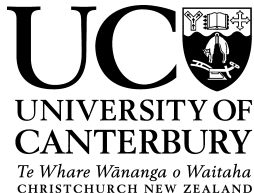


# EMTH211 Statistics

Fabian Dunker

University of Canterbury

2019



## Lecturer

- Fabian Dunker
- [fabian.dunker@canterbury.ac.nz](mailto:fabian.dunker@canterbury.ac.nz)
- Office hour Friday 11pm - 12pm  
and by appointment.
- Erskine 703

# Outline

## Outline

- 1 Descriptive statistics
- 2 Simple linear regression
- 3 Multiple linear regression

## Lecture notes

- Lecture notes by Richard Vale
- Slides

are available on Learn.

# Random variables

For example: measurement with measurement errors and perturbations, or events that occur with uncertainty.

## Definition

- We call an experiment or measurement whose outcome cannot be predicted a **random experiment**.
- A map that assigns real numbers to the outcomes of a random experiment is called **random variable**.

### Remark

You can think of a random variable as a real valued variable whose value is random, i.e. it can be different every time we observe the variable.

## Examples

- rolling a dice, possible outcomes are 1,2,3,4,5, or 6
- tossing a coin, we assign tail = 0, head = 1 to get a real valued variable.
- polls and surveys
- daily mean temperature
- electricity consumption
- life time of a machine in sec, min, hours, or years
- precision of a machine (often influenced by run time)
- measurement with measurement errors:  
temperature, pressure, velocity, voltage, current, ...
- bugs in a software
- Randomised software testing has a random variable as input.  
Hence, the output is also a random variable.

## Data

- Observing a random variable generates **data**.
- A collection of one or more observations is called a **sample**.
- The number of observation in a sample is called the **sample size** and is usually denoted by  $n$ .
- Statistical methods try to learn properties of the random variable by analysing properties of the sample.
- However, properties of a sample can differ from the properties of the random variable due to random perturbations. We need to be careful!

## Example

The average of a sample indicates what the average outcome of a random variable might be.

If we happen to get tail three times in a row when tossing a coin, the average of this sample will be misleading.

- Fair coin 50% head or tail

- 3 times tail =  $\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$

$$= 12.5\%$$

## Data vector

We store each sample in a vector, e.g.

$$\mathbf{x} = (x_1, x_2, \dots, x_n)^T,$$

where  $x_1, x_2, \dots, x_n$  are the individual observations. The sample size is  $n$ .

## Notation

- $\mathbf{x}$ ,  $\mathbf{y}$ ,  $\mathbf{b}$  vectors on the slides.
- $\underline{x}$ ,  $\underline{y}$ ,  $\underline{b}$  vectors on the white board.
- $\bar{x}$ ,  $\bar{y}$ ,  $\bar{b}$  mean of a data vector (next lecture).



• Vectors on the slides bold letters

(←)

$\tilde{x}, x$

$$x = (x_1, x_2, \dots, x_n)^T$$

mean on the board  $\tilde{x}$

## Two ore more random variables

Often two or more random variables are observed simultaneously.

**Example** Temperature and pressure in a boiler.

Temp (°C)	Pressure (kPa)
0	91
10	95
20	100
30	101
40	107
50	112

$$\mathbf{x} = (0, 10, 20, 30, 40, 50)^{\top}$$

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

**Goals** Describe and quantify by using linear algebra:

- probabilistic properties of one (or more) random variable,
- how variations in one random variable is linked to another random variable.

**Example** for random variables that depend on each other

- *temperature* ( $T$ ) and *pressure* ( $P$ )

Ideal gas law:  $P = cT$  with some constant  $c$ .

- *voltage* ( ~~$V$~~ ) and *current* ( $I$ )

~~$V$~~  =  $R I$  with resistance  $R$ .

- *daily mean temperature* and *electricity consumption* in Christchurch

The lower the temperature the higher the electricity consumption.

(This will be different in other places, e.g. New York)

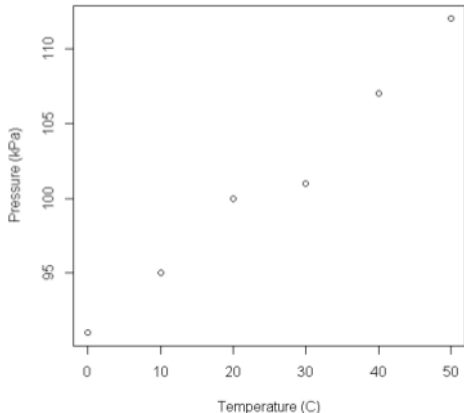
- *live time* of a machine and *production costs* of a machine  
A more durable machine is more expensive to produce.
- *Number of bugs* in a software and *number of randomised software tests*.

## Scatter plot

Plotting the values of two random variables against each other is called a **scatter plot**. Matlab command: `scatter(x,y)`

A scatter plot can give a rough idea about the data.

Temp (°C)	Pressure (kPa)
0	91
10	95
20	100
30	101
40	107
50	112



## 1.2 Sample mean and centering

**Sample mean** # Want to find the average outcome

- An important characteristic of a random variable is its average outcome.
- We will use the average of a sample to estimate the average of the random variable.

