

Notation

- Random variables X, Y, E
- observations / sample x, y, e
- Sample vector $\underset{i=1 \dots n}{x}, \underset{i=1 \dots n}{y}, \underset{i=1 \dots n}{e}$

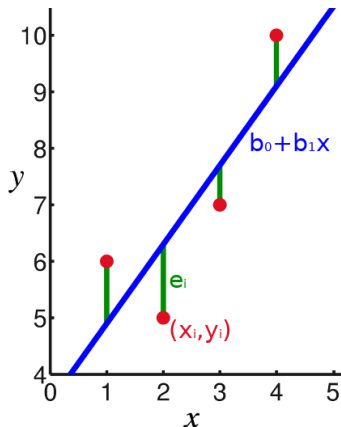
Covariance and correlation give some idea of direction and strength of linear dependencies. **Simple linear regression** tries to give an explicit linear function that explains one random variable Y by another random variable X *random variables*

$Y = b_0 + b_1 X + E.$

Here E are random variations in Y that cannot be explained by X .

- Random variables are denoted by capital letters, e.g. X, Y, E .
- Individual observations are denoted by lower-case letters often with an index, e.g. x_i, y_i, e_i .
- Samples are given by vectors, e.g. $\mathbf{x}, \mathbf{y}, \mathbf{e}$.

•



For random variables X, Y, E

$$Y = b_0 + b_1X + E$$

For **samples** x_j, y_j, e_j

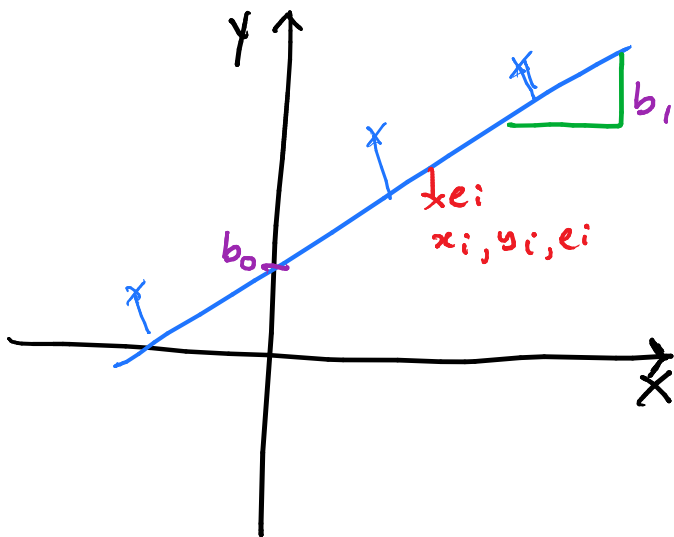
$$y_i = b_0 + b_1 x_i + e_i, \quad i = 1, 2, \dots, n$$

Applications

- b_0 and b_1 can contain interesting information about the data generating process, e.g. constant in the ideal gas law.
- Predict new values of Y given potential values of X , e.g. predict electricity consumption for the next day.

Notation

- () () () ()



Observations: x_i, y_i

unknown: b_0, b_1, e_i

$$\underline{e} = \underline{y} - X \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}$$

$$e_i = y_i - b_0 - b_1 x_i \quad i = 1 \dots n$$

Idea: minimise $\|\underline{e}\|^2$

$$\begin{aligned} \min_{b_0, b_1} \|\underline{e}\|^2 &= \|\underline{y} - b_0 \underline{1} - b_1 \underline{x}\|^2 \\ &= \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2 \end{aligned}$$

Given two data vectors \mathbf{x} , \mathbf{y} we have the regression equation for the samples

or $\mathbf{y} = b_0 \mathbf{1} + b_1 \mathbf{x} + \mathbf{e}$

with $\mathbf{e} = (e_1, e_2, \dots, e_n)^\top$.

Equivalently,

$$y = X \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} + e$$

with the $n \times 2$ matrix

$$\mathbf{X} = [\mathbf{1} \ \mathbf{x}] = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \text{matrix}$$

Notation

$$\hat{b}_0, \quad \hat{b}_1, \quad \hat{\mathbf{e}}$$

Idea

- $$\|\hat{\mathbf{e}}\|^2 = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{b}_0 - \hat{b}_1 x_i)^2.$$

Remember : $\tilde{x} \cdot \perp = 0$ (orthogonal)

$$\begin{aligned}\text{span} \{ \tilde{x}, \perp \} &= \text{span} \{ \tilde{x} + \bar{x} \perp, \perp \} \\ &= \text{span} \{ x, \perp \}\end{aligned}$$

$\frac{\tilde{x}}{\|\tilde{x}\|}, \frac{\perp}{\|\perp\|}$ orthonormal basis

Question

Because minimising $\sum_{i=1}^n \hat{e}_i^2$ leads to a least squares problem and this is "easy" to compute.

Since

$\underline{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$ orthonormal basis of $\text{span}\{\underline{1}, \underline{x}\}$

$\mathbf{e} = \mathbf{y} - \mathbf{X} \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}$ with $\mathbf{X} = [\underline{1} \ \underline{x}]$,

we are looking for the least squares solution to

$$\mathbf{y} = \mathbf{X} \begin{pmatrix} b_0 \\ b_1 \end{pmatrix}$$

This can be found by projecting \mathbf{y} on an orthonormal basis for the span of the columns of \mathbf{X} .

last lecture

• Regression model

$$Y = b_0 + b_1 X + E$$

↑ regressor / dependent variable

↑ regressor / independent variable

← error term / unobserved variable

- Observations: $y_i, x_i \quad i = 1 \dots n$

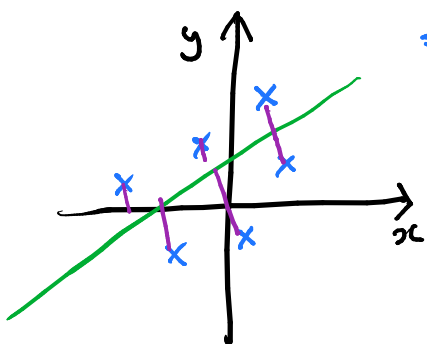
unobserved: $b_0, b_1, e_i \quad i = 1 \dots n$

$$y_i = b_0 + b_1 x_i + e_i \quad i = 1 \dots n$$

$$\underline{y} = b_0 + b_1 \underline{x} + \underline{e}$$

$$\underline{y} = \underline{X} \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} + \underline{e} \quad \underline{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

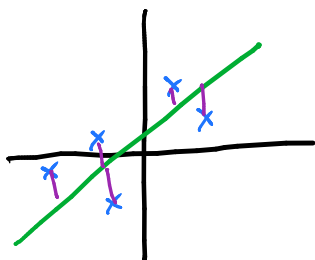
- Estimation: Find $\hat{b}_0, \hat{b}_1, \hat{e}$



