate?
2) Find out the rate.

Moore's law:

Gordon Moore predicted in 1965 that the number of transistors per chip would double every 18 months.

Based on the real data, check:
1) Whether is the chip capacity doubling at a fixed rate?
2) Find out the rate.

**Table 11.2.5**

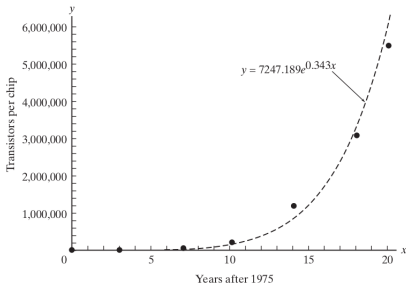| Chip | Year | Years after 1975, $x$ | Transistors per Chip, $y$ |
|---|---|---|---|
| 8080 | 1975 | 0 | 4,500 |
| 8086 | 1978 | 3 | 29,000 |
| 80286 | 1982 | 7 | 90,000 |
| 80386 | 1985 | 10 | 229,000 |
| 80486 | 1989 | 14 | 1,200,000 |
| Pentium | 1993 | 18 | 3,100,000 |
| Pentium Pro | 1995 | 20 | 5,500,000 |

*Source*: en.wikipedia.org/wiki/Transistor−count.

Sol. To check whether chip capacity doubles in a fixed rate, one needs to carry out exponential regression:

$$\implies \quad b = \cdots = 0.342810, \quad a = \cdots = e^{\ln a} = e^{8.89} = 7247.189.$$

**Sol.** To check whether chip capacity doubles in a fixed rate, one needs to carry out exponential regression:

**Table 11.2.6**

| Years after 1975, $x_i$ | $x_i^2$ | Transistors per Chip, $y_i$ | $\ln y_i$ | $x_i \cdot \ln y_i$ |
|---|---|---|---|---|
| 0 | 0 | 4,500 | 8.41183 | 0 |
| 3 | 9 | 29,000 | 10.27505 | 30.82515 |
| 7 | 49 | 90,000 | 11.40756 | 79.85292 |
| 10 | 100 | 229,000 | 12.34148 | 123.41480 |
| 14 | 196 | 1,200,000 | 13.99783 | 195.96962 |
| 18 | 324 | 3,100,000 | 14.94691 | 269.04438 |
| 20 | 400 | 5,500,000 | 15.52026 | 310.40520 |
| 72 | 1078 | | 86.90093 | 1009.51207 |

$$\implies \quad b = \cdots = 0.342810, \quad a = \cdots = e^{\ln a} = e^{8.89} = 7247.189.$$



$y = 7247.189 e^{0.343x}$

Finally, to find out the rate:

$$e^{0.343x} = e^{\ln 2 \times \frac{0.343}{\ln 2} x} = 2^{\frac{0.343}{\ln 2} x}$$

$$\frac{0.343}{\ln 2} x = 1 \quad \implies \quad x = \frac{\ln 2}{0.343} = 2.020837.$$

$\square$

# Other curvilinear models

**Table 11.2.10**

**a.** If $y = ae^{bx}$, then $\ln y$ is linear with $x$.

**b.** If $y = ax^b$, then $\log y$ is linear with $\log x$.

**c.** If $y = L/(1 + e^{a+bx})$, then $\ln\left(\dfrac{L-y}{y}\right)$ is linear with $x$.

**d.** If $y = \dfrac{1}{a+bx}$, then $\dfrac{1}{y}$ is linear with $x$.

**e.** If $y = \dfrac{x}{a+bx}$, then $\dfrac{1}{y}$ is linear with $\dfrac{1}{x}$.

**f.** If $y = 1 - e^{-x^b/a}$, then $\ln\ln\left(\dfrac{1}{1-y}\right)$ is linear with $\ln x$.

# Plan

# Chapter 11. Regression

# § 11.3 The Linear Model

Recall For any two random variables $X$ and $Y$, the regression curve of $Y$ on $X$, namely,

$$f(x) = \mathbb{E}[Y|X = x].$$

minimizes the squared error

$$\mathbb{E}\left[(Y - f(X))^2\right]$$

Difficulties The regression curve $y = \mathbb{E}[Y|x]$ is complicated and hard to obtain.

Compromise Assume that $f(x) = a + bx$ (i.e., the first order approximation)

# § 11.3 The Linear Model

Recall For any two random variables $X$ and $Y$, the regression curve of $Y$ on $X$, namely,
$$f(x) = \mathbb{E}[Y|X = x].$$
minimizes the squared error
$$\mathbb{E}\left[(Y - f(X))^2\right]$$

Difficulties The regression curve $y = \mathbb{E}[Y|x]$ is complicated and hard to obtain.

Compromise Assume that $f(x) = a + bx$ (i.e., the first order approximation)

# § 11.3 The Linear Model

Recall For any two random variables $X$ and $Y$, the regression curve of $Y$ on $X$, namely,

$$f(x) = \mathbb{E}[Y|X = x].$$

minimizes the squared error

$$\mathbb{E}\left[(Y - f(X))^2\right]$$

Difficulties The regression curve $y = \mathbb{E}[Y|x]$ is complicated and hard to obtain.

Compromise Assume that $f(x) = a + bx$   (i.e., the first order approximation)

# Simple linear model

**(Simple) linear model**:

1. $f_{Y|x}(y)$ is a normal pdf for any $x$ given.

2. The standard deviation, $\sigma$, of $Y|x$ is the same for all $x$, i.e.,

$$\sigma^2 \equiv \mathbb{E}[Y^2|x] - \mathbb{E}[Y|x]^2.$$

3. The mean of $Y|x$ is collinear, i.e.,

$$y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x.$$

4. All of the conditional distributions represnt indep. random variables.

Summary Let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \text{ are indep. and } \sim N(0, \sigma^2).$$

# Simple linear model

**(Simple) linear model**:

1. $f_{Y|x}(y)$ is a normal pdf for any $x$ given.

2. The standard deviation, $\sigma$, of $Y|x$ is the same for all $x$, i.e.,

$$\sigma^2 \equiv \mathbb{E}[Y^2|x] - \mathbb{E}[Y|x]^2.$$

3. The mean of $Y|x$ is collinear, i.e.,

$$y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x.$$

4. All of the conditional distributions represnt indep. random variables.

Summary Let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \text{ are indep. and } \sim N(0, \sigma^2).$$

# Simple linear model

**(Simple) linear model**:

1. $f_{Y|x}(y)$ is a normal pdf for any $x$ given.

2. The standard deviation, $\sigma$, of $Y|x$ is the same for all $x$, i.e.,

$$\sigma^2 \equiv \mathbb{E}[Y^2|x] - \mathbb{E}[Y|x]^2.$$

3. The mean of $Y|x$ is collinear, i.e.,

$$y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x.$$

4. All of the conditional distributions represent indep. random variables.

Summary  Let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \text{ are indep. and } \sim N(0, \sigma^2).$$

# Simple linear model

**(Simple) linear model**:

1. $f_{Y|x}(y)$ is a normal pdf for any $x$ given.

2. The standard deviation, $\sigma$, of $Y|x$ is the same for all $x$, i.e.,

$$\sigma^2 \equiv \mathbb{E}[Y^2|x] - \mathbb{E}[Y|x]^2.$$

3. The mean of $Y|x$ is collinear, i.e.,

$$y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x.$$

4. All of the conditional distributions represnt indep. random variables.

Summary Let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \text{ are indep. and } \sim N(0, \sigma^2).$$

# Simple linear model

**(Simple) linear model**:

1. $f_{Y|x}(y)$ is a normal pdf for any $x$ given.

2. The standard deviation, $\sigma$, of $Y|x$ is the same for all $x$, i.e.,

$$\sigma^2 \equiv \mathbb{E}[Y^2|x] - \mathbb{E}[Y|x]^2.$$
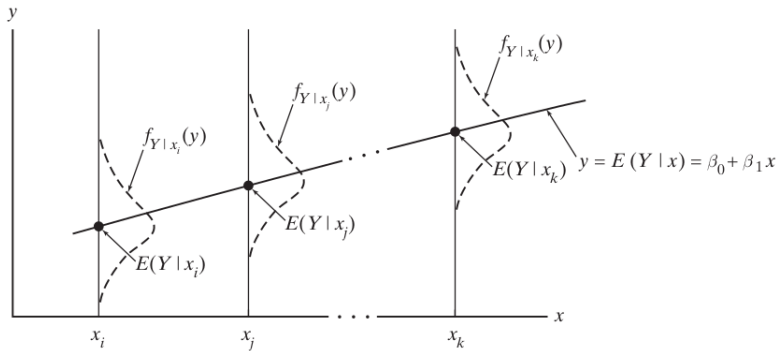
3. The mean of $Y|x$ is collinear, i.e.,

$$y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x.$$

4. All of the conditional distributions represnt indep. random variables.

Summary Let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.

$$\Updownarrow$$

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \text{ are indep. and} \sim N(0, \sigma^2).$$

# MLE for linear model

Thm. Let $(x_1, Y_1), \cdots, (x_n, Y_n)$ be a set of points satisfying the linear model, $\mathbb{E}[Y|x] = \beta_0 + \beta_1 x$.

($\Longleftrightarrow$ let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.)

The maximum likelihood estimators for $\beta_0$, $\beta_1$ and $\sigma^2$ are given by

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i Y_i - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n Y_i\right)}{n\left(\sum_{i=1}^n x_i^2\right) - \left(\sum_{i=1}^n x_i\right)^2}$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n Y_i - b \sum_{i=1}^n x_i}{n} = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left(Y_i - \widehat{Y}_i\right)^2, \qquad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

# MLE for linear model

Thm. Let $(x_1, Y_1), \cdots, (x_n, Y_n)$ be a set of points satisfying the linear model, $\mathbb{E}[Y|x] = \beta_0 + \beta_1 x$.

($\iff$ let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$ with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.)

The maximum likelihood estimators for $\beta_0$, $\beta_1$ and $\sigma^2$ are given by

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^n x_i Y_i - \left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n Y_i\right)}{n\left(\sum_{i=1}^n x_i^2\right) - \left(\sum_{i=1}^n x_i\right)^2}$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^n Y_i - b \sum_{i=1}^n x_i}{n} = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \left(Y_i - \hat{Y}_i\right)^2, \qquad \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

# MLE for linear model

Thm. Let $(x_1, Y_1), \cdots, (x_n, Y_n)$ be a set of points satisfying the linear model,
$\mathbb{E}[Y|x] = \beta_0 + \beta_1 x$.

