

# Introduction to statistical thinking

Catalina Vallejos (Alan Turing Institute), Aaron Lun (CRUK-CI)

2016-11-17

# Why statistics?

## DATA: BY THE NUMBERS



[www.phdcomics.com](http://www.phdcomics.com)

JORGE CHAM © 2004

# Why statistics?

**Statistics** is at the core of modern research

**What is statistics?**

# Why statistics?

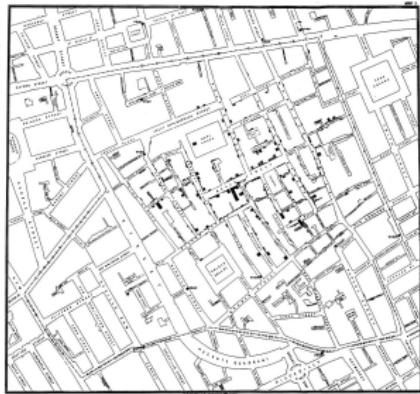
**Statistics** is at the core of modern research

## What is statistics?



→ everything from experimental design to figure preparation!

# What can we do with statistics?

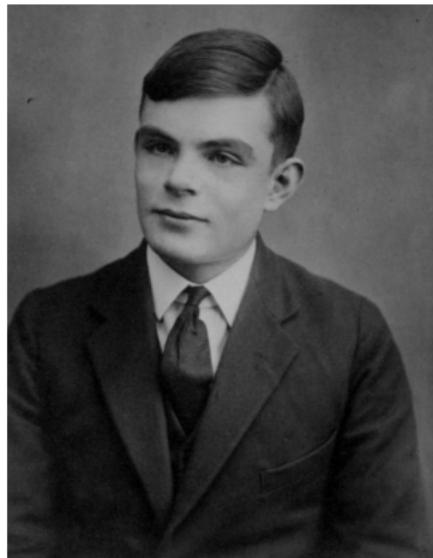


Public domain image

During the Soho cholera outbreak in 1854, John Snow used statistics to find an **association** between the quality of the water source and cholera cases

John Snow is considered one of the fathers of modern epidemiology

# What can we do with statistics?



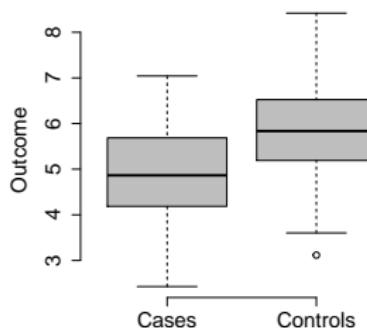
Public domain image

During WWII, Alan Turing used advanced statistics to **predict** the meaning of Enigma messages

Alan Turing is considered one of the fathers of modern computers

# What can we do with statistics?

## Comparisons between two groups



For example,

- ▶ Case/control studies
- ▶ Differential expression studies

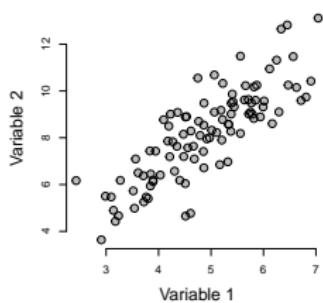
Is this difference *statistically significant*?

# What can we do with statistics?

## Finding associations between variables

For example,

- ▶ Dose level vs survival time
- ▶ Association studies (e.g. GWAS)
- ▶ Co-expression networks



Is there an association between  $X$  and  $Y$ ?

Which associations are *statistically significant*?

# Outline for the course

- ▶ **Introduction to statistics**
  - ▶ Types of data
  - ▶ Descriptive statistics
  - ▶ Practical: descriptive statistics in Rmarkdown
  - ▶ Statistical inference
- ▶ **Hypothesis testing (in R)**
- ▶ **Linear regression (in R)**

## Introduction to statistics

## Types of data

Prior to any analysis, it is important to know what type of data is available

# Types of data

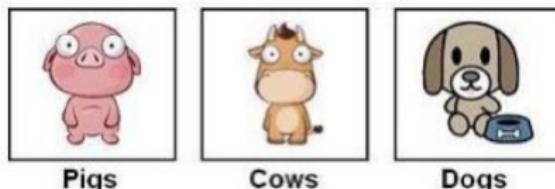
Prior to any analysis, it is important to know what type of data is available

Data can be classified into different **types of variables**:

- ▶ Categorical (nominal)
- ▶ Categorical (ordinal)
- ▶ Discrete
- ▶ Continuous

## Types of data: categorical (nominal)

To describe categories that cannot be ordered



Pigs

Cows

Dogs

Source: <http://www.restore.ac.uk>

Requirements:

- ▶ Same value assigned to all the members of level
- ▶ Same value not assigned to different levels
- ▶ Each observation only assigned to one level

Examples: gender, yes/no answers, surgery type, cancer type, eye colour, dead/alive, ethnicity, etc.

## Types of data: categorical (ordinal)

To describe categories that have a logical order



Requirements:

- ▶ Mutually exclusive fixed categories
- ▶ Implicit order
- ▶ Can say one category higher than another (but not how much)

Examples: stress level (1 = low, . . . , 7 = high), pain level  
(low/medium/high), education level (primary, secondary, . . . ), etc.

## Types of data: discrete

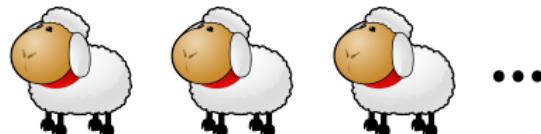


Image modified from Free Stock Photos

### Requirements:

- ▶ Fixed categories, can only take certain values
- ▶ Like ordinal but with well-defined distances
- ▶ Basically, anything counted (cardinal): how much?

Examples: number of tumours, hospital admissions, etc.

Sometimes treated as continuous if range is large!

## Types of data: continuous



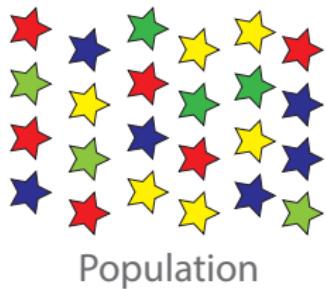
Requirements:

- ▶ Observations can take any value
- ▶ May have finite or infinite range
- ▶ Given any two values, one fits between

Examples: height, weight, blood pressure, temperature etc.