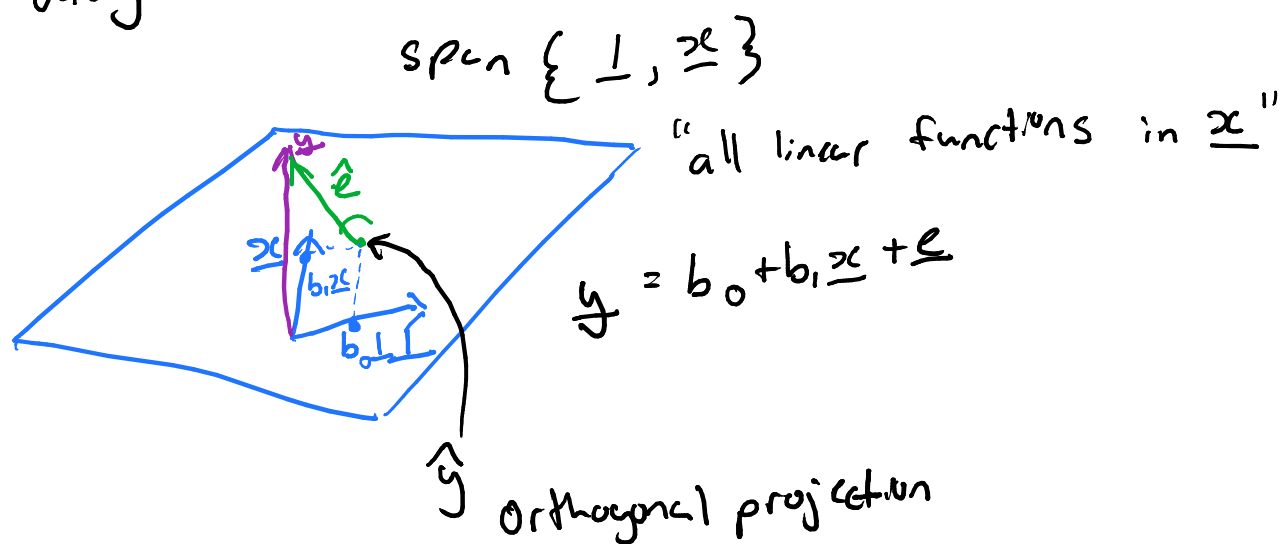
scatter plots

which one to choose

$\|\hat{e}\|^2$ is small



Vector diagram



\hat{b}_0, \hat{b}_1 is the solution to the least squares problem

$$\hat{y} = X \begin{pmatrix} \hat{b}_0 \\ \hat{b}_1 \end{pmatrix} \text{ because } e = y - \hat{y}$$

$$\min \|e\|^2 = \min \|y - X \begin{pmatrix} \hat{b}_0 \\ \hat{b}_1 \end{pmatrix}\|^2$$

• option 1 = for least squares

$$X^T X \begin{pmatrix} \hat{b}_0 \\ \hat{b}_1 \end{pmatrix} = X^T y \text{ normal equation}$$

• option 2 = orthogonal projection of y on an orthonormal basis of the columns of X

$$\text{span} \left\{ \frac{1}{\|1\|}, \frac{\tilde{x}}{\|\tilde{x}\|} \right\} = \text{span} \{1, x\}$$

orthonormal columns X

\hat{y} orthogonal projection of y on $\frac{1}{\|1\|}, \frac{\tilde{x}}{\|\tilde{x}\|}$

$$\begin{aligned} \hat{y} &= \frac{y \cdot \tilde{x}}{\|\tilde{x}\|} \cdot \frac{\tilde{x}}{\|\tilde{x}\|} + \frac{y \cdot 1}{\|1\|} \cdot \frac{1}{\|1\|} = \left\| \begin{bmatrix} 1 \\ \tilde{x} \end{bmatrix} \right\| \\ &= \frac{\tilde{x}}{\|\tilde{x}\|^2} (x - \bar{x} \cdot 1) + \frac{y \cdot 1}{\|1\|^2} \cdot 1 = \sqrt{\sum_{i=1}^n 1} = \sqrt{n} \end{aligned}$$

$\frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$ $y \cdot 1 = \sum_{i=1}^n y_i$

$$= \frac{y \cdot \tilde{x}}{\|\tilde{x}\|^2} \cdot x + \left(\bar{y} - \bar{x} \frac{y \cdot \tilde{x}}{\|\tilde{x}\|^2} \right) \cdot 1$$

$$= \hat{b}_0 \cdot 1 + \hat{b}_1 x$$

$$\hat{b}_1 = \frac{y \cdot \tilde{x}}{\|\tilde{x}\|^2} \quad \hat{b}_0 = (\bar{y} - \bar{x} \hat{b}_1)$$

Estimation

Least squares

best estimation

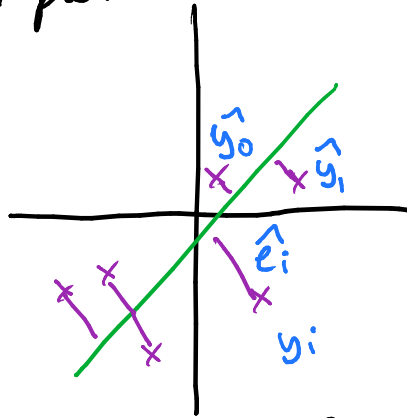
The two vectors

$$\frac{\mathbf{1}}{\|\mathbf{1}\|}, \quad \frac{\tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|} = \frac{\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}}{\|\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}\|}$$

are orthonormal and have the same span as the columns of $\mathbf{X} = [\mathbf{1} \ \mathbf{x}]$. Let $\hat{\mathbf{y}}$ be the orthogonal projection onto this basis.

$$\begin{aligned} \hat{\mathbf{y}} &= \frac{\mathbf{y} \cdot (\mathbf{x} - \bar{\mathbf{x}}\mathbf{1})}{\|\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}\|} \frac{\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}}{\|\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}\|} + \frac{\mathbf{y} \cdot \mathbf{1}}{\|\mathbf{1}\|} \frac{\mathbf{1}}{\|\mathbf{1}\|} \\ &= \frac{\mathbf{y} \cdot (\mathbf{x} - \bar{\mathbf{x}}\mathbf{1})}{\|\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}\|^2} \mathbf{x} + \left(\frac{\mathbf{y} \cdot \mathbf{1}}{\|\mathbf{1}\|^2} - \bar{\mathbf{x}} \frac{\mathbf{y} \cdot (\mathbf{x} - \bar{\mathbf{x}}\mathbf{1})}{\|\mathbf{x} - \bar{\mathbf{x}}\mathbf{1}\|^2} \right) \mathbf{1} \\ &= \underbrace{\frac{\mathbf{y} \cdot \tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|^2}}_{\hat{b}_1} \mathbf{x} + \underbrace{\left(\bar{\mathbf{y}} - \bar{\mathbf{x}} \frac{\mathbf{y} \cdot \tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|^2} \right)}_{\bar{\mathbf{y}} - \hat{b}_1 \bar{\mathbf{x}} = \hat{b}_0} \mathbf{1} \end{aligned}$$