($\Longleftrightarrow$ let $Y_1, \cdots, Y_n$ be independent r.v.'s where $Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$
with $x_i$ are known and $\beta_0$, $\beta_1$ and $\sigma^2$ are unknown.)

The maximum likelihood estimators for $\beta_0$, $\beta_1$ and $\sigma^2$ are given by

$$\hat{\beta}_1 = \frac{n \sum_{i=1}^{n} x_i Y_i - \left( \sum_{i=1}^{n} x_i \right) \left( \sum_{i=1}^{n} Y_i \right)}{n \left( \sum_{i=1}^{n} x_i^2 \right) - \left( \sum_{i=1}^{n} x_i \right)^2}$$

$$\hat{\beta}_0 = \frac{\sum_{i=1}^{n} Y_i - b \sum_{i=1}^{n} x_i}{n} = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \qquad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

**Proof** Since each $Y_i$ is assumed to be normally distributed with mean equal to $\beta_0 + \beta_1 x_i$ and variance equal to $\sigma^2$, the sample's likelihood function, $L$, is the product

$$L = \prod_{i=1}^{n} f_{Y|x_i}(y_i)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{1}{2}\left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma}\right)^2}$$

The maximum of $L$ occurs when the partial derivatives with respect to $\beta_0$, $\beta_1$, and $\sigma^2$ all vanish.

It will be easier, computationally, to differentiate $-2\ln L$, and the latter will be minimized for the same parameter values that maximize $L$. Here,

$$-2\ln L = n \cdot \ln(2\pi) + n \cdot \ln(\sigma^2) + \frac{1}{\sigma^2} \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)^2$$

Setting the three partial derivatives equal to 0 gives

$$\frac{\partial(-2\ln L)}{\partial \beta_0} = \frac{2}{\sigma^2} \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)(-1) = 0$$

$$\frac{\partial(-2\ln L)}{\partial \beta_1} = \frac{2}{\sigma^2} \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)(-x_i) = 0$$

$$\frac{\partial(-2\ln L)}{\partial \sigma^2} = \frac{n}{\sigma^2} - \frac{2}{(\sigma^2)^2} \sum_{i=1}^{n}(y_i - \beta_0 - \beta_1 x_i)^2 = 0$$

The first two equations depend only on $\beta_0$ and $\beta_1$, and the resulting solutions for $\hat{\beta}_0$ and $\hat{\beta}_1$ have the same forms that are given in the statement of the theorem. Substituting the solutions from the first two equations into the third gives the expression for $\hat{\sigma}^2$. $\qquad\square$

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\implies \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n - 2$ degrees of freedom. $\implies \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\implies \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n-2$ degrees of freedom. $\implies \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\implies \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n - 2$ degrees of freedom. $\implies \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\mathsf{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x - \bar{x})^2}$$

$$\mathsf{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\implies \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n - 2$ degrees of freedom. $\implies \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\qquad \Longrightarrow \quad \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n-2$ degrees of freedom. $\quad \Longrightarrow \quad \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Properties of linear model estimators

**Theorem:**

1. $\hat{\beta}_0$ and $\hat{\beta}_1$ are both normally distributed.
2. $\hat{\beta}_0$ and $\hat{\beta}_1$ are unbiased: $\mathbb{E}[\hat{\beta}_0] = \beta_0$ and $\mathbb{E}[\hat{\beta}_1] = \beta_1$.
3. Variances are eqal to

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{n \sum_{i=1}^{n}(x_i - \bar{x})^2} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2} \right]$$

4. $\hat{\beta}_1$, $\overline{Y}$ and $\hat{\sigma}^2$ are mutually independent. $\implies \widehat{Y}_i \perp \hat{\sigma}^2$
5. $\frac{n\hat{\sigma}^2}{\sigma^2} \sim$ Chi Square with $n-2$ degrees of freedom. $\implies \mathbb{E}[\sigma^2] = \frac{n-2}{n}\sigma^2$

Proof. ...

# Estimating $\sigma^2$

1. MLE:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2 = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$
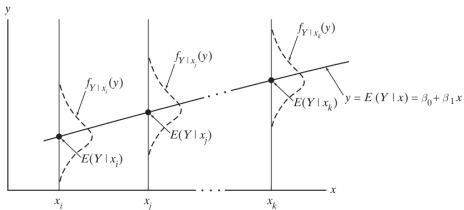
2. The unbiased estimator:

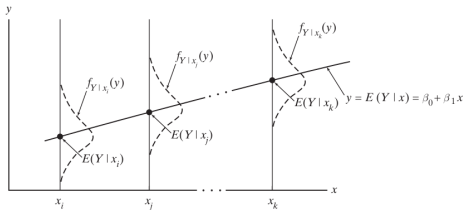$$MSE = S^2 = \frac{n}{n-2} \hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^{n} \left( Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$

# Estimating $\sigma^2$

1. MLE:

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2 = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$

2. The unbiased estimator:

$$MSE = S^2 = \frac{n}{n-2} \hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^{n} \left( Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$

Notation

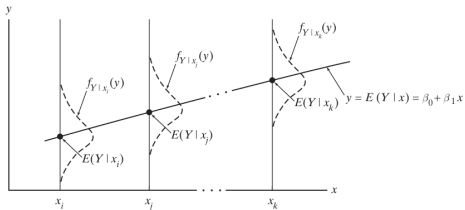| Parameter | Estimator | Estimate |
|-----------|-----------|----------|
| $\beta_1$ | $\hat{\beta}_1$ | $\beta_{1e}$ |
| $\beta_0$ | $\hat{\beta}_0$ | $\beta_{0e}$ |
| $\sigma$ | $S$ | $s$ |
| $\sigma^2$ | $S^2$ | $s^2$ |
| $\sigma^2$ | $\hat{\sigma}^2$ | $\sigma_e^2$ |
| | $\overline{Y}$ | $\bar{y}$ |
| | $\widehat{Y}_i$ | $\hat{y}_i = \beta_{0e} + \beta_{1e}x_i$ |

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations
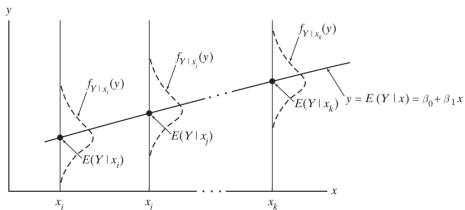
6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations

6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations
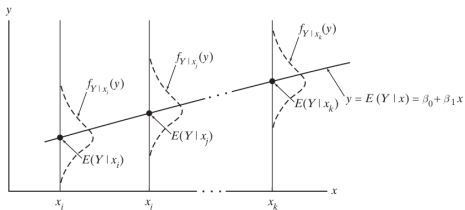
6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations
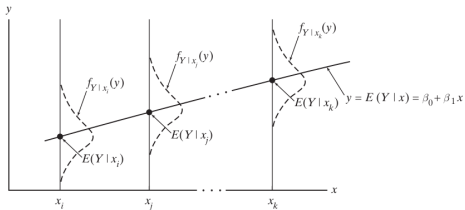
6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations

6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations

6. testing two slopes.

Drawing inferences on

1. the slope $\beta_1$

2. the intercept $\beta_0$

3. shape parameter $\sigma^2$

4. the regresion line itself
   $y = \mathbb{E}[Y|x] = \beta_0 + \beta_1 x$

5. the future observations

6. testing two slopes.

# 1. Drawing inferences on $\beta_1$

Thm. $T_{n-2} = \dfrac{\hat{\beta}_1 - \beta_1}{S \Big/ \sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}} \sim$ Student t distribution with df $= n - 2$.

1. Hypothesis test $H_0 : \beta_1 = \beta_1'$ vs. ....

2. C.I. for $\beta_1$:    $\beta_{1e} \pm t_{\alpha/2, n-2} \dfrac{s}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})}}$

# 1. Drawing inferences on $\beta_1$

Thm. $T_{n-2} = \dfrac{\hat{\beta}_1 - \beta_1}{S \Big/ \sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}} \sim$ Student t distribution with df $= n - 2$.

1. Hypothesis test $H_0 : \beta_1 = \beta_1'$ vs. ....

2. C.I. for $\beta_1$:    $\beta_{1e} \pm t_{\alpha/2, n-2} \dfrac{s}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})}}$

# 1. Drawing inferences on $\beta_1$

Thm. $T_{n-2} = \dfrac{\hat{\beta}_1 - \beta_1}{S \Big/ \sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}} \sim$ Student t distribution with df $= n - 2$.

1. Hypothesis test $H_0 : \beta_1 = \beta_1'$ vs. ....

2. C.I. for $\beta_1$: $\quad \beta_{1e} \pm t_{\alpha/2, n-2} \frac{s}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})}}$

# 2. Drawing inferences on $\beta_0$

The GLRT procedure for assessing the credibility of $H_0 : \beta_0 = \beta_{0_o}$ is based on a Student $t$ random variable with $n - 2$ degrees of freedom:

$$T_{n-2} = \frac{(\hat{\beta}_0 - \beta_{0_o}) \sqrt{n} \sqrt{\sum_{i=i}^{n} (x_i - \bar{x})^2}}{S \sqrt{\sum_{i=1}^{n} x_i^2}} = \frac{\hat{\beta}_0 - \beta_{0_o}}{\sqrt{\widehat{\text{Var}}(\hat{\beta}_0)}} \quad (11.3.6)$$

"Inverting" Equation 11.3.6 (recall the proof of Theorem 11.3.6) yields

$$\left[ \hat{\beta}_0 - t_{\alpha/2, n-2} \cdot \frac{s \sqrt{\sum_{i=1}^{n} x_i^2}}{\sqrt{n} \sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2}}, \, \hat{\beta}_0 + t_{\alpha/2, n-2} \cdot \frac{s \sqrt{\sum_{i=1}^{n} x_i^2}}{\sqrt{n} \sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2}} \right]$$

as the formula for a $100(1 - \alpha)\%$ confidence interval for $\beta_0$.

