

# Statistics of RNA-seq analysis

## Authors/Contributors:

- Dominique Laurent Couturier, CRUK-CI
- Zeynep Kalender-Atak, CRUK-CI
- Hugo Tavares, Bioinformatics Training Facility, University of Cambridge
- Abbi Edwards, CRUK-CI

```
> dds <- DESeqDataSetFromMatrix(cnts, DataFrame(cond), ~ cond)
> dds <- DESeq(dds)
> results(dds)
```

log2 fold change (MLE): cond 2 vs 1

Wald test p-value: cond 2 vs 1

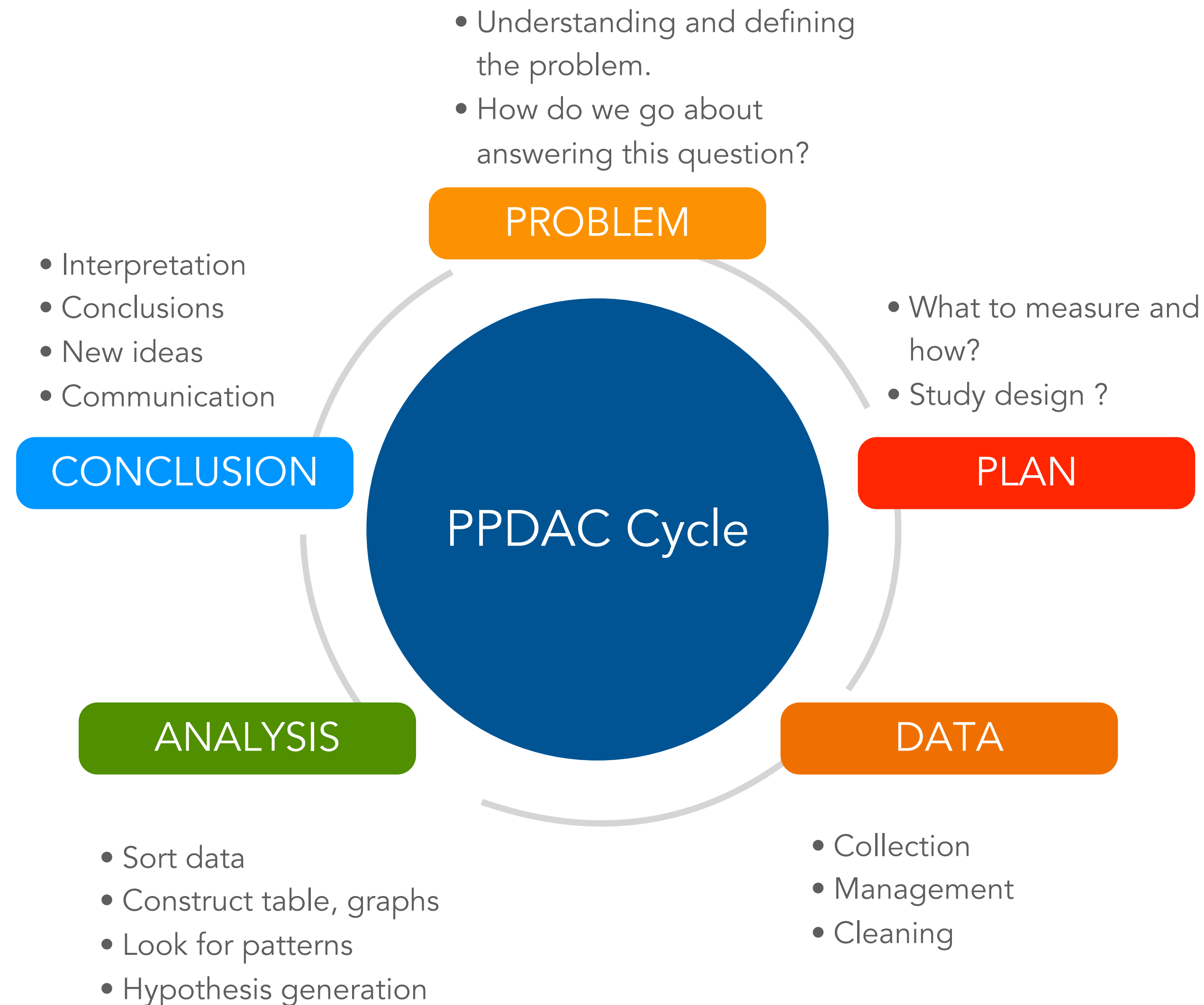
DataFrame with 1000 rows and 6 columns

	baseMean	log2FoldChange	lfcSE	stat	pvalue	padj
	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>
1	97.3140	-0.682067	0.344525	-1.979730	0.0477339	0.745842
2	109.9860	-0.228819	0.450720	-0.507676	0.6116808	0.944354
3	98.8111	0.104291	0.462113	0.225683	0.8214483	0.978382
4	103.2615	0.306400	0.297682	1.029284	0.3033460	0.944354
5	97.9406	0.316338	0.357242	0.885501	0.3758864	0.944354
...	...	...	...	...	...	...
996	86.8057	0.0467703	0.287042	0.162939	0.8705668	0.980044
997	101.4437	-0.2070806	0.339886	-0.609264	0.5423495	0.944354
998	78.1356	-0.6372790	0.369515	-1.724637	0.0845930	0.824310
999	89.2920	0.7554725	0.306192	2.467314	0.0136131	0.614613
1000	103.5569	-0.0728875	0.348655	-0.209053	0.8344065	0.978382

# Data literacy

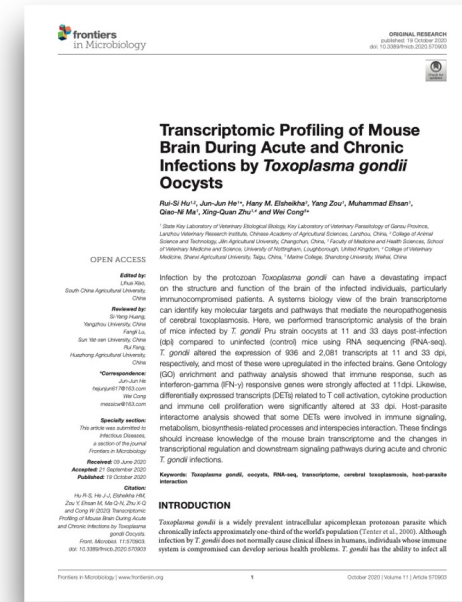
The ability to not only carry out statistical analysis on real-world problems, but also to understand and critique any conclusions drawn by others on the basis of statistics.

# Statistics as an investigative process of problem-solving and decision-making





# Statistics as an investigative process of problem-solving and decision-making



- IFN- $\gamma$  response increases as infection progresses
- Calcium response pathways are downregulated

- Understanding and defining the problem.
- How do we go about answering this question?

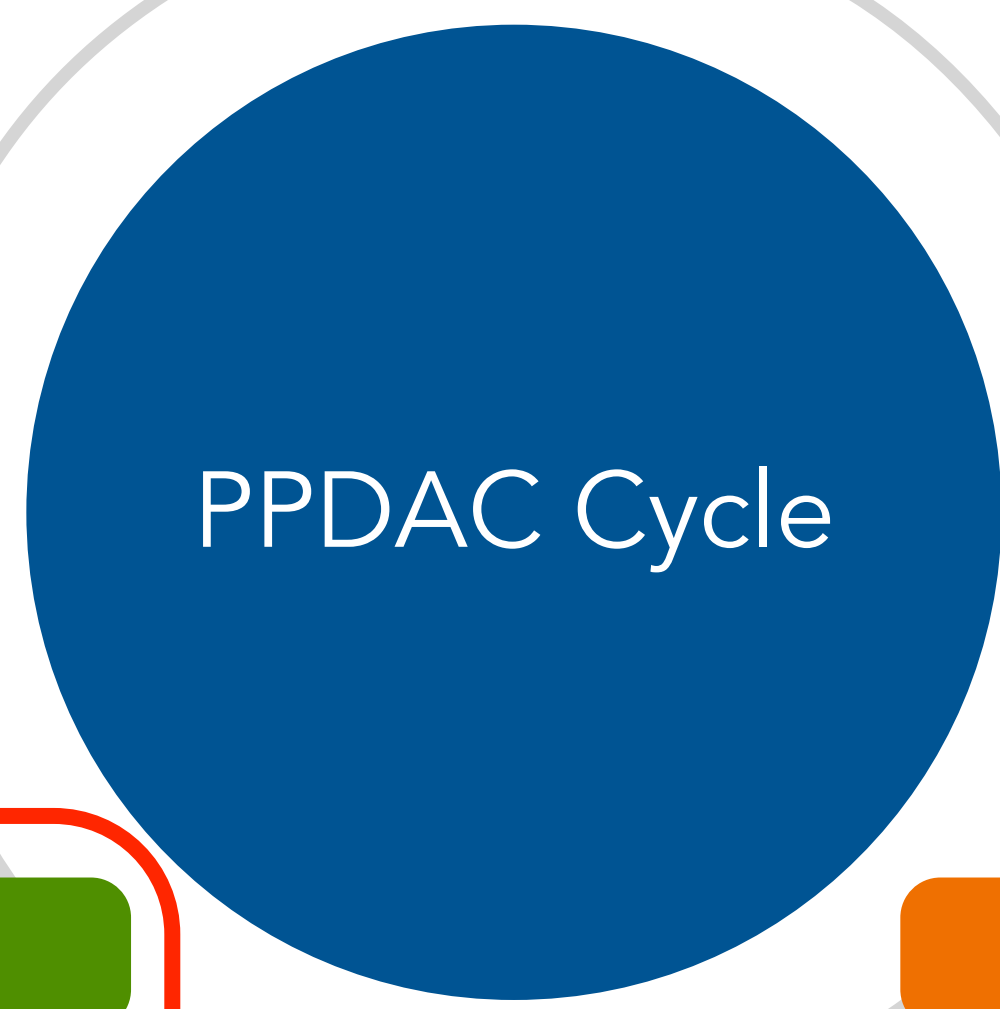
Toxoplasma gondii infection causes a host of severe neurological disorders. Our understanding of the molecular mechanisms associated with infection is incomplete.

We want to study the effect of Toxoplasma gondii infection (chronic and acute) in mouse brain

- Interpretation
- Conclusions
- New ideas
- Communication

CONCLUSION

PROBLEM



- What to measure and how?
- Study design ?

PLAN

- Total gene expression profile of the brain in infection versus no-infection
- A two-factor study with three biological replicates in each group with matched controls

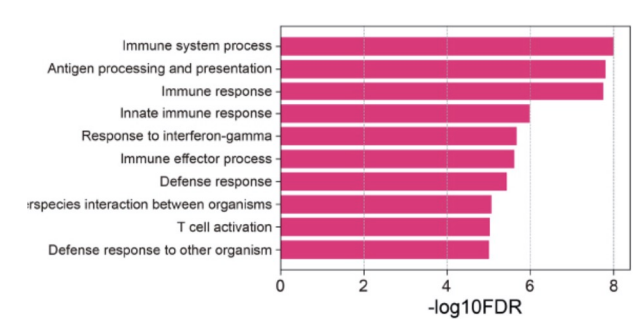
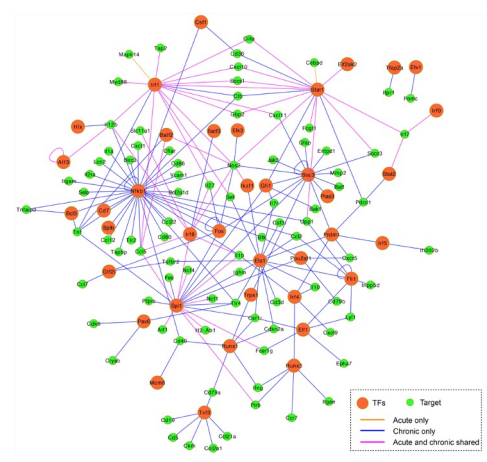
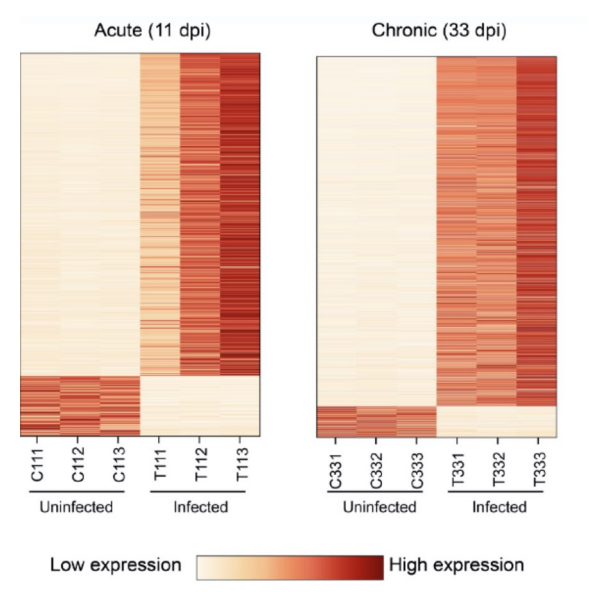
DATA

- Collection
- Management
- Cleaning

Profiling the total transcriptome with RNA-seq  
Preprocessing and quality control

ANALYSIS

- Sort data
- Construct table, graphs
- Look for patterns
- Hypothesis generation



# Outline

- Experimental Design
- General Statistical Concepts
- Statistical aspects specific to bulk RNA-seq analysis

# Outline

- Experimental Design
- Statistical Concepts - Bite size statistics
- Statistical aspects of bulk RNA-seq analysis

# Consequences of Poor Experimental Design

Inability to answer the questions we would like to answer

- **Cost** of experimentation.
- **Limited & Precious** material, esp. clinical samples.
- **Immortalization** of data sets in public databases and methods in the literature. Our bad science begets more bad science.
- **Ethical concerns** of experimentation: animals and clinical samples.



# A Well-Designed Experiment

## Should have

- Clear objectives
- Focus and simplicity
- Sufficient power
- Randomised comparisons

## And be

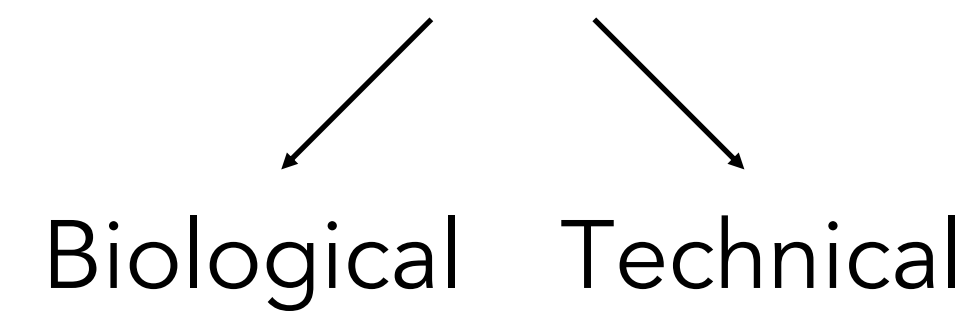
- Precise
- Unbiased
- Amenable to statistical analysis
- Reproducible

# Experimental Factors

- **Factors**: aspects of experiment that change and influence the outcome of the experiment
  - e.g. time, weight, drug, gender, ethnicity, country, plate, cage etc.
- **Variable type** depends on type of measurement:
  - Categorical (nominal) , e.g. sex
  - Categorical with ordering (ordinal), e.g. tumour grade
  - Discrete, e.g. shoe size, number of cells
  - Continuous, e.g. body weight in kg, height in cm
- **Independent and Dependent variables**
  - Independent variable (IV): what you change
  - Dependent variable (DV): what changes due to IV
  - “If (independent variable), then (dependent variable)”

# Sources of Variation

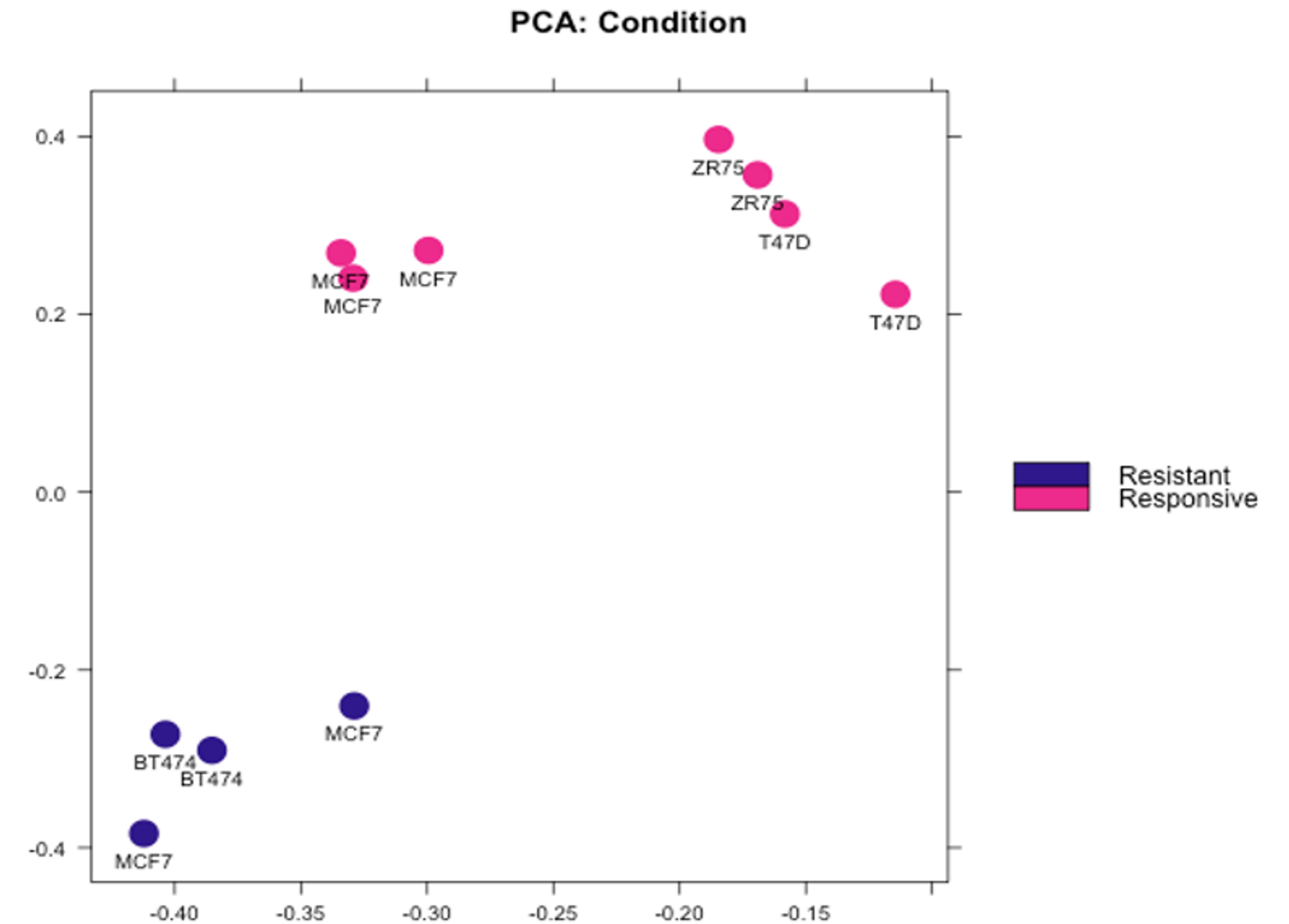
dependent variable = f ( independent variable ) + noise



- Biological "noise"
  - Biological processes are inherently stochastic
  - Single cells, cell populations, individuals, organs, species....
  - Timepoints, cell cycle, synchronized vs. unsynchronized
- Technical noise
  - Reagents, antibodies, temperatures, pollution
  - Platforms, runs, operators
- Replication is required to capture variance

# Types of Replication

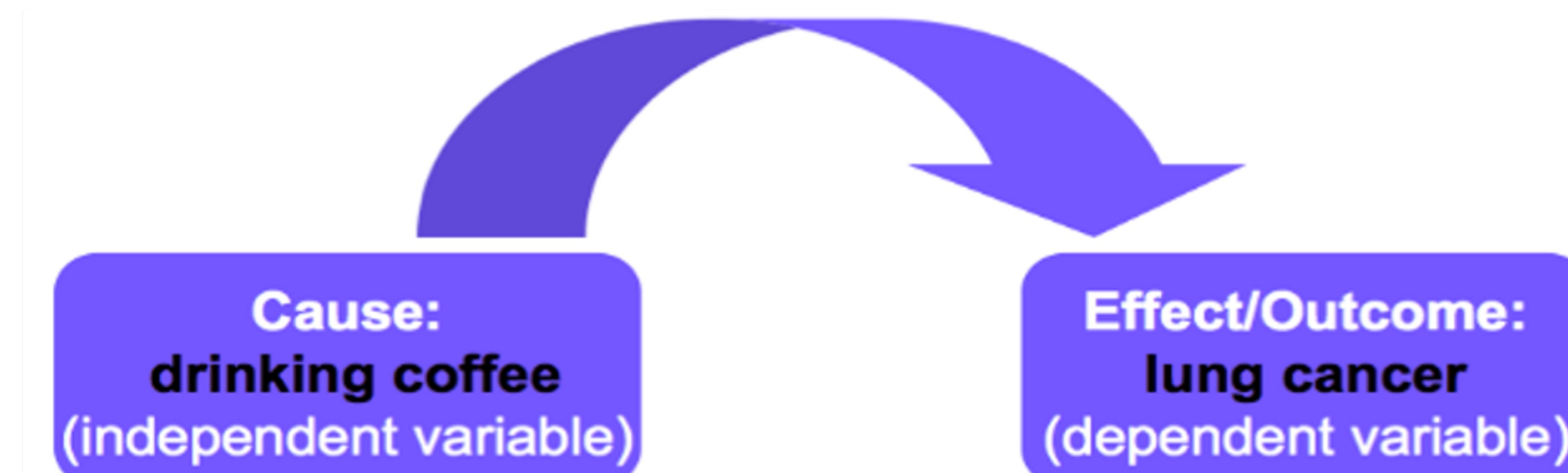
- Biological replication:
  - In vivo:
    - Patients
    - Mice
  - In vitro:
    - Different cell lines
    - Re-growing cells (passages)
- Technical replication:
  - Experimental protocol
  - Measurement platform (i.e. sequencer)