### Definition

For a data vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)^\top$  the **sample mean** is defined as

$$\bar{\mathbf{x}} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{1}{n} \sum_{i=1}^n x_i.$$

### Note

- The **average outcome** of the random variable is called **expectation** or **population mean**. It is not the same as the sample mean.
- Often the **sample mean** and the **expectation** are both just called the **mean**. Do not confuse these two concepts!

$$\underline{x} = (x_1 \ x_2 \ \dots \ x_n)^T \quad (\text{proofs})$$

$$\underline{y} = (y_1 \ y_2 \ \dots \ y_n)^T$$

(1):

$$\begin{aligned} \widetilde{x+y} &= \frac{1}{n} \sum_{i=1}^n (x_i + y_i) = \left( \frac{1}{n} \sum_{i=1}^n x_i \right) + \left( \frac{1}{n} \sum_{j=1}^n y_j \right) \\ &= \widetilde{x} + \widetilde{y} \end{aligned}$$

(2):

$$c \in \mathbb{R} \quad (c \text{ is linear})$$

$$\begin{aligned} \widetilde{cx} &= \frac{1}{n} \sum_{i=1}^n cx_i = c \cdot \frac{1}{n} \sum_{i=1}^n x_i \\ &= c \widetilde{x} \end{aligned}$$

(3)

$$\underline{1} = \left( 1 \cdot 1 \ \dots \ 1 \right)^T \Big\} n \text{ times}$$

$$\begin{aligned} \underline{x} + c \times \underline{1} &= (x_1 + c) + (x_2 + c) \ \dots + (x_n + c)^T \\ &= \underline{x} + c \end{aligned}$$

$$T = \frac{1}{n} \sum_{i=1}^n \cdot 1 = \frac{n}{n} = 1$$

## Sample mean

### Properties

Let  $\mathbf{x}, \mathbf{y}$  be two vectors of the same length  $n$  and  $c \in \mathbb{R}$ .

(1) •  $\overline{\mathbf{x} + \mathbf{y}} = \bar{\mathbf{x}} + \bar{\mathbf{y}}$

(2) •  $\overline{c\mathbf{x}} = c\bar{\mathbf{x}}$

(3) • Hence, the sample mean is a linear map from  $\mathbb{R}^n$  to  $\mathbb{R}$ .

### Example

- $\mathbf{x}$  temperature values in degree Celsius
- $\mathbf{x} + 237.5 \times \mathbf{1}$  temperature values in degree Kelvin
- Hence, the mean temperature in Kelvin is

$$\overline{\mathbf{x} + 237.5 \times \mathbf{1}} = \bar{\mathbf{x}} + \overline{237.5 \times \mathbf{1}} = \bar{\mathbf{x}} + 237.5 \times \bar{\mathbf{1}} = \bar{\mathbf{x}} + 237.5$$

### Notation

$\mathbf{1} = (1, 1, \dots, 1)^\top$  with length  $n$  (sample size).

## Sample mean

### Remark

The sample mean is not necessarily a good guess for future observations.

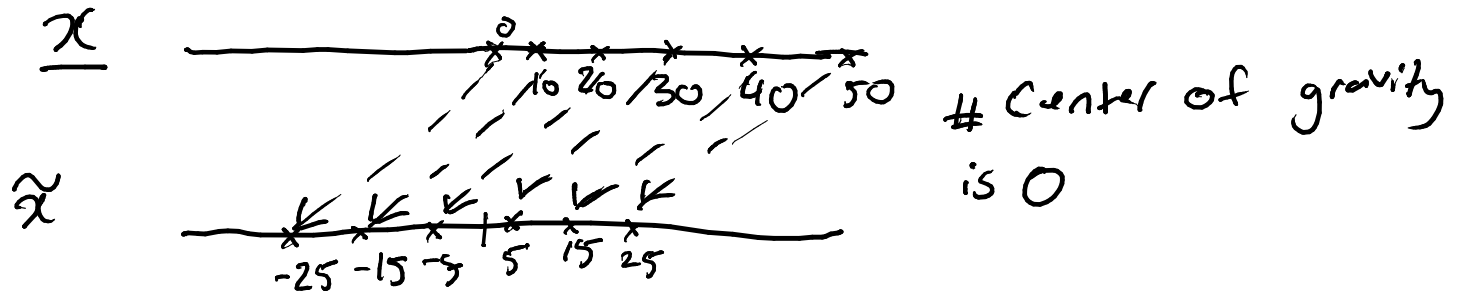
### Example

- The sample mean for a fair dice will be about 3.5. This is an impossible result for future observations.
- However, the average daily max temperature in Christchurch is  $16.8^{\circ}\text{C}$ . We have many days with a similar max temperature.



# # Centering

$$\tilde{x} = x - \bar{x} \mathbf{1}$$



$$\begin{aligned} \overline{\tilde{x}} &= \overline{x - \bar{x} \mathbf{1}} = \bar{x} - \bar{\bar{x} \mathbf{1}} \\ &= \bar{x} - \bar{x} \uparrow \mathbf{1} = \bar{x} - \bar{x} = 0 \end{aligned}$$

# • Last Lecture

- random variables

- sample  $x_1, x_2, \dots, x_n$

- data vector  $\underline{x} = (x_1, \dots, x_n)^T$

- sample mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

- centred data vector

$$\tilde{x} = \underline{x} - \bar{x} \underline{1} \quad \# \quad \underline{1} = (1, 1, \dots, 1)^T$$

$$c \in \mathbb{R} \quad \underline{x} \in \mathbb{R}^n$$

$$\underline{x} + c = \underline{x} + c \underline{1}$$



we want to

keep track of the

fact that this is a vector

## Centering

It is sometimes more convenient to work with a sample that has mean 0.

### Definition

We call

$$\tilde{\mathbf{x}} = \mathbf{x} - \bar{\mathbf{x}}\mathbf{1} = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})^\top$$

the **centred** data vector.