# 3. Drawing inferences on $\sigma^2$

Since $(n-2)S^2/\sigma^2$ has a $\chi^2$ pdf with $n-2$ df (if the $n$ observations satisfy the stipulations implicit in the simple linear model), it follows that

$$P\left[\chi^2_{\alpha/2,n-2} \leq \frac{(n-2)S^2}{\sigma^2} \leq \chi^2_{1-\alpha/2,n-2}\right] = 1-\alpha$$

Equivalently,

$$P\left[\frac{(n-2)S^2}{\chi^2_{1-\alpha/2,n-2}} \leq \sigma^2 \leq \frac{(n-2)S^2}{\chi^2_{\alpha/2,n-2}}\right] = 1-\alpha$$

in which case

$$\left[\frac{(n-2)s^2}{\chi^2_{1-\alpha/2,n-2}}, \frac{(n-2)s^2}{\chi^2_{\alpha/2,n-2}}\right]$$

becomes the *100(1 − α)% confidence interval for* $\sigma^2$ (recall Theorem 7.5.1). Testing $H_0 : \sigma^2 = \sigma_o^2$ is done by calculating the ratio

$$\chi^2 = \frac{(n-2)s^2}{\sigma_o^2}$$

which has a $\chi^2$ distribution with $n-2$ df when the null hypothesis is true. Except for the degrees of freedom ($n-2$ rather than $n-1$), the appropriate decision rules for one-sided and two-sided $H_1$'s are similar to those given in Theorem 7.5.2.

# 4. Drawing inference on the regression line

Intuition tells us that a reasonable point estimator for $E(Y \mid x)$ is the height of the regression line at $x$ — that is, $\hat{Y} = \hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_1 x$. By Theorem 11.3.2, the latter is unbiased:

$$E(\hat{Y}) = E(\hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_1 x) = E(\hat{\boldsymbol{\beta}}_0) + x\, E(\hat{\boldsymbol{\beta}}_1) = \beta_0 + \beta_1 x$$
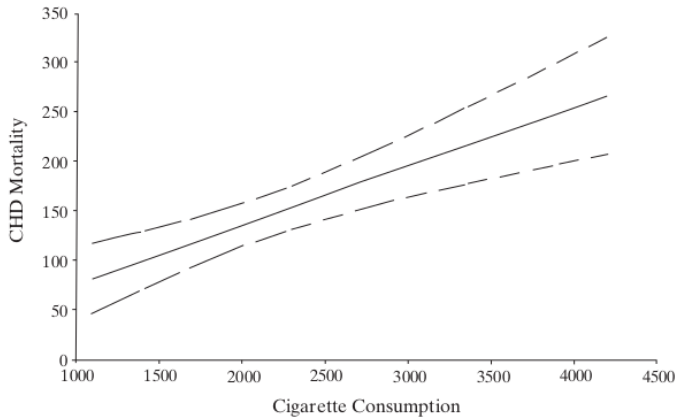
Of course, to use $\hat{Y}$ in any inference procedure requires that we know its variance. But

$$
\begin{aligned}
\mathrm{Var}(\hat{Y}) &= \mathrm{Var}(\hat{\boldsymbol{\beta}}_0 + \hat{\boldsymbol{\beta}}_1 x) = \mathrm{Var}(\bar{Y} - \hat{\boldsymbol{\beta}}_1 \bar{x} + \hat{\boldsymbol{\beta}}_1 x) \\
&= \mathrm{Var}[\bar{Y} + \hat{\boldsymbol{\beta}}_1 (x - \bar{x})] \\
&= \mathrm{Var}(\bar{Y}) + (x - \bar{x})^2 \mathrm{Var}(\hat{\boldsymbol{\beta}}_1) \quad \text{(why?)} \\
&= \frac{1}{n}\sigma^2 + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2}\sigma^2 \\
&= \sigma^2 \left[ \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2} \right]
\end{aligned}
$$

An application of Definition 7.3.3, then, allows us to construct a Student $t$ random variable based on $\hat{Y}$. Specifically,

$$T_{n-2} = \frac{\hat{Y} - (\beta_0 + \beta_1 x)}{\sigma \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n} (x_i - \bar{x})^2}}} \Bigg/ \sqrt{\frac{(n-2)S^2}{\frac{\sigma^2}{n-2}}} = \frac{\hat{Y} - (\beta_0 + \beta_1 x)}{S \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n} (x_i - \bar{x})^2}}}$$

has a Student $t$ distribution with $n - 2$ degrees of freedom. Isolating $\beta_0 + \beta_1 x = E(Y \mid x)$ in the center of the inequalities $P(-t_{\alpha/2, n-2} \leq T_{n-2} \leq t_{\alpha/2, n-2}) = 1 - \alpha$ produces a $100(1 - \alpha)\%$ confidence interval for $E(Y \mid x)$.

Figure 11.3.4

# 5. Drawing inference on future observations

Let $(x_1, Y_1), (x_2, Y_2), \ldots, (x_n, Y_n)$ be a set of $n$ points that satisfy the assumptions of the simple linear model, and let $(x, Y)$ be a hypothetical future observation, where $Y$ is independent of the $n$ $Y_i$'s. A *prediction interval* is a range of numbers that contains $Y$ with a specified probability.

Consider the difference $\hat{Y} - Y$. Clearly,

$$E(\hat{Y} - Y) = E(\hat{Y}) - E(Y) = (\beta_0 + \beta_1 x) - (\beta_0 + \beta_1 x) = 0$$

and

$$\mathrm{Var}(\hat{Y} - Y) = \mathrm{Var}(\hat{Y}) + \mathrm{Var}(Y)$$

$$= \sigma^2 \left[ \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n} (x_i - \bar{x})^2} \right] + \sigma^2$$

$$= \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum\limits_{i=1}^{n} (x_i - \bar{x})^2} \right]$$

Following exactly the same steps that were taken in the derivation of Theorem 11.3.7, a Student $t$ random variable with $n-2$ degrees of freedom can be constructed from $\hat{Y} - Y$ (using Definition 7.3.3). Inverting the equation $P(-t_{\alpha/2,n-2} \leq T_{n-2} \leq t_{\alpha/2,n-2}) = 1 - \alpha$ will then yield the prediction interval $(\hat{y} - w, \hat{y} + w)$ given in Theorem 11.3.8.

**Theorem 11.3.8**   *Let $(x_1, Y_1), (x_2, Y_2), \ldots,$ and $(x_n, Y_n)$ be a set of $n$ points that satisfy the assumptions of the simple linear model. A $100(1 - \alpha)\%$ prediction interval for $Y$ at the fixed value $x$ is given by $(\hat{y} - w, \hat{y} + w)$, where*

$$w = t_{\alpha/2,n-2} \cdot s \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\displaystyle\sum_{i=1}^{n} (x_i - \bar{x})^2}}$$

*and $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$.*   □

Does smoking contribute to coronary heat disease?

1) Test $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 > 0$ at $\alpha = 0.05$.

2) Find C.I. for $\beta_1$ with the same $\alpha$.

E.g. 1 Does smoking contribute to coronary heat disease?

| Country | Cigarette Consumption per Adult per Year, $x$ | CHD Mortality per 100,000 (ages 35–64), $y$ |
|---|---|---|
| United States | 3900 | 256.9 |
| Canada | 3350 | 211.6 |
| Australia | 3220 | 238.1 |
| New Zealand | 3220 | 211.8 |
| United Kingdom | 2790 | 194.1 |
| Switzerland | 2780 | 124.5 |
| Ireland | 2770 | 187.3 |
| Iceland | 2290 | 110.5 |
| Finland | 2160 | 233.1 |
| West Germany | 1890 | 150.3 |
| Netherlands | 1810 | 124.7 |
| Greece | 1800 | 41.2 |
| Austria | 1770 | 182.1 |
| Belgium | 1700 | 118.1 |
| Mexico | 1680 | 31.9 |
| Italy | 1510 | 114.3 |
| Denmark | 1500 | 144.9 |
| France | 1410 | 59.7 |
| Sweden | 1270 | 126.9 |
| Spain | 1200 | 43.9 |
| Norway | 1090 | 136.3 |

**Table 11.3.1**

1) Test $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 > 0$ at $\alpha = 0.05$.

2) Find C.I. for $\beta_1$ with the same $\alpha$.

**Table 11.3.1**

| Country | Cigarette Consumption per Adult per Year, $x$ | CHD Mortality per 100,000 (ages 35–64), $y$ |
|---|---|---|
| United States | 3900 | 256.9 |
| Canada | 3350 | 211.6 |
| Australia | 3220 | 238.1 |
| New Zealand | 3220 | 211.8 |
| United Kingdom | 2790 | 194.1 |
| Switzerland | 2780 | 124.5 |
| Ireland | 2770 | 187.3 |
| Iceland | 2290 | 110.5 |
| Finland | 2160 | 233.1 |
| West Germany | 1890 | 150.3 |
| Netherlands | 1810 | 124.7 |
| Greece | 1800 | 41.2 |
| Austria | 1770 | 182.1 |
| Belgium | 1700 | 118.1 |
| Mexico | 1680 | 31.9 |
| Italy | 1510 | 114.3 |
| Denmark | 1500 | 144.9 |
| France | 1410 | 59.7 |
| Sweden | 1270 | 126.9 |
| Spain | 1200 | 43.9 |
| Norway | 1090 | 136.3 |

1) Test $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 > 0$ at $\alpha = 0.05$.

2) Find C.I. for $\beta_1$ with the same $\alpha$.

Sol. http://r-statistics.co/Linear-Regression.html
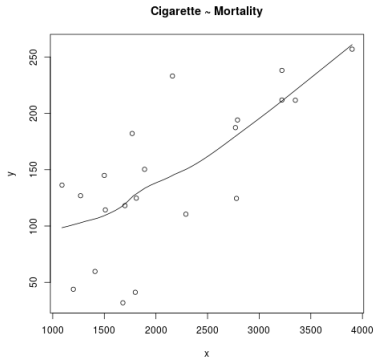
1. Let's first take of look of the data by scatter plot:

```
1  scatter.smooth(x=x, y=y, main="Cigarette ~ Mortality")
```

Suggests a linearly increasing relationship between *x* and *y*.

1. Let's first take of look of the data by scatter plot:

```
1  scatter.smooth(x=x, y=y, main="Cigarette ~ Mortality")
```

Suggests a linearly increasing relationship between *x* and *y*.
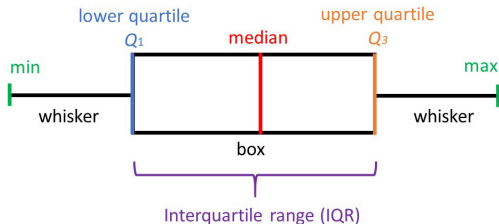
1. Let's first take of look of the data by scatter plot:

```
1  scatter.smooth(x=x, y=y, main="Cigarette ~ Mortality")
```

**Cigarette ~ Mortality**

Suggests a linearly increasing relationship between *x* and *y*.