scatter plot



$$\hat{b}_0 + \hat{b}_1 x$$

why is

$$\begin{aligned} y \cdot \tilde{x} &= \tilde{y} \cdot \tilde{x} \\ \tilde{y} \cdot \tilde{x} &= (y - \bar{y} \mathbf{1}) \cdot \tilde{x} \quad \text{True because} \\ &= y \cdot \tilde{x} - \underbrace{\bar{y} \mathbf{1} \cdot \tilde{x}}_{=0} \end{aligned}$$

$$\Rightarrow \hat{b}_1 = \frac{y \cdot \tilde{x}}{\|\tilde{x}\|^2} = \frac{\tilde{y} \cdot \tilde{x}}{\|\tilde{x}\|^2} = \text{Corr}(\dots)$$

$$= \text{Corr}(x, y) \cdot \frac{\text{sd}(y)}{\text{sd}(x)}$$

$$= \frac{\text{cov}(x, y)}{\text{var}(x)}$$

Result of least squares

Result of least squares

$$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i \quad i = 1, 2, \dots, n$$

$$y_i = \hat{b}_0 + \hat{b}_1 x_i + \hat{e}_i \quad i = 1, 2, \dots, n$$

with

- $\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$

- $\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}}$

A scatter plot illustrating linear regression. The x-axis is labeled x and ranges from 0 to 5. The y-axis is labeled y and ranges from 4 to 10. A blue line represents the fitted regression line. Five data points are plotted as red dots. Vertical green lines connect each data point to the regression line, representing the residuals. The equation $\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i$ is shown in blue. The residual for the point $(2, 5)$ is labeled \hat{e}_i in green. The point $(2, 5)$ is labeled (x_i, y_i) in red.

Notation We call $\hat{e}_i = \text{residual}$

- \hat{b}_0, \hat{b}_1 - ordinary least squares (OLS) estimators
- \hat{y}_i - fitted values
- $\hat{\mathbf{y}}$ - vector of fitted values
- \hat{e}_i - residuals
- $\hat{\mathbf{e}}$ - vector of residuals
- $\hat{y} = \hat{b}_0 + \hat{b}_1 x$, with $x \in \mathbb{R}$
- regression line.

The regression coefficients have a straight forward interpretation.

- For every increase of x by one unit a change of y by \hat{b}_1 units can be expected on average.
- If $x = 0$, then on average $y = \hat{b}_0$.

last lecture

• OLS estimations

$$\hat{b}_1 = \frac{y \cdot \hat{x}}{\|\hat{x}\|^2}$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

• regression line

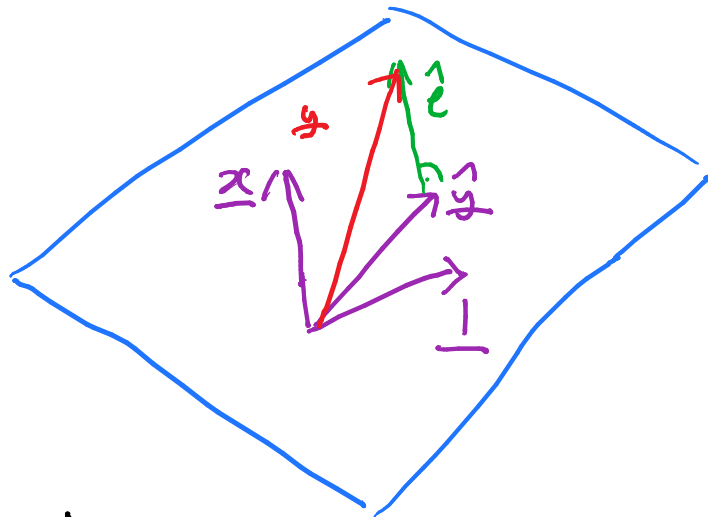
$$\hat{b}_0 + \hat{b}_1 x \quad x \in \mathbb{R}$$

• Fitted values

$$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_i$$

• residuals $\hat{e}_i = y_i - \hat{y}_i$

$$\text{Hence, } y_i = \hat{b}_0 + \hat{b}_1 x_i + \hat{e}_i$$



By construction

$$\bullet \underline{\hat{e}} \perp \underline{1}$$

$$\bullet \underline{\hat{e}} \perp \underline{x}$$

$$\bullet \underline{\hat{e}} \perp \underline{\hat{x}} \quad \text{are orthogonal}$$

$$\bullet \underline{\hat{e}} \perp \underline{\hat{y}}$$

$$\bullet \bar{\underline{\hat{e}}} = \frac{1}{n} \sum_{i=1}^n \underline{\hat{e}}_i = \frac{1}{n} \underbrace{\underline{\hat{e}} \cdot \underline{1}}_{=0} = 0$$

$$\bullet \bar{\hat{y}} = \bar{y} \quad ? \quad \text{Tutorial}$$

$$\bullet \text{cov}(\underline{x}, \underline{\hat{e}}) = \frac{1}{n-1} \underbrace{\underline{x} \cdot \underline{\hat{e}}}_{=0} = 0$$

$$= \text{corr}(\underline{x}, \underline{\hat{e}})$$

$$\bullet \text{cov}(\underline{\hat{y}}, \underline{\hat{e}}) = \frac{1}{n-1} \underline{\hat{y}} \cdot \underline{\hat{e}} = 0$$

$$= \text{corr}(\underline{\hat{y}}, \underline{\hat{e}})$$

$$\underline{\tilde{\hat{e}}} = \underline{\hat{e}}$$

$$\min \|\underline{\hat{e}}\|^2 = \min \|\underline{\tilde{\hat{e}}}\|^2 = \min n \times \text{var}(\underline{\hat{e}})$$

Estimation

Properties

Matlab Function 'regress'

- $\widehat{\mathbf{e}} = \mathbf{0}$
- $\widehat{\mathbf{y}} = \bar{\mathbf{y}}$
- $\text{cov}(\mathbf{x}, \widehat{\mathbf{e}}) = 0$, i.e. $\widetilde{\mathbf{x}} \perp \widehat{\mathbf{e}}$
- $\text{cov}(\widehat{\mathbf{y}}, \widehat{\mathbf{e}}) = 0$, i.e. $\widetilde{\widehat{\mathbf{y}}} \perp \widehat{\mathbf{e}}$

Hence, the residuals of the OLS estimator ...

- contain no information that could be explained by a linear function of \mathbf{x} .
- are on average 0.
- have the smallest variance possible.