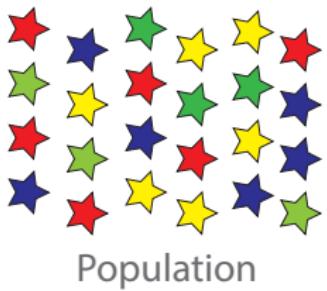
# Population versus sample

**Census** records information from every subject within a population



# Population versus sample

**Census** records information from every subject within a population

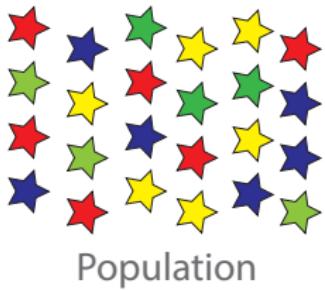


**Scientific studies** typically rely on a subset of a population



# Population versus sample

**Census** records information from every subject within a population



Population

**Scientific studies** typically rely on a subset of a population



Population

A good sample must be a good representative of the population!

# Descriptive statistics versus statistical inference

## Descriptive statistics

summarizes the information  
available in the data

6 x  +  
6 x  +  
6 x  +  
6 x  +

# Descriptive statistics versus statistical inference

## Descriptive statistics

summarizes the information available in the data

6 x  +  
6 x  +  
6 x  +  
6 x  +

## Statistical inference

extrapolates sample information to population level

$$P(\text{Red Star}) = 1/4$$

$$P(\text{Green Star}) = 1/4$$

$$P(\text{Blue Star}) = 1/4$$

$$P(\text{Yellow Star}) = 1/4$$

# Descriptive statistics versus statistical inference

## Descriptive statistics

summarizes the information available in the data

6 x  +  
6 x  +  
6 x  +  
6 x  +

## Statistical inference

extrapolates sample information to population level

$$P(\text{Red Star}) = 1/4$$

$$P(\text{Green Star}) = 1/4$$

$$P(\text{Blue Star}) = 1/4$$

$$P(\text{Yellow Star}) = 1/4$$

A good sample must be a good representative of the population!

## Descriptive analyses

## Descriptive analyses

There are two basic forms of data summary

- ▶ Numerical: frequency tables, summary measures, etc
- ▶ Graphical: histograms, scatter-plots, boxplots, etc

We need to choose the summary depending on the data type!

# Descriptive analyses: categorical variables

## Numerical summary

Category	Frequency	Cumulative frequency	Relative frequency	Cum. rel. frequency
★	6	6	0.25	0.25
★	6	12	0.25	0.50
★	6	18	0.25	0.75
★	6	24	0.25	1.00

# Descriptive analyses: categorical variables

## Numerical summary

Category	Frequency	Cumulative frequency	Relative frequency	Cum. rel. frequency
★	6	6	0.25	0.25
★	6	12	0.25	0.50
★	6	18	0.25	0.75
★	6	24	0.25	1.00

## Graphical summary



## Descriptive analysis: discrete/continuous variables

**Frequency tables** are not useful  
unless we *categorize* the data

## Descriptive analysis: discrete/continuous variables

**Frequency tables** are not useful  
unless we *categorize* the data

Example:

Birth weight	Rel. freq
Very low (<1.5 kg.)	0.05
Low (1.5-2.5 kg.)	0.10
Normal ( $\geq 2.5$ kg.)	0.85

## Descriptive analysis: discrete/continuous variables

**Frequency tables** are not useful unless we *categorize* the data

Instead, we typically prefer *summary measures*

Example:

Birth weight	Rel. freq
Very low (<1.5 kg.)	0.05
Low (1.5-2.5 kg.)	0.10
Normal ( $\geq 2.5$ kg.)	0.85

- ▶ Mean
- ▶ Median
- ▶ Mode
- ▶ Variance
- ▶ Quantiles

Which one?

## Descriptive analysis: discrete/continuous variables

Firstly, we can use measures to describe the **location** of the data

## Descriptive analysis: discrete/continuous variables

Firstly, we can use measures to describe the **location** of the data

### Mean ( $\bar{x}$ )

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

### Median ( $x_{50\%}$ )

50% of observations lie below  $x_{50\%}$

# Descriptive analysis: discrete/continuous variables

Firstly, we can use measures to describe the **location** of the data

When the data is normal ...

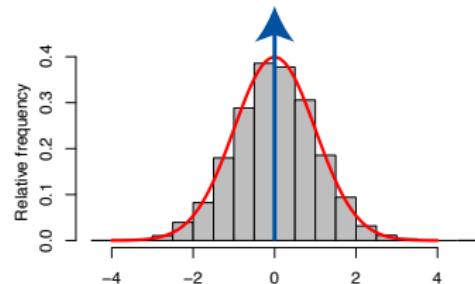
**Mean ( $\bar{x}$ )**

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

**Median ( $x_{50\%}$ )**

50% of observations lie below  $x_{50\%}$

Mean = Median



... but that is not always true

## Descriptive analysis: example

Suppose we record the number of Facebook friends for 7 colleagues:

311, 345, 270, 310, 243, 5300, 11

**Mean** ( $\bar{x}$ )

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_7}{7} = 970$$

**Median** ( $x_{50\%}$ )

$$11, 243, 270, \textcolor{red}{310}, 311, 345, 5300 \Rightarrow x_{50\%} = 310$$

Which one provides a better description for the location of the data?

## Descriptive analysis: example

Now suppose the data is slightly different:

311, 345, 270, 310, 243, 530, 11

**Mean** ( $\bar{x}$ )

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_7}{7} = 289$$

**Median** ( $x_{50\%}$ )

$$11, 243, 270, 310, 311, 345, 530 \Rightarrow x_{50\%} = 310$$

What happened?