2. Check outliers using boxplot.

Any datapoint that lies outside the $r \times$IQR is considered an outlier.

Generally, $r = 1.5$.

2. Check outliers using boxplot.



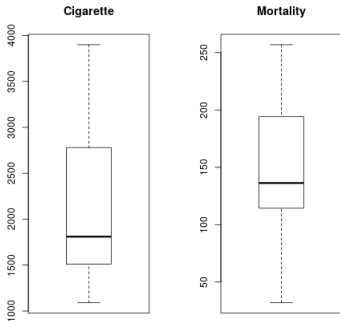Any datapoint that lies outside the $r \times$IQR is considered an outlier.

Generally, $r = 1.5$.

```
1  r <− 1.5
2  par(mfrow=c(1, 2))   # divide  graph area in 2 columns
3  boxplot(x, main="Cigarette", range=r, sub=paste("Outlier rows: ", boxplot.stats(x, coef=r)$out))
          # box plot for 'Cigarette'
4  boxplot(y, main="Mortality", range=r, sub=paste("Outlier rows: ", boxplot.stats(y, coef=r)$out))
          # box plot for 'Mortality'
```
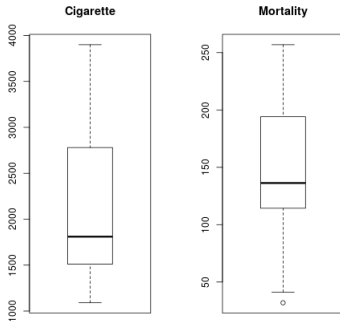


$r = 1.5$                         $r = 1$

# 3. Compute kernel density estimates

```r
library (e1071)
plot(density(x), main="Density Plot: Cigarette", ylab="Frequency",
     sub=paste("Skewness:", round(e1071::skewness(x), 2))) # density plot for 'Cigarette'
polygon(density(x), col="red")
plot(density(y), main="Density Plot: Mortality", ylab="Frequency",
     sub=paste("Skewness:", round(e1071::skewness(y), 2))) # density plot for 'Mortality'
polygon(density(y), col="red")
```

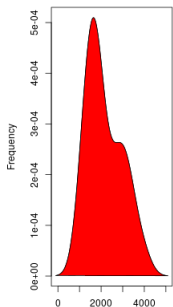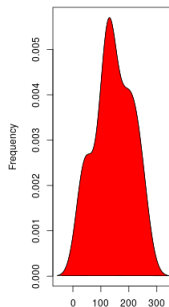# 3. Compute kernel density estimates

```
1  library (e1071)
2  plot(density(x), main="Density Plot: Cigarette", ylab="Frequency",
3       sub=paste("Skewness:", round(e1071::skewness(x), 2))) # density plot for 'Cigarette'
4  polygon(density(x), col="red")
5  plot(density(y), main="Density Plot: Mortality ", ylab="Frequency",
6       sub=paste("Skewness:", round(e1071::skewness(y), 2))) # density plot for 'Mortality '
7  polygon(density(y), col="red")
```

4. Compute correlation coeficient.

Correlation is a statistical measure with values in $[-1, 1]$ that suggests
the level of linear dependence between two variables.

A value closer to 0 suggests a weak relationship between the variables.
A low correlation $(-0.2, 0.2)$ probably suggests that much of variation of
the response variable $Y$ is unexplained by the predictor $X$, in which
case, we should probably look for better explanatory variables.

```
> cor(x,y)
[1]  0.7295154
```

4. Compute correlation coeficient.

Correlation is a statistical measure with values in $[-1, 1]$ that suggests the level of linear dependence between two variables.

A value closer to 0 suggests a weak relationship between the variables. A low correlation $(-0.2, 0.2)$ probably suggests that much of variation of the response variable $Y$ is unexplained by the predictor $X$, in which case, we should probably look for better explanatory variables.

```
1  > cor(x,y)
2  [1]  0.7295154
```

4. Compute correlation coeficient.

Correlation is a statistical measure with values in $[-1, 1]$ that suggests the level of linear dependence between two variables.

A value closer to 0 suggests a weak relationship between the variables. A low correlation $(-0.2, 0.2)$ probably suggests that much of variation of the response variable $Y$ is unexplained by the predictor $X$, in which case, we should probably look for better explanatory variables.

```
> cor(x,y)
[1]  0.7295154
```

5. Compute linear regression.

```
1  > CigMort <- data.frame("Cigarette" = x, "Mortality" = y) # Build the data frame
2  > linearMod <- lm(Mortality ~ Cigarette, data=CigMort) # linear regression
3  > print (linearMod) # Print out the result
4
5  Call:
6  lm(formula = Mortality ~ Cigarette, data = CigMort)
7
8  Coefficients:
9  (Intercept)    Cigarette
10     15.7711       0.0601
```

$$y = 15.7711 + 0.0601x$$

5. Compute linear regression.

```
1  > CigMort <- data.frame("Cigarette" = x, "Mortality" = y) # Build the data frame
2  > linearMod <- lm(Mortality ~ Cigarette, data=CigMort) # linear regression
3  > print(linearMod) # Print out the result
4
5  Call:
6  lm(formula = Mortality ~ Cigarette, data = CigMort)
7
8  Coefficients:
9  (Intercept)    Cigarette
10     15.7711       0.0601
```

$$y = 15.7711 + 0.0601x$$

5. Compute linear regression.

```
1  > CigMort <- data.frame("Cigarette" = x, "Mortality " = y)  # Build the data frame
2  > linearMod <- lm(Mortality ~ Cigarette, data=CigMort) # linear regression
3  > print(linearMod) # Print out the result
4
5  Call :
6  lm(formula = Mortality ~ Cigarette, data = CigMort)
7
8  Coefficients :
9  (Intercept)    Cigarette
10     15.7711       0.0601
```

$$y = 15.7711 + 0.0601x$$

## 6. Check statistical significance of the linear model

```
1  > summary(linearMod)
2
3  Call:
4  lm(formula = Mortality ~ Cigarette, data = CigMort)
5
6  Residuals:
7      Min     1Q  Median     3Q     Max
8  −84.835 −40.809 5.058 28.814  87.518
9
10 Coefficients:
11             Estimate Std. Error  t value Pr(>|t|)
12 (Intercept) 15.77115  29.57889   0.533  0.600085
13 Cigarette    0.06010   0.01293   4.649  0.000175 ∗∗∗
14 −−−
15 Signif. codes:  0 "∗∗∗" 0.001 "∗∗" 0.01 "∗" 0.05 "." 0.1 " " 1
16
17 Residual standard error: 46.71 on 19 degrees of freedom
18 Multiple R−squared: 0.5322, Adjusted R−squared: 0.5076
19 F−statistic: 21.62 on 1 and 19 DF, p−value: 0.0001749
```

0.1 By default, p-values are computed for $H_0 : \beta_i = 0$ vs. $H_1 : \beta_i \neq 0$, $i = 0, 1$.
0.2 The more stars by the variable's p-Value, the more significant the variable.

6. Check statistical significance of the linear model

```
1  > summary(linearMod)
2
3  Call:
4  lm(formula = Mortality ~ Cigarette, data = CigMort)
5
6  Residuals:
7      Min     1Q  Median    3Q     Max
8  −84.835 −40.809 5.058 28.814 87.518
9
10  Coefficients:
11              Estimate Std. Error  t value Pr(>|t|)
12  (Intercept) 15.77115  29.57889   0.533  0.600085
13  Cigarette    0.06010   0.01293   4.649  0.000175 ***
14  −−−
15  Signif. codes:  0 "***" 0.001 "**" 0.01 "*" 0.05 "." 0.1 " " 1
16
17  Residual standard error: 46.71 on 19 degrees of freedom
18  Multiple R−squared: 0.5322, Adjusted R−squared: 0.5076
19  F−statistic:  21.62 on 1 and 19 DF, p−value: 0.0001749
```

0.1 By default, p-values are computed for $H_0 : \beta_i = 0$ vs. $H_1 : \beta_i \neq 0$, $i = 0, 1$.
0.2 The more stars by the variable's p-Value, the more significant the variable.

# 6. Check statistical significance of the linear model

```
1  > summary(linearMod)
2
3  Call:
4  lm(formula = Mortality ~ Cigarette, data = CigMort)
5
6  Residuals:
7      Min      1Q  Median      3Q     Max
8  −84.835 −40.809  5.058  28.814  87.518
9
10 Coefficients:
11              Estimate Std. Error  t value Pr(>|t|)
12 (Intercept) 15.77115   29.57889   0.533  0.600085
13 Cigarette    0.06010    0.01293   4.649  0.000175 ∗∗∗
14 −−−
15 Signif. codes:  0 "∗∗∗" 0.001 "∗∗" 0.01 "∗" 0.05 "." 0.1 " " 1
16
17 Residual standard error: 46.71 on 19 degrees of freedom
18 Multiple R−squared: 0.5322, Adjusted R−squared: 0.5076
19 F−statistic: 21.62 on 1 and 19 DF, p−value: 0.0001749
```

0.1 By default, p-values are computed for $H_0 : \beta_i = 0$ vs. $H_1 : \beta_i \neq 0$, $i = 0, 1$.

0.2 The more stars by the variable's p-Value, the more significant the variable.

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value= 0.000175

Conclusion: reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value= 0.600

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +   coef <- summary(linearMod)$coefficient
5  +   df <- linearMod$df.residual
6  +   lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +   ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +   print(paste("95% C.I. for the slope is beta_",i,
9  +               " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s.
$H_1 : \beta_1 \neq 0$

*t*-score is 4.4649.

*p*-value= 0.000175

Conclusion: reject at
$\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s.
$H_1 : \beta_0 \neq 0$

*t*-score is 0.533.

*p*-value= 0.600

Conclusion: fail to reject at
$\alpha = 0.05$.