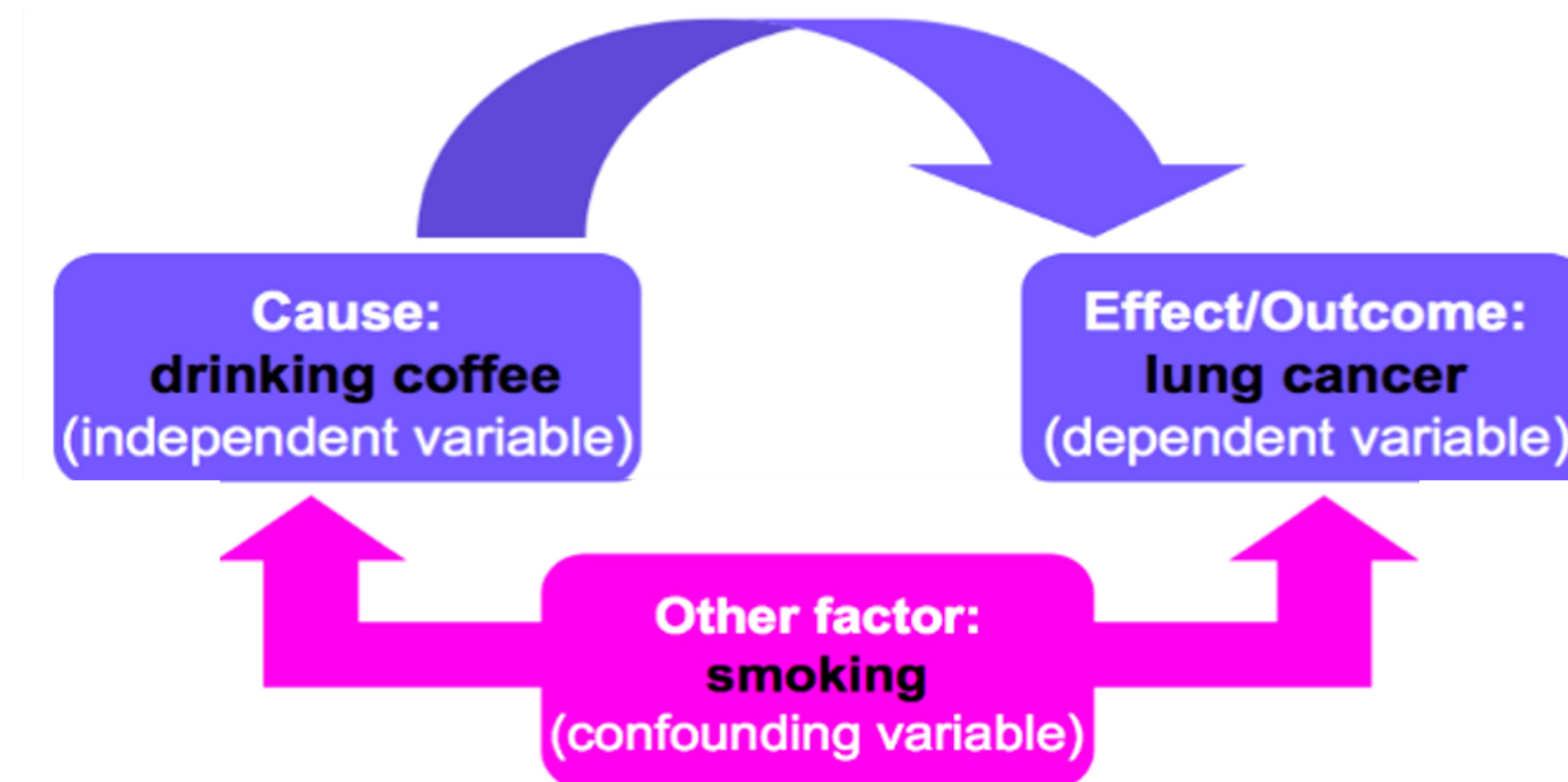
# Confounding Factors

- Also known as extraneous, hidden, lurking or masking factors, or the third variable or mediator variable.
- May mask an actual association or falsely demonstrate an apparent association between the independent & dependent variables.
- Hypothetical Example would be a study of coffee drinking and lung cancer.



# Confounding Factors

- Also known as extraneous, hidden, lurking or masking factors, or the third variable or mediator variable.
- May mask an actual association or falsely demonstrate an apparent association between the independent & dependent variables.
- Hypothetical Example would be a study of coffee drinking and lung cancer.

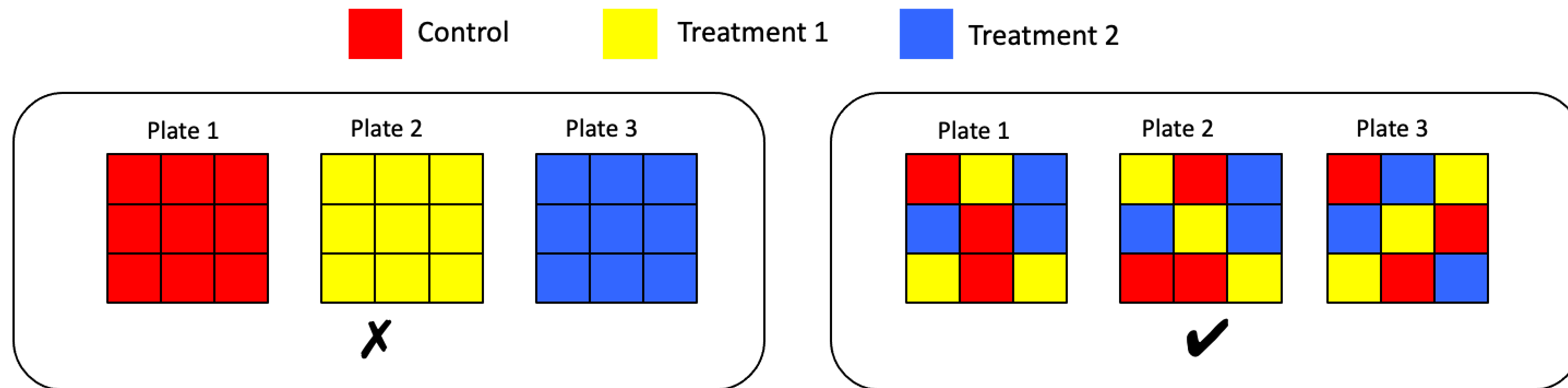


# Solutions

- Write it all down!!!!!!!!!!
- Controlling technical effects:
  - Randomisation
    - Statistical analyses assume randomised comparisons
    - May not see issues caused by non-randomised comparisons
    - Make every decision random not arbitrary
    - Caveat: over-randomization can increase error
  - Blinding
    - Especially important where subjective measurements are taken
    - Potentially multiple degrees of blinding (eg. double-blinding)

# Randomised Block Design

- Blocking is the arranging of experimental units in groups (blocks) that are similar to one another.



- Each plate contains spatially randomised equal proportions of:
  - Control
  - Treatment 1
  - Treatment 2
- controlling plate effects.



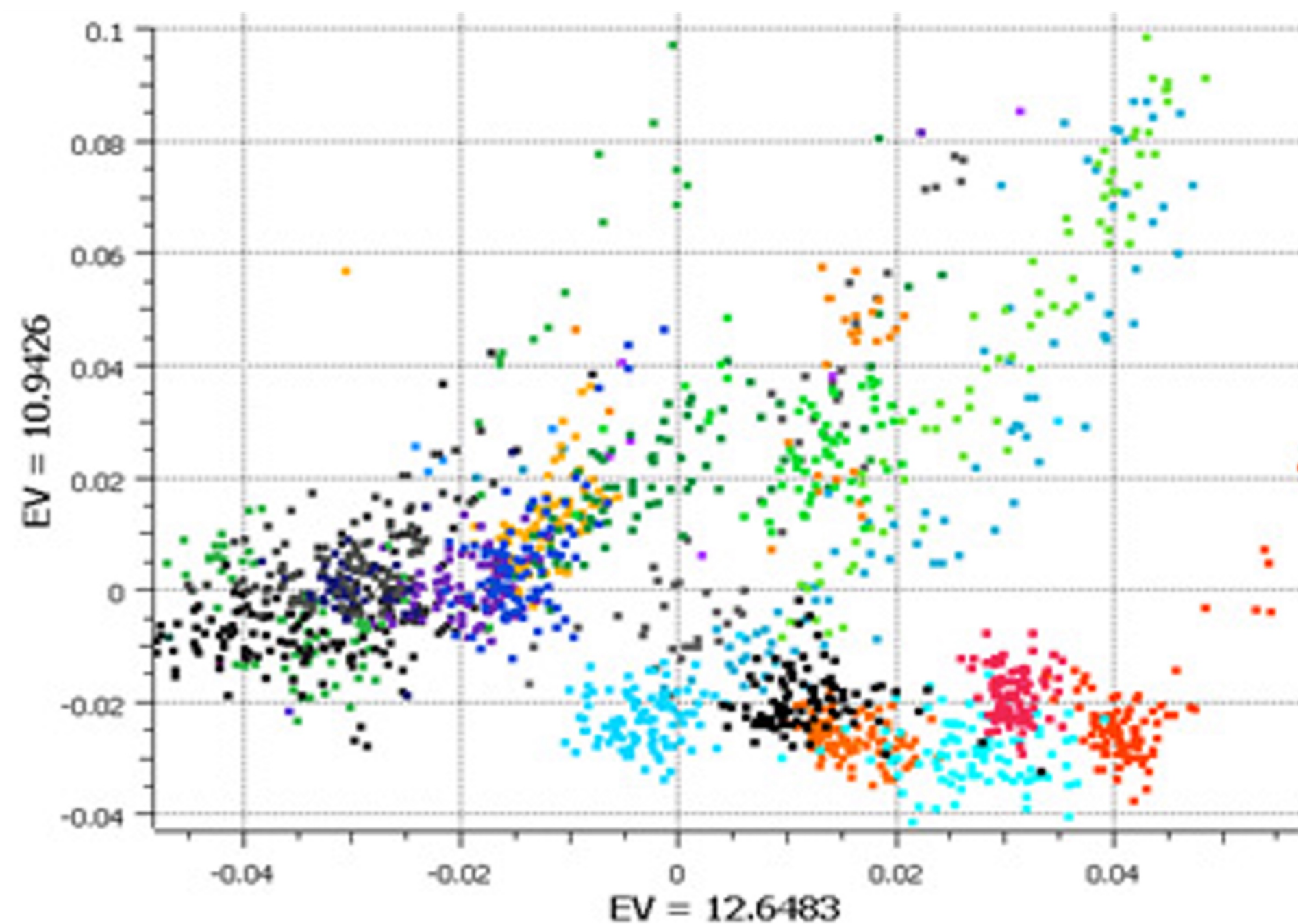
# Randomised Block Design

- Good design example: Alzheimer's study from GlaxoSmithKline

Plate effects by plate

Left PCA plot show large plate effects.

Each colour corresponds to a different plate





# Randomised Block Design

- Good design example: Alzheimer's study from GlaxoSmithKline

Plate effects by plate

Left PCA plot show large plate effects.  
Each colour corresponds to a different plate

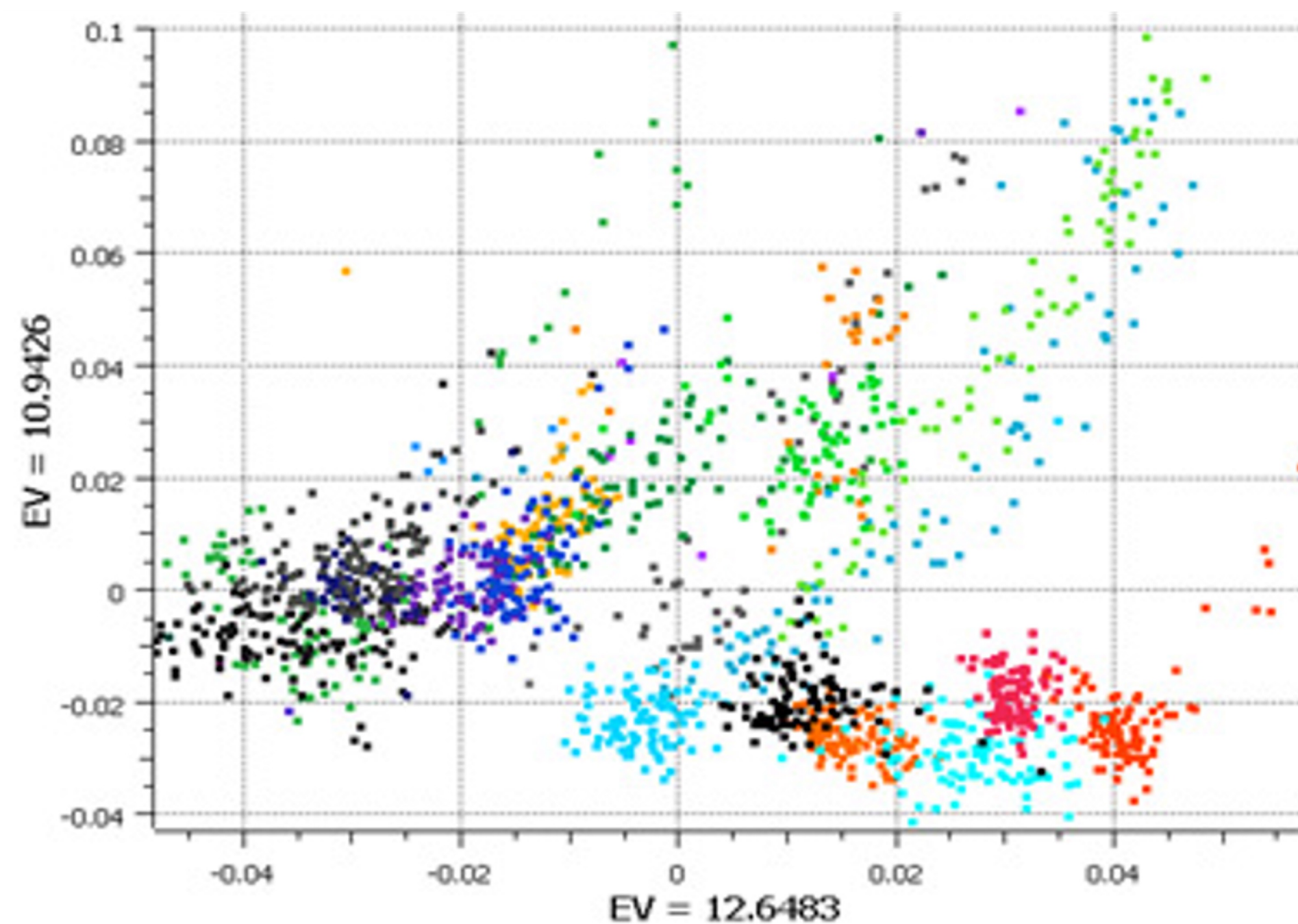
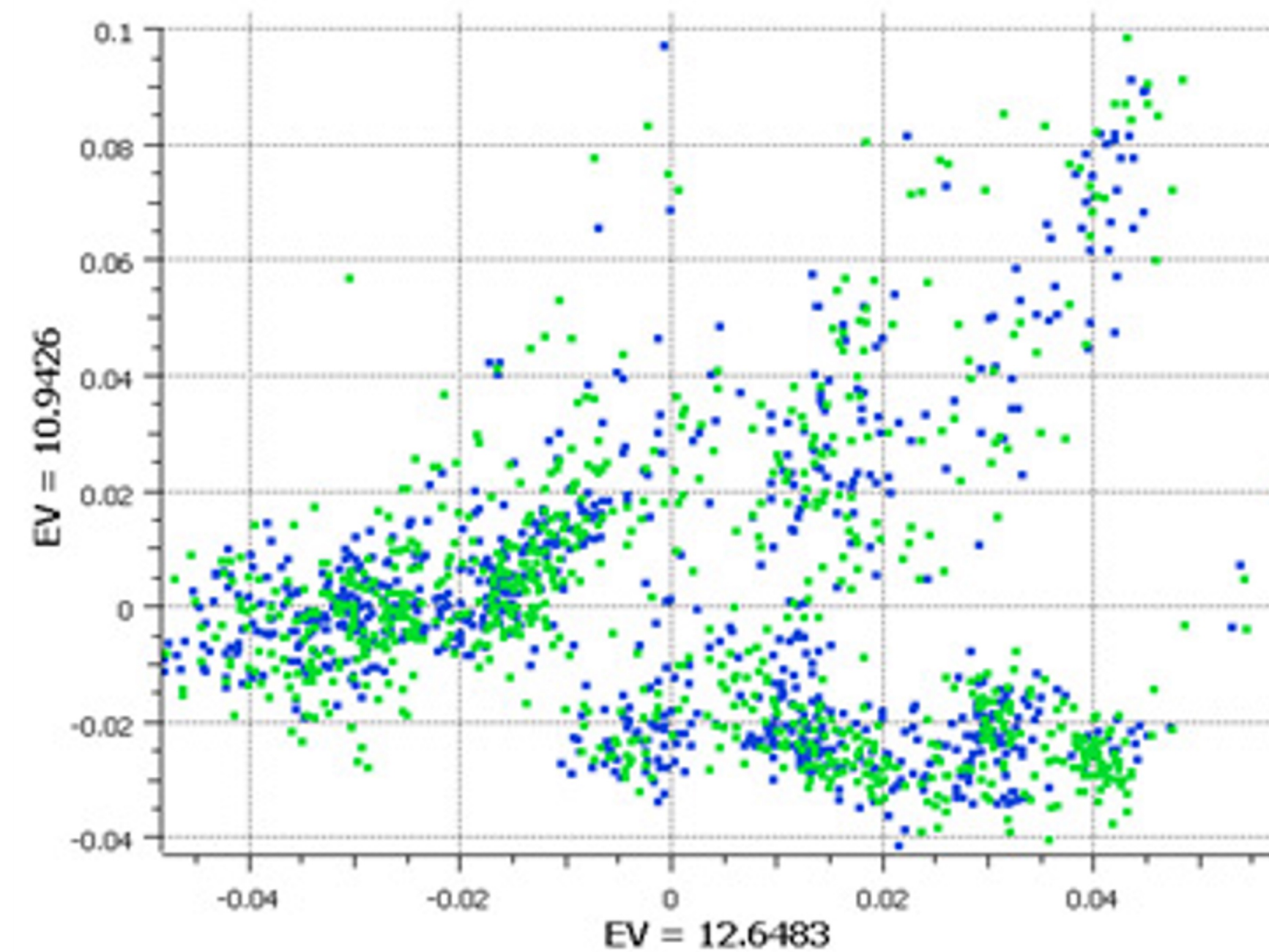


Plate effects by case/control

Right PCA plot shows each plate cluster contains  
equal proportions of cases (blue) and controls (green).