## Descriptive analysis: discrete/continuous variables

It is also important to summarize the **spread** of the data

**Standard deviation** ( $\text{sd}(x)$ )

$$\text{sd}(x) = \sqrt{\frac{(x_1 - \bar{x})^2 + \cdots + (x_n - \bar{x})^2}{n}}$$

**Interquartile range** ( $\text{IQR}(x)$ )

$$\text{IQR}(x) = x_{75\%} - x_{25\%}$$

## Descriptive analysis: example

In the first example, we have

**Standard deviation** ( $\text{sd}(x)$ )

$$\text{sd}(x) = \sqrt{\frac{(x_1 - \bar{x})^2 + \dots + (x_7 - \bar{x})^2}{7}} = 1912.57$$

**Interquartile range** ( $\text{IQR}(x)$ )

$$11, \textcolor{red}{243}, 270, \textcolor{red}{310}, 311, \textcolor{red}{345}, 5300 \Rightarrow \text{IQR}(x) = 345 - 243$$

## Descriptive analysis: example

In the first example, we have

### Standard deviation ( $\text{sd}(x)$ )

$$\text{sd}(x) = \sqrt{\frac{(x_1 - \bar{x})^2 + \cdots + (x_7 - \bar{x})^2}{7}} = 1912.57$$

### Interquartile range ( $\text{IQR}(x)$ )

$$11, \textcolor{red}{243}, 270, \textcolor{red}{310}, 311, \textcolor{red}{345}, 5300 \Rightarrow \text{IQR}(x) = 345 - 243$$

In the second example, we have

### Standard deviation ( $\text{sd}(x)$ )

$$\text{sd}(x) = \sqrt{\frac{(x_1 - \bar{x})^2 + \cdots + (x_7 - \bar{x})^2}{7}} = 153.79$$

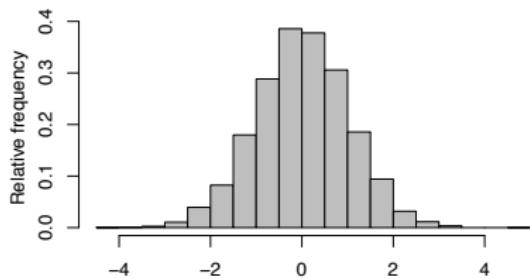
### Interquartile range ( $\text{IQR}(x)$ )

$$11, \textcolor{red}{243}, 270, \textcolor{red}{310}, 311, \textcolor{red}{345}, 530 \Rightarrow \text{IQR}(x) = 345 - 243$$

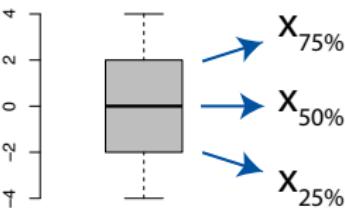
# Descriptive analysis: discrete/continuous variables

We can also summarize the distribution of the data **graphically**

**Histogram**



**Boxplot**



**\*\*Link to 'exploration' shiny app\*\***

# Before we continue: the normal distribution

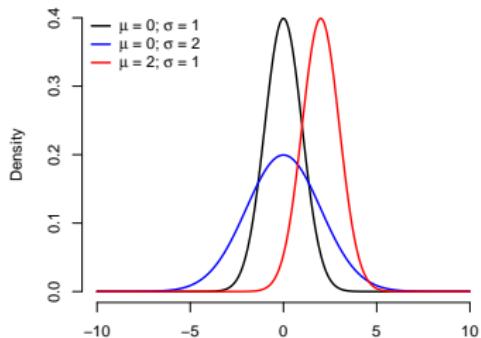
The normal (or Gaussian) distribution is very popular!

Indexed by two parameters:

- ▶  $\mu \rightarrow$  location
- ▶  $\sigma$  (or  $\sigma^2$ )  $\rightarrow$  spread

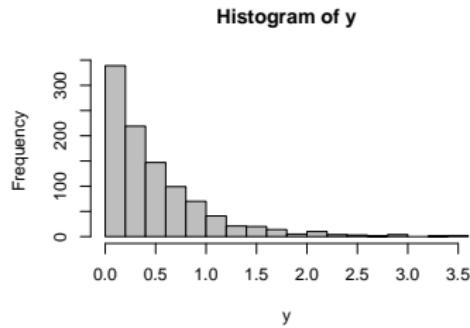
Rule of thumb:

- ▶  $\sim 95\%$  of the data lies between  $(\mu - 2\sigma, \mu + 2\sigma)$
- ▶  $\sim 99\%$  of the data lies between  $(\mu - 3\sigma, \mu + 3\sigma)$

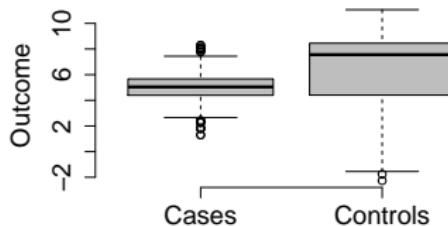
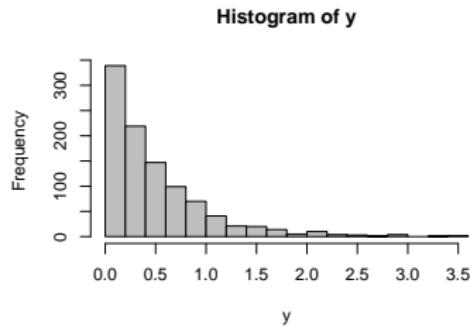


The distribution is symmetric  
around  $\mu$ !

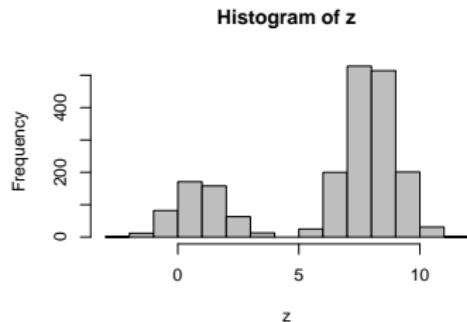
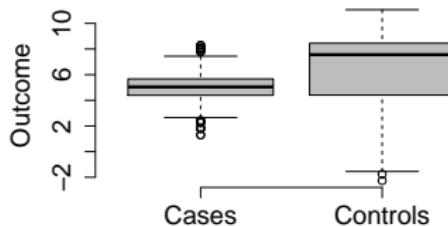
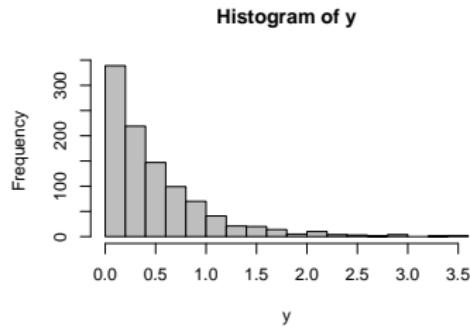
# Descriptive analysis: is the data normally distributed?