• Example (Temperature, pressure)

$$\underline{x} = [0, 10, 20, 30, 40, 50]^T$$

$$\bar{x} = 25$$

$$\tilde{x} = x - 25 = [-25, -15, -5, 5, 10, 25]^T$$

$$\begin{aligned} \|\tilde{x}\|^2 &= 25^2 + 15^2 + 5^2 + 5^2 + 15^2 + 25^2 \\ &= 1750 \end{aligned}$$

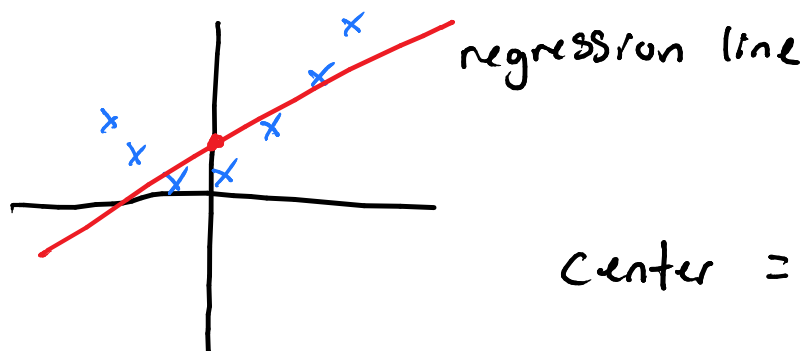
$$\underline{y} = [91, 95, 100, 101, 107, 112]^T$$

$$\bar{y} = 101$$

$$\begin{aligned} \hat{b}_1 = \underline{y} \cdot \tilde{x} &= \frac{91(-25) + 95(-15) + 100(-5) + 101(5) + 107(15) + 112(25)}{1750} \end{aligned}$$

$$\begin{aligned} \hat{b}_0 &= \bar{y} - \hat{b}_1 \bar{x} \approx 101 - 0.406 \times 25 \\ &\approx 90.9 \end{aligned}$$

$$\hat{y}_i = 90.9 + 0.406 \times x_i$$



Center = temp - mean(temp)

Estimation

Example Temperature and pressure

$$\tilde{\mathbf{x}} = (-25, -15, -5, 5, 15, 25)^\top$$

$$\mathbf{y} = (91, 95, 100, 101, 107, 112)^\top$$

$$\bar{y} = 101 \quad \text{and} \quad \bar{x} = 25$$

$$\begin{aligned} \mathbf{y} \cdot \tilde{\mathbf{x}} &= \tilde{\mathbf{y}} \cdot \tilde{\mathbf{x}} = \left((-25)91 + (-15)95 + (-5)100 \right. \\ &\quad \left. + (5)101 + (15)107 + (25)112 \right) \\ &= 710 \end{aligned}$$

$$\|\tilde{\mathbf{x}}\|^2 = 1750$$

$$\hat{b}_1 = \frac{\mathbf{y} \cdot \tilde{\mathbf{x}}}{\|\tilde{\mathbf{x}}\|^2} = \frac{710}{1750} \approx 0.406$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x} = 101 - \frac{710}{1750} 25 \approx 90.9.$$

Hence, $\hat{y} = 90.9 + 0.406x$.

Estimation

Example

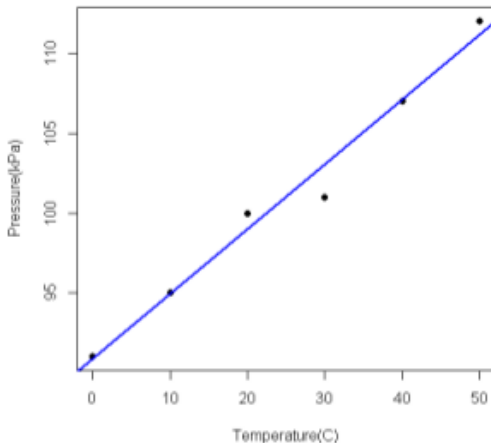


Figure 3: Running example: pressure versus temperature in a boiler.

Estimation

Example

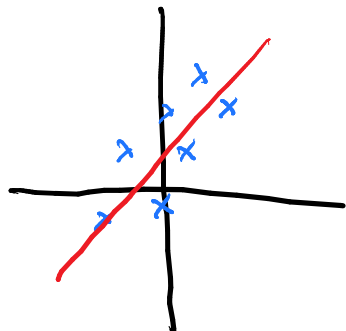
Fitted values

Temp ($^{\circ}\text{C}$) x_i	Pressure (kPa) y_i	$\hat{b}_0 + \hat{b}_1 x_i$
0	91	90.9
10	95	94.9
20	100	99.0
30	101	103
40	107	107
50	112	111

Table 1: Temperature, pressure, and predicted pressure

Interpretation

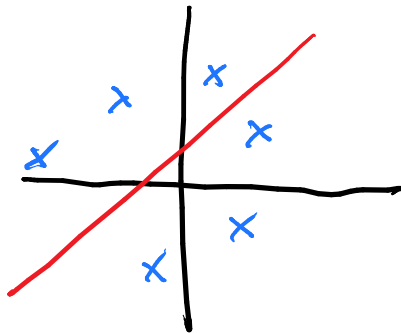
- The regression line is $\hat{y} = 90.9 + 0.406x$.
- For every increase (decrease) of temperature by one degree Celsius an increase (decrease) of pressure by $b_1 = 0.406$ kPa can be expected.
- The pressure that we can expect at 0°C is $b_0 = 90.9$ kPa.



good fit

$$\text{var}(\hat{y}) = \frac{1}{n-1} \|\tilde{y}\|^2$$

$$TSS = (n-1) \text{var}(\hat{y})$$



bad fit

2.3 Goodness-of-fit

Notation

- **Total sum of squares (TSS)**

$$TSS = \|\tilde{\mathbf{y}}\|^2 = \sum_{i=1}^n (y_i - \bar{y})^2$$

This quantifies all variation in the sample \mathbf{y} .

- **Explained sum of squares (ESS)**

$$ESS = \|\hat{\tilde{\mathbf{y}}}\|^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

The variation in \mathbf{y} that is explained by the regression line.

- **Residual sum of squares (RSS)**

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

The variation in \mathbf{y} that is not explained by the regression line.

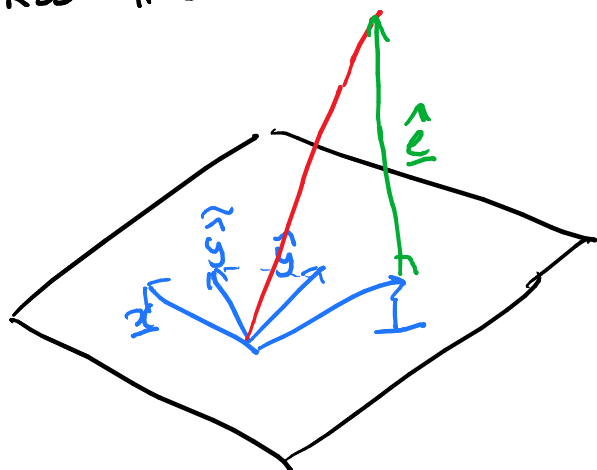
last lecture

• properties of OLS $\hat{b}_0, \hat{b}_1, \hat{e}$

• examples

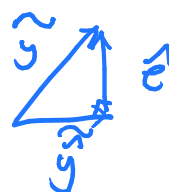
• $TSS = \|\tilde{y}\|^2$, $ESS = \|\hat{y}\|^2$