### Example

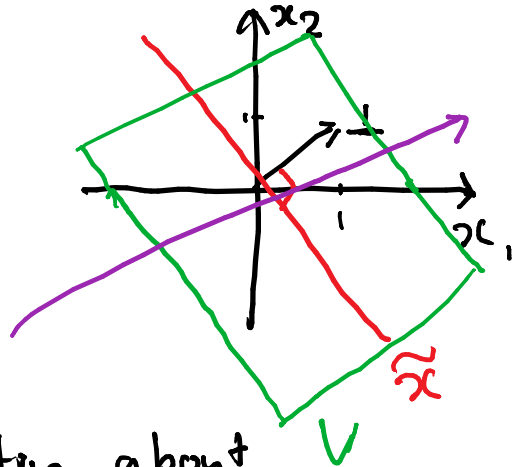
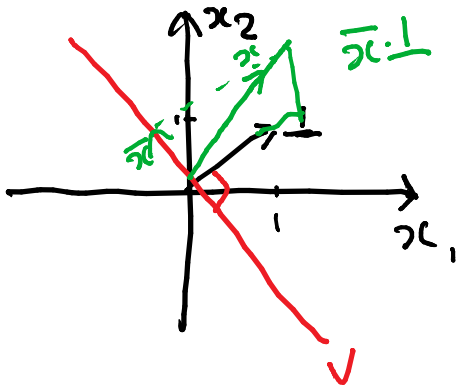
Temperature  $\mathbf{x} = (0, 10, 20, 30, 40, 50)^\top$

$$\begin{aligned}\bar{x} &= \frac{1}{6}(0 + 10 + 20 + 30 + 40 + 50) \\ &= \frac{150}{6} = 25\end{aligned}$$

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

$$\begin{aligned}
 \tilde{x} \cdot \perp &= (x - \bar{x} \perp) \cdot \perp \\
 &= \sum_{i=1}^n x_i - \bar{x} \\
 &= n \cdot \frac{1}{n} \sum_{i=1}^n x_i - \bar{x} \\
 &= n \bar{x} - n \bar{x} = 0
 \end{aligned}$$

• Vector diagram



- $\bar{x}$  contains information about the mean.
- $\tilde{x}$  contains all other information.

## Centering

A centred vector has two important properties.

### Properties

- 1 Centred vectors have 0 mean

$$\widetilde{\mathbf{x}} = 0.$$

- 2 Centred vectors are orthogonal to  $\mathbf{1}$

$$\widetilde{\mathbf{x}} \cdot \mathbf{1} = 0.$$

## Orthogonal decomposition

Centering and mean decompose  $\mathbb{R}^n$  into two orthogonal subspaces

$$V = \left\{ \mathbf{x} : \sum_{i=1}^n x_i = 0 \right\} \quad \text{and} \quad V^\perp = \text{span}\{\mathbf{1}\}$$

with  $\dim(V) = n - 1$  and  $\dim(V^\perp) = 1$ .

Hence, every  $\mathbf{x} \in \mathbb{R}^n$  can be written as

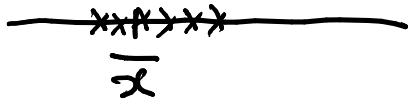
$$\mathbf{x} = \underbrace{\mathbf{x} - \bar{x}\mathbf{1}}_{=\tilde{\mathbf{x}} \in V} + \underbrace{\bar{x}\mathbf{1}}_{\in V^\perp}.$$

## Degrees of freedom

- All centred vectors of length  $n$  form a  $n - 1$  dimensional space.
- A centred vector of length  $n$  contains the same amount of information as a non-centred vector of length  $n - 1$ .
- We say  $\tilde{\mathbf{x}}$  has  $n - 1$  **degrees of freedom**

- Spread

2c



Small spread



Large spread

## 1.3 Spread

The **mean** gives a **rough** idea where the **data** are. Now we want to describe how much the data vary between measurements, i.e. how they spread around the sample mean.

### Sample variance

We measure the spread of  $\mathbf{x}$  around the sample mean  $\bar{\mathbf{x}}$  by the squared Euclidean norm of  $\mathbf{x} - \bar{\mathbf{x}} = \tilde{\mathbf{x}}$  normalised by its *degrees of freedom*.

**Remember**  $\tilde{\mathbf{x}} \in V$  and  $\dim(V) = n - 1$ .

**Definition**

The **sample variance** of  $\mathbf{x}$  is

variance = average distance from  $\bar{\mathbf{x}}$

$\bar{\mathbf{x}}$  = sample mean

$$\text{var}(\mathbf{x}) = \frac{1}{n-1} \|\tilde{\mathbf{x}}\|^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{\mathbf{x}})^2.$$



## Sample variance

### Notation

- We will often just say *variance* instead of *sample variance*.
- A similar concept exists for random variables. It called the *population variance*.