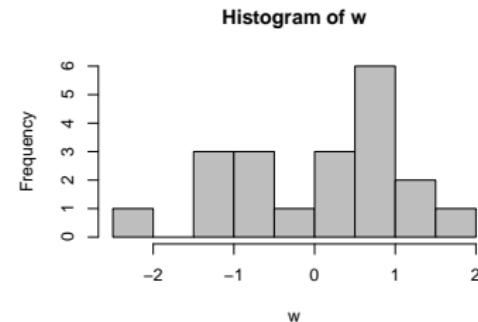
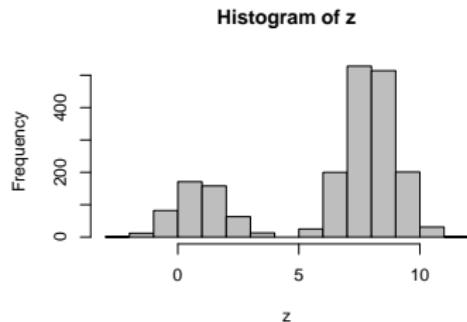
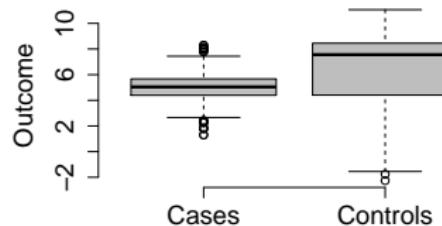
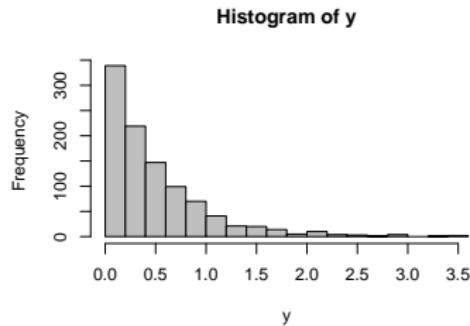
# Descriptive analysis: is the data normally distributed?



# Descriptive analysis: is the data normally distributed?



# Descriptive analysis: is the data normally distributed?

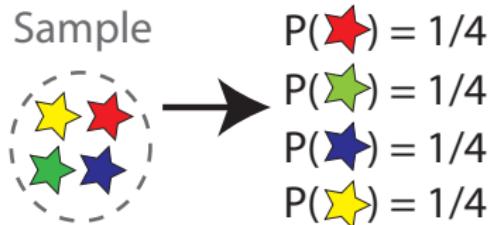


## Statistical inference

# Statistical inference: estimation

## Statistical inference

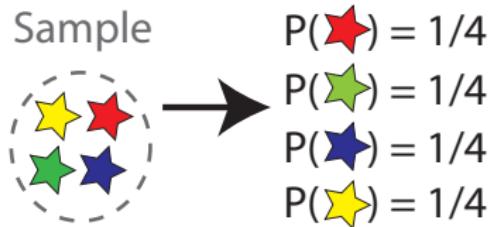
extrapolates sample information  
to population level



# Statistical inference: estimation

## Statistical inference

extrapolates sample information  
to population level

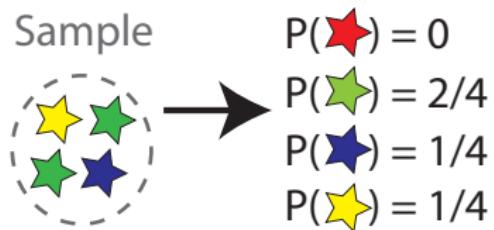


In this example, we want to **estimate** the proportion of red, green, blue and yellow stars in the population based on the observed data

# Statistical inference: estimation

## Statistical inference

extrapolates sample information  
to population level

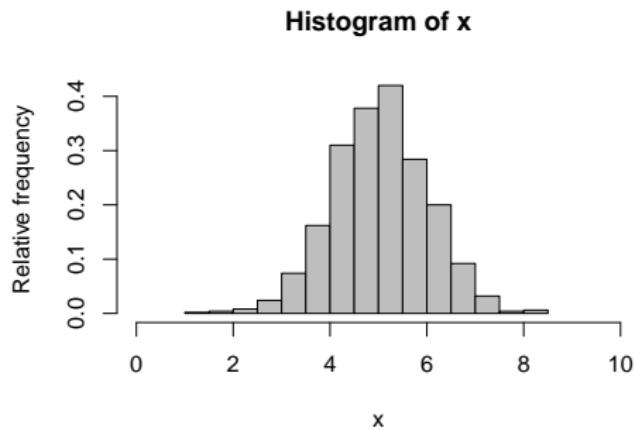


If we observe a different sample,  
we obtain different results

Note: in this sample we didn't  
observe red stars!

## Statistical inference: estimation

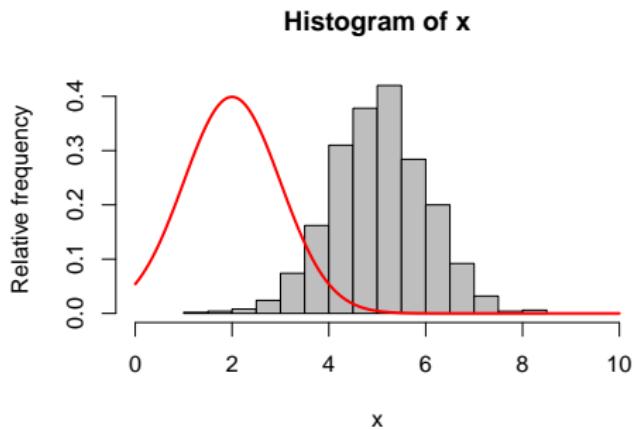
To model the data (e.g. expression of a gene across samples) as being normally distributed with mean  $\mu$  and standard deviation 1.



The aim here is to estimate the value of  $\mu$ !