# Experimental Controls

- Ideal : Everything is identical across conditions except the variable you are testing
- Controlling errors
  - Type I: False Positives
    - Negative controls: should have minimal or no effect
  - Type II: False Negatives
    - Positive controls: known effect
- Technical controls
  - Detect/correct technical biases
  - Normalise measurements (quantification)

# Examples of Experimental Controls

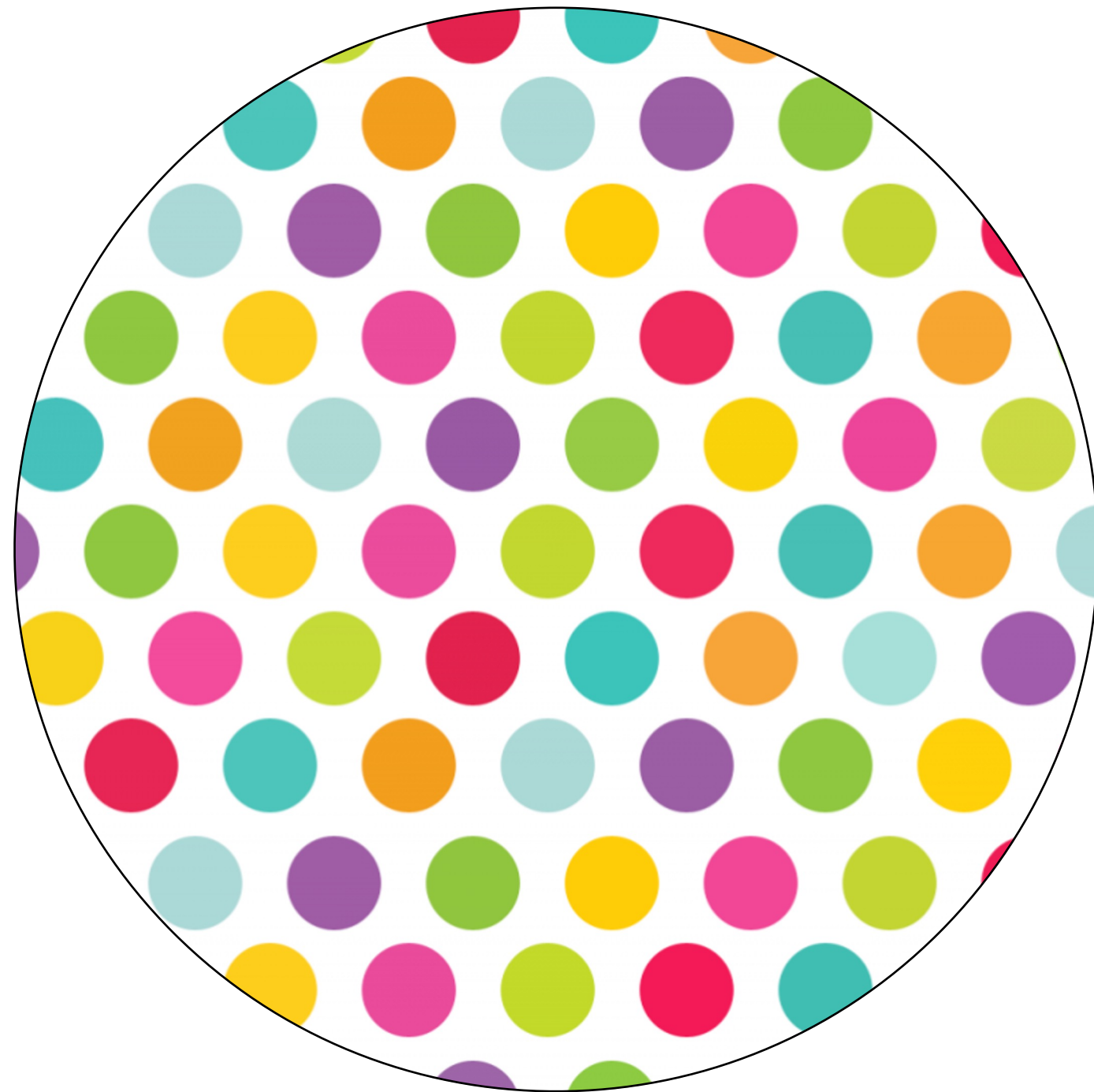
- Wild-type organism (knockouts)
- Inactive siRNA (silencing)
- Vehicle (treatments)
- Spike-ins (quantification/normalisation)
- “Gold standard” data points
- Multi-level controls
- e.g. contrast Vehicle/Input vs. Treatment/Input

# Outline

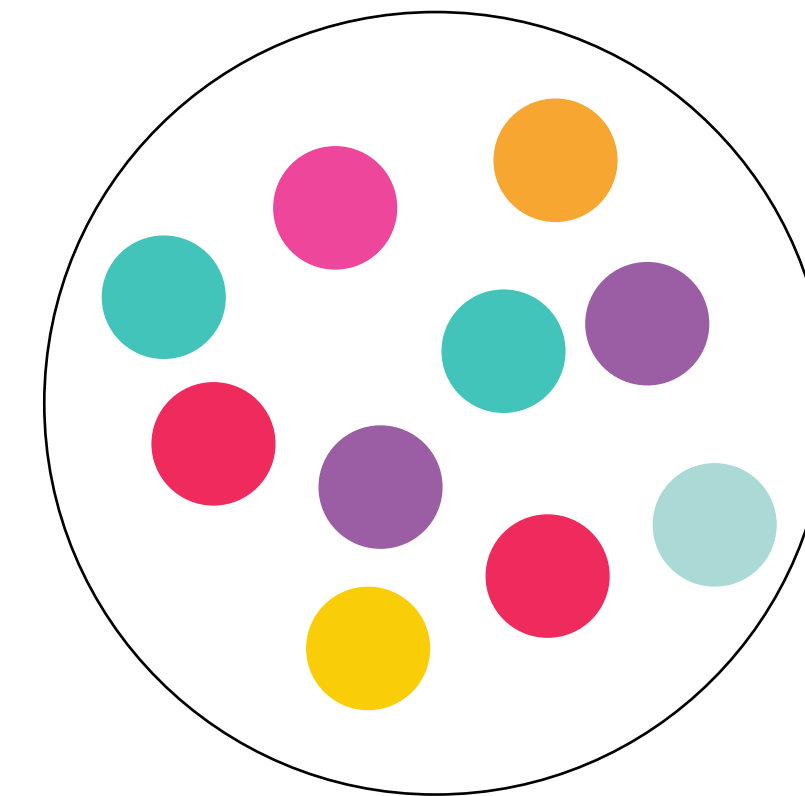
- Experimental Design
- Statistical Concepts
- Statistical aspects of bulk RNA-seq analysis



# Basics on inferential statistics and hypothesis testing



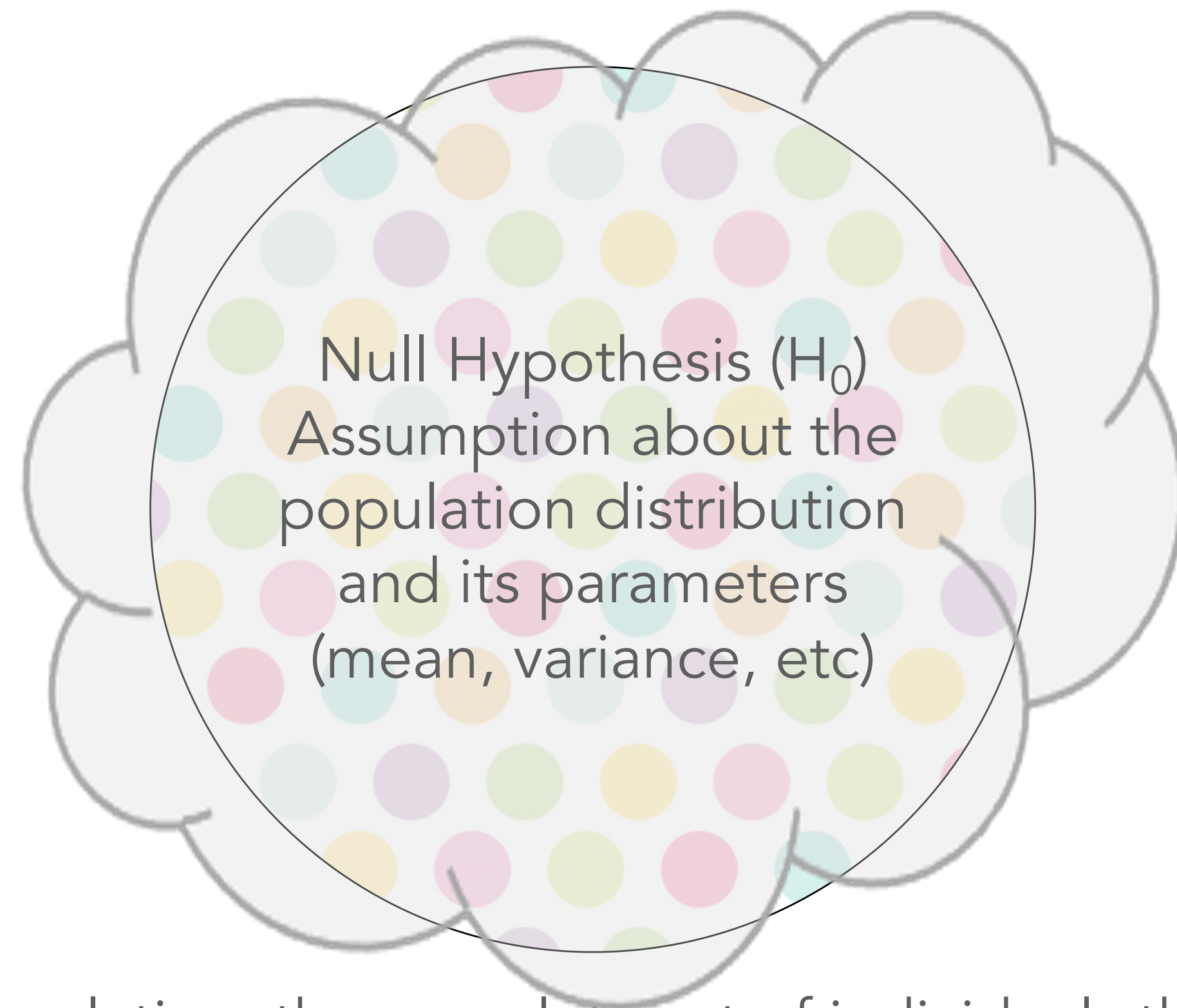
Population: the complete set of individuals that we are interested in



Sample: smaller set of individuals that is representative of the population

Variable: what we are interested in measuring

# Basics on inferential statistics and hypothesis testing



Population: the complete set of individuals that we are interested in

Inference means two things:

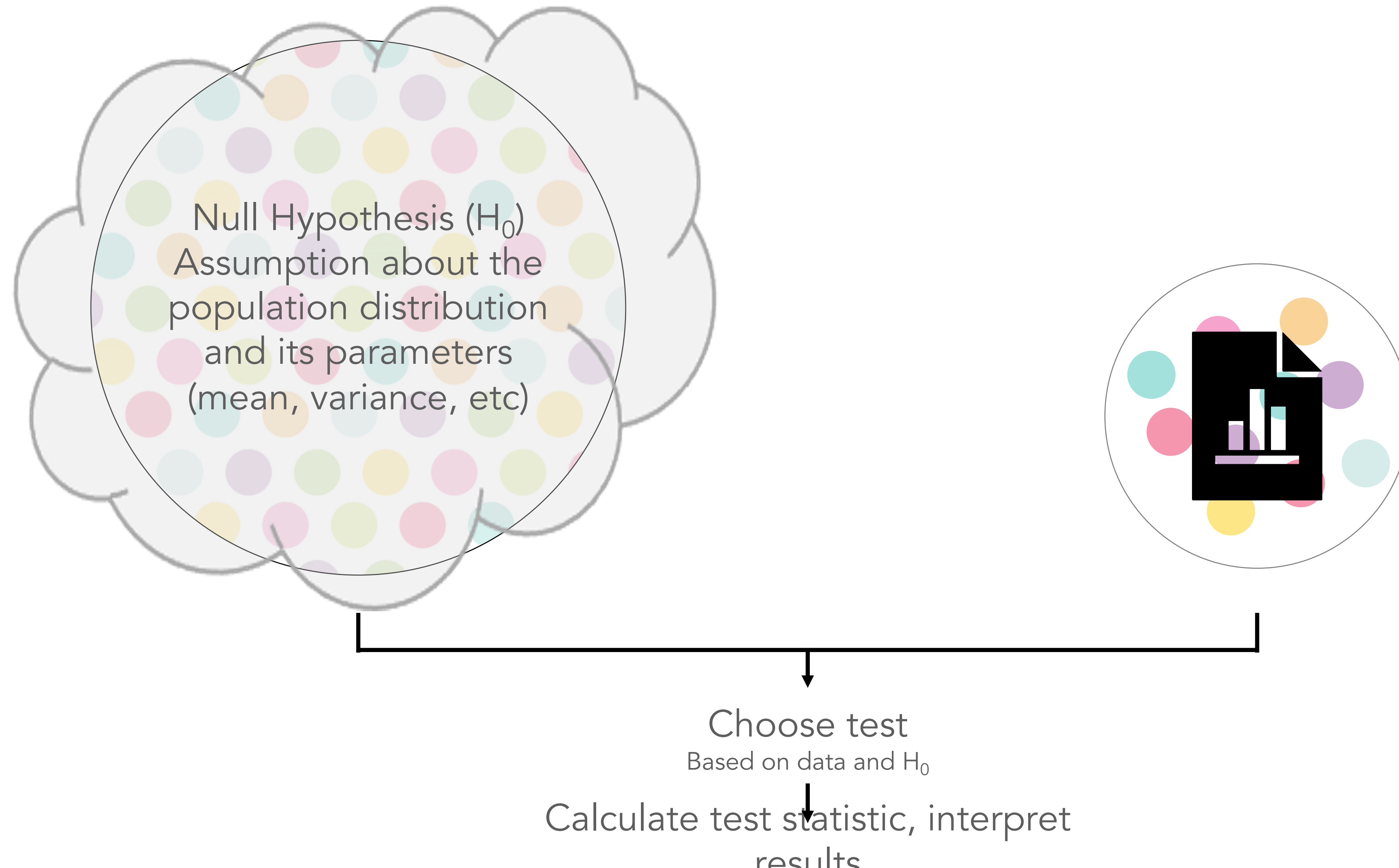
1. Estimating population parameters
2. Testing hypothesis regarding the population distribution



Sample: smaller set of individuals that is representative of the population

Variable: what we are interested in measuring

# Basics on inferential statistics and hypothesis testing



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$H_0$ : Drug has no effect on response time

$H_1$ : Drug has an effect on response time

# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

S

$$H_1: \mu \neq 1.2$$

S



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

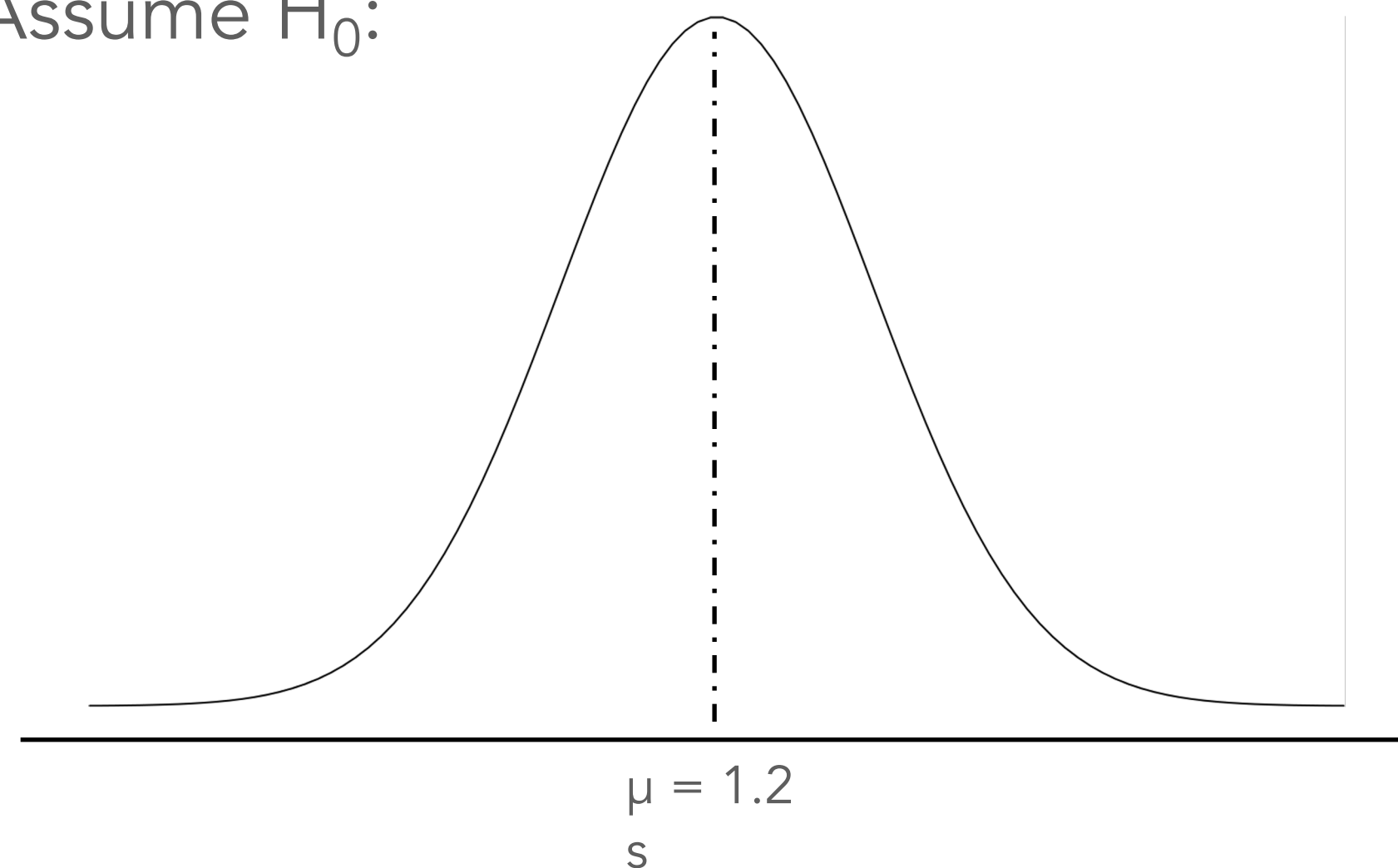
$$H_0: \mu = 1.2$$

s

$$H_1: \mu \neq 1.2$$

s

Assume  $H_0$ :



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

s

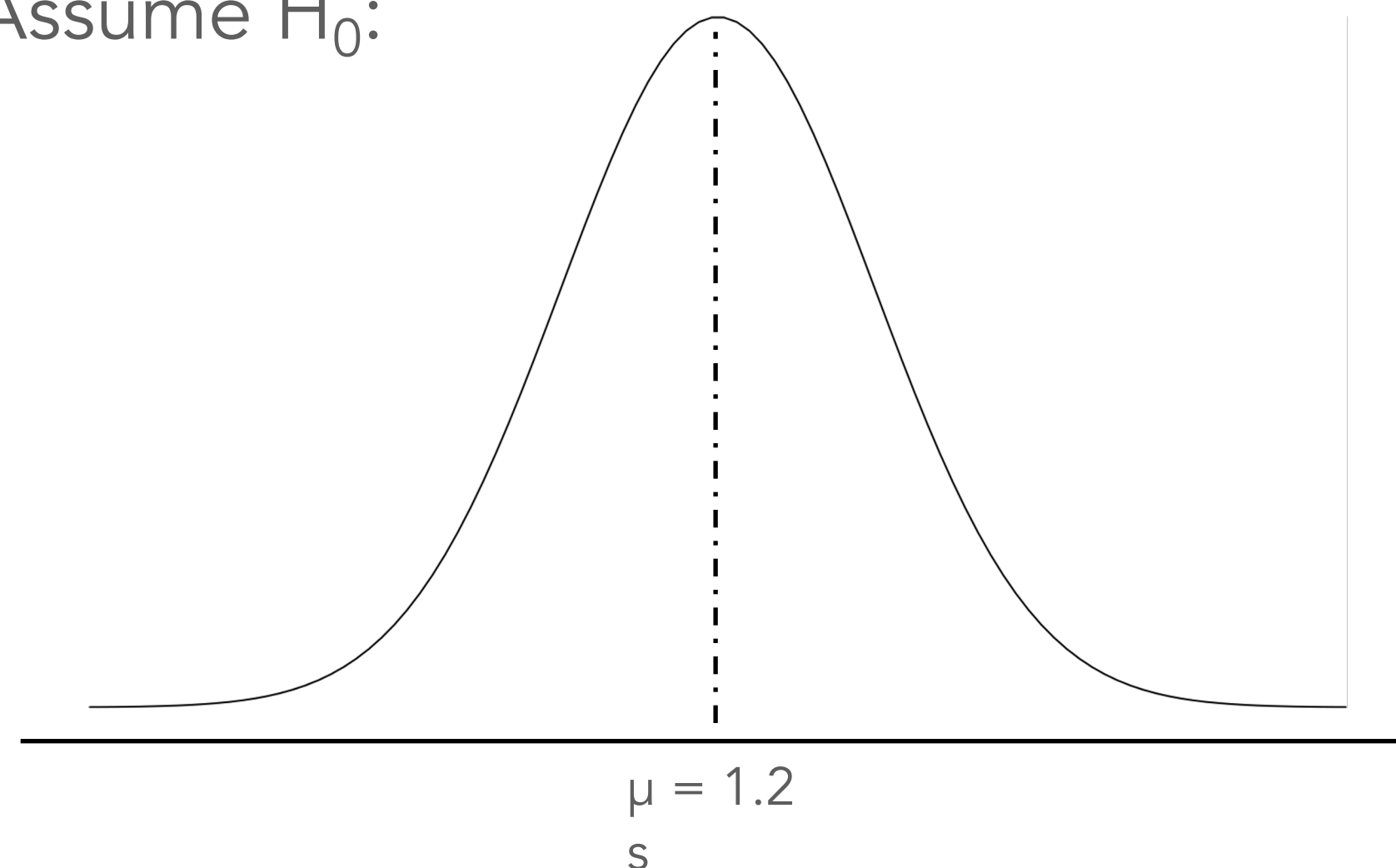
$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\bar{m} - \mu}{s / \sqrt{n}}$$

Assume  $H_0$ :