95% C.I. for $\beta_0$:

```r
> # 95% C.I. for slope parameter beta_1
> alpha <- 0.05
> for (i in c(1,0)) {
+    coef <- summary(linearMod)$coefficient
+    df <- linearMod$df.residual
+    lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
+    ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
+    print(paste("95% C.I. for the slope is beta_",i,
+                " is (", round(lbd,3), ",", round(ubd,3),")"))
+ }
[1] "95% C.I. for the slope is beta_ 1 is ( 0.049 , 0.071 )"
[1] "95% C.I. for the slope is beta_ 0 is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +    coef <- summary(linearMod)$coefficient
5  +    df <- linearMod$df.residual
6  +    lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +    ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +    print(paste("95% C.I. for the slope is beta_",i,
9  +                 " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s.
$H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at
$\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s.
$H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at
$\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +    coef <- summary(linearMod)$coefficient
5  +    df <- linearMod$df.residual
6  +    lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +    ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +    print(paste("95% C.I. for the slope is beta_",i,
9  +                " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s.
$H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at
$\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s.
$H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at
$\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +    coef <- summary(linearMod)$coefficient
5  +    df <- linearMod$df.residual
6  +    lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +    ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +    print(paste("95% C.I. for the slope is beta_",i,
9  +              " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at $\alpha = 0.05$.

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +    coef <- summary(linearMod)$coefficient
5  +    df <- linearMod$df.residual
6  +    lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +    ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +    print(paste("95% C.I. for the slope is beta_",i,
9  +                " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s.
$H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at
$\alpha = 0.05$.

Testing $H_0 : \beta_0 = 0$ v.s.
$H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at
$\alpha = 0.05$.

95% C.I. for $\beta_1$:

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <− 0.05
3  > for (i in c(1,0)) {
4  +   coef <− summary(linearMod)$coefficient
5  +   df <− linearMod$df.residual
6  +   lbd <− coef[i+1,1] − pt(1−alpha/2,df) ∗ coef[i+1,2]
7  +   ubd <− coef[i+1,1] + pt(1−alpha/2,df) ∗ coef[i+1,2]
8  +   print(paste("95% C.I. for the slope is beta_",i,
9  +               " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1 is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0 is ( −8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at $\alpha = 0.05$.

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

95% C.I. for $\beta_0$:

```
> # 95% C.I. for slope parameter beta_1
> alpha <- 0.05
> for (i in c(1,0)) {
+     coef <- summary(linearMod)$coefficient
+     df <- linearMod$df.residual
+     lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
+     ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
+     print(paste("95% C.I. for the slope is beta_",i,
+              " is (", round(lbd,3), ",", round(ubd,3),")"))
+ }
[1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
[1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +   coef <- summary(linearMod)$coefficient
5  +   df <- linearMod$df.residual
6  +   lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +   ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +   print(paste("95% C.I. for the slope is beta_",i,
9  +               " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1  is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0  is ( -8.753 , 40.295 )"
```

Testing $H_0 : \beta_1 = 0$ v.s. $H_1 : \beta_1 \neq 0$

$t$-score is 4.4649.

$p$-value$= 0.000175$

Conclusion: reject at $\alpha = 0.05$.

95% C.I. for $\beta_1$:

Testing $H_0 : \beta_0 = 0$ v.s. $H_1 : \beta_0 \neq 0$

$t$-score is 0.533.

$p$-value$= 0.600$

Conclusion: fail to reject at $\alpha = 0.05$.

95% C.I. for $\beta_0$:

```
1  > # 95% C.I. for slope parameter beta_1
2  > alpha <- 0.05
3  > for (i in c(1,0)) {
4  +     coef <- summary(linearMod)$coefficient
5  +     df <- linearMod$df.residual
6  +     lbd <- coef[i+1,1] - pt(1-alpha/2,df) * coef[i+1,2]
7  +     ubd <- coef[i+1,1] + pt(1-alpha/2,df) * coef[i+1,2]
8  +     print(paste("95% C.I. for the slope is beta_",i,
9  +                 " is (", round(lbd,3), ",", round(ubd,3),")"))
10 + }
11 [1] "95% C.I. for the slope is beta_ 1 is ( 0.049 , 0.071 )"
12 [1] "95% C.I. for the slope is beta_ 0 is ( -8.753 , 40.295 )"
```

7. Compute R-Squared and the adjusted R-Squared.

$$R^2 = 1 - \frac{SSE}{SST} \quad \text{and} \quad R^2_{adj} = 1 - \frac{MSE}{MST}$$

```
1  > names(summary(linearMod))
2  [1] "call"            "terms"          "residuals"      "coefficients"
3  [5] "aliased"         "sigma"          "df"             "r.squared"
4  [9] "adj.r.squared"   "fstatistic"     "cov.unscaled"
5  > summary(linearMod)$r.squared
6  [1] 0.5321927
7  > summary(linearMod)$adj.r.squared
8  [1] 0.5075712
```

The large $r^2$ or $r^2_{adj}$ the better, the more powerful or expressive is the L.M.

7. Compute R-Squared and the adjusted R-Squared.

$$R^2 = 1 - \frac{SSE}{SST} \quad \text{and} \quad R^2_{adj} = 1 - \frac{MSE}{MST}$$

```
> names(summary(linearMod))
 [1] "call"            "terms"            "residuals"       "coefficients"
 [5] "aliased"         "sigma"            "df"              "r.squared"
 [9] "adj.r.squared"   "fstatistic"       "cov.unscaled"
> summary(linearMod)$r.squared
[1] 0.5321927
> summary(linearMod)$adj.r.squared
[1] 0.5075712
```

The large $r^2$ or $r^2_{adj}$ the better, the more powerful or expressive is the L.M.

7. Compute R-Squared and the adjusted R-Squared.

$$R^2 = 1 - \frac{SSE}{SST} \quad \text{and} \quad R^2_{adj} = 1 - \frac{MSE}{MST}$$

```
1  > names(summary(linearMod))
2  [1] "call"           "terms"         "residuals"      "coefficients"
3  [5] "aliased"        "sigma"         "df"             "r.squared"
4  [9] "adj.r.squared"  "fstatistic"    "cov.unscaled"
5  > summary(linearMod)$r.squared
6  [1] 0.5321927
7  > summary(linearMod)$adj.r.squared
8  [1] 0.5075712
```

The large $r^2$ or $r^2_{adj}$ the better, the more powerful or expressive is the L.M.

7. Compute R-Squared and the adjusted R-Squared.

$$R^2 = 1 - \frac{SSE}{SST} \quad \text{and} \quad R^2_{adj} = 1 - \frac{MSE}{MST}$$

```
1  > names(summary(linearMod))
2    [1] "call"           "terms"          "residuals"     "coefficients"
3    [5] "aliased"        "sigma"          "df"            "r.squared"
4    [9] "adj.r.squared"  "fstatistic"     "cov.unscaled"
5  > summary(linearMod)$r.squared
6    [1] 0.5321927
7  > summary(linearMod)$adj.r.squared
8    [1] 0.5075712
```

The large $r^2$ or $r^2_{adj}$ the better, the more powerful or expressive is the L.M.

## 8. Residue standard error and *F*-statistic

$$\text{Residue standard error} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-q}}$$

$$F = \frac{MSR}{MSE} = \frac{SSR/(q-1)}{SSE/(n-q)} \sim \text{F-distribution } (df_1 = q-1, df_2 = n-q)$$

```
> names(summary(linearMod))
 [1] "call"            "terms"          "residuals"      "coefficients"
 [5] "aliased"         "sigma"          "df"             "r.squared"
 [9] "adj.r.squared"   "fstatistic"     "cov.unscaled"
> summary(linearMod)$sigma
[1] 46.70826
> summary(linearMod)$fstatistic
   value   numdf   dendf
21.61501 1.00000 19.00000
> f <- summary(linearMod)$fstatistic
> pf( f[1], f[2], f[3], lower=FALSE)
       value
0.0001748805
```

8. Residue standard error and *F*-statistic

$$\text{Residue standard error} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-q}}$$

$$F = \frac{MSR}{MSE} = \frac{SSR/(q-1)}{SSE/(n-q)} \sim \text{F-distribution } (df_1 = q-1, df_2 = n-q)$$

```
> names(summary(linearMod))
 [1] "call"           "terms"           "residuals"      "coefficients"
 [5] "aliased"        "sigma"           "df"             "r.squared"
 [9] "adj.r.squared" "fstatistic"     "cov.unscaled"
> summary(linearMod)$sigma
[1] 46.70826
> summary(linearMod)$fstatistic
   value    numdf   dendf
21.61501 1.00000 19.00000
> f <- summary(linearMod)$fstatistic
> pf(f[1], f[2], f[3], lower=FALSE)
       value
0.0001748805
```

8. Residue standard error and *F*-statistic

$$\text{Residue standard error} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-q}}$$

$$F = \frac{MSR}{MSE} = \frac{SSR/(q-1)}{SSE/(n-q)} \sim \text{F-distribution } (df_1 = q-1, df_2 = n-q)$$

```
> names(summary(linearMod))
 [1] "call"            "terms"          "residuals"      "coefficients"
 [5] "aliased"         "sigma"          "df"             "r.squared"
 [9] "adj.r.squared"  "fstatistic"     "cov.unscaled"
> summary(linearMod)$sigma
[1] 46.70826
> summary(linearMod)$fstatistic
   value    numdf    dendf
21.61501 1.00000 19.00000
> f <- summary(linearMod)$fstatistic
> pf( f [1], f [2], f [3], lower=FALSE)
       value
0.0001748805
```

8. Residue standard error and *F*-statistic

$$\text{Residue standard error} = \sqrt{MSE} = \sqrt{\frac{SSE}{n-q}}$$

$$F = \frac{MSR}{MSE} = \frac{SSR/(q-1)}{SSE/(n-q)} \sim \text{F-distribution } (df_1 = q - 1, df_2 = n - q)$$

```
1  > names(summary(linearMod))
2  [1] "call"              "terms"           "residuals"     "coefficients "
3  [5] "aliased"           "sigma"           "df"            "r.squared"
4  [9] "adj.r.squared" " fstatistic "       "cov.unscaled"
5  > summary(linearMod)$sigma
6  [1] 46.70826
7  > summary(linearMod)$fstatistic
8     value    numdf    dendf
9  21.61501 1.00000 19.00000
10 > f <- summary(linearMod)$fstatistic
11 > pf( f [1], f [2], f [3], lower=FALSE)
12        value
13 0.0001748805
```

9. Model selection:

Akaike's information criterion — AIC (Akaike, 1974)

$$AIC = -2\ln(\widehat{L}) + 2q$$

Bayesian information criterion — BIC (Schwarz, 1978)

$$BIC = -2\ln(\widehat{L}) + q\ln(n)$$

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
1  > AIC(linearMod)
2  [1] 224.9383
3  > BIC(linearMod)
4  [1] 228.0719
```

The lower the better!