## Statistical inference: estimation

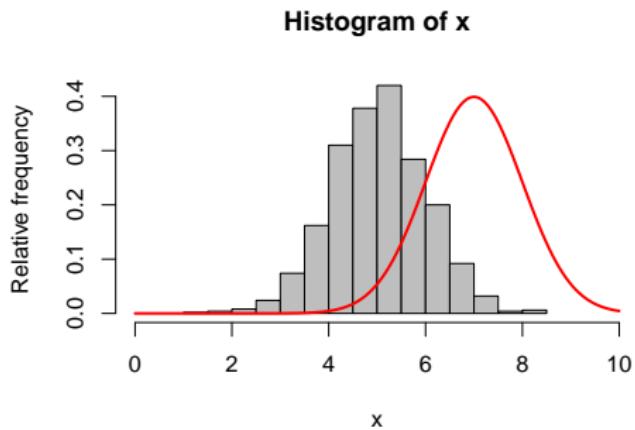
To model the data (e.g. expression of a gene across samples) as being normally distributed with mean  $\mu$  and standard deviation 1.



$$\mu = 2?$$

## Statistical inference: estimation

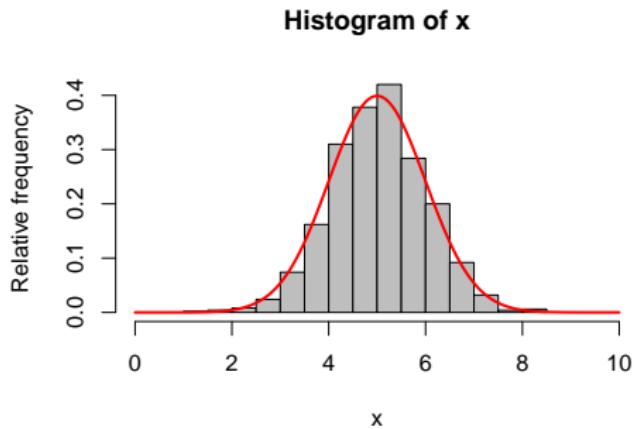
To model the data (e.g. expression of a gene across samples) as being normally distributed with mean  $\mu$  and standard deviation 1.



$$\mu = ?$$

## Statistical inference: estimation

To model the data (e.g. expression of a gene across samples) as being normally distributed with mean  $\mu$  and standard deviation 1.



$$\mu = 5?$$

## Statistical inference: estimation

We can estimate  $\mu$  using the observed data through

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

**\*\*Link to 'estimation' shiny app\*\***

## Statistical inference: estimation

We can estimate  $\mu$  using the observed data through

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

[\\*\\*Link to 'estimation' shiny app\\*\\*](#)

Estimates of  $\mu$  that are based on a **finite sample** are not exact:

- ▶ The larger the sample size,  $\bar{x}$  gets closer to  $\mu$
- ▶ The smaller  $\sigma$ , the faster  $\bar{x}$  gets closer to  $\mu$

## Statistical inference: estimation

We can estimate  $\mu$  using the observed data through

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

[\\*\\*Link to 'estimation' shiny app\\*\\*](#)

Estimates of  $\mu$  that are based on a **finite sample** are not exact:

- ▶ The larger the sample size,  $\bar{x}$  gets closer to  $\mu$
- ▶ The smaller  $\sigma$ , the faster  $\bar{x}$  gets closer to  $\mu$

Variability in the estimation quantified through the **standard error**:

$$\text{S.E.} = \sigma / \sqrt{n}$$

## Statistical inference: estimation

We can estimate  $\mu$  using the observed data through

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

[\\*\\*Link to 'estimation' shiny app\\*\\*](#)

Estimates of  $\mu$  that are based on a **finite sample** are not exact:

- ▶ The larger the sample size,  $\bar{x}$  gets closer to  $\mu$
- ▶ The smaller  $\sigma$ , the faster  $\bar{x}$  gets closer to  $\mu$

Variability in the estimation quantified through the **standard error**:

$$\text{S.E.} = \sigma / \sqrt{n}$$

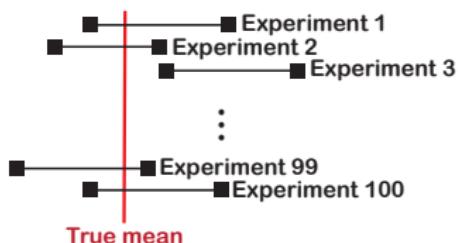
**We need to provide a measure of uncertainty every time we provide an estimation!**

# Statistical inference: confidence intervals

A confidence interval (CI) is a *random* interval

In repeated experiments ...

95% of the time CI covers the *true* mean

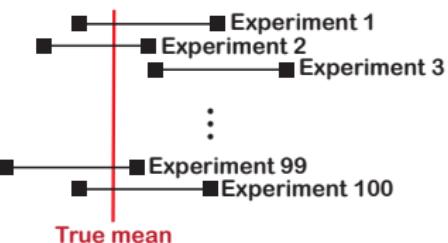


# Statistical inference: confidence intervals

A confidence interval (CI) is a *random* interval

In repeated experiments ...

95% of the time CI covers the *true* mean



For the mean, the 95% CI is given by

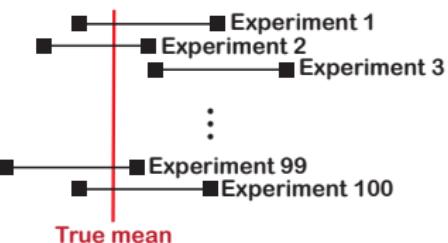
$$(\bar{x} - 1.96 \text{ S.E.}, \bar{x} + 1.96 \text{ S.E.})$$

# Statistical inference: confidence intervals

A confidence interval (CI) is a *random* interval

In repeated experiments ...

95% of the time CI covers the *true* mean



For the mean, the 95% CI is given by

$$(\bar{x} - 1.96 \text{ S.E.}, \bar{x} + 1.96 \text{ S.E.})$$

This assumes the data is **normally distributed**. What if not?

## Statistical inference: central Limit Theorem

If the data is not normally distributed, the **Central Limit Theorem** guarantees this result is still valid, provided the sample size is large

[\\*\\*Link to 'clt' shiny app\\*\\*](#)

# Statistical inference: hypothesis testing

Often, the aim of the analysis is not only **estimation**.

For example,

- ▶ *Is the treatment effective?*
- ▶ *Is a gene differentially expressed?*
- ▶ *Is there an association between genotype and phenotype?*



These questions can be translated as **hypothesis testing** problems

Comic taken from Significance Magazine, December 2008.

# Statistical inference: hypothesis testing



Am I pregnant?

**Disclaimer:** this image was  
downloaded from the internet and  
does not reflect the life of the  
instructors!

# Statistical inference: hypothesis testing

## Basic setup



Am I pregnant?

Disclaimer: this image was downloaded from the internet and does not reflect the life of the instructors!

Formulate the ‘null’ and alternative hypothesis

Calculate a “test statistic” from the data

Desicion rule

# Statistical inference: hypothesis testing

## Basic setup (example)



Am I pregnant?

Disclaimer: this image was downloaded from the internet and does not reflect the life of the instructors!