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

s

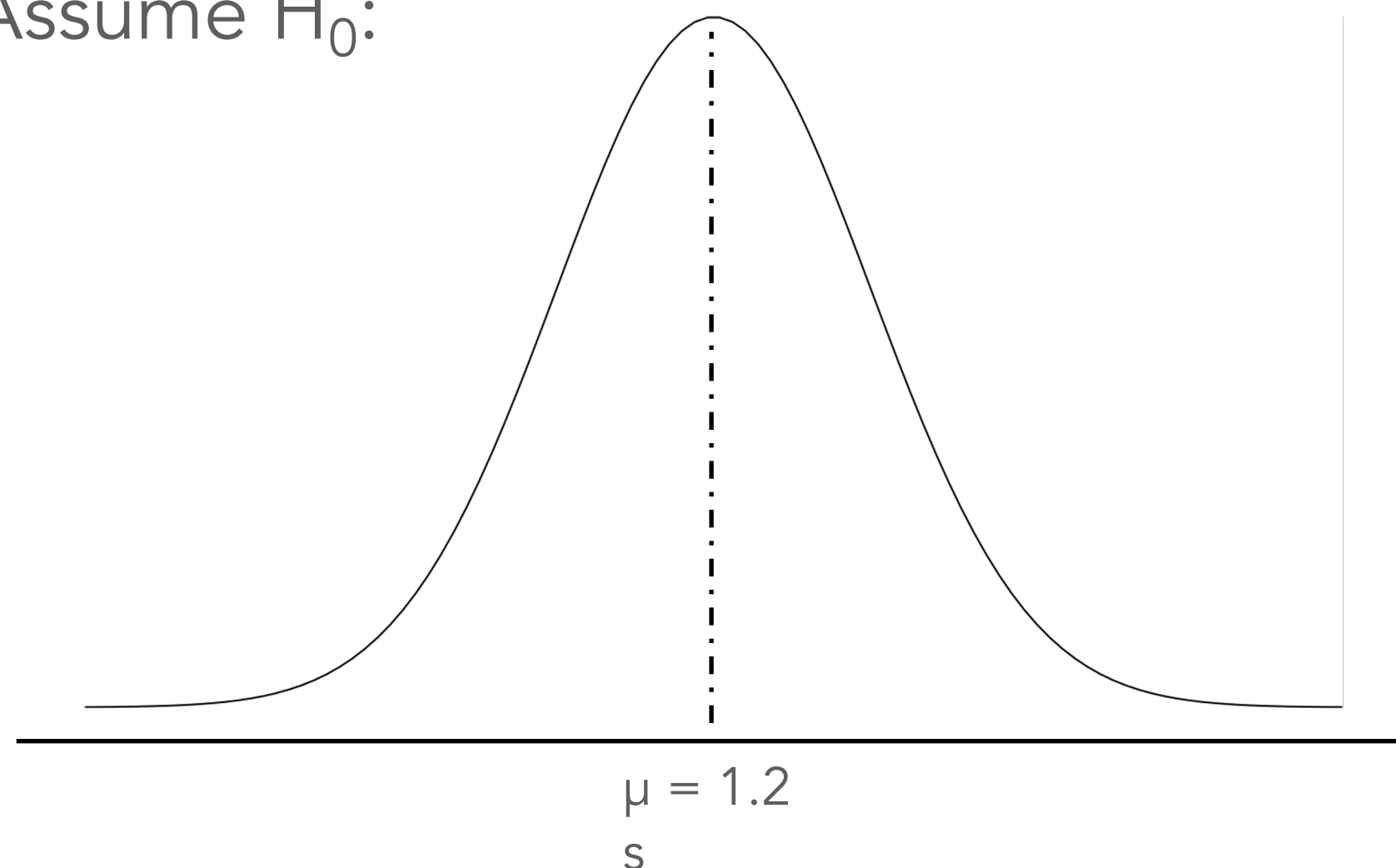
$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\overset{1.05}{\underset{\curvearrowright}{m}} - \overset{1.2}{\underset{\curvearrowright}{\mu}}}{\underset{0.5}{\underset{\curvearrowright}{s}} / \underset{100}{\underset{\curvearrowright}{\sqrt{n}}}}$$

Assume  $H_0$ :



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

s

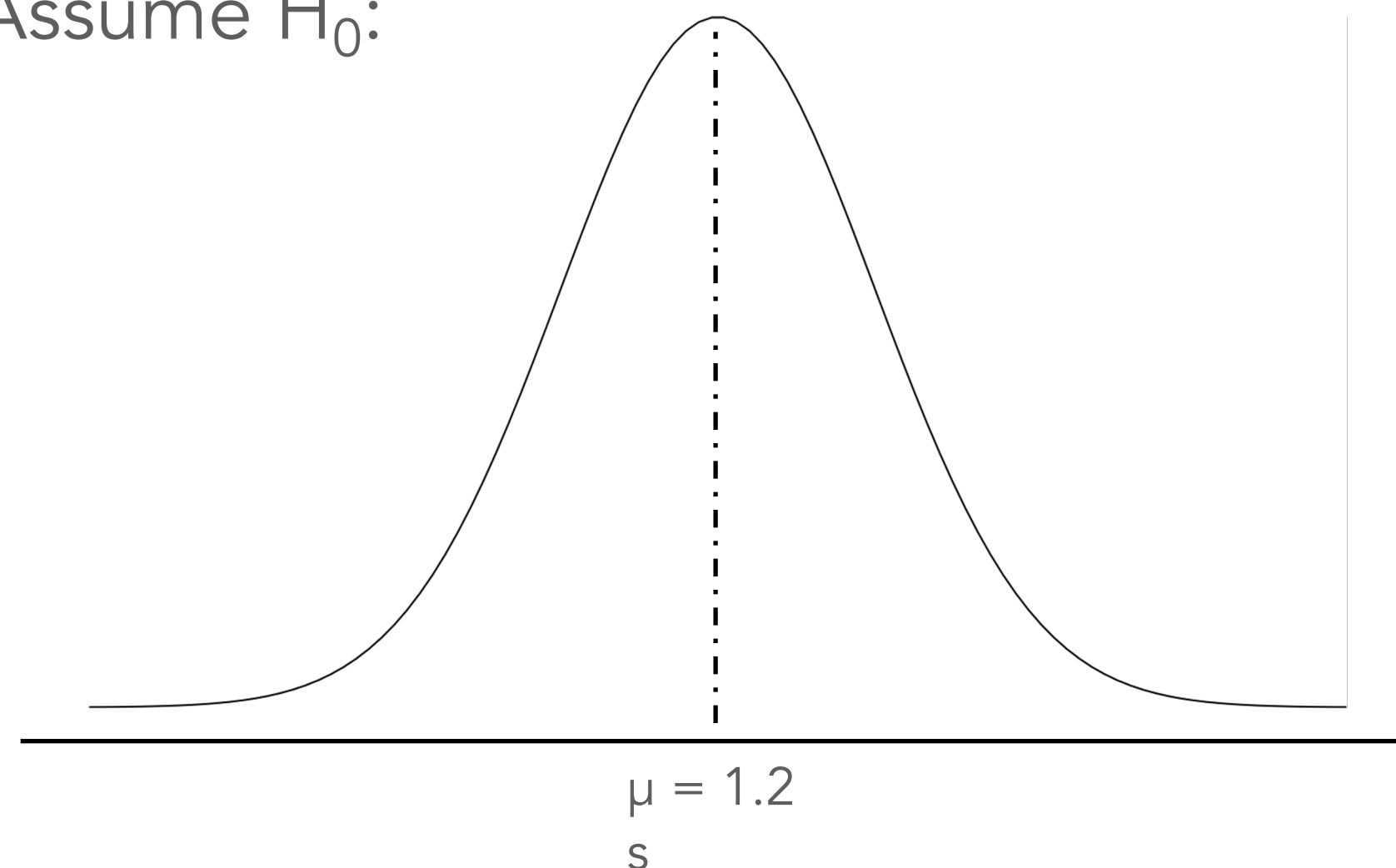
$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\bar{m} - \mu}{s / \sqrt{n}} = -3$$

Assume  $H_0$ :



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

s

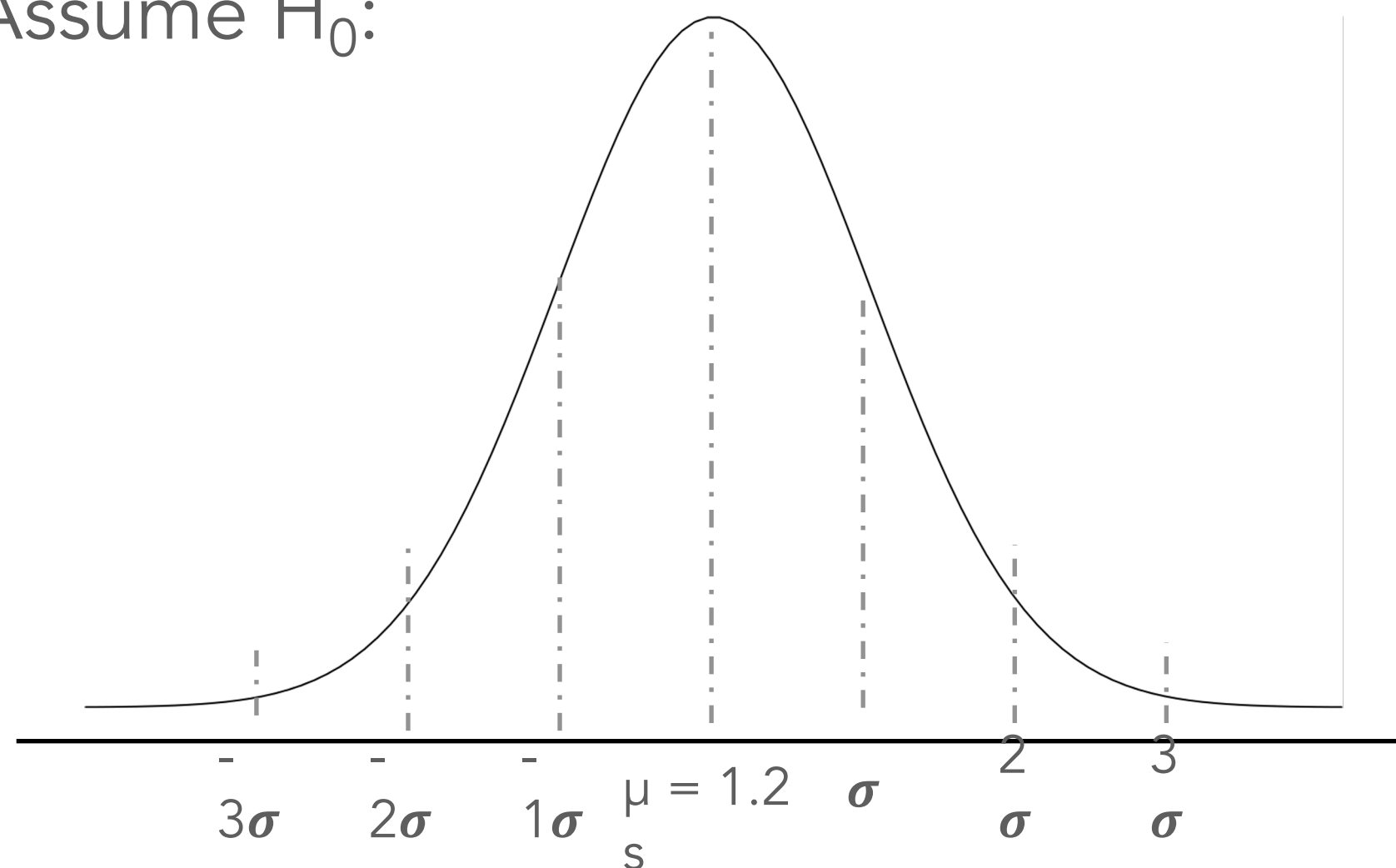
$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\bar{m} - \mu}{s/\sqrt{n}} = -3$$

Assume  $H_0$ :



This means that the sample mean (1.05) is 3 standard deviations away from the mean



# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

s

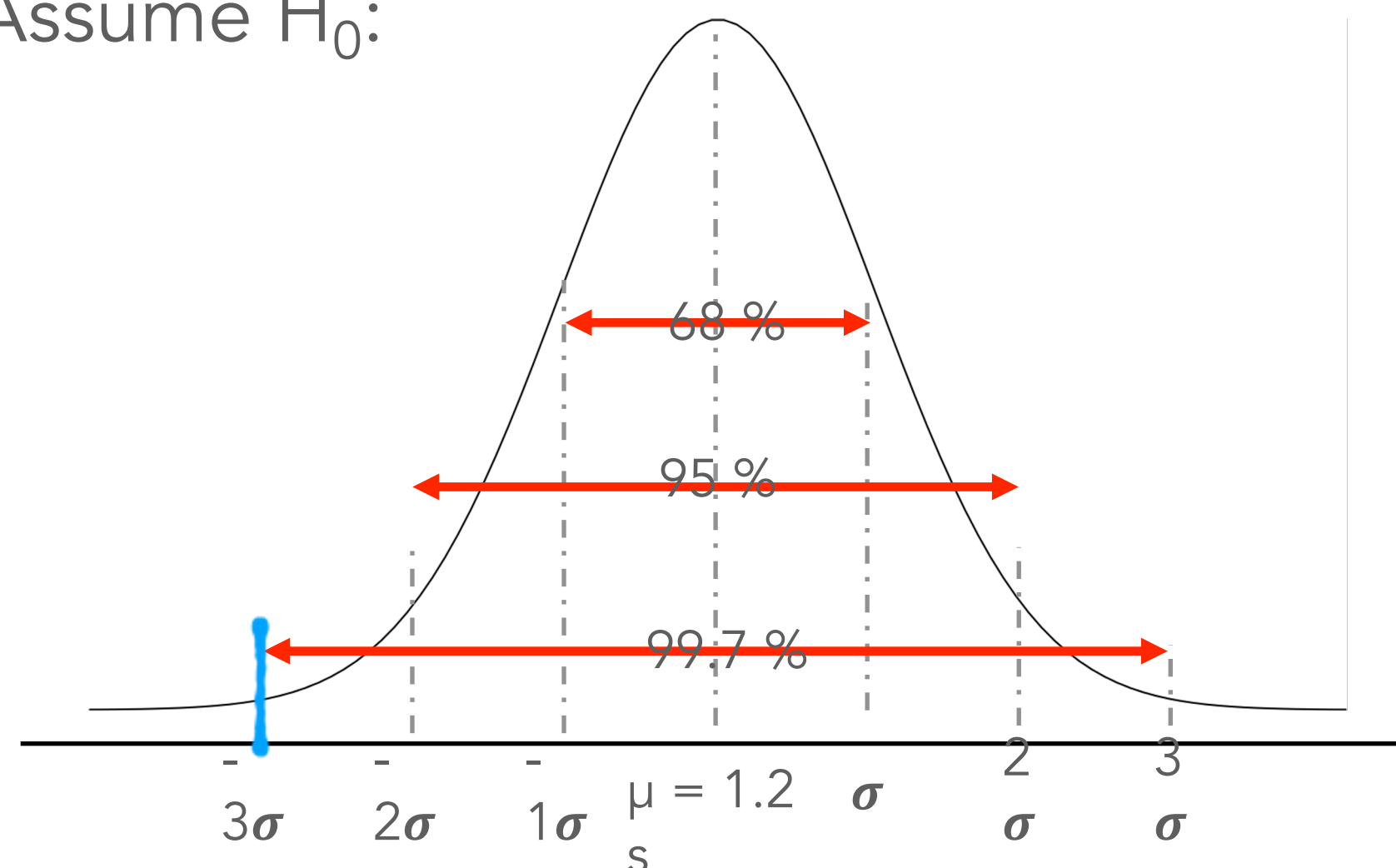
$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\bar{m} - \mu}{s/\sqrt{n}} = -3$$

Assume  $H_0$ :



This means that the sample mean (1.05) is 3 standard deviations away from the mean

What is the probability of observing a test statistic as extreme as 1.05?

# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

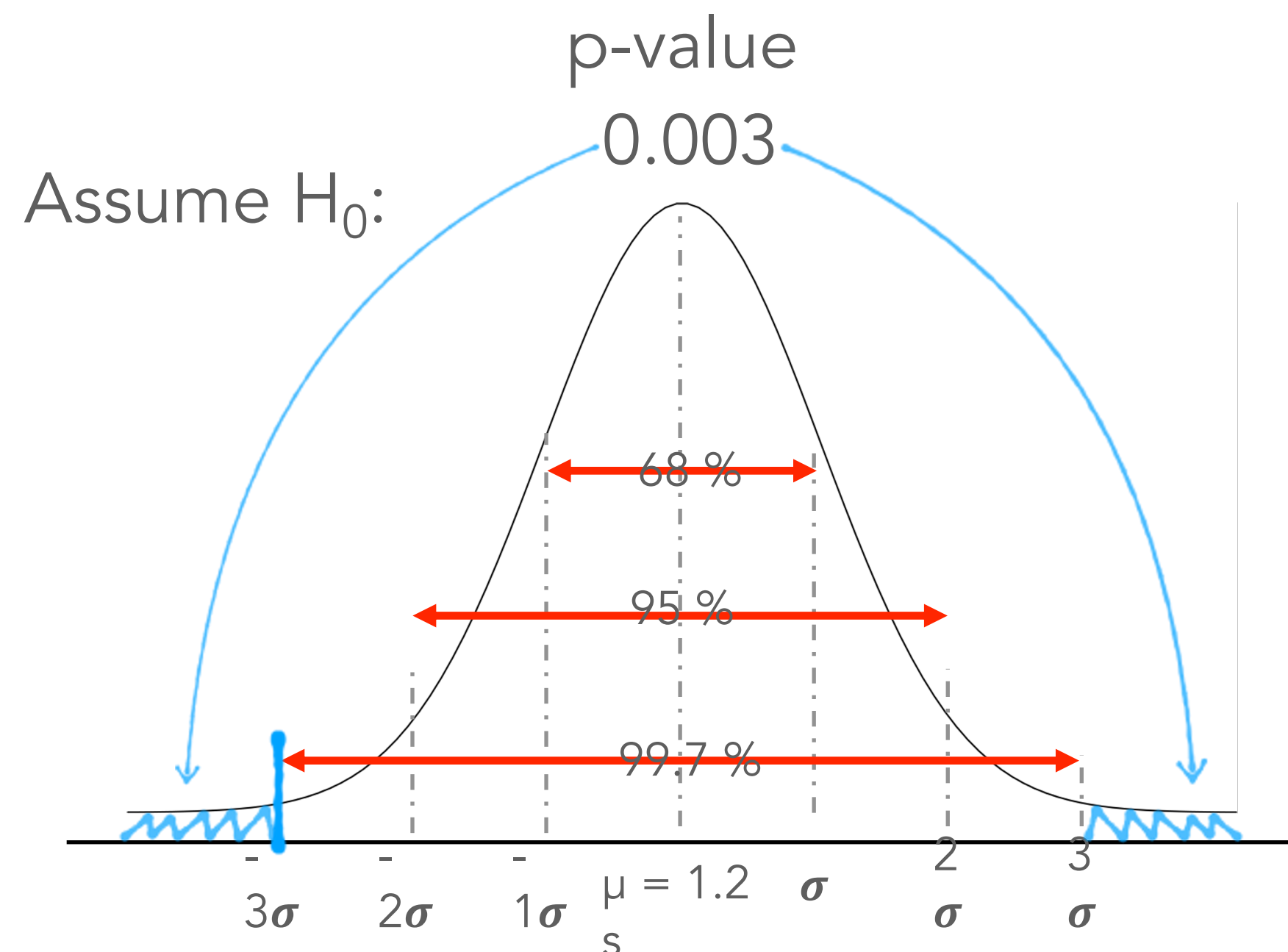
s

$$H_1: \mu \neq 1.2$$

s

Calculate test  
statistic

$$t = \frac{\bar{m} - \mu}{s / \sqrt{n}} = -3$$



This means that the sample mean (1.05) is 3 standard deviations away from the mean

What is the probability of observing a test statistic as extreme as 1.05?

$$\text{p-value} = 2 \min[P(t \leq t_{\text{obs}} | H_0), P(t \geq t_{\text{obs}} | H_0)]$$

# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

$$H_0: \mu = 1.2$$

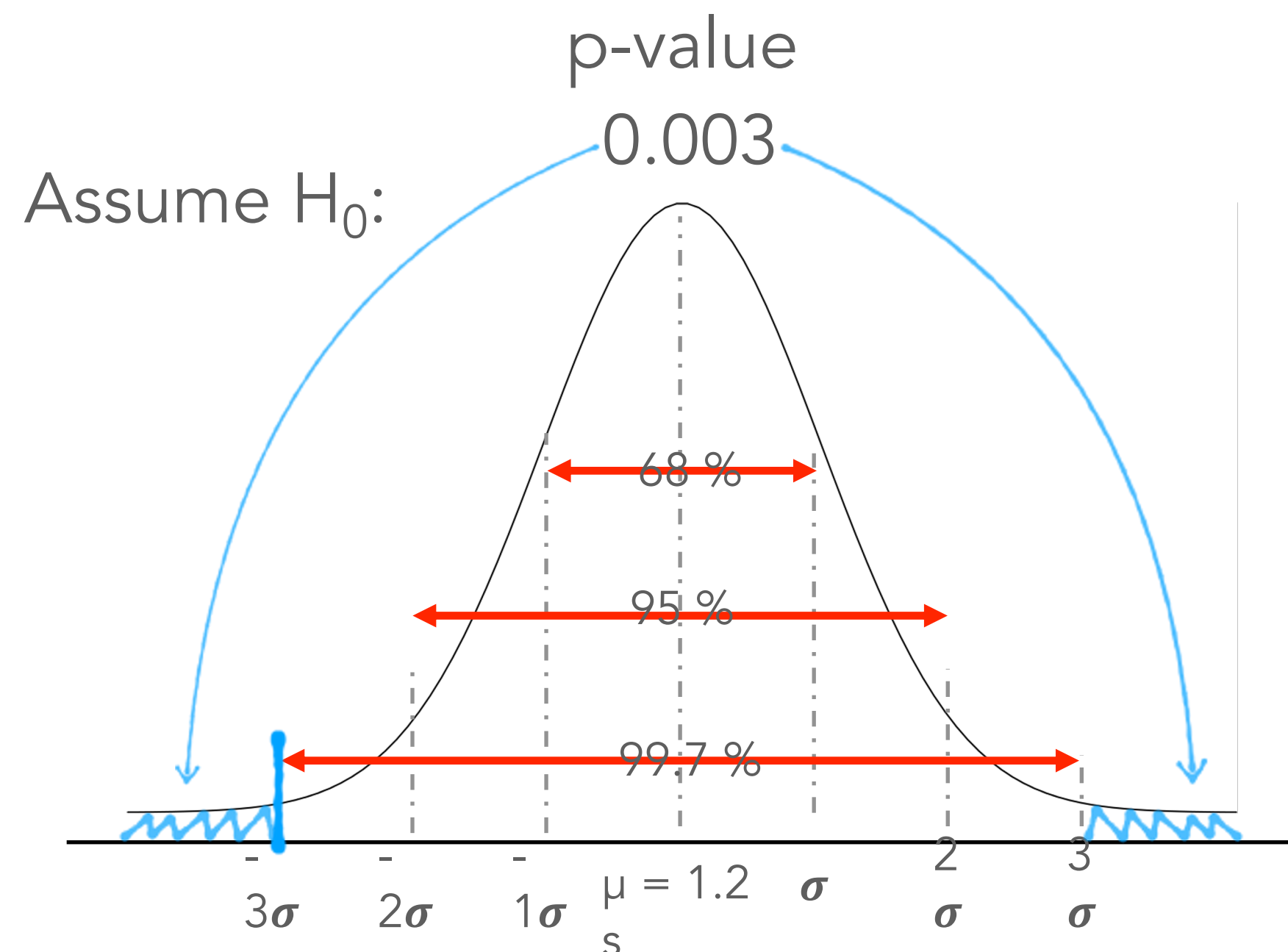
s

$$H_1: \mu \neq 1.2$$

s

Calculate test statistic

$$t = \frac{\bar{m} - \mu}{s/\sqrt{n}} = -3$$



This means that the sample mean (1.05) is 3 standard deviations away from the mean

What is the probability of observing a test statistic as extreme as 1.05?

$$\text{p-value} = 2 \min[P(t \leq t_{\text{obs}} | H_0), P(t \geq t_{\text{obs}} | H_0)]$$

We reject the null hypothesis!