$RSS = \|\hat{e}\|^2$



(i) \hat{e} and \hat{y} are orthogonal

(ii) $\tilde{y} = \hat{y} + \hat{e}$
 $\tilde{y} = \underbrace{\tilde{y}}_{\hat{y}} + \underbrace{\hat{e}}_{\hat{e}} = \hat{y} + \hat{e} = \hat{y} + \hat{e}$



$$\|\tilde{y}\|^2 = \|\hat{y}\|^2 + \|\hat{e}\|^2$$

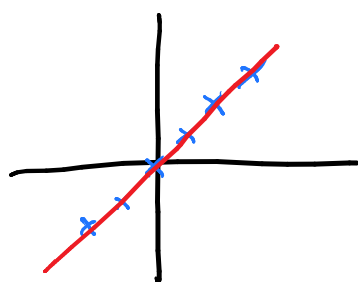
$$TSS = ESS + RSS$$

$$\Rightarrow R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

$\underbrace{\hspace{10em}}_{\leq 1}$

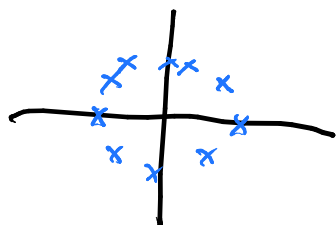
$0 \leq R^2 \leq 1$

$R^2 = 1$



perfect fit

$R^2 = 0$



regression is useless

• Calculating R^2

$$R^2 = \frac{\| \hat{y} \|^2}{\| y \|^2} = \frac{\text{var}(\hat{y})}{\text{var}(y)}$$

$$\hat{y} = X \begin{pmatrix} \hat{b}_0 \\ \hat{b}_1 \end{pmatrix}$$

Goodness-of-fit

We want to measure how good the regression line fits the data by taking the ratio of explained and total sum of squares.

A large value indicates good fit, a small value indicates a less good fit.

Definition The **R-squared** of the regression is

$$R^2 = \frac{ESS}{TSS}$$

Remember $\tilde{\mathbf{y}} \perp \hat{\mathbf{e}}$

Pythagorean theorem

$$\|\mathbf{y} - \bar{y}\mathbf{1}\|^2 = \|\bar{y}\mathbf{1} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \mathbf{y}\|^2, \quad \text{i.e.} \quad \|\tilde{\mathbf{y}}\|^2 = \|\hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{e}}\|^2$$

Hence, $TSS = ESS + RSS$ and

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}.$$

Goodness-of-fit

Properties of the R-squared

- $0 \leq R^2 \leq 1$
- If $R^2 = 1$, the fit is perfect. All data points are on the regression line and all residuals equal 0.
- If $R^2 = 0$, the regression line does not fit the data and is useless.
- If R^2 is close to 1, it is very likely that the regression found an actual relation in the data and will probably work well for prediction.
- If R^2 is small, it is unclear whether we found an actual relation or not. Prediction will very likely not work well.

2.4 Confidence intervals

Remember

We assume the regression equation with the true b_0 and b_1 for the random variables X and Y

$$Y = b_0 + b_1X + E.$$

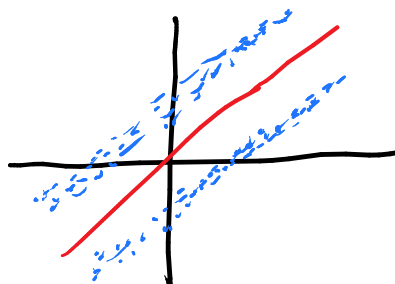
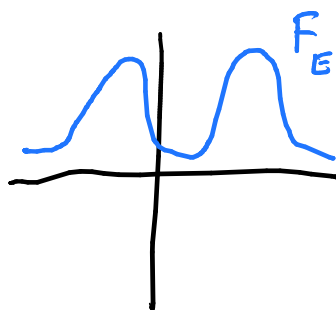
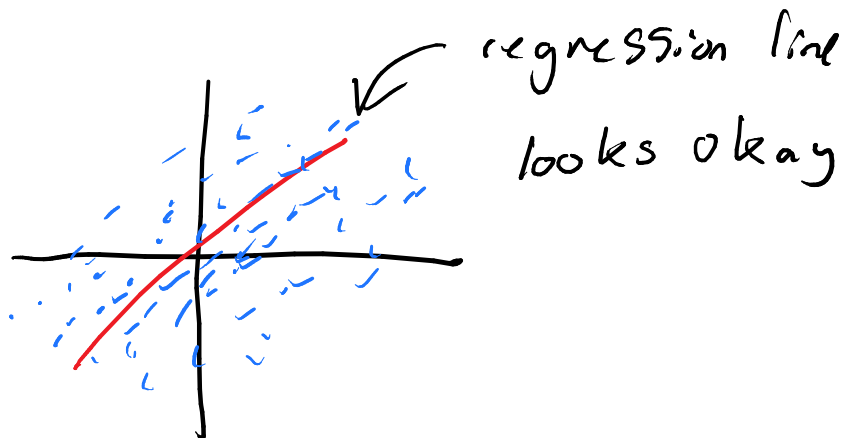
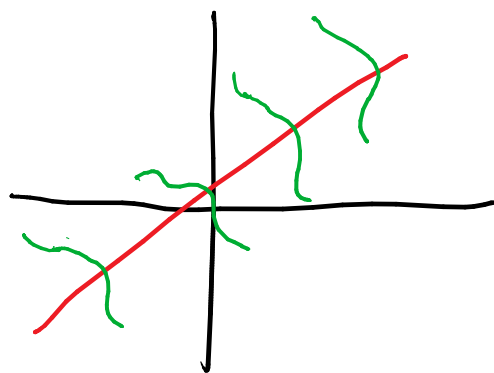
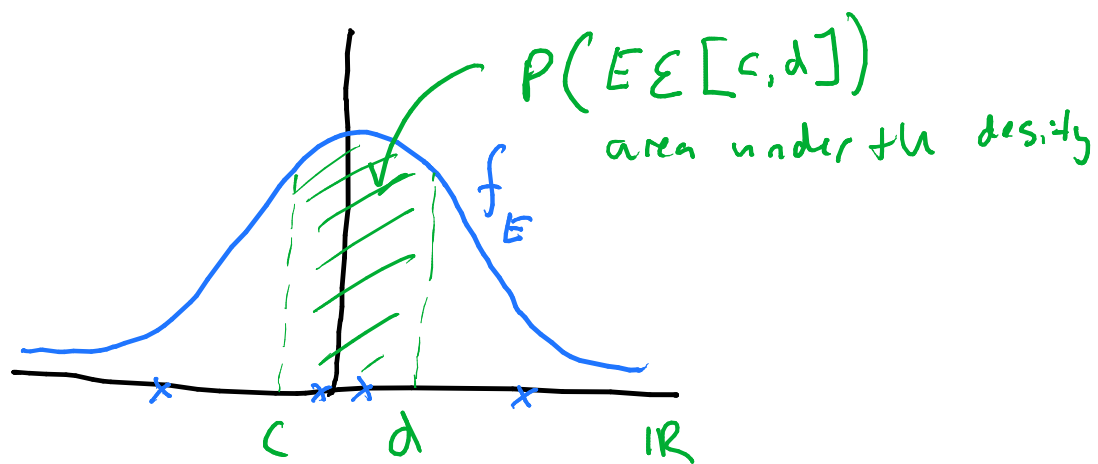
With samples \mathbf{x} and \mathbf{y} we estimate \hat{b}_0 and \hat{b}_1 such that