9. Model selection:

Akaike's information criterion — AIC (Akaike, 1974)

$$AIC = -2\ln(\widehat{L}) + 2q$$

Bayesian information criterion — BIC (Schwarz, 1978)

$$BIC = -2\ln(\widehat{L}) + q\ln(n)$$

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
1  > AIC(linearMod)
2  [1] 224.9383
3  > BIC(linearMod)
4  [1] 228.0719
```

The lower the better!

9. Model selection:

| Akaike's information criterion | Bayesian information criterion |
|---|---|
| — AIC (Akaike, 1974) | — BIC (Schwarz, 1978) |

$$AIC = -2\ln(\widehat{L}) + 2q \qquad\qquad BIC = -2\ln(\widehat{L}) + q\ln(n)$$

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
1  > AIC(linearMod)
2  [1] 224.9383
3  > BIC(linearMod)
4  [1] 228.0719
```

The lower the better!

9. Model selection:

<table>
<tr><td>Akaike's information criterion<br>— AIC (Akaike, 1974)</td><td>Bayesian information criterion<br>— BIC (Schwarz, 1978)</td></tr>
<tr><td>$AIC = -2\ln(\widehat{L}) + 2q$</td><td>$BIC = -2\ln(\widehat{L}) + q\ln(n)$</td></tr>
</table>

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
1  > AIC(linearMod)
2  [1] 224.9383
3  > BIC(linearMod)
4  [1] 228.0719
```

The lower the better!

9. Model selection:

| Akaike's information criterion — AIC (Akaike, 1974) | Bayesian information criterion — BIC (Schwarz, 1978) |
|---|---|
| $AIC = -2\ln(\widehat{L}) + 2q$ | $BIC = -2\ln(\widehat{L}) + q\ln(n)$ |

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
> AIC(linearMod)
[1] 224.9383
> BIC(linearMod)
[1] 228.0719
```

The lower the better!

9. Model selection:

| Akaike's information criterion | Bayesian information criterion |
|---|---|
| — AIC (Akaike, 1974) | — BIC (Schwarz, 1978) |

$$AIC = -2\ln(\widehat{L}) + 2q \qquad BIC = -2\ln(\widehat{L}) + q\ln(n)$$

$\widehat{L}$: the maximum of likelihood.

$q$: the number of parameters in the model.

$n$: the sample size.

```
1  > AIC(linearMod)
2  [1] 224.9383
3  > BIC(linearMod)
4  [1] 228.0719
```

The lower the better!

10. Does L.M. fit our model?

| Statistic | criterion | our case |
|-----------|-----------|----------|
| $R^2$ | Higher the better (>0.7) | 0.53 |
| $R^2_{adj}$ | Higher the better | 0.51 |
| AIC | Lower the better | 225 |
| BIC | Lower the better | 228 |
| $\vdots$ | $\vdots$ | $\vdots$ |

## 11. Drawing inference on $\mathbb{E}(Y|x)$

Find 95% C.I. for $Y$ at $x = 4200$.

11. Drawing inference on $\mathbb{E}(Y|x)$

   Find 95% C.I. for $Y$ at $x = 4200$.

11. Drawing inference on $\mathbb{E}(Y|x)$

   Find 95% C.I. for $Y$ at $x = 4200$.

11. Drawing inference on $\mathbb{E}(Y|x)$

   Find 95% C.I. for $Y$ at $x = 4200$.

   Here, $n = 21$, $t_{.025, 19} = 2.0930$, $\sum_{i=1}^{21} (x_i - \bar{x})^2 = 13{,}056{,}523.81$, $s = 46.707$, $\hat{\beta}_0 =$ 15.7661, $\hat{\beta}_1 = 0.0601$, and $\bar{x} = 2148.095$. From Theorem 11.3.7, then,

   $$\hat{y} = 15.7661 + 0.0601(4200) = 268.1861$$

   $$w = 2.0930(46.707)\sqrt{\frac{1}{21} + \frac{(4200 - 2148.095)^2}{13{,}056{,}523.81}} = 59.4714$$

   and the 95% confidence interval for $E(Y|4200)$ is

   $$(268.1861 - 59.4714, \ 268.1861 - 59.4714)$$

   which rounded to two decimal places is

   $$(208.71, 327.66)$$

```
1  s <- summary(linearMod)$sigma
2  beta <- linearMod$coefficients
3  z <- seq(1000,4500,1)
4  hatY <- beta[1]+beta[2]*z
5  w <- qt(0.975,19) * s * sqrt(1/21+(z-mean(x))^2/(sum((x-mean(x))^2)))
6  matplot(z,cbind(hatY,hatY+w,hatY-w),type = c("l","l","l"),lwd=c(3,4,4))
7  points(x, y, pch = 19)
8  abline(v=4200,col = "blue",  lty  = 4)
9  abline(h=208.71,col = "blue",  lty  = 4)
10 abline(h=327.66,col = "blue",  lty  = 4)
11 text(4200,50,4200)
12 text(1200,203,208.71)
13 text(1200,331,327.66)
```

12. Drawing inference on future observations.

Find 95% prediction interval for $Y$ at $x = 4200$.

12. Drawing inference on future observations.

Find 95% prediction interval for $Y$ at $x = 4200$.

12. Drawing inference on future observations.

   Find 95% prediction interval for $Y$ at $x = 4200$.

12. Drawing inference on future observations.

Find 95% prediction interval for $Y$ at $x = 4200$.

When $x = 4200$, $\hat{y} = 268.1861$ for both intervals. From Theorem 11.3.8, the width of the 95% prediction interval for $Y$ is:

$$w = 2.0930(46.707)\sqrt{1 + \frac{1}{21} + \frac{(4200 - 2148.095)^2}{13{,}056{,}523.81}} = 114.4725$$

The 95% prediction interval, then, is

$$(268.1861 - 114.4725, \ 268.1861 + 114.4725)$$

which rounded to two decimal places is

$$(153.76, \ 382.61)$$

which makes it 92% wider than the 95% confidence interval for $E(Y|4200)$. ∎

382.61

327.66

208.71

153.76

4200

cbind(hatY, hatY + w, hatY - w, hatY + f, hatY - f)

z

```r
1  s <- summary(linearMod)$sigma
2  beta <- linearMod$coefficients
3  z <- seq(1000,4500,1)
4  hatY <- beta[1]+beta[2]*z
5  w <- qt(0.975,19) * s * sqrt(1/21+(z-mean(x))^2/(sum((x-mean(x))^2)))
6  f <- qt(0.975,19) * s * sqrt(1+1/21+(z-mean(x))^2/(sum((x-mean(x))^2)))
7  matplot(z,cbind(hatY,hatY+w,hatY-w,hatY+f,hatY-f),
8          type = c("l","l","l","l","l"),lwd=c(3,4,4,4,4))
9  points(x, y, pch = 19)
10 abline(v=4200,col = "blue", lty = 4)
11 abline(h=208.71,col = "blue", lty = 4)
12 abline(h=327.66,col = "blue", lty = 4)
13 text(4200,50,4200)
14 text(1200,208.71-5,208.71)
15 text(1200,327.66+5,327.66)
16 abline(h=153.76,col = "red", lty = 4)
17 abline(h=382.61,col = "red", lty = 4)
18 text(4200,50,4200)
19 text(1200,153.76-5,153.76)
20 text(1200,382.61+5,382.61)
```
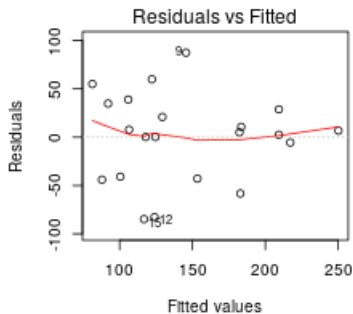
13. More about diagnozing the linear model:

```
1  # diagnostic  plots
2  layout(matrix(c(1,2,3,4),2,2))  # optional  4 graphs/page
3  plot(linearMod)
```

13. More about diagnozing the linear model:

```
1  # diagnostic  plots
2  layout(matrix(c(1,2,3,4) ,2,2) )  # optional  4 graphs/page
3  plot (linearMod)
```

E.g. 2  Find 95% C.I. for the amount of increas year-by-year in the cost of Toyota Camry sedan.

E.g. 2 Find 95% C.I. for the amount of increas year-by-year in the cost of Toyota Camry sedan.

| Table 11.3.2 | | |
|---|---|---|
| Year | Year after 2005 | Suggested Retail Price ($) |
| 2005 | 0 | 7,935 |
| 2006 | 1 | 8,495 |
| 2007 | 2 | 10,160 |
| 2008 | 3 | 10,817 |
| 2009 | 4 | 11,078 |
| 2010 | 5 | 11,967 |
| 2011 | 6 | 12,658 |
| 2012 | 7 | 13,844 |
| 2013 | 8 | 13,982 |
| 2014 | 9 | 14,629 |

Data from: kbb.com

Sol. We first find the regression:

Sol. We first find the regression:



**Figure 11.3.3**

The slope of the line, $\hat{\beta}_1$, represents the amount of increase year-by-year in the cost of an older model. Often a range of values is better than a single estimate, so a good way to provide this is using a confidence interval for the true value $\beta_1$.

Here,

$$\sqrt{\sum_{i=0}^{9}(x_i - \bar{x})^2} = \sqrt{82.5} = 9.083$$

and from Equation 11.3.5, $s^2 = \frac{1}{10-2}\left(\sum_{i=0}^{9}y_i^2 - \hat{\beta}_0\sum_{i=0}^{9}y_i - \hat{\beta}_1\sum_{i=0}^{9}x_iy_i\right)$

$$\frac{1}{8}[1,382,678,777 - (8188.67)(115,565) - (748.41)(581,786)] = 117,727.98$$

so $s = \sqrt{117,727.98} = 343.11$.

Using $t_{.025,8} = 2.3060$, the expression given in Theorem 11.3.6 reduces to

$$\left(748.41 - 2.3060\frac{343.11}{9.083},\ 748.41 + 2.3060\frac{343.11}{9.083}\right) = (\$661.30, \$835.52)$$

# 7. Testing the equality of two slopes

**Table 11.3.3**

| Date | Day no., $x(=x^*)$ | Strain $A$ pop$^n$, $y$ | Strain $B$ pop$^n$, $y^*$ |
|------|------|------|------|
| Feb 2 | 0 | 100 | 100 |
| May 13 | 100 | 250 | 203 |
| Aug 21 | 200 | 304 | 214 |
| Nov 29 | 300 | 403 | 295 |
| Mar 8 | 400 | 446 | 330 |
| Jun 16 | 500 | 482 | 324 |



**Figure 11.3.5**

Do you believe that $\beta_1 = \beta_1^*$?