# A simple example

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug subjecting each to neurological stimulus and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats response times is 1.05 seconds with the sample standard deviation of 0.5 seconds. Do you think that the drug has an effect on response time ?

Constructed the null and alternative hypothesis about the population

$$H_0: \mu = 1.2$$

s

$$H_1: \mu \neq 1.2$$

s

Calculate test statistic

$$t = \frac{\bar{m} - \mu}{s / \sqrt{n}} = -3$$

Calculated test statistic

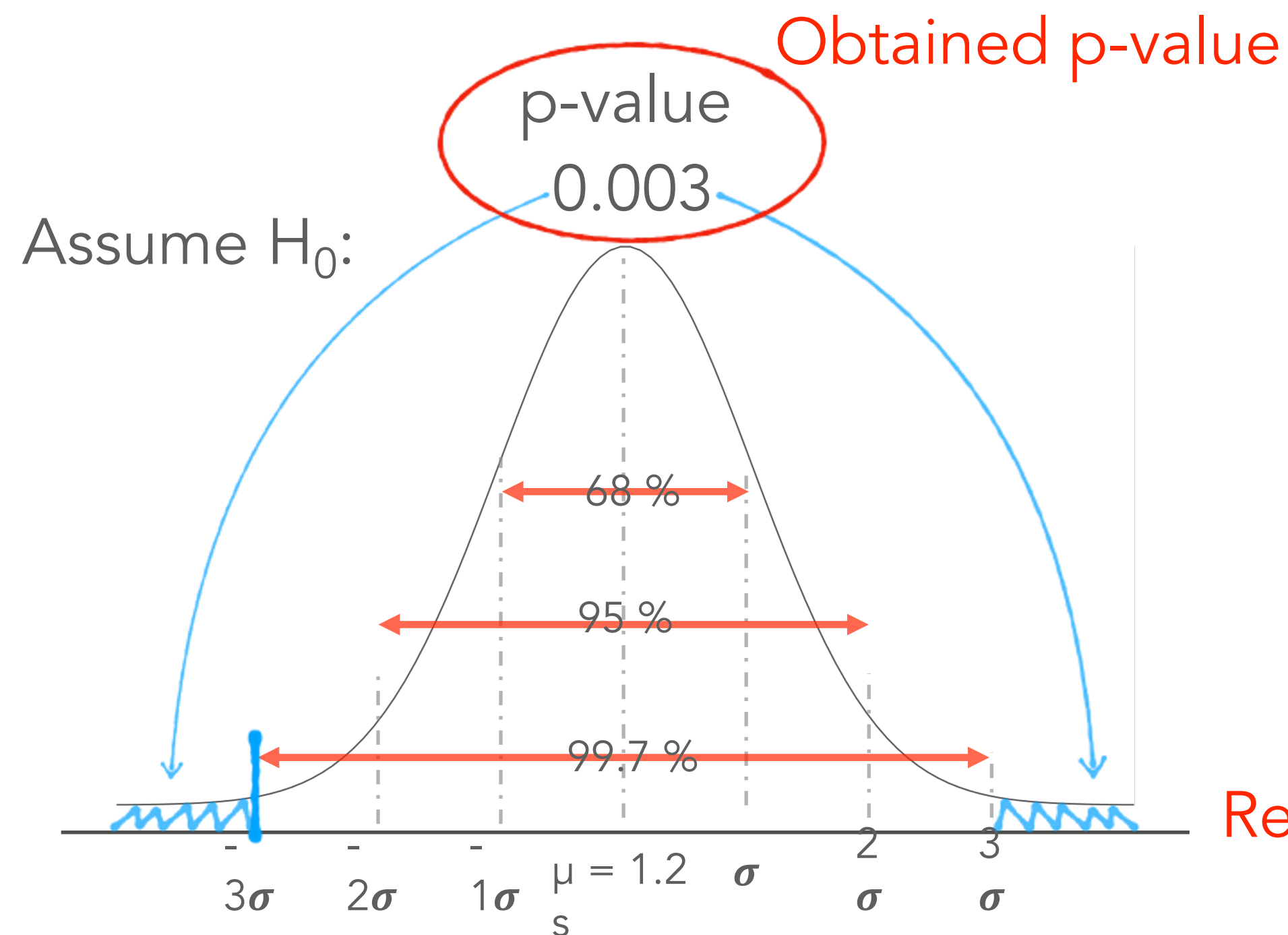
This means that the sample mean (1.05) is 3 standard deviations away from the mean

What is the probability of observing a test statistic as extreme as 1.05?

$$p\text{-value} = 2 \min[P(t \leq t_{\text{obs}} | H_0), P(t \geq t_{\text{obs}} | H_0)]$$

Reached a conclusion

We reject the null hypothesis!



# Key Concepts - Hypothesis Testing

- All statistical tests are based on assumptions!
- All statistics can be wrong
- Statistical tests are probabilistic in nature
- There is always a chance that the result is wrong (even when all assumptions met perfectly):
  - Either significant result when no difference (Type I),
  - Or insignificant results when there is an actual difference (Type II)



# Type I and Type II Errors

- All hypothesis tests involve making a decision:

Is this result significant or not?

- This decision can be wrong in two ways:

Type I error or False positive  
This is when you reject the null hypothesis when it is true

"You're pregnant !"



Type II error or False negative  
This is when you fail to reject the null hypothesis when it isn't true

"You're not pregnant"



# Type I and Type II Errors

$H_0: \mu = 1.2$

$s$

$H_1: \mu \neq 1.2$

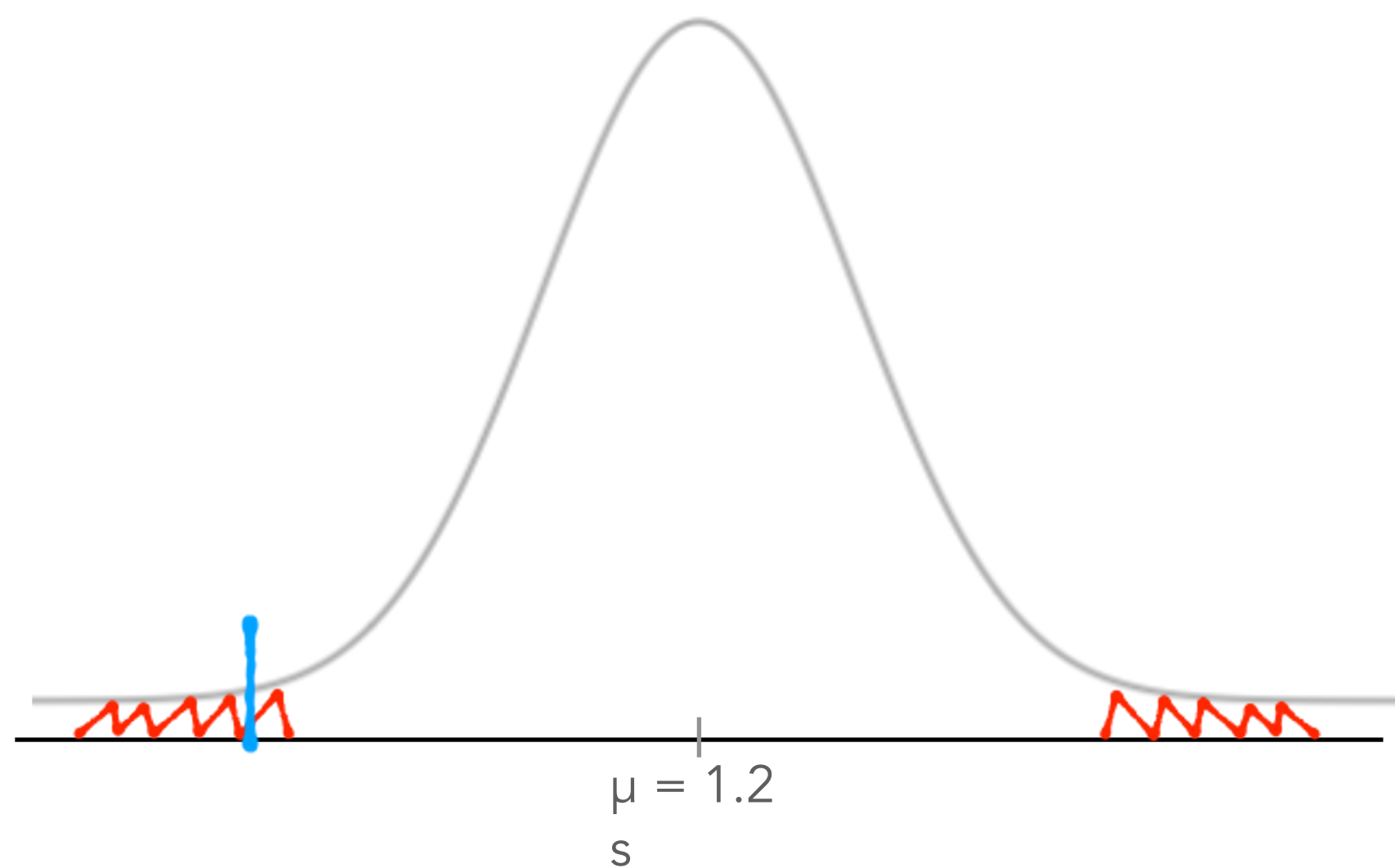
$s$

if p-value  $> \alpha \rightarrow$  do not reject  $H_0$

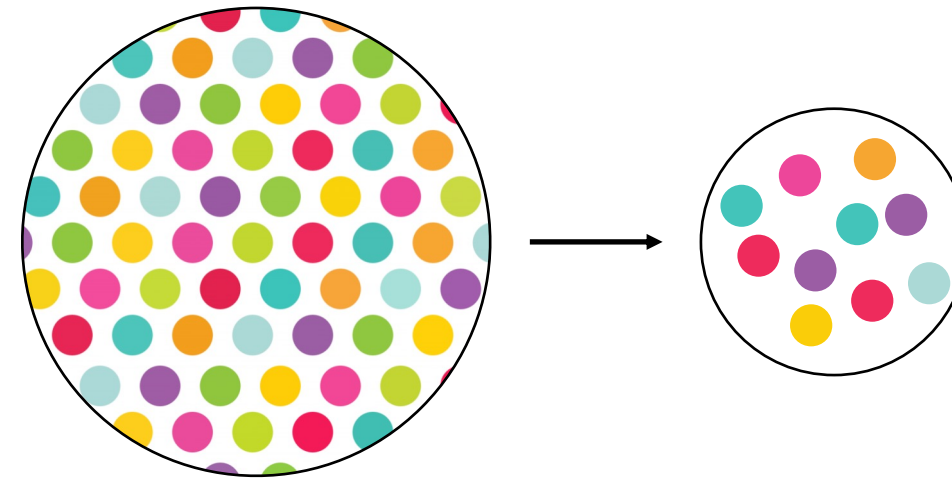
if p-value  $< \alpha \rightarrow$  reject  $H_0$  in favour of  $H_1$

$\alpha=0.05 \rightarrow$  the type I error, the probability of rejecting

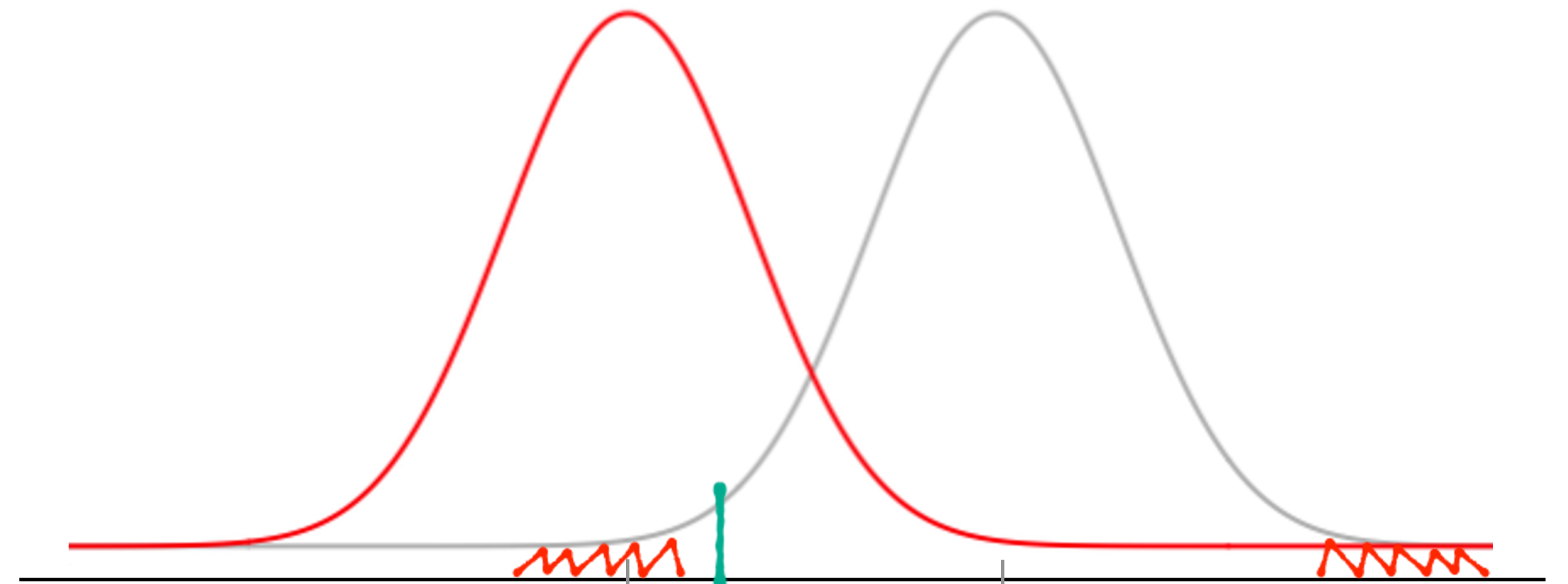
$H_0$  when  $H_0$  is correct



Suppose  $H_1$  true:



Depending on your sampling, you might fail to reject  $H_0$



# Type I and Type II Errors

$$H_0: \mu = 1.2$$

s

$$H_1: \mu \neq 1.2$$

s

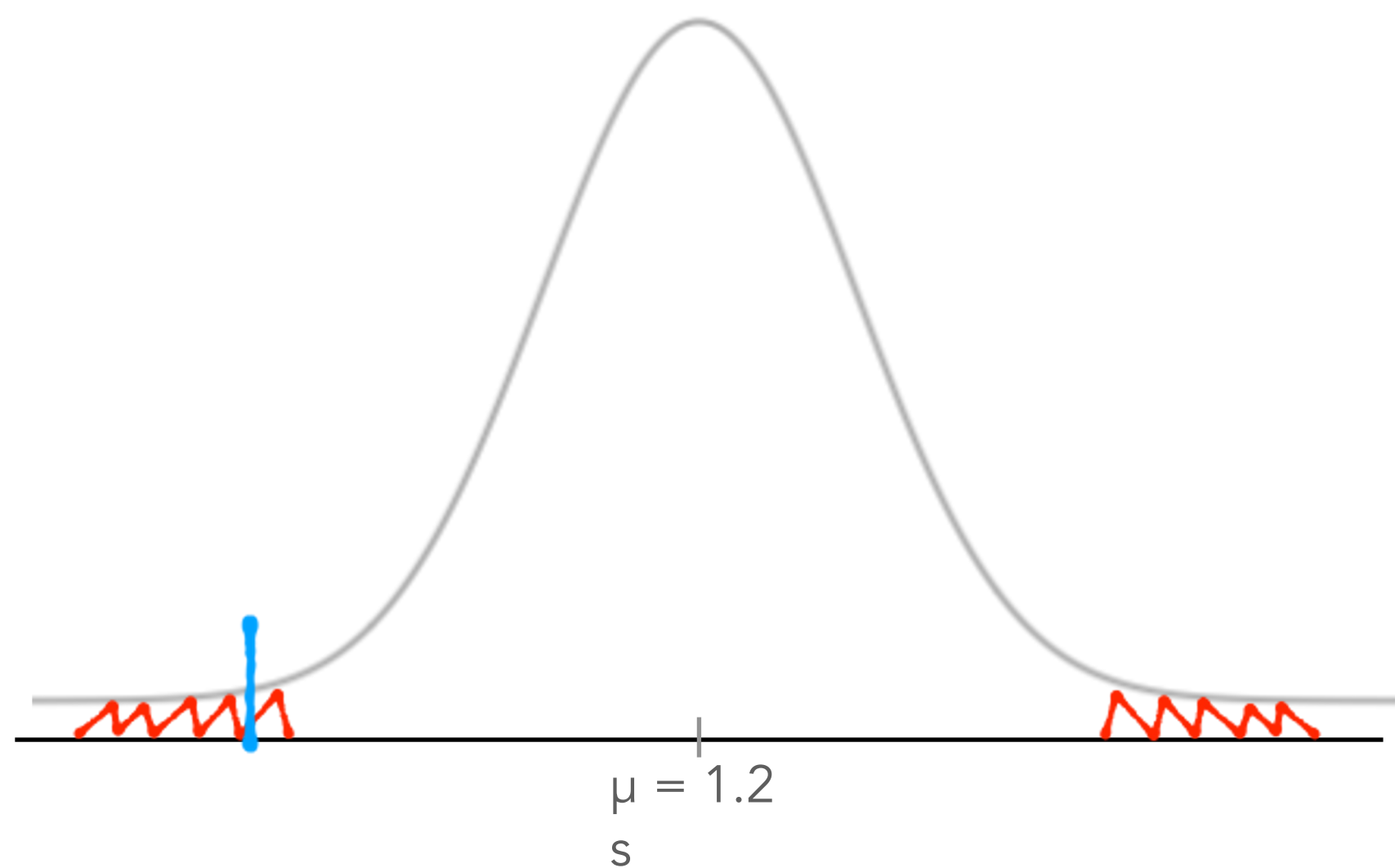
if p-value  $> \alpha \rightarrow$  do not reject  $H_0$

if p-value  $< \alpha \rightarrow$  reject  $H_0$  in favour of

$H_1$

$\alpha=0.05 \rightarrow$  the type I error, the probability of rejecting

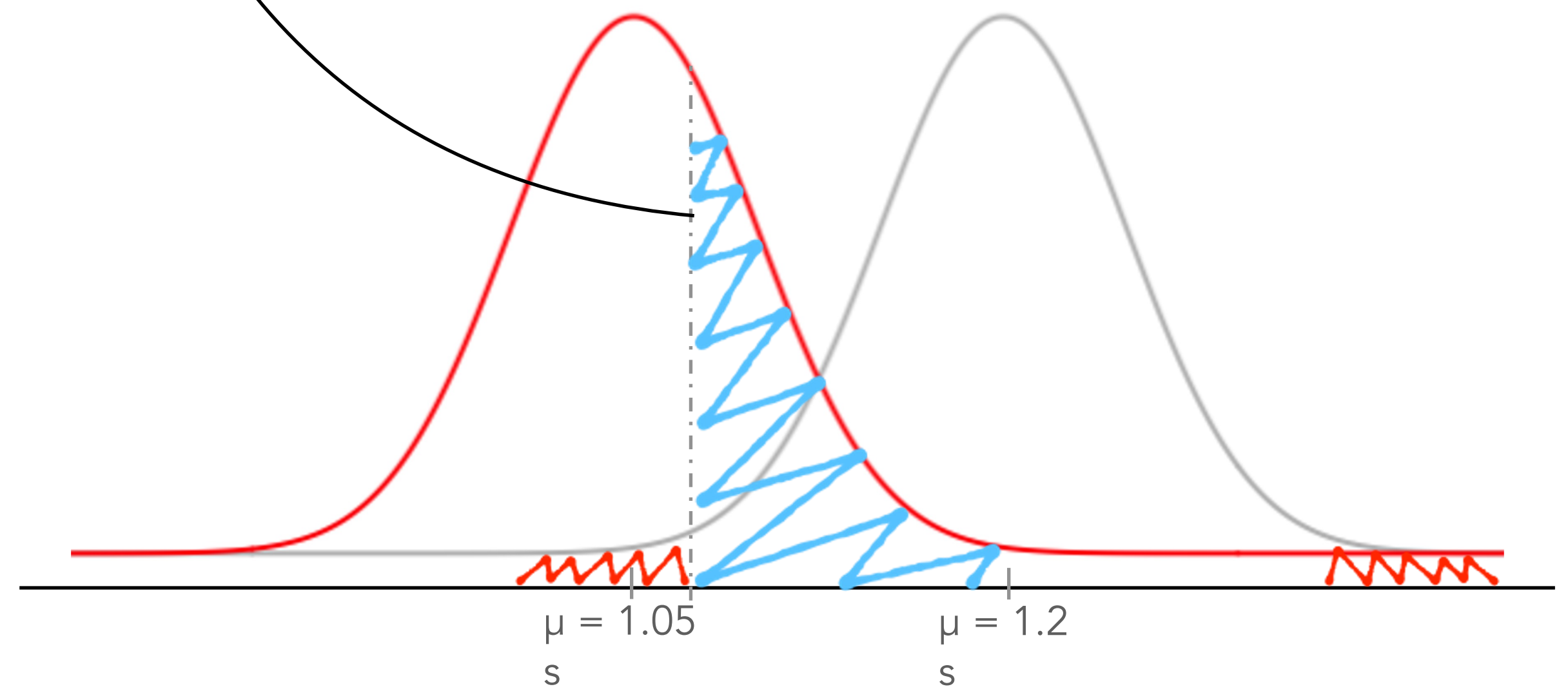
$H_0$  when  $H_0$  is correct



Suppose  $H_1$  true:

$\theta \rightarrow$  the type II error, the probability of not rejecting

$H_0$  when  $H_1$  is correct



# Type I and Type II Errors

$H_0: \mu = 1.2$

s

$H_1: \mu \neq 1.2$

s

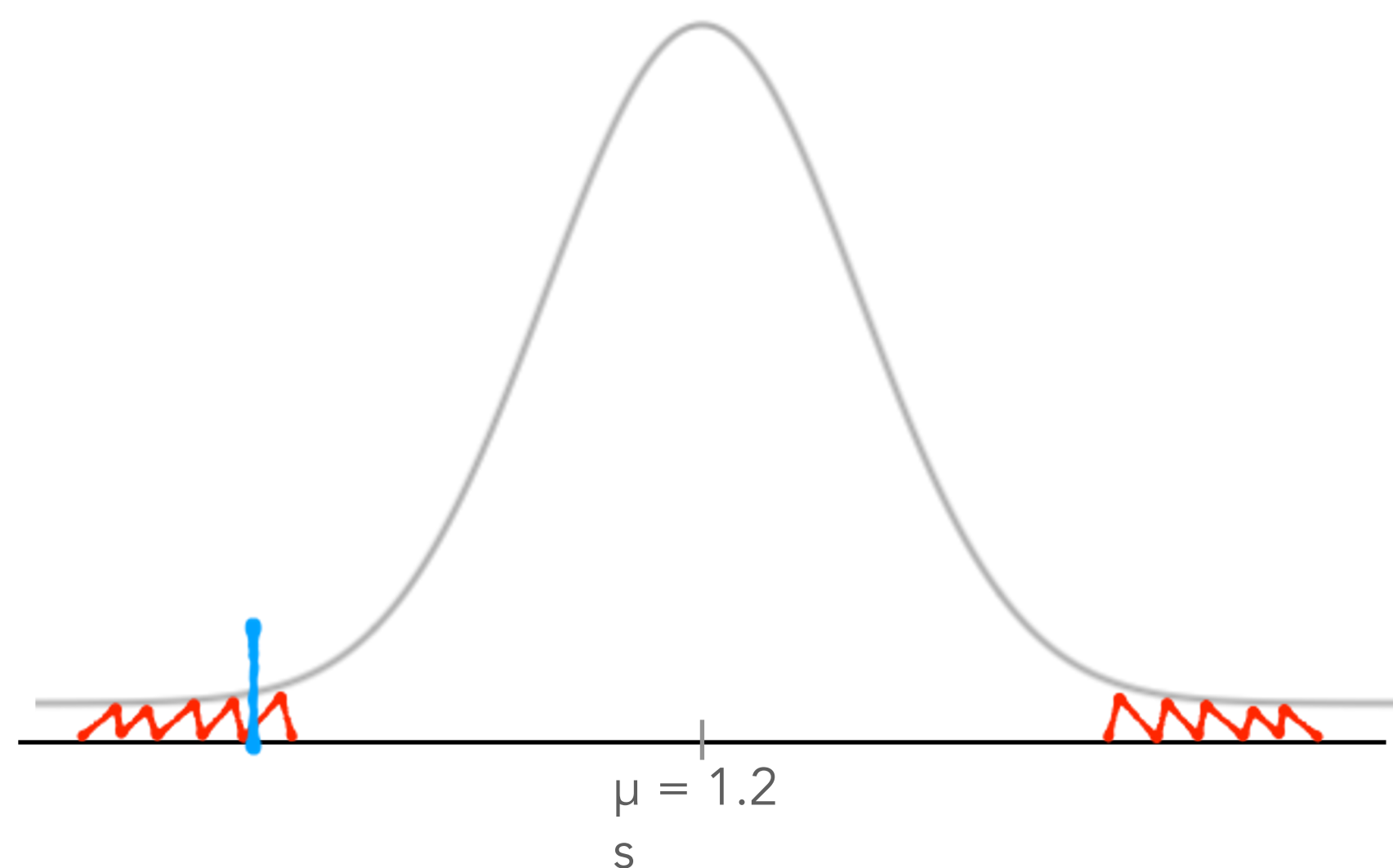
if p-value  $> \alpha \rightarrow$  do not reject  $H_0$

if p-value  $< \alpha \rightarrow$  reject  $H_0$  in favour of

$H_1$

$\alpha=0.05 \rightarrow$  the type I error, the probability of rejecting

$H_0$  when  $H_0$  is correct

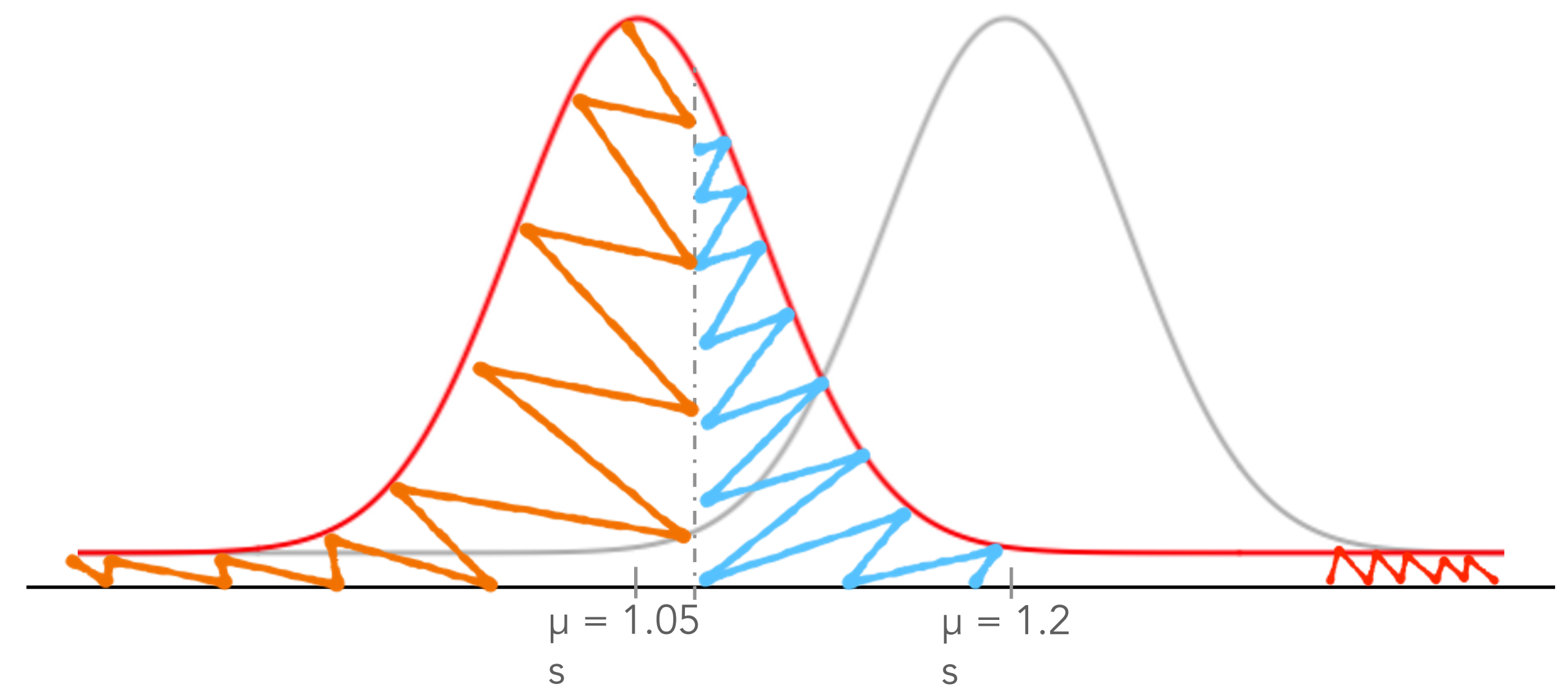


Suppose  $H_1$  true:

$\theta \rightarrow$  the type II error, the probability of not rejecting

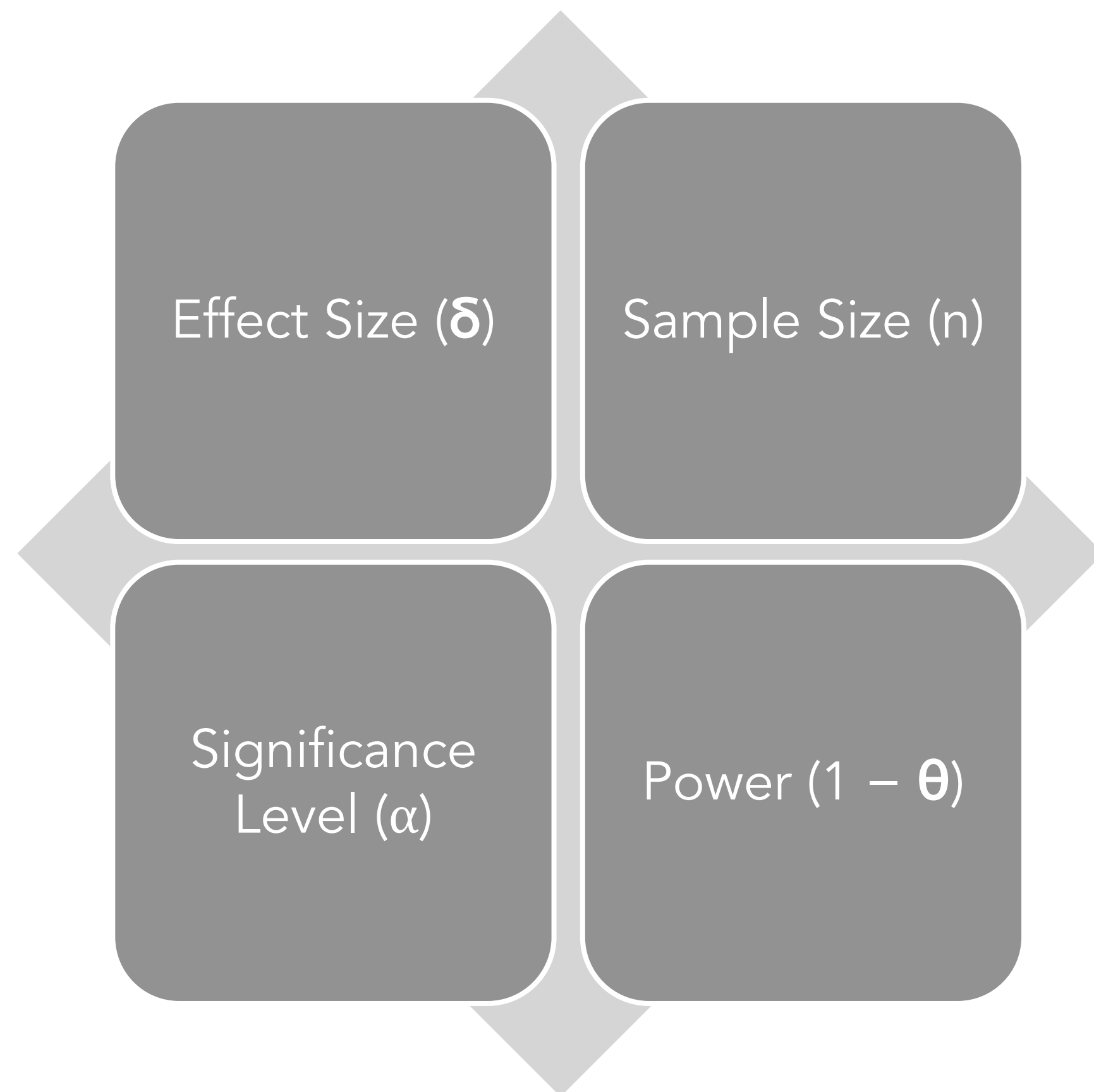
$H_0$  when  $H_1$  is correct

$1 - \theta \rightarrow$  Power is the probability that we actually detect an effect that exists



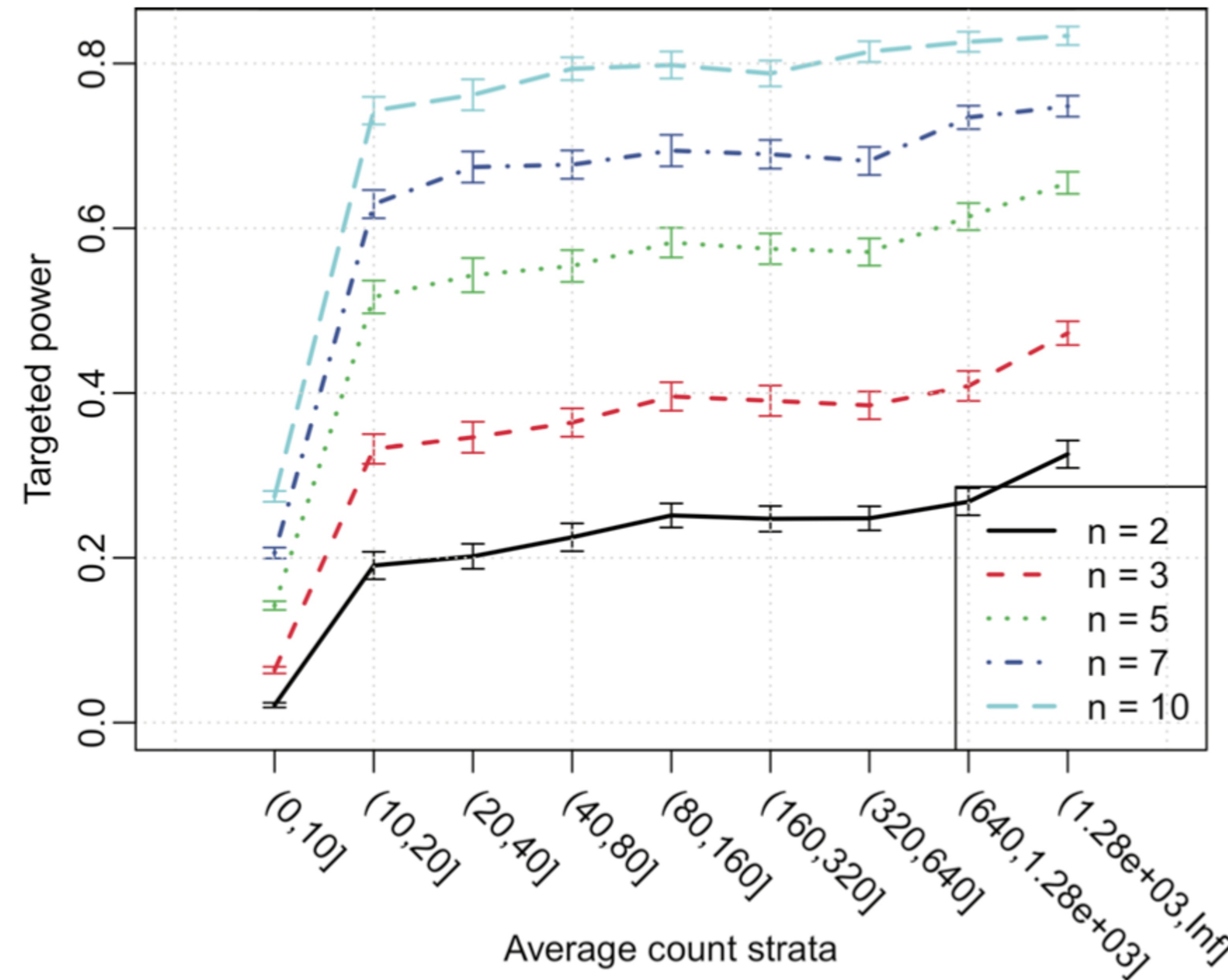


# Power Analysis



- The four concepts are linked
- If we know three, we can work out the fourth
- **Power calculation:** Aim is to define the probability ( $1 - \theta$ ) to detect an effect size of interest ( $\delta$ ) at the  $\alpha$  level with a sample size of  $n$  biological replicates
- **Sample size calculation:** Aim is to define the sample size ( $n$ ) allowing to detect an effect size of interest ( $\delta$ ) at the  $\alpha$  level with a given probability ( $1 - \theta$ ).