$$\mathbf{y} = \hat{b}_0\mathbf{1} + \hat{b}_1\mathbf{x} + \hat{\mathbf{e}}.$$

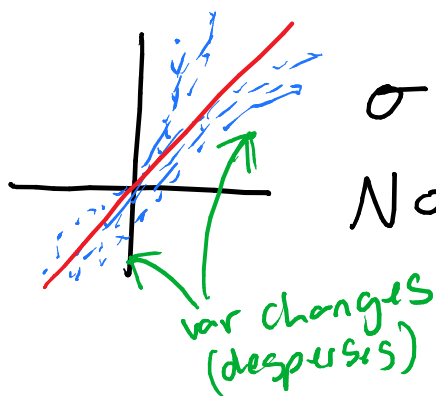
Estimation errors

- Every time we collect data for the random variables X and Y , we will get different data vectors \mathbf{x} and \mathbf{y} and different estimates \hat{b}_0 and \hat{b}_1 .
- The estimators \hat{b}_0 and \hat{b}_1 can be understood as random variables.
- The estimation errors $\hat{b}_0 - b_0$ and $\hat{b}_1 - b_1$ will be random.
- How much can we trust \hat{b}_0 and \hat{b}_1 ?

Density is a non negative function



Not okay



σ depends on X
Not okay

The assumption holds in many applications, e.g. measurement errors.

Confidence intervals

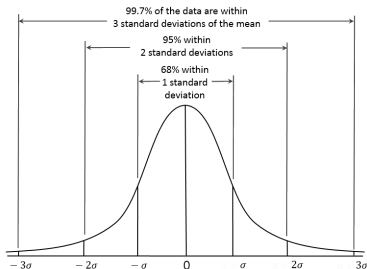
We will answer this question only for b_1 and only under the following assumption.

Assumption

The relation $Y = b_0 + b_1X + E$ holds with $E \sim \mathcal{N}(0, \sigma^2)$, i.e. normally distributed with mean 0 and variance σ^2 .

Interpretation

- When Y is on average equal to $b_0 + b_1X$, E will have mean 0.
- E has the probability density $f_E(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{x^2}{2\sigma^2}\right)$



68% chance that $E \in [-\sigma, \sigma]$

95% chance that
 $E \in [-1.96\sigma, 1.96\sigma]$

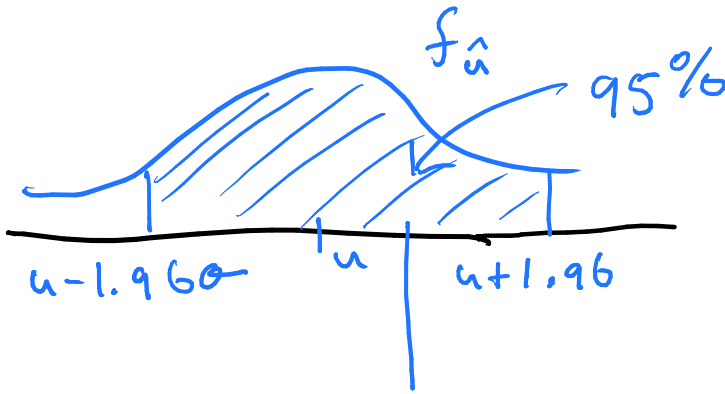
99.7% chance that $E \in [-3\sigma, 3\sigma]$

$$\alpha = 0.05 \rightarrow 95\%$$

$$\alpha = 0.1 \rightarrow 90\%$$

$$\alpha = 0.01 \rightarrow 99\%$$

depends on the engineering field



with probability 0.95

$$\hat{u} \in [u - 1.96\sigma, u + 1.96\sigma]$$

$$\Rightarrow |u - \hat{u}| \leq 1.96\sigma$$

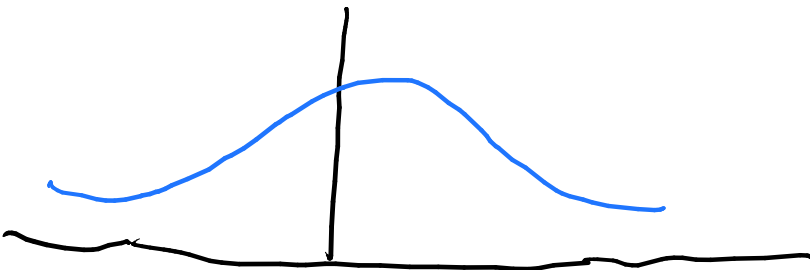
$$\Rightarrow u \in [\hat{u} - 1.96\sigma, \hat{u} + 1.96\sigma]$$

95% confidence interval

for u .

$$\hat{b}_1 \sim t(n-2)$$

t -distribution
with $n-2$ df



Confidence intervals

Definition

A $100(1 - \alpha)\%$ **confidence interval** for b_1 is an interval $[c, d]$ constructed from the data vectors \mathbf{x} , \mathbf{y} such that the chance for b_1 to be in $[c, d]$ is $100(1 - \alpha)\%$.

In other words

If we collect a lot of different samples and construct a confidence interval for each sample

$$\mathbf{x}_1, \mathbf{y}_1 \rightsquigarrow [c_1, d_1]$$

$$\mathbf{x}_2, \mathbf{y}_2 \rightsquigarrow [c_2, d_2]$$

$$\vdots$$

then $b_1 \in [a_j, b_j]$ can be expected in $100(1 - \alpha)\%$ of the cases.

Remark

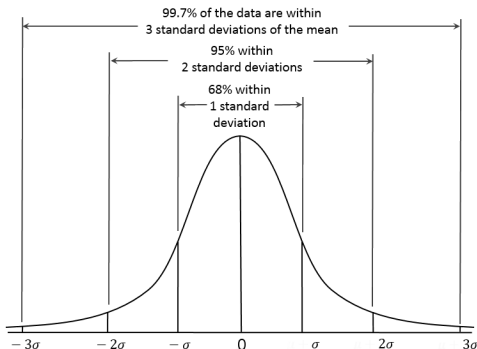
Typically, we chose $\alpha = 0.05$ or $\alpha = 0.1$ and construct 95% or a 90% confidence interval respectively.

Confidence intervals

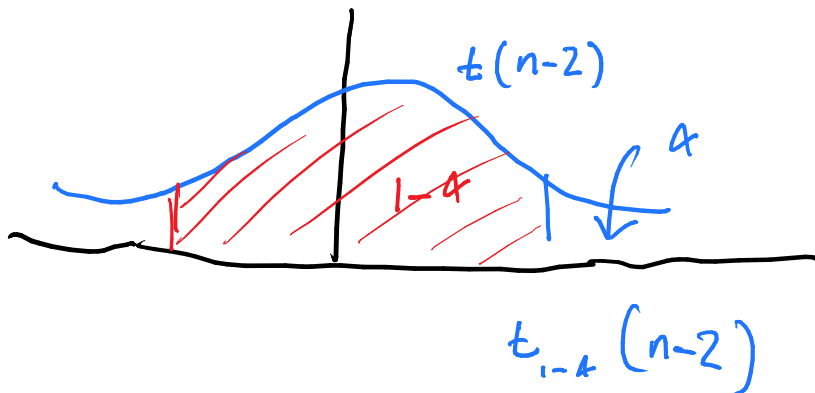
Example

Let \hat{u} be an estimator for some value u with $\hat{u} \sim \mathcal{N}(u, \sigma^2)$.