Formulate the ‘null’ and alternative hypothesis

e.g.  $H_0$ : Not pregnant vs  $H_1$ : Pregnant

Calculate a “test statistic” from the data

e.g. summary measure based on hormonal levels

Decision rule

e.g. is the data more extreme than what is expected by chance (for a non-pregnant woman)?

# Statistical inference: hypothesis testing

No test is exact:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	Correct True positive	Wrong False positive
Do not reject null hypothesis	Wrong False negative	Correct True negative

# Statistical inference: hypothesis testing

No test is exact:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	Correct True positive	Wrong False positive
Do not reject null hypothesis	Wrong False negative	Correct True negative

Error rates depend on e.g. sample size ... what else?

# Statistical inference: hypothesis testing

No test is exact:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	Correct True positive	Wrong False positive
Do not reject null hypothesis	Wrong False negative	Correct True negative

Error rates depend on e.g. sample size ... what else?

Beware of multiple testing correction issues!

# Statistical inference: hypothesis testing

Suppose that 100 women take the test

- ▶ 70 of them were truly pregnant
- ▶ the test was positive for 75 women
- ▶ the result was wrong for 10 non-pregnant women



Am I pregnant?

Complete the table:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis		
Do not reject null hypothesis		

# Statistical inference: hypothesis testing

Suppose that 100 women take the test

- ▶ 70 of them were truly pregnant
- ▶ the test was positive for 75 women
- ▶ the result was wrong for 10 non-pregnant women



Am I pregnant?

Complete the table:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis		10/30
Do not reject null hypothesis		

# Statistical inference: hypothesis testing

Suppose that 100 women take the test

- ▶ 70 of them were truly pregnant
- ▶ the test was positive for 75 women
- ▶ the result was wrong for 10 non-pregnant women



Am I pregnant?

Complete the table:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis		10/30
Do not reject null hypothesis		20/30

# Statistical inference: hypothesis testing

Suppose that 100 women take the test

- ▶ 70 of them were truly pregnant
- ▶ the test was positive for 75 women
- ▶ the result was wrong for 10 non-pregnant women



Am I pregnant?

Complete the table:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis		20/30

# Statistical inference: hypothesis testing

Suppose that 100 women take the test

- ▶ 70 of them were truly pregnant
- ▶ the test was positive for 75 women
- ▶ the result was wrong for 10 non-pregnant women



Am I pregnant?

Complete the table:

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

## Statistical inference: hypothesis testing

Typically, two types of error are of main interest:

- ▶ Type I error

$$\alpha = p(\text{Reject } H_0 \text{ when } H_0 \text{ is true})$$

- ▶ Type II error

$$\beta = p(\text{Do not reject } H_0 \text{ when } H_0 \text{ is false})$$

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

## Statistical inference: hypothesis testing

Typically, two types of error are of main interest:

- ▶ Type I error

$$\alpha = p(\text{Reject } H_0 \text{ when } H_0 \text{ is true})$$

- ▶ Type II error

$$\beta = p(\text{Do not reject } H_0 \text{ when } H_0 \text{ is false})$$

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

Most hypothesis tests are designed to control type I error!

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **significance level** of the test?
- ▶ What is the **power** of the test?

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **significance level** of the test?
  - ▶ Type I error ( $\alpha$ ), i.e.  $10/30 = 0.33$
- ▶ What is the **power** of the test?

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **significance level** of the test?
  - ▶ Type I error ( $\alpha$ ), i.e.  $10/30 = 0.33$
- ▶ What is the **power** of the test?
  - ▶  $1 - \text{type II error } (1 - \beta)$ , i.e.  $65/70 = 0.93$

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **sensitivity** of the test?
- ▶ What is the **specificity** of the test?

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **sensitivity** of the test?
  - ▶ True positive rate, i.e.  $65/70 = 0.93$
- ▶ What is the **specificity** of the test?

# Statistical inference: hypothesis testing

In terms of this table

	Null hypothesis does not hold	Null hypothesis holds
Reject null hypothesis	65/70	10/30
Do not reject null hypothesis	5/70	20/30

- ▶ What is the **sensitivity** of the test?
  - ▶ True positive rate, i.e.  $65/70 = 0.93$
- ▶ What is the **specificity** of the test?
  - ▶ True negative rate, i.e.  $20/30 = 0.67$

# Statistical inference: hypothesis testing

## Recall: basic setup

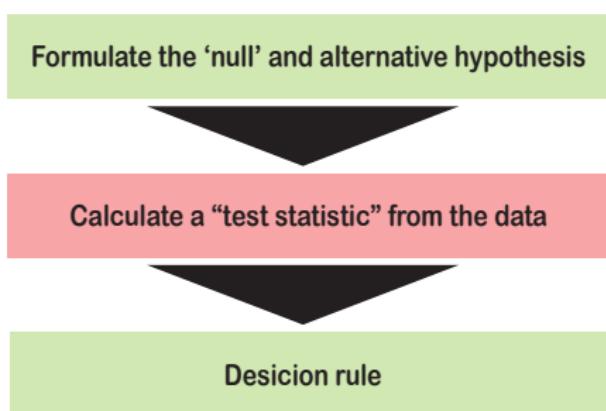
Formulate the ‘null’ and alternative hypothesis

Calculate a “test statistic” from the data

Decision rule

# Statistical inference: hypothesis testing

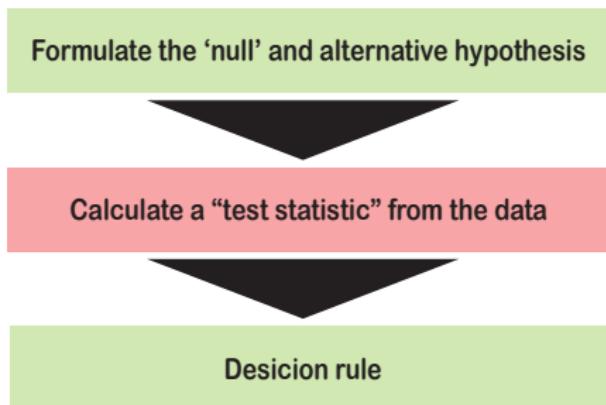
## Recall: basic setup



A **test statistic** is a summary calculated from the data whose distribution is known (if  $H_0$  is true)

# Statistical inference: hypothesis testing

## Recall: basic setup



A **test statistic** is a summary calculated from the data whose distribution is known (if  $H_0$  is true)