Or is $\beta_1 > \beta_1^*$ statisically significantly?

**Theorem**
**11.3.9**

Let $(x_1, Y_1), (x_2, Y_2), \ldots, (x_n, Y_n)$ and $(x_1^*, Y_1^*), (x_2^*, Y_2^*), \ldots, (x_m^*, Y_m^*)$ be two indepen-
dent sets of points, each satisfying the assumptions of the simple linear model—that is,
$E(Y \mid x) = \beta_0 + \beta_1 x$ and $E(Y^* \mid x^*) = \beta_0^* + \beta_1^* x^*$.

**a.** Let

$$T = \frac{\hat{\beta}_1 - \hat{\beta}_1^* - (\beta_1 - \beta_1^*)}{S \sqrt{\dfrac{1}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2} + \dfrac{1}{\sum\limits_{i=1}^{m}(x_i^* - \bar{x}^*)^2}}}$$

where

$$S = \sqrt{\frac{\sum\limits_{i=1}^{n}[Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2 + \sum\limits_{i=1}^{m}[Y_i^* - (\hat{\beta}_0^* + \hat{\beta}_1^* x_i)]^2}{n + m - 4}}$$

Then T has a Student t distribution with $n + m - 4$ degrees of freedom.

**b.** To test $H_0 : \beta_1 = \beta_1^*$ versus $H_1 : \beta_1 \neq \beta_1^*$ at the $\alpha$ level of significance, reject $H_0$ if t
is either $(1) \leq -t_{\alpha/2, n+m-4}$ or $(2) \geq t_{\alpha/2, n+m-4}$, where

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_1^*}{s \sqrt{\dfrac{1}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2} + \dfrac{1}{\sum\limits_{i=1}^{m}(x_i^* - \bar{x}^*)^2}}}$$

(One-sided tests are defined in the usual way by replacing $\pm t_{\alpha/2, n+m-4}$ with either
$t_{\alpha, n+m-4}$ or $-t_{\alpha, n+m-4}$.)

$S^2 = SSE$ and $q = 4$.

**Theorem 11.3.9**

Let $(x_1, Y_1), (x_2, Y_2), \ldots, (x_n, Y_n)$ and $(x_1^*, Y_1^*), (x_2^*, Y_2^*), \ldots, (x_m^*, Y_m^*)$ be two independent sets of points, each satisfying the assumptions of the simple linear model—that is, $E(Y \mid x) = \beta_0 + \beta_1 x$ and $E(Y^* \mid x^*) = \beta_0^* + \beta_1^* x^*$.

**a.** Let

$$T = \frac{\hat{\beta}_1 - \hat{\beta}_1^* - (\beta_1 - \beta_1^*)}{S\sqrt{\dfrac{1}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2} + \dfrac{1}{\sum\limits_{i=1}^{m}(x_i^* - \bar{x}^*)^2}}}$$

where

$$S = \sqrt{\frac{\sum\limits_{i=1}^{n}[Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2 + \sum\limits_{i=1}^{m}[Y_i^* - (\hat{\beta}_0^* + \hat{\beta}_1^* x_i)]^2}{n + m - 4}}$$

Then $T$ has a Student $t$ distribution with $n + m - 4$ degrees of freedom.

**b.** To test $H_0 : \beta_1 = \beta_1^*$ versus $H_1 : \beta_1 \neq \beta_1^*$ at the $\alpha$ level of significance, reject $H_0$ if $t$ is either $(1) \leq -t_{\alpha/2, n+m-4}$ or $(2) \geq t_{\alpha/2, n+m-4}$, where

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_1^*}{s\sqrt{\dfrac{1}{\sum\limits_{i=1}^{n}(x_i - \bar{x})^2} + \dfrac{1}{\sum\limits_{i=1}^{m}(x_i^* - \bar{x}^*)^2}}}$$

(One-sided tests are defined in the usual way by replacing $\pm t_{\alpha/2, n+m-4}$ with either $t_{\alpha, n+m-4}$ or $-t_{\alpha, n+m-4}$.)

$$S^2 = SSE \text{ and } q = 4.$$

Sol. Test
$$H_0 : \beta_1 = \beta_1^* \quad \text{v.s.} \quad H_1 : \beta_1 > \beta_1^*.$$

Long computations ... $t = 2.50$.

Critical region: $t > t_{0.05,8} = 1.8595$.

Reject.

Test

$$H_0 : \beta_1 = \beta_1^* \quad \text{v.s.} \quad H_1 : \beta_1 > \beta_1^*.$$

Long computations ... $t = 2.50$.

Critical region: $t > t_{0.05,8} = 1.8595$.

Reject. □

Sol. Test

$$H_0 : \beta_1 = \beta_1^* \quad \text{v.s.} \quad H_1 : \beta_1 > \beta_1^*.$$

Long computations ... $t = 2.50$.

Critical region: $t > t_{0.05,8} = 1.8595$.

Reject. □

Sol. Test

$$H_0 : \beta_1 = \beta_1^* \quad \text{v.s.} \quad H_1 : \beta_1 > \beta_1^*.$$

Long computations ... $t = 2.50$.

Critical region: $t > t_{0.05,8} = 1.8595$.

Reject. □

```
> # Example 11.3.4
> # Read data first
> Input <- ("
+ x      yA    yB
+ 0       100   100
+ 100    250   203
+ 200    304   214
+ 300    403   295
+ 400    446   330
+ 500    482   324
+ ")
> Data = read.table(textConnection(Input),
+                      header=TRUE)
> Data
    x   yA   yB
1    0 100 100
2 100 250 203
3 200 304 214
4 300 403 295
5 400 446 330
6 500 482 324
```

```
 1  > # fit  the  first  model ...
 2  > DataA <− data.frame(x = Data$x,yA = Data$yA)
 3  > fitA  <− lm(yA~x, DataA)
 4  > summary(fitA)
 5
 6  Call :
 7  lm(formula = yA ~ x,  data = DataA)
 8
 9  Residuals:
10       1       2       3       4       5       6
11  −45.333 30.467 10.267 35.067   3.867 −34.333
12
13  Coefficients :
14              Estimate Std. Error  t  value Pr(>|t |)
15  (Intercept) 145.33333 26.86684  5.409 0.00566 ∗∗
16  x             0.74200   0.08874  8.362 0.00112 ∗∗
17  −−−
18  Signif . codes:  0 '∗∗∗' 0.001 '∗∗' 0.01 '∗' 0.05 '.' 0.1 ' ' 1
19
20  Residual standard error: 37.12 on 4 degrees of freedom
21  Multiple R−squared: 0.9459, Adjusted R−squared: 0.9324
22  F−statistic : 69.92 on 1 and 4 DF, p−value: 0.001119
```
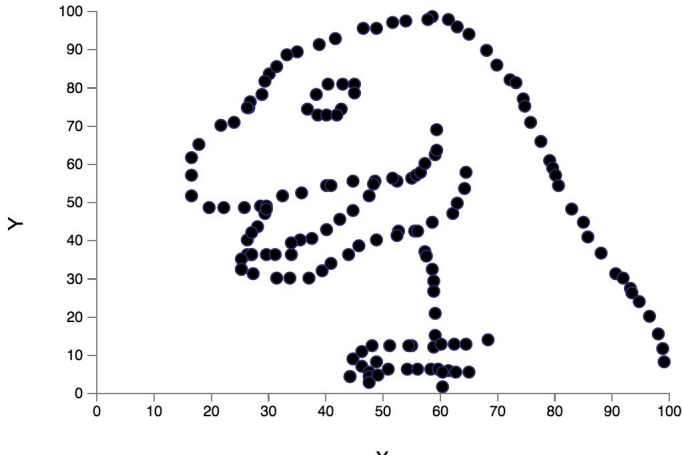
```
 1  > # fit  the second model ...
 2  > DataB <− data.frame(x = Data$x,yB = Data$yB)
 3  > fitB  <− lm(yB~x, DataB)
 4  > summary(fitB)
 5
 6  Call :
 7  lm(formula = yB ~ x,  data = DataB)
 8
 9  Residuals:
10       1        2       3        4       5        6
11  −31.333 26.467 −7.733 28.067 17.867 −33.333
12
13  Coefficients :
14               Estimate Std. Error  t  value Pr(>|t|)
15  (Intercept) 131.33333 22.77255  5.767 0.00449 ∗∗
16  x             0.45200   0.07522  6.009 0.00386 ∗∗
17  − − −
18  Signif . codes:  0 '∗∗∗' 0.001 '∗∗' 0.01 '∗' 0.05 '.' 0.1 ' ' 1
19
20  Residual standard error: 31.46 on 4 degrees of freedom
21  Multiple R−squared: 0.9003, Adjusted R−squared: 0.8754
22  F−statistic : 36.11 on 1 and 4 DF, p−value: 0.00386
```

```
1   > # Now compute t−score and p−value
2   > sA <− summary(fitA)$coefficients
3   > sA
4               Estimate  Std. Error   t value    Pr(>|t|)
5   (Intercept) 145.3333  26.86683800  5.409395   0.005656733
6   x             0.7420   0.08873825  8.361671   0.001118570
7   > sB <− summary(fitB)$coefficients
8   > sB
9               Estimate  Std. Error   t value    Pr(>|t|)
10  (Intercept) 131.3333  22.77254682  5.767178   0.004486443
11  x             0.4520   0.07521525  6.009420   0.003860274
12  > db <− (sA[2,1]−sB[2,1]) # difference of beta_1's
13  > db
14  [1] 0.29
15  > sd <− sqrt(sB[2,2]^2+sA[2,2]^2) # standard deviation
16  > sd
17  [1] 0.1163263
18  > df <− (fitA$df.residual+fitB$df.residual) # degrees of freedom
19  > df
20  [1] 8
21  > td <− db/sd # t−score
22  > pv <− 2*pt(−abs(td), df) # two−sided p−value
23  > print(paste("t−score is ", round(td,3),
24  +             "and p−value is", round(pv,3)))
25  [1] "t−score is  2.493 and p−value is 0.037"
```

⚠️

You should always visualize your data
before any analysis

N = 157 ; X mean = 50.7333 ; X SD = 19.5661 ; Y mean = 46.495 ; Y SD = 27.2828 ;
Pearson correlation = -0.1772