### Example

Temperature  $\mathbf{x} = (0, 10, 20, 30, 40, 50)^\top$

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

$$\|\tilde{\mathbf{x}}\|^2 = 25^2 + 15^2 + 5^2 + 5^2 + 15^2 + 25^2 = 1750$$

$$\text{var}(\mathbf{x}) = \frac{1}{6-1} \|\tilde{\mathbf{x}}\|^2 = \frac{1750}{5} = 350.$$

• check

$$\text{Var}(\underline{x} + c\underline{1}) = \frac{1}{n-1} \|\widetilde{\underline{x} + c\underline{1}}\|^2$$

$$\begin{aligned}\widetilde{\underline{x} + c\underline{1}} &= \underline{x} + c\underline{1} - (\overline{\underline{x} + c\underline{1}})\underline{1} \\ &= \underline{x} + c\underline{1} - (\bar{x} + c)\underline{1} \\ &= \underline{x} - \bar{x}\underline{1} = \hat{\underline{x}}\end{aligned}$$

$$\text{var}(\underline{x} + c\underline{1}) = \text{var}(\underline{x})$$

$$\begin{aligned}\text{var}(c\underline{x}) &= \frac{1}{n-1} \|c\underline{x}\|^2 \\ &= \frac{1}{n-1} c^2 \|\underline{x}\|^2 \\ &= c^2 \text{var}(\underline{x})\end{aligned}$$

## Sample variance

### Properties

Let  $c \in \mathbb{R}$ .

- $\text{var}(\mathbf{x} + c\mathbf{1}) = \text{var}(\mathbf{x})$  translation invariance
- $\text{var}(c\mathbf{x}) = c^2 \text{var}(\mathbf{x})$
- The smaller the variance the closer are  $x_1, x_2, \dots, x_n$  to the sample mean  $\bar{x}$ .

### Example

Assume  $\mathbf{x}$  are some distances measured in metre and  $\mathbf{y}$  are the same distances measured in centimetre.

- $\bar{y} = 100\bar{x}$

# fixed by standard deviation

- $\text{var}(\mathbf{y}) = 10000 \text{var}(\mathbf{x})$

$$\begin{aligned}
 \text{sd}(x) &= \sqrt{\text{Var}(Cx)} = \sqrt{\frac{1}{n-1} \|Cx\|^2} \\
 &= \sqrt{\frac{1}{n-1} C^2 \|x\|^2} \\
 &= \sqrt{C^2 \text{Var}(x)}
 \end{aligned}$$

$$\text{Var}(x) = \sigma^2_x \quad \sigma = \text{sigma}$$

$$\text{sd}(x) = \sigma_x$$

## Standard deviation

The quadratic change of the variance when changing the units of measurements is counterintuitive and does sometimes cause problems.

*# just square root the variance*

**Definition** The **standard deviation** of  $\mathbf{x}$  is

$$sd(\mathbf{x}) = \frac{1}{\sqrt{n-1}} \|\tilde{\mathbf{x}}\| = \sqrt{var(\mathbf{x})}.$$

**Properties** Let  $c \in \mathbb{R}$ .

- $sd(\mathbf{x} + c\mathbf{1}) = sd(\mathbf{x})$  translation invariance
- $sd(c\mathbf{x}) = |c|sd(\mathbf{x})$

The standard deviation is in the same units of measurement as the data  $\mathbf{x}$ . Changing from metre to centimetre results in multiplication by 100.

## Sample variance and standard deviation

### Notation

Some textbooks use the following notation

- $s_x^2 = \text{var}(\mathbf{x})$
- $s_x = \text{sd}(\mathbf{x})$ .

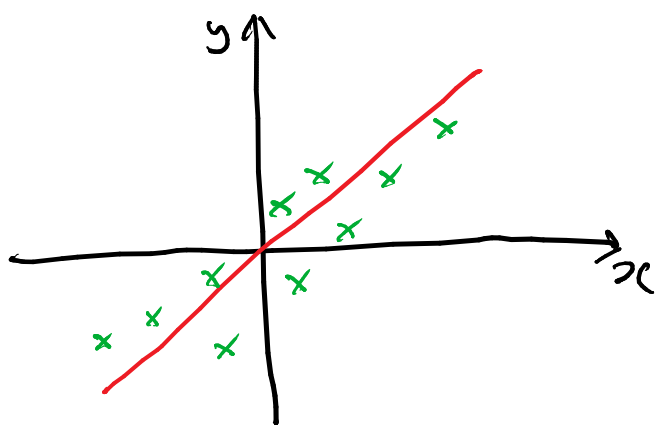
### Remark

Other measures for the spread such as

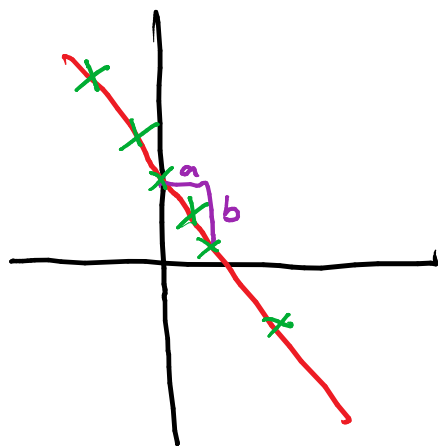
$$\frac{1}{n-1} \sum_{i=1}^n |x_i - \bar{\mathbf{x}}| = \frac{1}{n-1} \|\mathbf{x} - \bar{\mathbf{x}}\|_1 = \frac{1}{n-1} \|\tilde{\mathbf{x}}\|_1$$

are possible. But they would lead to complicated formulas next week! In most fields of statistics *var* and *sd* are used for convenience. However, in image processing and big data other measures of the spread are sometimes useful.