- Hence, with a 95% chance $|\hat{u} - u| \leq 1.96\sigma$.
- Thus, with a 95% chance $u \in [\hat{u} - 1.96\sigma, \hat{u} + 1.96\sigma]$.
- $[\hat{u} - 1.96\sigma, \hat{u} + 1.96\sigma]$ is a 95% confidence interval.



$\hat{b}_1 \sim t(n-2)$
 t -distribution
 with $n-2$ df



$$\begin{array}{ccc}
 \alpha/2 & & \alpha/2 \\
 -b_1 - \frac{\alpha}{2}(n-1) & & t_{1-\frac{\alpha}{2}}(n-2)
 \end{array}$$

Confidence intervals

t-distribution

Confidence intervals for b_1 work similarly, except that we cannot work with the normal distribution.

- Instead we have to work with the so called ***t*-distribution** with $n - 2$ **degrees of freedom**.
- Replace 1.96 by the so called **critical value** $t_{1-\frac{\alpha}{2}}(n-2)$.
- The critical value can be looked up in a table or use the Matlab command `tinv(1- α /2,n-2)`.

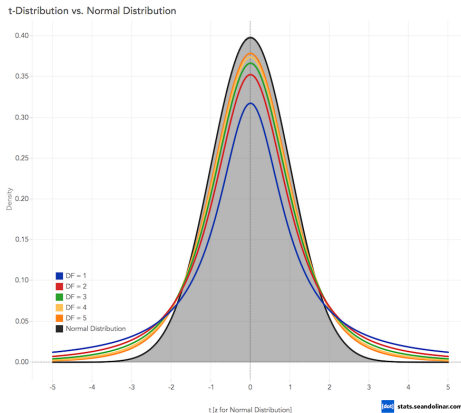
Theorem With probability $100(1 - \alpha)\%$

$$b_1 \in \left[\hat{b}_1 - t_{1-\frac{\alpha}{2}}(n-2)se(\hat{b}_1), \hat{b}_1 + t_{1-\frac{\alpha}{2}}(n-2)se(\hat{b}_1) \right]$$

where $se(\hat{b}_1) = \sqrt{\frac{1}{n-2} \frac{RSS}{\|\tilde{\mathbf{x}}\|^2}}$ is the so called **standard error**.

Confidence intervals

t -distribution



- It is similar to the normal distribution but has heavier tails.
- Often called **Student's t -distribution**.
- First published 1908 by William Sealy Gosset under the pseudonym “Student” while working for the Guinness Brewery.

cum. prob	f. ₅₀	f. ₇₅	f. ₈₀	f. ₈₅	f. ₉₀	f. ₉₅	f. ₉₇₅	f. ₉₉	f. ₉₉₅	f. ₉₉₉	f. ₉₉₉₅
one-tail	0.50	0.25	0.20	0.15	0.10	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	1.00	0.50	0.40	0.30	0.20	0.10	0.05	0.02	0.01	0.002	0.001
df											
1	0.000	1.000	1.376	1.963	3.078	6.314	12.71	31.82	63.66	318.31	636.62
2	0.000	0.816	1.061	1.386	1.886	2.920	4.303	6.965	9.925	22.327	31.599
3	0.000	0.765	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215	12.924
4	0.000	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.173	8.610
5	0.000	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893	6.869
6	0.000	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208	5.959
7	0.000	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	0.000	0.706	0.889	1.108	1.397	1.860	2.306	2.896	3.355	4.501	5.041
9	0.000	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250	4.297	4.781
10	0.000	0.700	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	0.000	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	0.000	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930	4.318
13	0.000	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852	4.221
14	0.000	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787	4.140
15	0.000	0.691	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073
16	0.000	0.690	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686	4.015
17	0.000	0.689	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646	3.965
18	0.000	0.688	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610	3.922
19	0.000	0.688	0.861	1.066	1.328	1.729	2.093	2.539	2.861	3.579	3.883
20	0.000	0.687	0.860	1.064	1.325	1.725	2.086	2.528	2.845	3.552	3.850
21	0.000	0.686	0.859	1.063	1.323	1.721	2.080	2.518	2.831	3.527	3.819
22	0.000	0.686	0.858	1.061	1.321	1.717	2.074	2.508	2.819	3.505	3.792
23	0.000	0.685	0.858	1.060	1.319	1.714	2.069	2.500	2.807	3.485	3.768
24	0.000	0.685	0.857	1.059	1.318	1.711	2.064	2.492	2.797	3.467	3.745
25	0.000	0.684	0.856	1.058	1.316	1.708	2.060	2.485	2.787	3.450	3.725
26	0.000	0.684	0.856	1.058	1.315	1.706	2.056	2.479	2.779	3.435	3.707
27	0.000	0.684	0.855	1.057	1.314	1.703	2.052	2.473	2.771	3.421	3.690
28	0.000	0.683	0.855	1.056	1.313	1.701	2.048	2.467	2.763	3.408	3.674
29	0.000	0.683	0.854	1.055	1.311	1.699	2.045	2.462	2.756	3.396	3.659
30	0.000	0.683	0.854	1.055	1.310	1.697	2.042	2.457	2.750	3.385	

Or look for the confidence % in the last row.

Confidence intervals

Interpretation

- A small confidence interval indicates that the estimator \hat{b}_1 is quite precise.
- A wide confidence interval indicates that the estimator \hat{b}_1 is less reliable.
- If 0 is not in the confidence interval, we have strong evidence that $b_1 \neq 0$ and that X has an impact on Y .
We say b_1 is **significant**.
- If we collect more data and increase the sample size n , the confidence interval becomes tighter.

Remember

$$b_1 \in \left[\hat{b}_1 - t_{1-\frac{\alpha}{2}}(n-2)se(\hat{b}_1), \hat{b}_1 + t_{1-\frac{\alpha}{2}}(n-2)se(\hat{b}_1) \right]$$

$$\text{with } se(\hat{b}_1) = \sqrt{\frac{1}{n-2} \frac{RSS}{\|\tilde{\mathbf{x}}\|^2}}$$

Confidence intervals

Example 95% confidence interval for temperature and pressure

$$\hat{b}_0 = 90.9, \hat{b}_1 = 0.406$$

$$\hat{y} = 90.9 + 0.406x$$

$$\mathbf{y} = (91, 95, 100, 101, 107, 112)^\top$$

$$\begin{aligned}\hat{\mathbf{e}} &= (91 - 90.9, 95 - 95, 100 - 99, 101 - 103, 107 - 107, 112 - 111)^\top \\ &= (0.1, 0, 1, -2, 1)^\top\end{aligned}$$

$$RSS = 0.1^2 + 0^2 + 1^2 + (-2)^2 + 1^2 = 6.01$$

Hence, the standard error is

$$se(\hat{b}_1) = \sqrt{\frac{1}{6-2} \frac{6.01}{1750}} \approx 0.029.$$

Confidence intervals

Example 95% confidence interval for temperature and pressure

We have $n - 2 = 6 - 2 = 4$ degrees of freedom and

$$1 - \frac{\alpha}{2} = 1 - \frac{0.05}{2} = 0.975.$$

$$t_{0.975}(4) \approx 2.78$$

$$t_{0.975}(4)se(\hat{b}_1) \approx 2.78 \times 0.029 \approx 0.081$$

This gives the 95% confidence interval for b_1

$$[0.406 - 0.081, 0.406 + 0.081] = [0.325, 0.487].$$

We can see that b_1 is significant.