# Plan
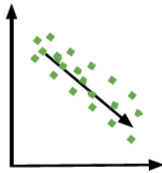
# Chapter 11. Regression

CORRELATION

Positive
Correlation

Zero
Correlation

Negative
Correlation

$$\text{Corr(X,Y)} = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y} \quad \text{— Covarianced normalized by Standard Deviation}$$

Correlation between X and Y

Standard deviation of X

Standard deviation of Y

Notation: $Corr(X, Y) = \rho(X, Y) = \rho_{XY}$

$Var(X) = \sigma_X^2$, $Var(Y) = \sigma_Y^2$, $Cov(X, Y) = \sigma_{XY}$

$$\Downarrow$$

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

$$\text{Corr(X,Y)} = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y} \Bigg] \text{Covarianced normalized by Standard Deviation}$$

Correlation between X and Y

Standard deviation of X

Standard deviation of Y

Notation: $Corr(X, Y) = \rho(X, Y) = \rho_{XY}$

$\text{Var}(X) = \sigma_X^2, \text{Var}(Y) = \sigma_Y^2, \text{Cov}(X, Y) = \sigma_{XY}$

$$\Downarrow$$

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

$$\text{Corr(X,Y)} = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y} \Bigg] \text{Covarianced normalized by Standard Deviation}$$

Correlation between X and Y

Standard deviation of X

Standard deviation of Y

Notation: $Corr(X, Y) = \rho(X, Y) = \rho_{XY}$

$$\text{Var}(X) = \sigma_X^2, \text{Var}(Y) = \sigma_Y^2, \text{Cov}(X, Y) = \sigma_{XY}$$

$$\Downarrow$$

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

For any two r.v.s $X$ and $Y$,

    a. $|\rho(X, Y)| \leq 1$

    b. $\rho(X, Y) = 1$ if and only if $Y = aX + b$ for some $a > 0$ and $b \in \mathbb{R}$;

      $\rho(X, Y) = -1$ if and only if $Y = aX + b$ for some $a < 0$ and $b \in \mathbb{R}$.

Proof. ......

Thm. For any two r.v.s $X$ and $Y$,

a. $|\rho(X, Y)| \leq 1$

b. $\rho(X, Y) = 1$ if and only if $Y = aX + b$ for some $a > 0$ and $b \in \mathbb{R}$;

$\rho(X, Y) = -1$ if and only if $Y = aX + b$ for some $a < 0$ and $b \in \mathbb{R}$.

Proof. ......  □

Thm. For any two r.v.s $X$ and $Y$,

   a. $|\rho(X, Y)| \leq 1$

   b. $\rho(X, Y) = 1$ if and only if $Y = aX + b$ for some $a > 0$ and $b \in \mathbb{R}$;

     $\rho(X, Y) = -1$ if and only if $Y = aX + b$ for some $a < 0$ and $b \in \mathbb{R}$.

Proof. ......

Thm. For any two r.v.s $X$ and $Y$,

a. $|\rho(X, Y)| \leq 1$

b. $\rho(X, Y) = 1$ if and only if $Y = aX + b$ for some $a > 0$ and $b \in \mathbb{R}$;

$\rho(X, Y) = -1$ if and only if $Y = aX + b$ for some $a < 0$ and $b \in \mathbb{R}$.

Proof. ......

Thm. For any two r.v.s $X$ and $Y$,

a. $|\rho(X, Y)| \leq 1$

b. $\rho(X, Y) = 1$ if and only if $Y = aX + b$ for some $a > 0$ and $b \in \mathbb{R}$;

   $\rho(X, Y) = -1$ if and only if $Y = aX + b$ for some $a < 0$ and $b \in \mathbb{R}$.

Proof. ......  $\qquad \square$

## Estimating $\rho(X, Y)$: Sample correlation coefficient

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

$$= \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]}\sqrt{\mathbb{E}[Y^2] - \mathbb{E}[Y]}}$$

$$\Downarrow$$

$$R = \frac{n\sum_{i=1}^{n} X_i Y_i - \left(\sum_{i=1}^{n} X_i\right)\left(\sum_{i=1}^{n} Y_i\right)}{\sqrt{n\sum_{i=1}^{n} X_i^2 - \left(\sum_{i=1}^{n} X_i\right)^2}\sqrt{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_i\right)^2}}$$

*Person product-moment correlation coefficent*
*or simply*
*Sample corelation coefficient*

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

$$= \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]}\sqrt{\mathbb{E}[Y^2] - \mathbb{E}[Y]}}$$

$$\Downarrow$$

$$R = \frac{n\sum_{i=1}^{n} X_i Y_i - \left(\sum_{i=1}^{n} X_i\right)\left(\sum_{i=1}^{n} Y_i\right)}{\sqrt{n\sum_{i=1}^{n} X_i^2 - \left(\sum_{i=1}^{n} X_i\right)^2}\sqrt{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_i\right)^2}}$$

*Person product-moment correlation coefficent*
*or simply*
*Sample corelation coefficient*

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

$$= \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]}\sqrt{\mathbb{E}[Y^2] - \mathbb{E}[Y]}}$$

$$\Downarrow$$

$$R = \frac{n\sum_{i=1}^{n} X_i Y_i - \left(\sum_{i=1}^{n} X_i\right)\left(\sum_{i=1}^{n} Y_i\right)}{\sqrt{n\sum_{i=1}^{n} X_i^2 - \left(\sum_{i=1}^{n} X_i\right)^2}\sqrt{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_i\right)^2}}$$

*Person product-moment correlation coefficent*
*or simply*
*Sample corelation coefficient*

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

$$= \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sqrt{\mathbb{E}[X^2] - \mathbb{E}[X]}\sqrt{\mathbb{E}[Y^2] - \mathbb{E}[Y]}}$$

$$\Downarrow$$

$$R = \frac{n\sum_{i=1}^{n} X_i Y_i - \left(\sum_{i=1}^{n} X_i\right)\left(\sum_{i=1}^{n} Y_i\right)}{\sqrt{n\sum_{i=1}^{n} X_i^2 - \left(\sum_{i=1}^{n} X_i\right)^2}\sqrt{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_i\right)^2}}$$

*Person product-moment correlation coefficent*
or simply
*Sample corelation coefficient*

Thm.

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

Thm.

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

Thm.

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

Thm.

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

$$R^2 = 1 - \frac{SSE}{SST} = \frac{SST - SSE}{SST} = \frac{SSTR}{SST}$$

where

$$SSE = \sum_{i=1}^{n} \left( Y_i - \widehat{Y}_i \right)^2, \quad \widehat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$SST = \sum_{i=1}^{n} \left( Y_i - \overline{Y}_i \right)^2, \quad \text{and} \quad SSTR = SST - SSE.$$

Remark   SSE: sum of square errors $\sim$ the variation in $y_i$'s not explained by L.M.

SST: Total sum of squares $\sim$ total variability.

SSTR: Treatment sum of sqrs. $\sim$ the variation in $y_i$'s explained by L.M.

$R^2$ (or $r^2$ when $X_i$ and $Y_i$ are replaced by $x_i$ and $y_i$) $\sim$ proportion of total variation in the $y_i$'s that can be attributed to L.M.

*Coefficient of determination* or simply *R squared*

Proof

$\square$

# Adjusted R-squared

Def. The adjusted R-squareed:

$$R^2_{adj} := 1 - \frac{MSE}{MST}$$

where

$$MSE = \frac{SSE}{n-q} \qquad \text{and} \qquad MST = \frac{SST}{n-1}$$

and $q$ is number of parameters in the model.

Relation:

$$R^2_{adj} = 1 - \left(1 - R^2\right) \frac{n-1}{n-q}$$

MSE: Mean squared error.

MST: Mean squared total.

MSR = MSTR: Mean square for treatment (or regresssion).

$$MSR = MSTR = \frac{SSTR}{q-1}$$

# Adjusted R-squared

Def. The adjusted R-squareed:

$$R^2_{adj} := 1 - \frac{MSE}{MST}$$

where

$$MSE = \frac{SSE}{n-q} \qquad \text{and} \qquad MST = \frac{SST}{n-1}$$

and $q$ is number of parameters in the model.

Relation:

$$R^2_{adj} = 1 - \left(1 - R^2\right) \frac{n-1}{n-q}$$

MSE: Mean squared error.

MST: Mean squared total.

MSR = MSTR: Mean square for treatment (or regresssion).

$$MSR = MSTR = \frac{SSTR}{q-1}$$

# Adjusted R-squared

Def. The adjusted R-squareed:

$$R^2_{adj} := 1 - \frac{MSE}{MST}$$

where

$$MSE = \frac{SSE}{n - q} \qquad \text{and} \qquad MST = \frac{SST}{n - 1}$$

and $q$ is number of parameters in the model.

Relation:

$$R^2_{adj} = 1 - \left(1 - R^2\right) \frac{n - 1}{n - q}$$

MSE: Mean squared error.

MST: Mean squared total.

MSR = MSTR: Mean square for treatment (or regresssion).

$$MSR = MSTR = \frac{SSTR}{q - 1}$$

# Adjusted R-squared

Def.  The adjusted R-squareed:

$$R_{adj}^2 := 1 - \frac{MSE}{MST}$$

where

$$MSE = \frac{SSE}{n-q} \qquad \text{and} \qquad MST = \frac{SST}{n-1}$$

and $q$ is number of parameters in the model.

Relation:

$$R_{adj}^2 = 1 - \left(1 - R^2\right) \frac{n-1}{n-q}$$

MSE: Mean squared error.

MST: Mean squared total.

MSR = MSTR: Mean square for treatment (or regresssion).

$$MSR = MSTR = \frac{SSTR}{q-1}$$

# Adjusted R-squared

Def. The adjusted R-squareed:

$$R^2_{adj} := 1 - \frac{MSE}{MST}$$

where

$$MSE = \frac{SSE}{n-q} \qquad \text{and} \qquad MST = \frac{SST}{n-1}$$

and $q$ is number of parameters in the model.

Relation:

$$R^2_{adj} = 1 - \left(1 - R^2\right) \frac{n-1}{n-q}$$

MSE: Mean squared error.

MST: Mean squared total.

MSR = MSTR: Mean square for treatment (or regresssion).

$$MSR = MSTR = \frac{SSTR}{q-1}$$

# Plan

# Chapter 11. Regression

# § 11.5 The Bivariate Normal Distribution