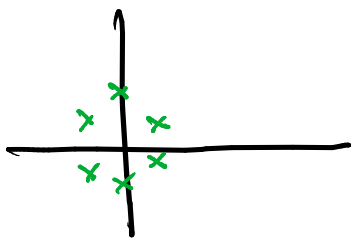
## Scatter plot



## linear relationship



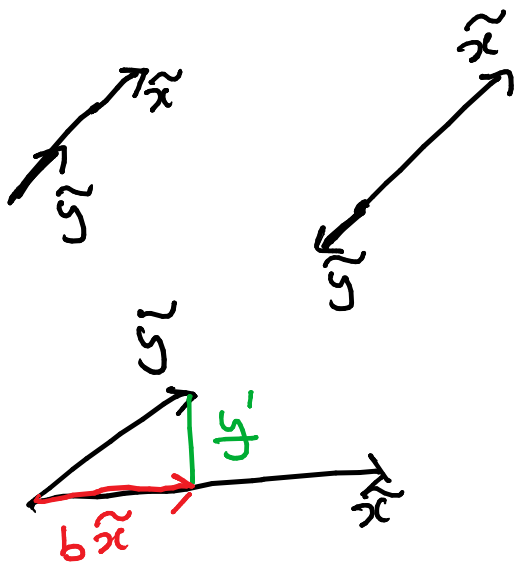
## non-linear relationships



• for centred vectors

$$\begin{aligned}\tilde{y} &= \underline{y} - \bar{y} \underline{1} = (b\underline{x} + a\underline{1}) - (\overline{b\underline{x} + a\underline{1}}) \underline{1} \\ &= b\underline{x} + a\underline{1} - b\bar{x}\underline{1} - a\underline{1} \\ &= b\underline{x} - b\bar{x} \\ &= b\tilde{x}\end{aligned}$$

## vector diagram



• perfect linear relation

## 1.4 Covariance and correlation

We want to quantify how two data vectors  $\mathbf{x}$  and  $\mathbf{y}$  of the same length  $n$  behave together. Can  $\mathbf{x}$  predict  $\mathbf{y}$  in a linear way?

$$\mathbf{y} = b\mathbf{x} + a\mathbf{1}, \quad a, b \in \mathbb{R}$$

It is easier to work with the centred vectors. If the linear equation holds, then

$$\tilde{\mathbf{y}} = b\tilde{\mathbf{x}}.$$

The  $a\mathbf{1}$  cancels out.

Such a perfect linear relation is often too much to ask for. But

$$\tilde{\mathbf{y}} = b\tilde{\mathbf{x}} + \mathbf{y}' \quad \text{with } \mathbf{y}' \in \text{span}\{\tilde{\mathbf{x}}\}^\perp$$

holds.

**Idea** The  $b$  quantifies the linear dependency while  $\mathbf{y}'$  is a component of  $\tilde{\mathbf{y}}$  which is linear independent of  $\tilde{\mathbf{x}}$ .



check

$$\text{cov}(\underline{x}, \underline{y}) = \frac{1}{n-1} \tilde{\underline{x}} \cdot \hat{\underline{y}} = \frac{1}{n-1} \tilde{\underline{x}} (b\tilde{\underline{x}} + \underline{y}')$$

with  $\underline{y}' \cdot \tilde{\underline{x}} = 0$

$$= \frac{1}{n-1} \tilde{\underline{x}} \cdot b\tilde{\underline{x}} + \underbrace{\tilde{\underline{x}} \cdot \underline{y}'}_{=0 \text{ (orthogonal)}}$$

$$= b \frac{1}{n-1} \tilde{\underline{x}} \cdot \tilde{\underline{x}}$$

$$= b \frac{1}{n-1} \|\tilde{\underline{x}}\|^2 = b \text{var}(\underline{x})$$

## Sample covariance

### Definition

The **sample covariance** of  $\mathbf{x}$  and  $\mathbf{y}$  is

# dot product  
- projection (orthogonal)

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{n-1} \tilde{\mathbf{x}} \cdot \tilde{\mathbf{y}} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

### Notation

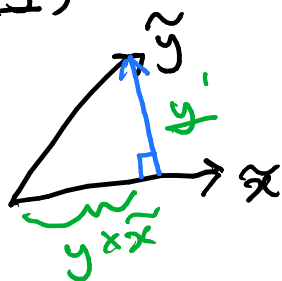
- We will often say *covariance* instead of *sample covariance*.
- Some textbooks use the notation  $c_{\mathbf{x}, \mathbf{y}} = \text{cov}(\mathbf{x}, \mathbf{y})$ .

**Properties** Let  $c \in \mathbb{R}$

- $\text{cov}(\mathbf{x}, \mathbf{y}) = b \times \text{var}(\mathbf{x})$
- $\text{cov}(\mathbf{x}, \mathbf{x}) = \text{var}(\mathbf{x})$
- $\text{cov}(\mathbf{x}, \mathbf{y}) = \text{cov}(\mathbf{y}, \mathbf{x})$       symmetry
- $\text{cov}(c\mathbf{x}, \mathbf{y}) = c \times \text{cov}(\mathbf{x}, \mathbf{y}) = \text{cov}(\mathbf{x}, c\mathbf{y})$
- $\text{cov}(\mathbf{x} + c\mathbf{1}, \mathbf{y}) = \text{cov}(\mathbf{x}, \mathbf{y}) = \text{cov}(\mathbf{x}, \mathbf{y} + c\mathbf{1})$   
translation invariance

# Last lecture

- $\tilde{x} \perp \bar{x}$  are orthogonal
- $x = \tilde{x} + \bar{x}$
- $\tilde{x}$  has  $n-1$  degrees of freedom
- $\text{var}(x) = \frac{1}{n-1} \|x\|^2$   
 $\text{var}(cx) = \frac{n-1}{c^2} \text{var}(x)$
- $\text{sd}(x) = \sqrt{\text{var}(x)}$   
 $\text{sd}(cx) = |c| \text{sd}(x)$
- $\text{cov}(x, y)$