We will explore some example test statistics later today

# Statistical inference: hypothesis testing

## Recall: basic setup

Formulate the ‘null’ and alternative hypothesis

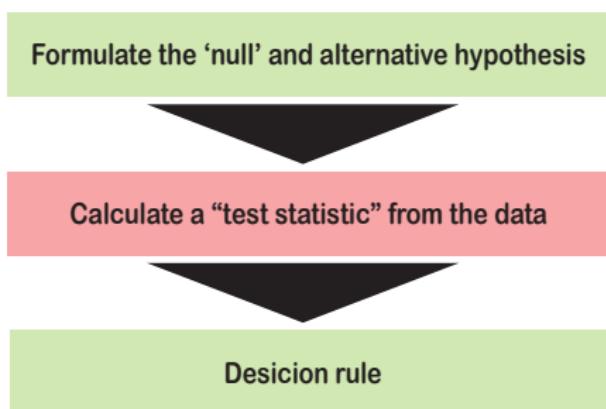
Calculate a “test statistic” from the data

Decision rule

We can also use a  
**p-value**

# Statistical inference: hypothesis testing

## Recall: basic setup



We can also use a  
**p-value**

**What is a p-value?**

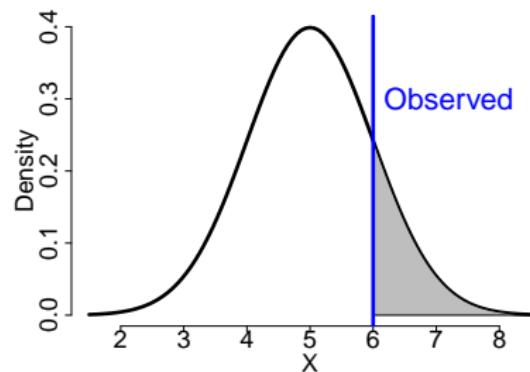
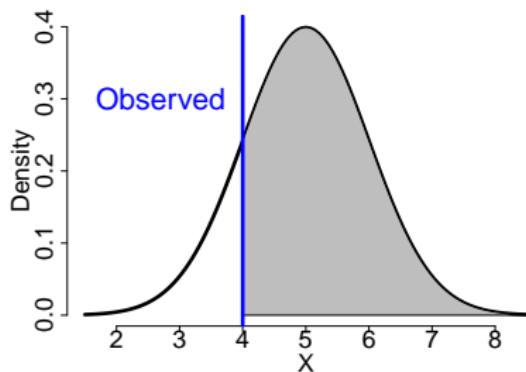
## Statistical inference: what is a p-value?

A **p-value** is the probability of observing the current data (or more extreme) given than  $H_0$  is true.

## Statistical inference: what is a p-value?

A **p-value** is the probability of observing the current data (or more extreme) given than  $H_0$  is true. For example

$$H_0 : \mu \leq 5 \text{ vs } H_1 : \mu > 5$$



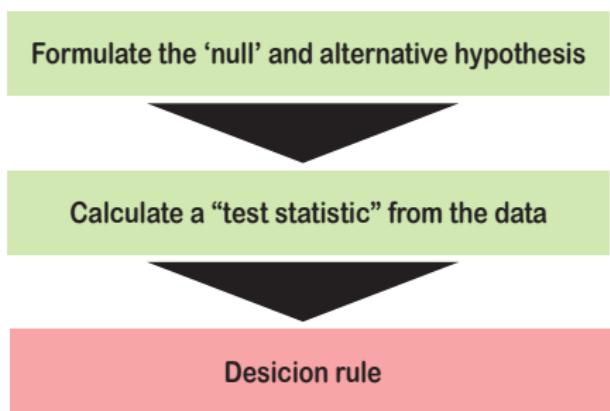
The p-value is equal to the gray area

The smaller p-value, the stronger the evidence against  $H_0$

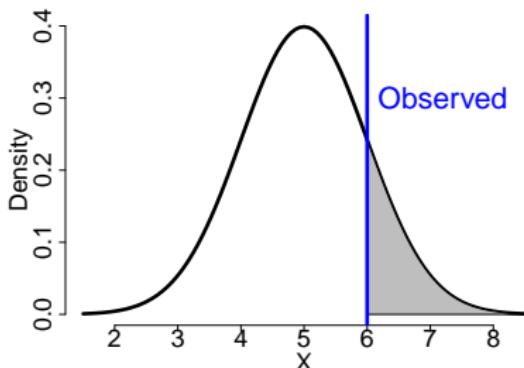
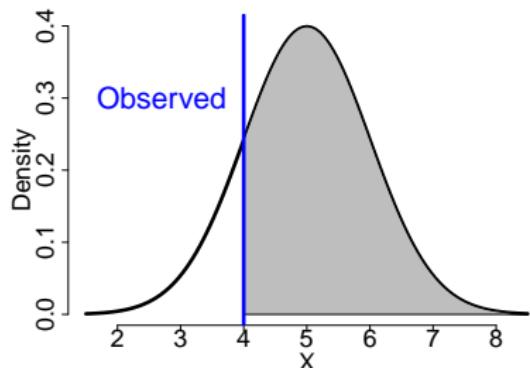
\*\*Link to 'pvalue' shiny app\*\*

# Statistical inference: decision rule of an hypothesis test

## Recall: basic setup



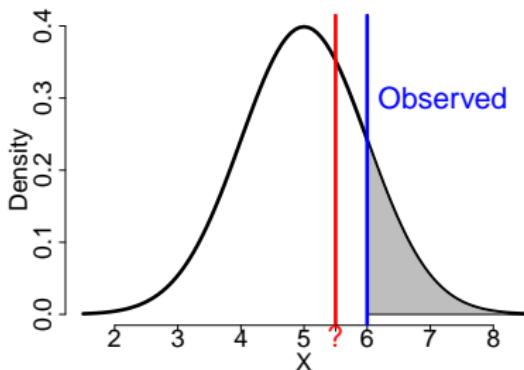
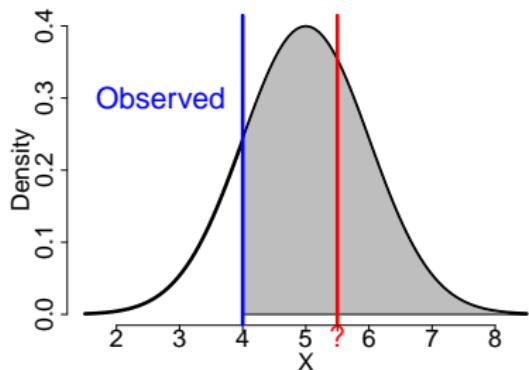
# Statistical inference: decision rule of an hypothesis test



- ▶ If the  $p$ -value is **large**, the data seems to support  $H_0$ .
- ▶ If the  $p$ -value is **small**, the data does not seem to support  $H_0$ .