2.5 t-test

Statisticians like to describe the statement that b_1 is significant in a different way. It is common to derive this statement from a statistical **hypothesis test**. Let us understand statistical testing step by step.

Hypothesis

- A statistical test is a tool that gives evidence **against** a hypothesis, which is called the **null hypothesis** H_0 .
- If you want evidence that b_1 is significant, you need to start with the opposite statement.

$$H_0 : b_1 = 0$$

- Sometime the **alternative hypothesis** H_1 is introduced as well. Although, it is not necessary to state it explicitly because it is always the opposite of H_0 .

$$H_1 : b_1 \neq 0$$

***t*-test**

Test statistic

- We want to find evidence against H_0 . However, we assume for a moment that H_0 is true.
- The next step is to find an estimator of which we know the distribution whenever H_0 is true. We know that

$$\frac{\hat{b}_1}{se(\hat{b}_1)} \sim t \quad \text{with } n - 2 \text{ degrees of freedom}$$

when $b_1 = 0$. This estimator is called **test statistic**.

- If $\hat{b}_1/se(\hat{b}_1)$ takes a value that is very unlikely for the t -distribution, we have evidence that H_0 was probably not true in the first place.

t-test

Critical values

- Critical values have the property that there is only a $100\alpha\%$ chance that

$$\hat{b}_1 / se(\hat{b}_1) < -t_{1-\frac{\alpha}{2}}(n-2) \quad \text{or} \quad t_{1-\frac{\alpha}{2}}(n-2) < \hat{b}_1 / se(\hat{b}_1)$$

when $\frac{\hat{b}_1}{se(\hat{b}_1)} \sim t$ with $n - 2$ degrees of freedom.

- Hence, when

$$|\hat{b}_1 / se(\hat{b}_1)| > t_{1-\frac{\alpha}{2}}(n-2)$$

we **reject** the null-hypothesis on the $100(1 - \alpha)\%$ **confidence level**.

- If

$$|\hat{b}_1 / se(\hat{b}_1)| \leq t_{1-\frac{\alpha}{2}}(n-2),$$

we **fail to reject** the null-hypothesis. We have neither evidence for nor against the null-hypothesis!

Note

- The test described above is called **two sided t -test**.
- The test rejects H_0 when 0 is not in the $100(1 - \alpha)\%$ confidence interval and fails to reject when 0 is in the interval.
- We cannot rule out a null-hypothesis with 100% certainty. We have to live with the fact that we can only reject on a certain confidence level, e.g. 90% or 95%.
- α is called the **significance level** of the test.

One sided t -test

- _____

0 1 1 1 1 1

- 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99

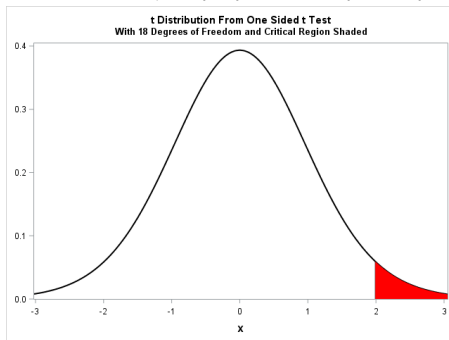
1. *Journal of the American Medical Association*, 1997; 278: 1039-1044.

- Fail to reject if $-t_{1-\alpha}(n-2) \leq \hat{b}_1/se(\hat{b}_1)$.

One sided t -test

- $$H_0 : b_1 \leq 0.$$

- Fail to reject if $\hat{b}_1/se(\hat{b}_1) \leq t_{1-\alpha}(n-2)$.



t-test

Example Temperature and pressure $\hat{y} = 90.9 + 0.406x$

Can we confirm that b_1 is positive on the 99% confidence level?

$$H_0 : b_1 \leq 0$$

The critical value is $t_{0.99}(4) = 3.747$.

Compute the test statistic:

$$\begin{aligned}\hat{b}_1 &= 0.406 \\ se(\hat{b}_1) &\approx 0.029 \\ \frac{\hat{b}_1}{se(\hat{b}_1)} &\approx \frac{0.406}{0.029} \approx 14.\end{aligned}$$

Since $14 > 3.747$ the null hypothesis can be rejected. We have strong evidence that $b_1 > 0$.

2.6 Prediction intervals

In many applications \hat{b}_0 and \hat{b}_1 are estimated and used to predict Y for a new value x^* of X by the corresponding fitted value \hat{y}^*

$$\hat{y}^* = \hat{b}_0 + \hat{b}_1 x^*.$$

How precise is this prediction compared to the actual value y^* ?

Theorem

A $100(1 - \alpha)\%$ confidence interval for y^* is $[\hat{y}^* - \tau, \hat{y}^* + \tau]$ with

$$\tau = t_{1-\frac{\alpha}{2}}(n-2) \sqrt{\frac{RSS}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\|\tilde{\mathbf{x}}\|^2}}.$$

Notation

This interval is called **prediction interval**.

$$\hat{y}^* \pm t_{1-\frac{\alpha}{2}}(n-2) \sqrt{\frac{RSS}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\|\tilde{\mathbf{x}}\|^2}}$$

1. **Introduction**

Prediction intervals are

- wide for small n and if x^* is far away from the mean \bar{x} .
- small for large n and if x^* is close to \bar{x} .

TABLE 1. Summary of the data

Interpolation works much better than extrapolation. Be careful if x^* is outside the range of the data \mathbf{x} .

Prediction intervals

$$\hat{y}^* \pm t_{1-\frac{\alpha}{2}}(n-2) \sqrt{\frac{RSS}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\|\tilde{\mathbf{x}}\|^2}}$$

Example Temperature and pressure $\hat{y} = 90.9 + 0.406x$

Let $x^* = 24^\circ\text{C}$. What is a 95% prediction interval?

The critical value is $t_{0.975}(4) = 2.78$. The prediction interval is given by

$$\begin{aligned} 90.9 + 0.406 \times 24 \pm 2.78 \sqrt{\frac{601}{4}} \sqrt{1 + \frac{1}{6} + \frac{(24 - 25)^2}{1750}} \\ \approx 100.84 \pm 3.75 \end{aligned}$$

or equivalently $[97.79, 104.59]$.