## check properties

$$\text{cov}(x, x) = \frac{1}{n-1} \tilde{x} \cdot \tilde{x} = \frac{1}{n-1} \|\tilde{x}\|^2 = \text{var}(x)$$

$$\text{cov}(x, y) = \frac{1}{n-1} \tilde{x} \cdot \tilde{y} = \frac{1}{n-1} \tilde{y} \cdot \tilde{x} = \text{cov}(y, x)$$

$$\text{cov}(cx, y) = \frac{1}{n-1} c \tilde{x} \cdot \tilde{y}$$

$$= c \frac{1}{n-1} \tilde{x} \cdot \tilde{y}$$

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

$$= \text{cov}(x, cy)$$

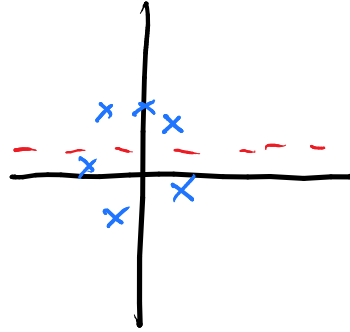
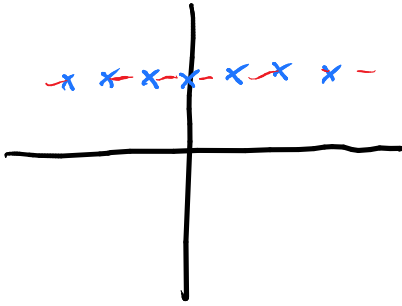
by symmetry

$$\text{cov}(x + c\bar{x}, y) = \text{cov}(x, y)$$

because  $\widetilde{x + c\bar{x}} = \tilde{x}$

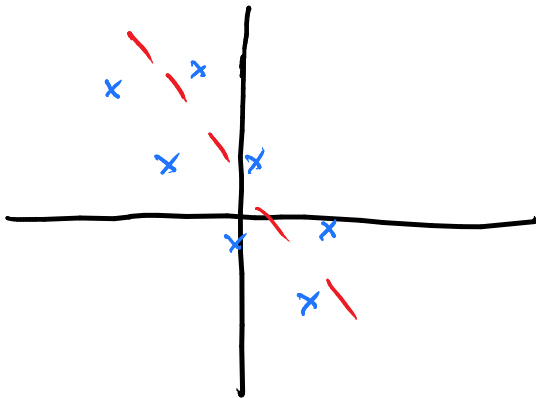
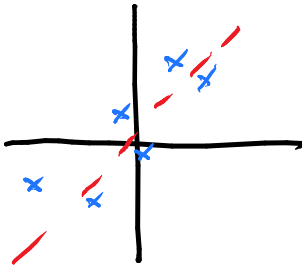
scatter plot

$$b = 0$$



$$\text{cov}(x, y) = 0$$

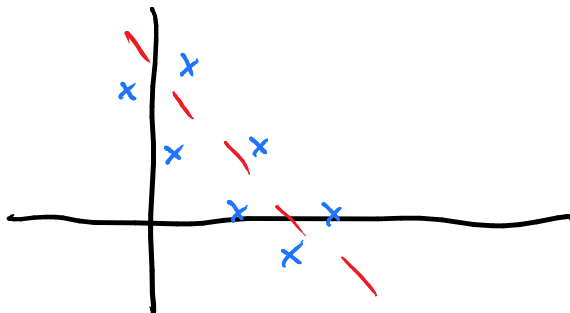
$$b > 0 \quad \text{cov}(x, y) > 0$$



$$b < 0 \quad \text{cov}(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\overset{\nearrow V}{\cancel{U}} = RI$$

$$\text{var}(\underline{u}, i) > 0$$



$$\text{cov}(\underline{t}, \underline{e}) < 0$$

## Sample covariance

### Interpretation

- The covariance quantifies the linear relation between  $\mathbf{x}$  and  $\mathbf{y}$ .
- If they do not even have an approximately linear relation,  $\text{cov}(\mathbf{x}, \mathbf{y}) = 0$ . In this case  $\mathbf{x} \cdot \mathbf{y} = 0$ .
- If  $\mathbf{x}$  and  $\mathbf{y}$  tend to be above their respective means at the same time and below their respective means at the same time, then  $\text{cov}(\mathbf{x}, \mathbf{y}) > 0$ .
- If  $\mathbf{x}$  and  $\mathbf{y}$  tend to go into opposite directions relative to their respective means, then  $\text{cov}(\mathbf{x}, \mathbf{y}) < 0$ .
- The absolute value of  $\text{cov}(\mathbf{x}, \mathbf{y})$  depends on  $\text{var}(\mathbf{x})$  and by symmetry also on  $\text{var}(\mathbf{y})$ .

### Example

- Voltage and current in a wire will have positive covariance.
- Electricity consumption and temperature in Christchurch will have negative covariance.

## Sample correlation

Let us remove the dependency on  $var(\mathbf{x})$  and  $var(\mathbf{y})$  from the covariance.

**Definition** The **sample correlation** of  $\mathbf{x}$  and  $\mathbf{y}$  is

$$corr(\mathbf{x}, \mathbf{y}) = \frac{cov(\mathbf{x}, \mathbf{y})}{sd(\mathbf{x})sd(\mathbf{y})} = \frac{\tilde{\mathbf{x}} \cdot \tilde{\mathbf{y}}}{\|\tilde{\mathbf{x}}\| \|\tilde{\mathbf{y}}\|} .$$

## Notation

- We will often just say *correlation* instead of *sample correlation*.
- This concept is also called **Pearson correlation**.
- Some textbooks use the notation  $r_{\mathbf{x}, \mathbf{y}} = corr(\mathbf{x}, \mathbf{y})$ .