# Power Analysis in Differential Expression Analysis



(Wu, Wang and Wu (2015))



# Outline

- Experimental Design
- Statistical Concepts - Bite size statistics
- Statistical aspects of bulk RNA-seq analysis

# Statistical Aspects of Differential Expression Analysis

## Linear Modeling

Model the expression of each gene as linear combination of explanatory factors (eg. treatment, age, sex, etc.)

$$Y \sim Normal(\mu, \sigma)$$

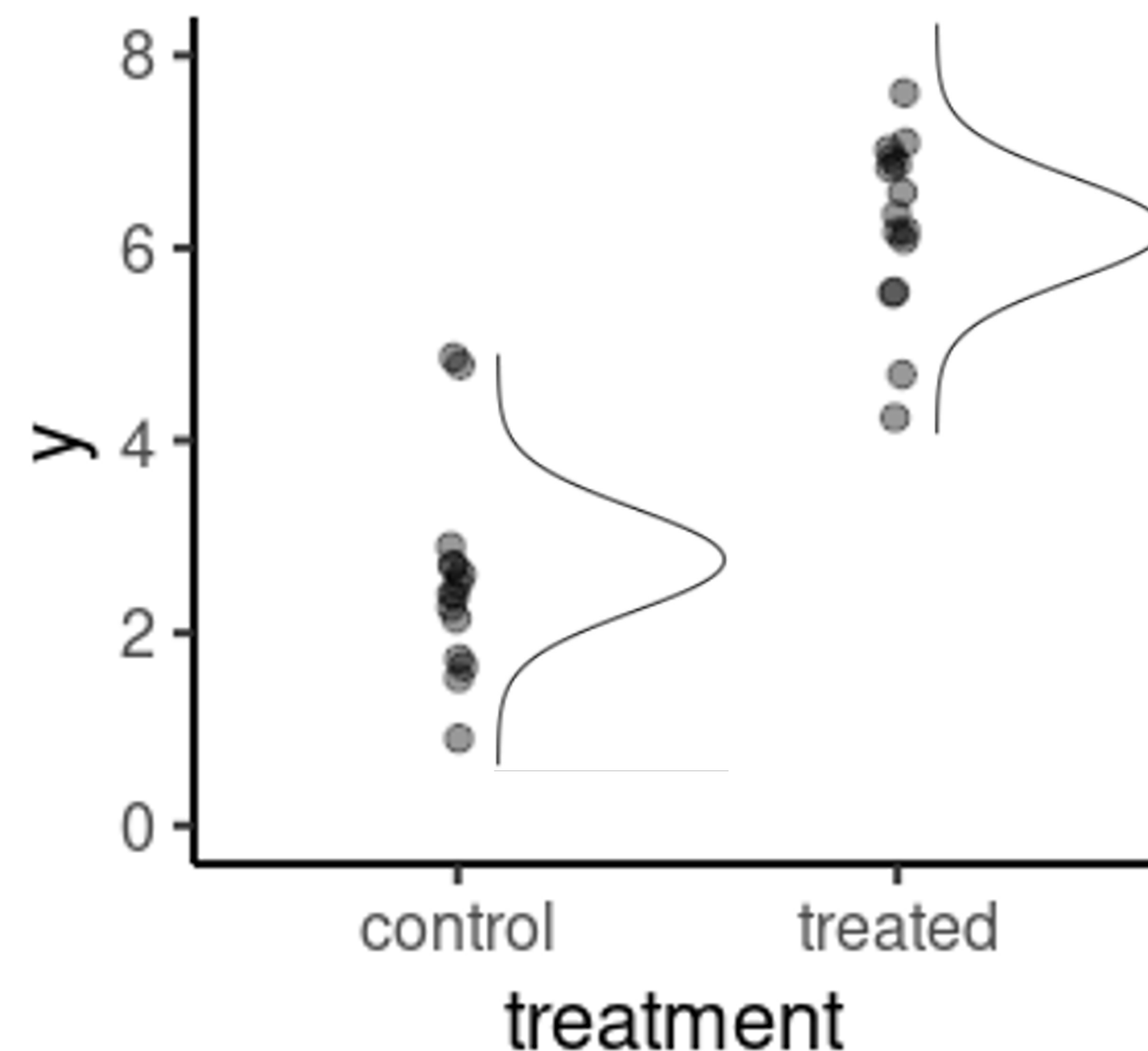
$$\mu = \beta_0 + \beta_1 * treatment + \beta_2 * age + \dots$$

$y$  - expression of the gene

$\beta_i$  - parameters we want to estimate from the data

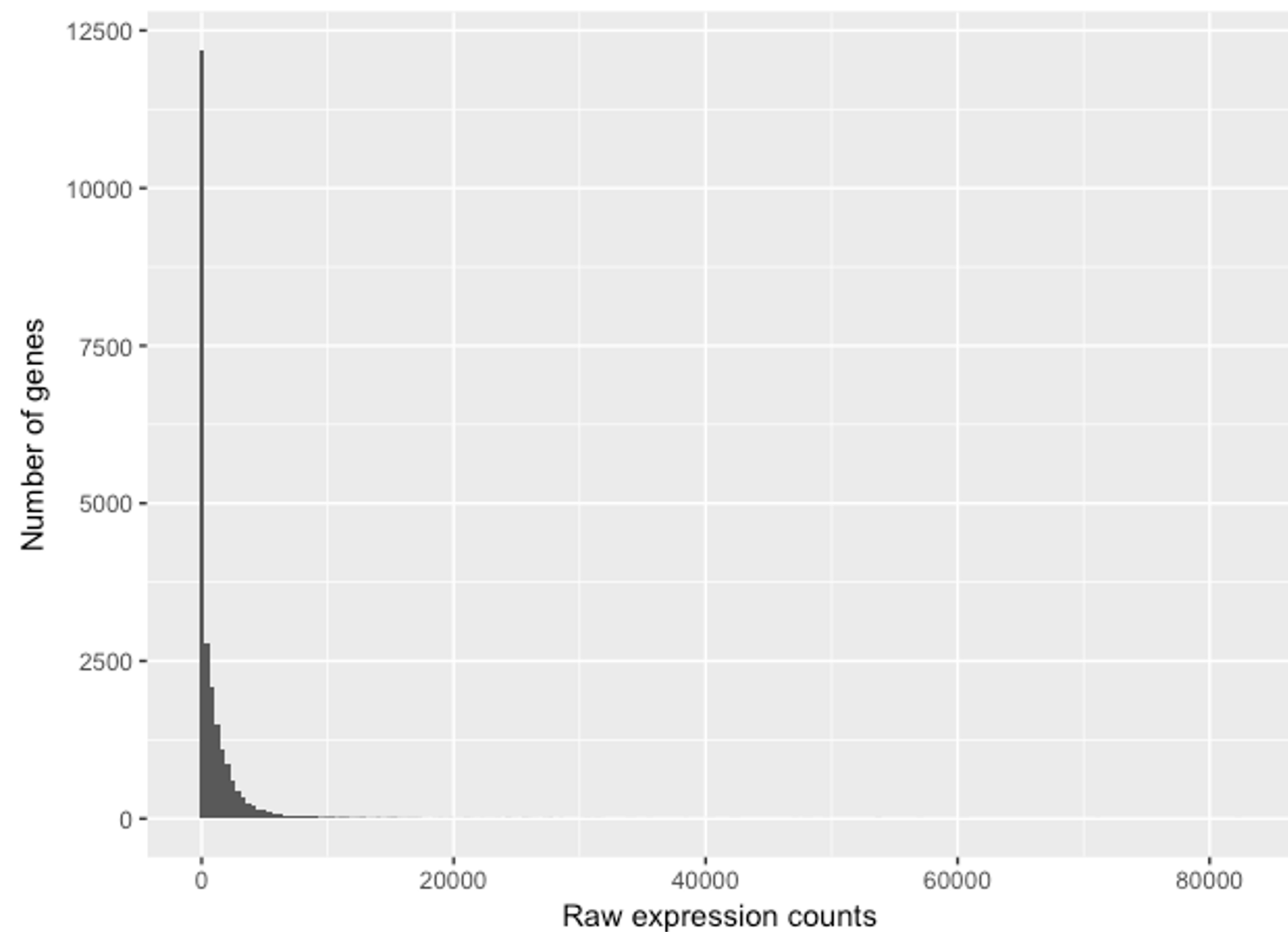
$\beta_0$  - the “intercept” (the value of expression when all other parameters are set at a reference level)

$\sigma$  - the standard deviation (uncertainty) of our model (also estimated from the data)



# Statistical Aspects of Differential Expression Analysis

## Characteristics of RNA-seq data



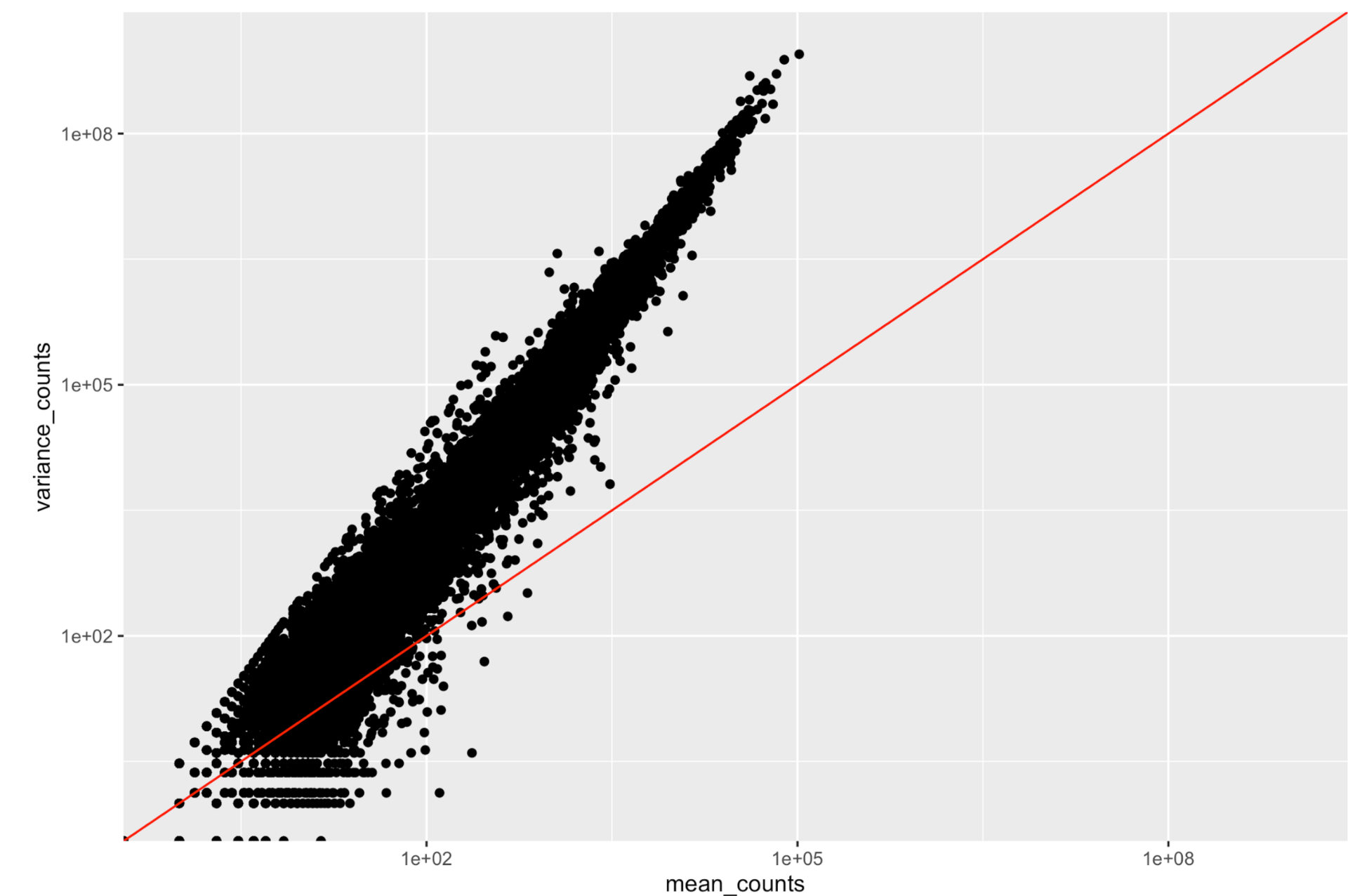
This plot illustrates some common features of RNA-seq count data:

- a low number of counts associated with a large proportion of genes
- a long right tail due to the lack of any upper limit for expression
- large dynamic range

Looking at the shape of the histogram, we see that it is not normally distributed.

# Statistical Aspects of Differential Expression Analysis

## Characteristics of RNA-seq data



To assess the properties of the data we are working with, we can look at the mean-variance relationship.

For the genes with high mean expression, the variance across replicates tends to be greater than the mean (scatter is above the red line).

Essentially, the Negative Binomial is a good approximation for data where the mean  $<$  variance, as is the case with RNA-Seq count data.

# Statistical Aspects of Differential Expression Analysis

## Linear Modeling

Model the expression of each gene as linear combination of explanatory factors (eg. treatment, age, sex, etc.)

$$\text{Counts} \sim NB(\mu, \phi)$$

$$\mu = sq$$

$$\log_2(q) = \beta_0 + \beta_1 * \text{treatment} + \beta_2 * \text{age} + \dots$$

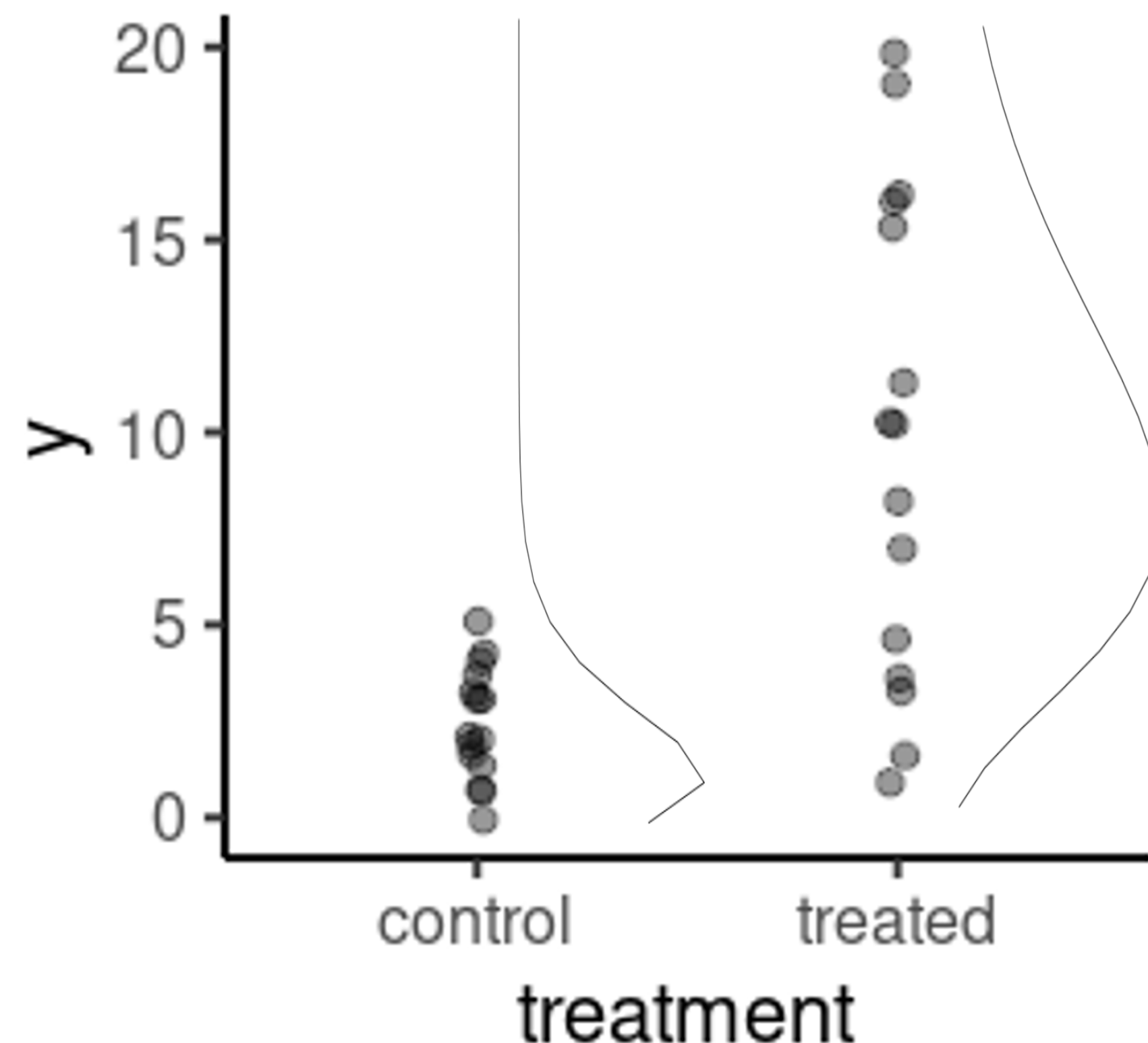
**counts** - expression of the gene

$\beta_i$  - parameters we want to estimate from the data

$\beta_0$  - the “intercept” (the value of expression when all other parameters are set at a reference level)

$\phi$  - the “dispersion” (uncertainty) of our model (also estimated from the data)

**s** - scaling factor (sequencing depth and transcript composition)





# Statistical Aspects of Differential Expression Analysis

## Linear Modeling

Model the expression of each gene as linear combination of explanatory factors (eg. treatment, age, sex, etc.)

$$\text{Counts} \sim NB(\mu, \phi)$$

$$\mu = sq$$

$$\log_2(q) = \beta_0 + \beta_1 * \text{treatment} + \beta_2 * \text{age} + \dots$$

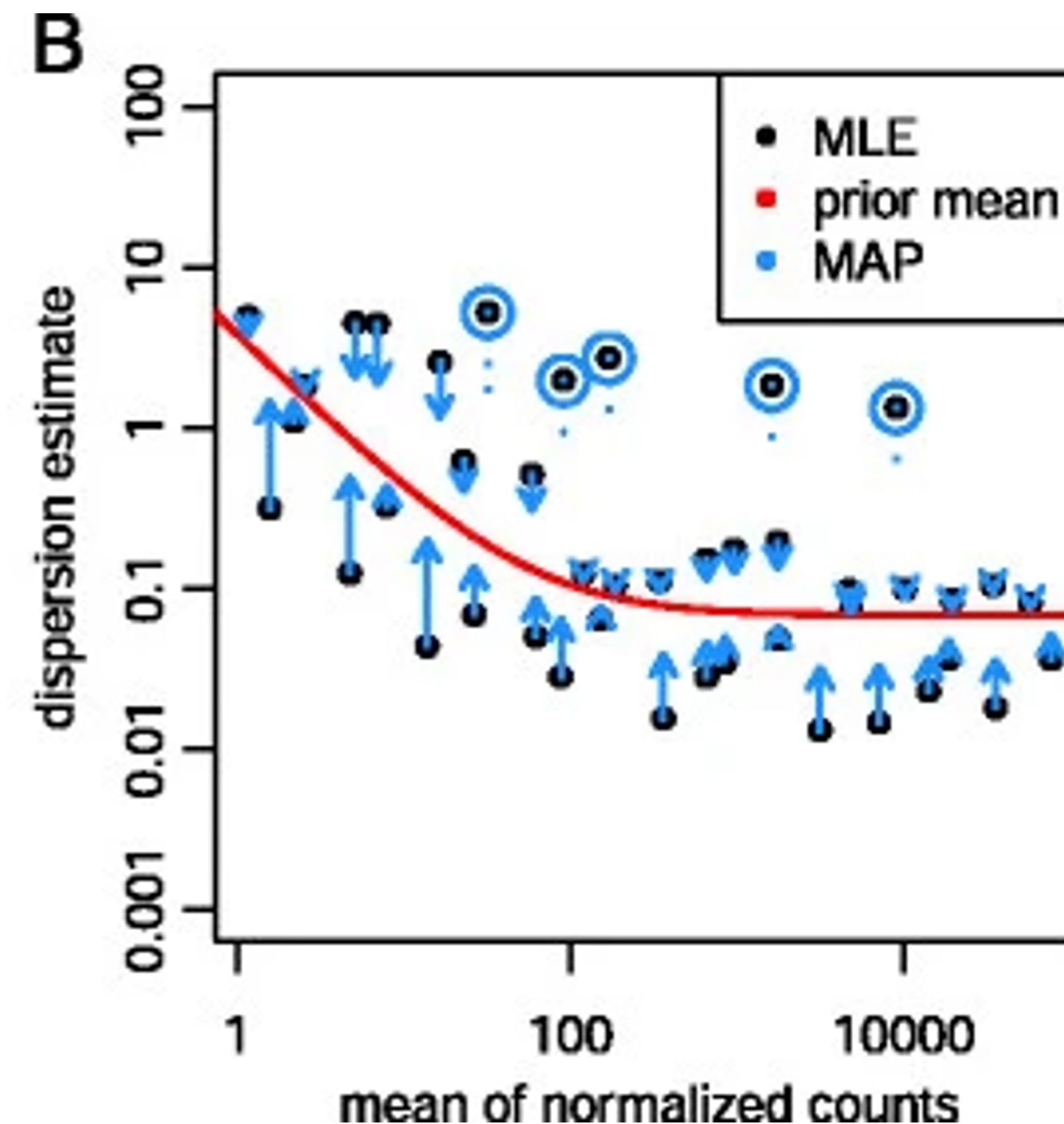
**counts** - expression of the gene

$\beta_i$  - parameters we want to estimate from the data

$\beta_0$  - the “intercept” (the value of expression when all other parameters are set at a reference level)

$\phi$  - the “dispersion” (uncertainty) of our model (also estimated from the data)

**s** - scaling factor (sequencing depth and transcript composition)



# Statistical Aspects of Differential Expression Analysis

## Linear Modeling

Model the expression of each gene as linear combination of explanatory factors (eg. treatment, age, sex, etc.)

$$\text{Counts} \sim NB(\mu, \phi)$$

$$\mu = sq$$

$$\log_2(q) = \beta_0 + \beta_1 * \text{treatment} + \beta_2 * \text{age} + \dots$$

**counts** - expression of the gene

$\beta_i$  - parameters we want to estimate from the data

$\beta_0$  - the “intercept” (the value of expression when all other parameters are set at a reference level)

$\phi$  - the “dispersion” (uncertainty) of our model (also estimated from the data)

**s** - scaling factor (sequencing depth and transcript composition)

**Coefficients** are estimated for each sample group along with their standard error.

The coefficients are the estimates for the **log2 fold-changes**, and will be used as input for hypothesis testing.