Where to place the cut-off?

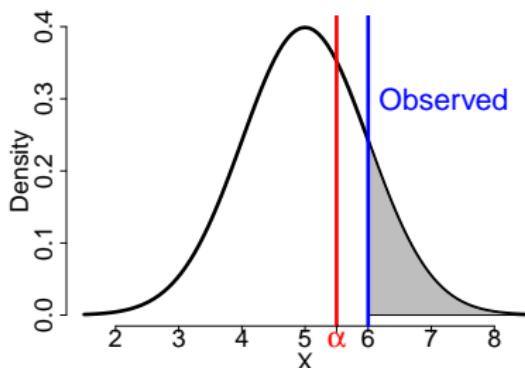
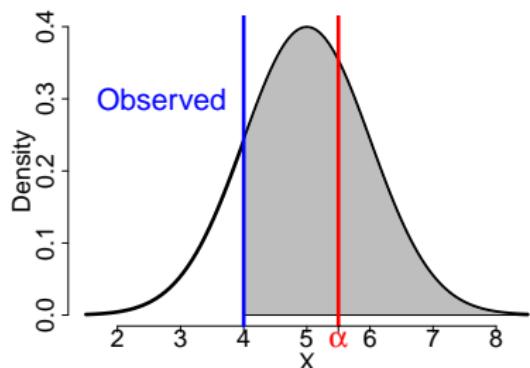
# Statistical inference: decision rule of an hypothesis test



- ▶ If the  $p$ -value is **large**, the data seems to support  $H_0$ .
- ▶ If the  $p$ -value is **small**, the data does not seem to support  $H_0$ .

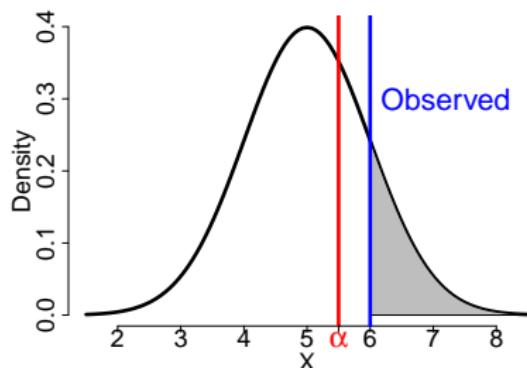
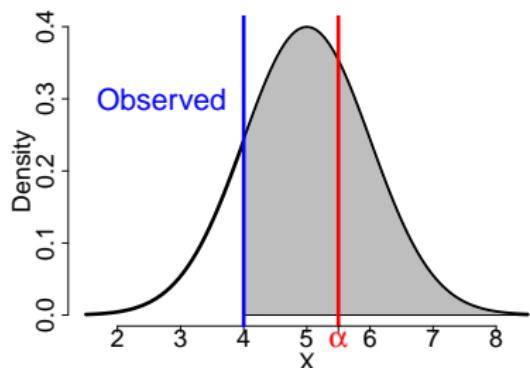
Where to place the cut-off?

## Statistical inference: decision rule of an hypothesis test



The cut-off is typically set to control **type I error**

## Statistical inference: decision rule of an hypothesis test

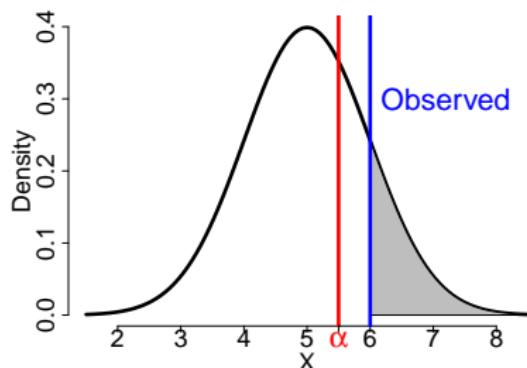
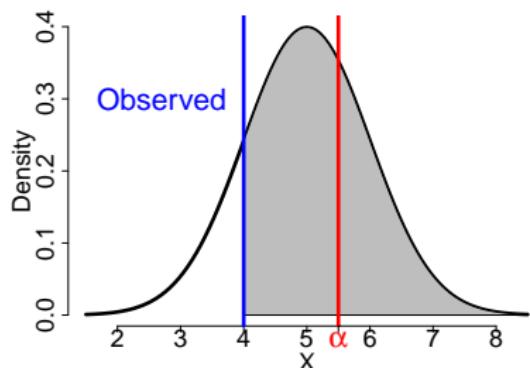


The cut-off is typically set to control **type I error**

$$\alpha = P(\text{Reject } H_0 \text{ when } H_0 \text{ is true})$$

**Popular choice:**  $\alpha = 0.05$

# Statistical inference: decision rule of an hypothesis test



The cut-off is typically set to control **type I error**

$$\alpha = P(\text{Reject } H_0 \text{ when } H_0 \text{ is true})$$

**Popular choice:**  $\alpha = 0.05$

**Cut-off must be set before running the test!**

## Statistical inference: what test to use?

Standard statistical softwares (e.g. R) have a wide range of hypothesis tests built-in ⇒ easy to get numbers back!

## Statistical inference: what test to use?

Standard statistical softwares (e.g. R) have a wide range of hypothesis tests built-in ⇒ easy to get numbers back!

How to choose?

# Statistical inference: what test to use?

Standard statistical softwares (e.g. R) have a wide range of hypothesis tests built-in  $\Rightarrow$  easy to get numbers back!

## How to choose?

Before running a test it is important to recognise:

- ▶ What type of data is available?
- ▶ What is the relevant hypothesis?
- ▶ What assumptions underlie the test?

# Statistical inference: what test to use?

Standard statistical softwares (e.g. R) have a wide range of hypothesis tests built-in  $\Rightarrow$  easy to get numbers back!

## How to choose?

Before running a test it is important to recognise:

- ▶ What type of data is available?
- ▶ What is the relevant hypothesis?
- ▶ What assumptions underlie the test?

We will explore some examples in detail this afternoon ...

## Acknowledgements

Part of the materials discussed today have been adapted from materials provided by

- ▶ Mark Dunning (CRUK-CI)
- ▶ Paul Kirk (MRC-BSU)