- $\text{cov}(\underline{x}, \underline{y}) \neq \text{corr}(\underline{x}, \underline{y})$

have the same sign because we divide by  $\text{sd}(\underline{x})\text{sd}(\underline{y}) > 0$

- $\text{corr}$  is symmetric because  $\text{cov}$  is symmetric.

- $\text{corr}(\underline{x} + c\underline{1}, \underline{y}) = \text{corr}(\underline{x}, \underline{y})$

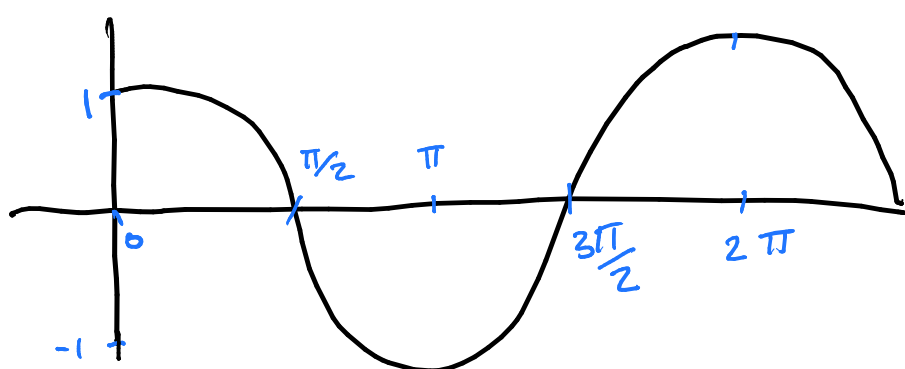
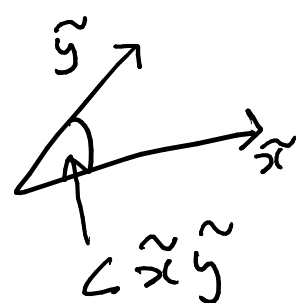
because  $\widetilde{\underline{x} + c\underline{1}} = \tilde{\underline{x}}$

$$\begin{aligned} \text{corr}(c\underline{x}, \underline{y}) &= \frac{c\underline{x} \cdot \underline{y}}{\|c\underline{x}\| \|\underline{y}\|} \\ &= \frac{c}{|c|} \frac{\tilde{\underline{x}} \cdot \underline{y}}{\|\tilde{\underline{x}}\| \|\underline{y}\|} \\ &= \text{Sign}(c) \end{aligned}$$

- $\underline{x} \cdot \underline{y} = \|\underline{x}\| \|\underline{y}\| \cos(\angle \underline{x} \underline{y})$

Hence,

$$\text{corr}(\underline{x}, \underline{y}) = \cos(\angle \underline{x} \underline{y})$$

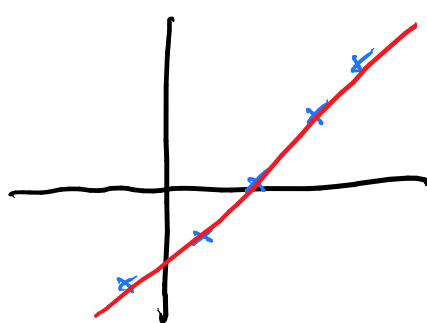


- $\text{corr}(\underline{x}, \underline{y}) = 1 \rightarrow \angle \underline{x}, \underline{y} = 0$



perfect linear relation

$$\underline{y} = b\underline{x} \quad b > 0$$

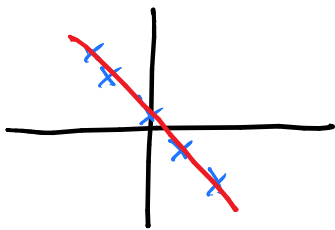
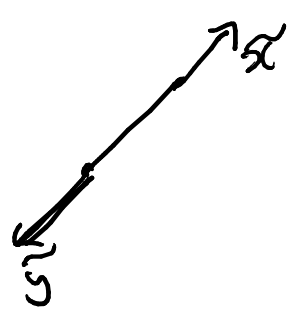


$$\angle \underline{x} \underline{y} = \pi$$

perfect linear relation

$$\underline{y} = b\underline{x} \quad b < 0$$

- $\text{corr}(\underline{x}, \underline{y}) = -1$





## Sample correlation

**Properties** Let  $c \in \mathbb{R}$ .

- $\text{corr}(\mathbf{x}, \mathbf{y})$  and  $\text{cov}(\mathbf{x}, \mathbf{y})$  have the same sign.  
 $\text{corr}(\mathbf{x}, \mathbf{y}) = 0$  if and only if  $\text{cov}(\mathbf{x}, \mathbf{y}) = 0$ .
- $\text{corr}(\mathbf{x}, \mathbf{y}) = \text{corr}(\mathbf{y}, \mathbf{x})$  symmetry
- $\text{corr}(\mathbf{x} + c\mathbf{1}, \mathbf{y}) = \text{corr}(\mathbf{x}, \mathbf{y}) = \text{corr}(\mathbf{x}, \mathbf{y} + c\mathbf{1})$   
translation invariance
- $\text{corr}(c\mathbf{x}, \mathbf{y}) = \text{sign}(c)\text{corr}(\mathbf{x}, \mathbf{y}) = \text{corr}(\mathbf{x}, c\mathbf{y})$   
scale invariance
- $\text{corr}(\mathbf{x}, \mathbf{y}) = \cos(\angle \tilde{\mathbf{x}}, \tilde{\mathbf{y}})$
- $-1 \leq \text{corr}(\mathbf{x}, \mathbf{y}) \leq 1$ .