# Statistical Aspects of Differential Expression Analysis

## Linear Modeling

$$Counts \sim NB(\mu, \phi)$$

$$\mu = sq$$

$$\log_2(q) = \beta_0 + \beta_1 * treatment + \beta_2 * age + \dots$$

**counts** - expression of the gene

$\beta_i$  - parameters we want to estimate from the data

$\beta_0$  - the “intercept” (the value of expression when all other parameters are set at a reference level)

$\phi$  - the “dispersion” (uncertainty) of our model (also estimated from the data)

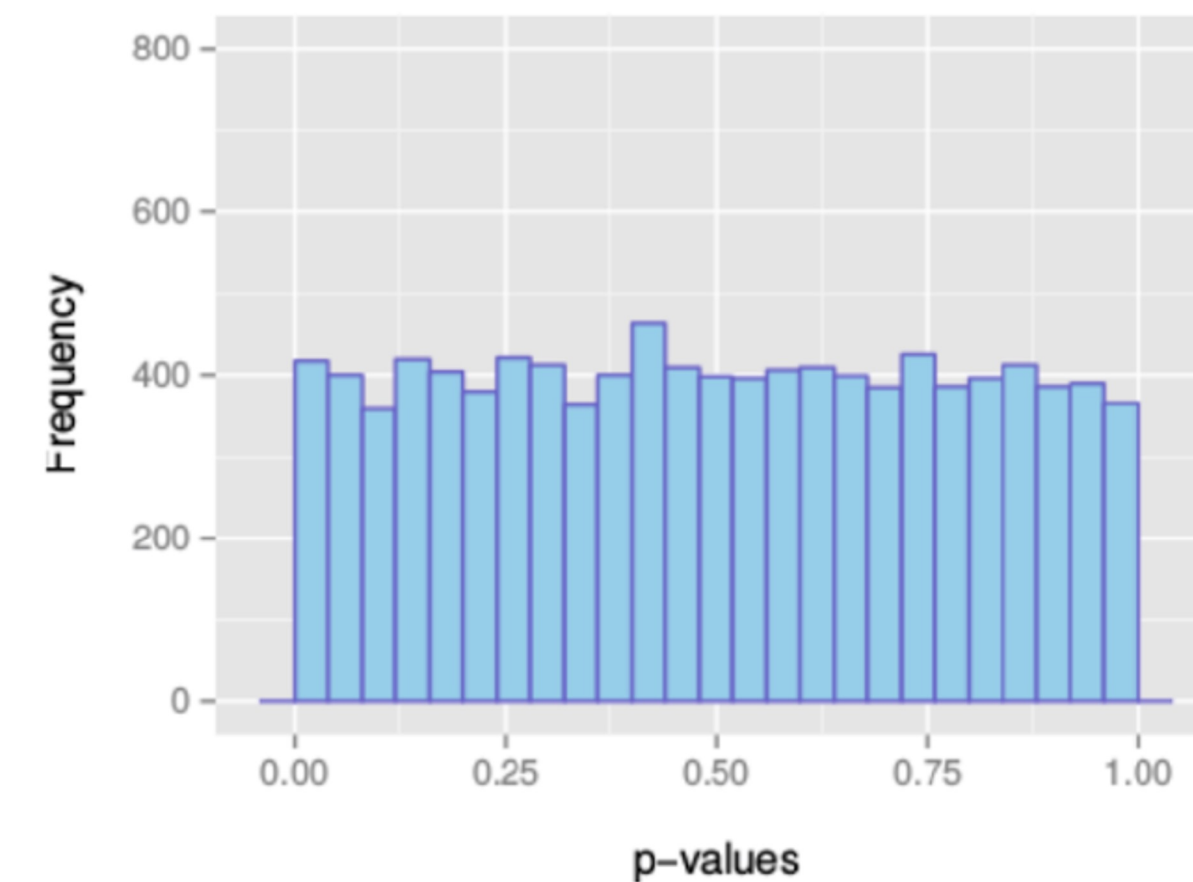
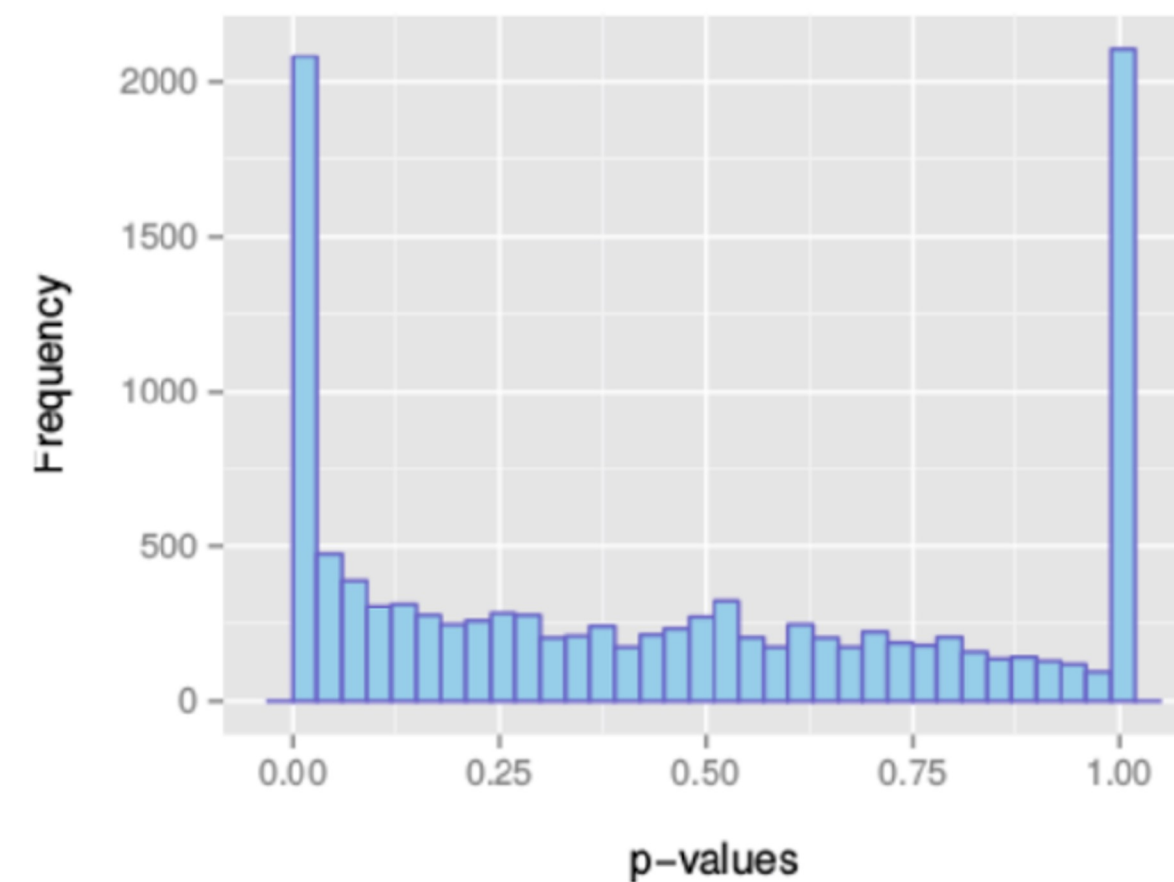
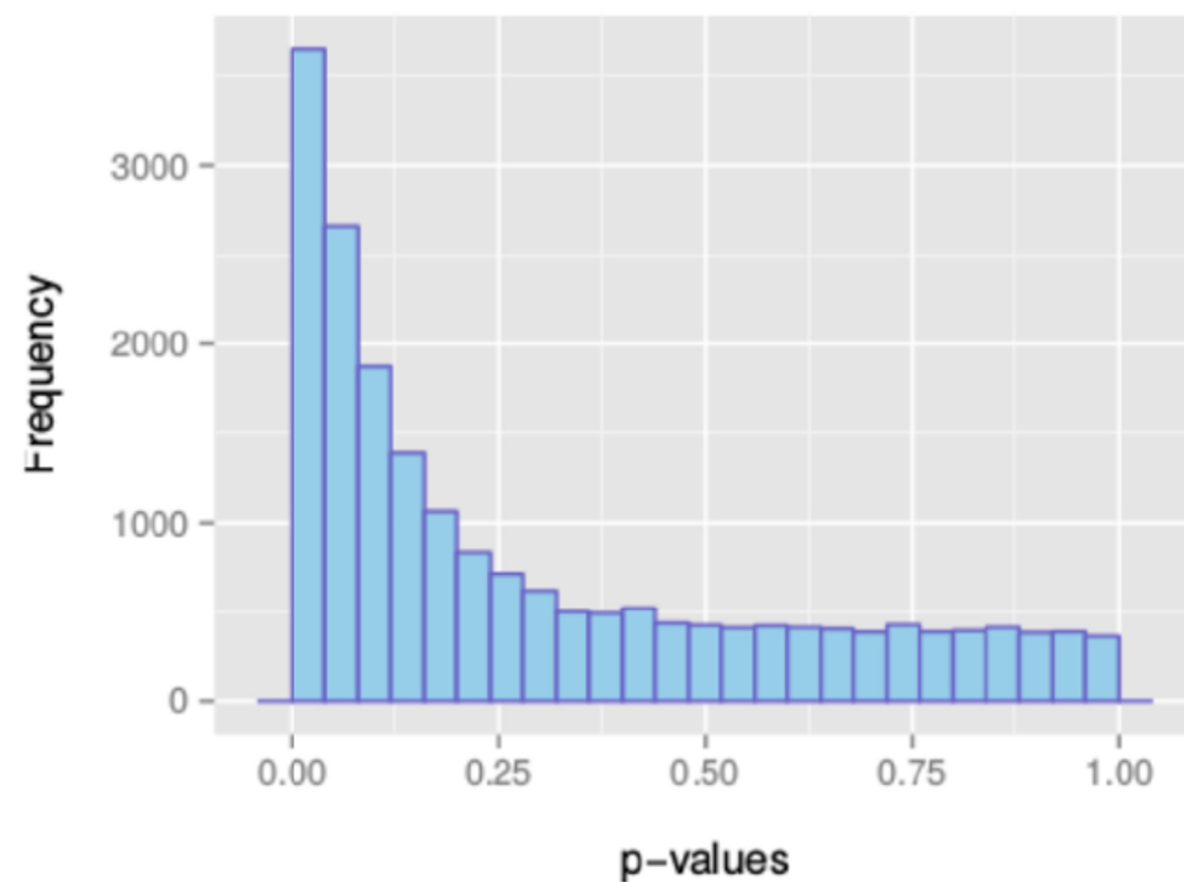
**s** - scaling factor (sequencing depth and transcript composition)

Summary:

- Use **negative binomial linear regression** to model gene expression in RNA-seq
- Calculate **size factors** for each sample to account for differences in sequencing depth and transcript composition between samples
- Estimate **dispersion** for each gene by “borrowing” information across genes for more precise estimates when sample sizes are small (as is typical in RNA-seq experiments)
- Estimate model **coefficients** which are used to define test hypothesis ( $\beta_i = 0$ )

# P-value Histograms

Examples of expected overall distribution



(a) : the most desirable shape

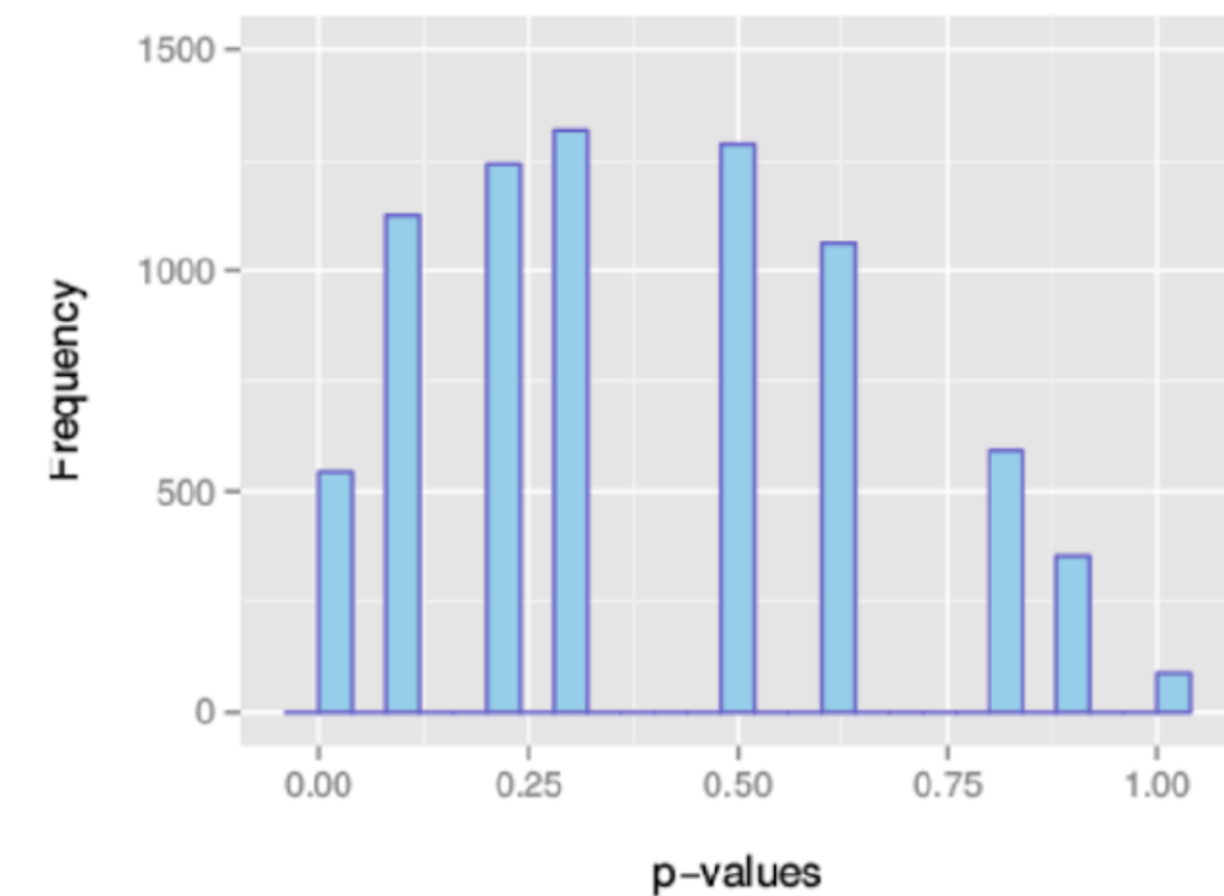
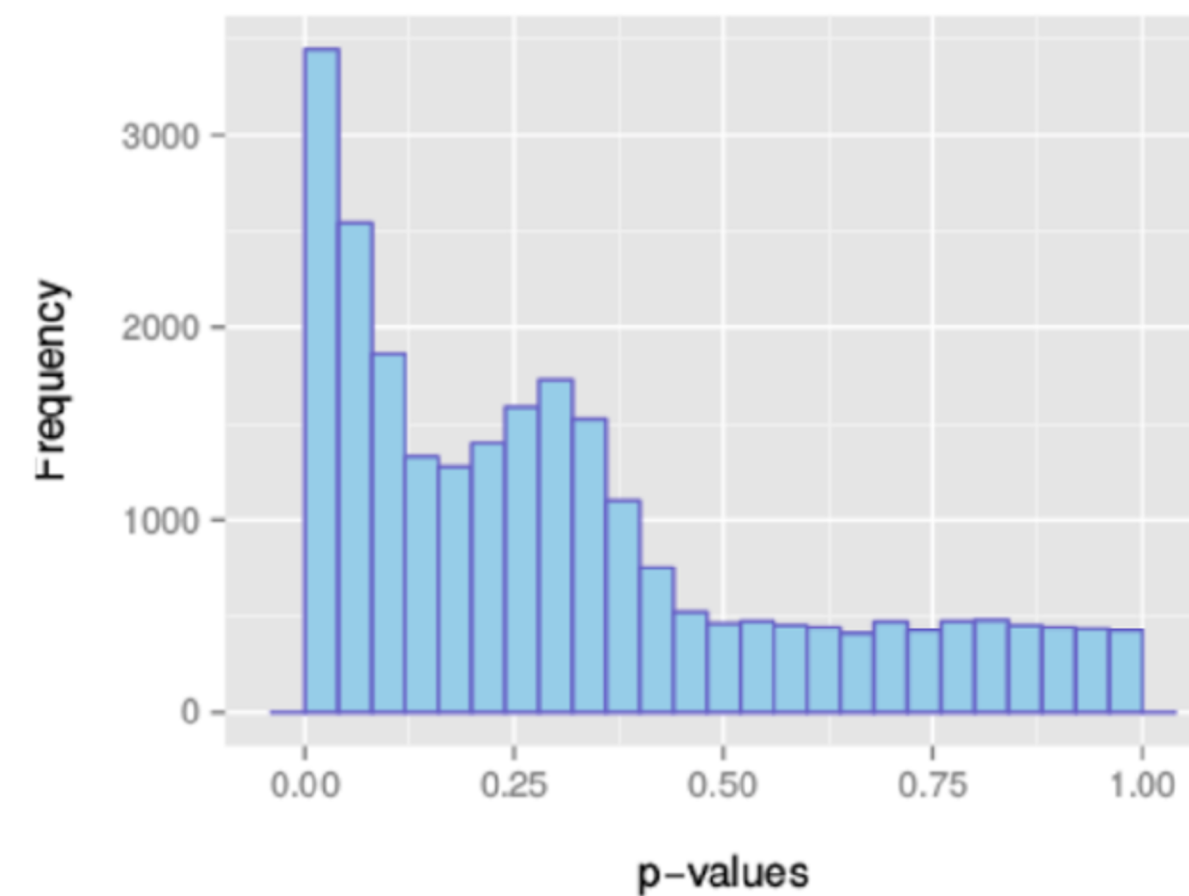
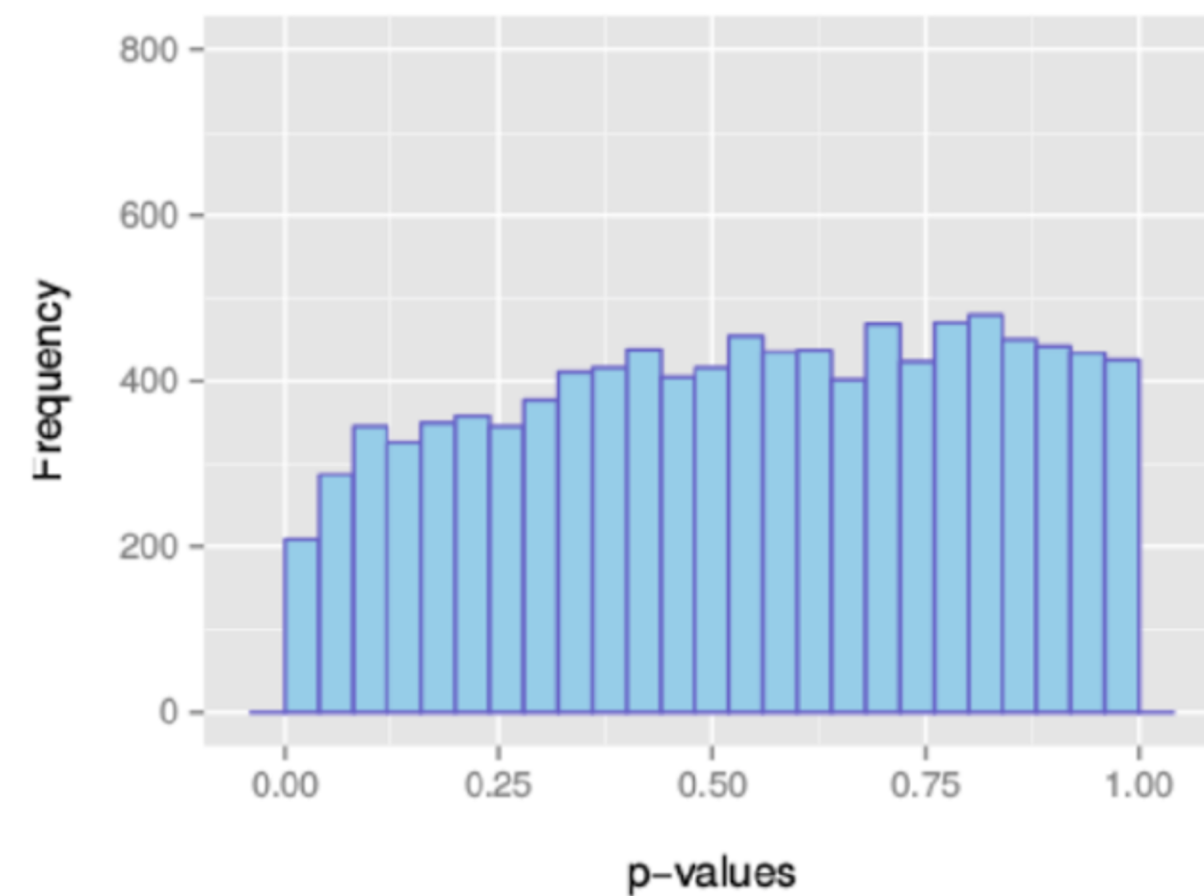
(b) : very low counts genes usually have large p-values

(c) : do not expect positive tests after correction



# P-value Histograms

Examples of unexpected overall distribution



- (a) : indicates a batch effect (confounding hidden variables)
- (b) : the test statistics may be inappropriate (due to strong correlation structure for instance)
- (c) : discrete distribution of p-values : unexpected



# Multiplicity Correction

- A gene with a significance cut-off of  $\alpha = 0.05$ , means there is a 5% chance it is a false positive.
- If we test for 20,000 genes for differential expression at  $\alpha = 0.05$ , we would expect to find 1,000 genes by chance
- If we found 3000 genes to be differentially expressed total, roughly one third of our genes are false positives!
- The more genes we test, the more we inflate the false positive rate. This is the multiple testing problem.

# Multiplicity Correction

- Bonferroni: The adjusted p-value is calculated by:  $\alpha^k$  (k = total number of tests). This is a very conservative approach
- FDR/Benjamini-Hochberg: Benjamini and Hochberg (1995) defined the concept of FDR and created an algorithm to control the expected FDR below a specified level given a list of independent p-values.

# Conclusions

- Assumptions assumptions assumptions