Here  $\angle \tilde{\mathbf{x}}, \tilde{\mathbf{y}}$  is the angle between  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$  and

$$\text{sign}(c) = \begin{cases} -1 & \text{if } c < 0 \\ 0 & \text{if } c = 0 \\ 1 & \text{if } c > 0. \end{cases}$$

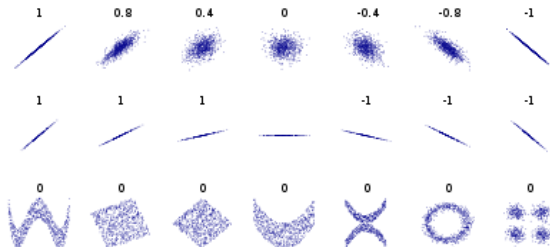
## Sample correlation

### Note

- $\text{corr}(\mathbf{x}, \mathbf{y}) = 1$  if  $\tilde{\mathbf{y}} = b\tilde{\mathbf{x}}$  with  $b > 0$ , i.e.  $\angle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} = 0$
- $\text{corr}(\mathbf{x}, \mathbf{y}) = -1$  if  $\tilde{\mathbf{y}} = b\tilde{\mathbf{x}}$  with  $b < 0$ , i.e.  $\angle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} = \pi$  ( $180^\circ$ ).

### Interpretation

The correlation measures how perfect the linear relation  $\tilde{\mathbf{y}} = b\tilde{\mathbf{x}} + \mathbf{y}'$  is.



Be careful when interpreting **cov** and **corr**! Only approximately linear relations are detected.

## Covariance and correlation

### Example

Temperature and pressure in a boiler.

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

$$\tilde{\mathbf{y}} = (-10, -6, -1, 0, 6, 11)^\top$$

$$\begin{aligned} \text{cov}(\mathbf{x}, \mathbf{y}) &= \frac{1}{6-1} \left( (-25)(-10) + (-15)(-6) + (-5)(-1) \right. \\ &\quad \left. + (5)(0) + (15)(6) + (25)(11) \right) \\ &= \frac{710}{5} = 142 \end{aligned}$$

$$\begin{aligned} \text{corr}(\mathbf{x}, \mathbf{y}) &= \frac{(-25)(-10) + (-15)(-6) + (-5)(-1) + (5)(0) + (15)(6) + (25)(11)}{\sqrt{25^2 + 15^2 + 5^2 + 5^2 + 15^2 + 25^2} \sqrt{10^2 + 6^2 + 1^2 + 0^2 + 6^2 + 11^2}} \\ &= \frac{710}{\sqrt{1750} \sqrt{294}} \approx 0.9898 \end{aligned}$$

# • Last Lecture

- $\text{cov}(\underline{x}, \underline{y}) = \frac{\tilde{\underline{x}} \cdot \tilde{\underline{y}}}{n-1}$

- translation invariance

- Symmetric

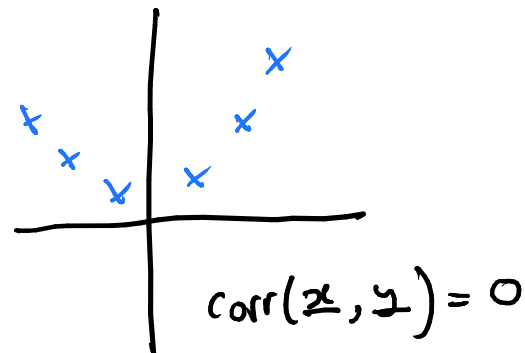
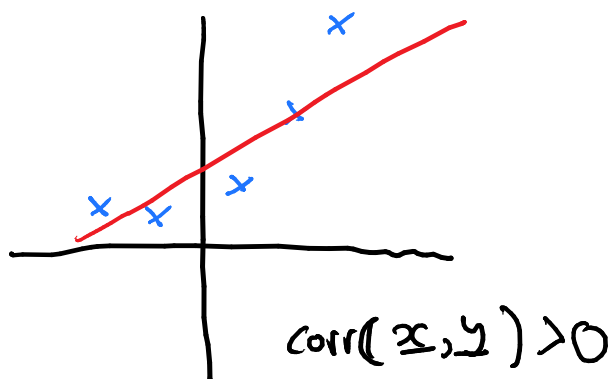
- $\text{cov}(c\underline{x}, \underline{y}) = c \text{cov}(\underline{x}, \underline{y})$

- $\text{corr}(\underline{x}, \underline{y}) = \frac{\tilde{\underline{x}} \cdot \tilde{\underline{y}}}{\|\tilde{\underline{x}}\| \|\tilde{\underline{y}}\|} = \cos(\angle \tilde{\underline{x}}, \tilde{\underline{y}})$

- translation invariance

- Scale invariance

- Symmetric



el-monthly electricity in kWh

temp-monthly mean temperature in °F

house in upstate NY

covariance matrix

$$\begin{pmatrix} \text{cov}(\text{el}, \text{temp}) & \text{cov}(\underline{x}, \underline{y}) \\ \text{cov}(\underline{x}, \underline{y}) & \text{var}(\underline{y}) \end{pmatrix}$$

$$^{\circ}\text{C} = \frac{5}{9} (^{\circ}\text{F} - 32) \quad \# \text{ scaling}$$