



CANCER
RESEARCH
UK

Cambridge
Institute

Together we are
beating cancer

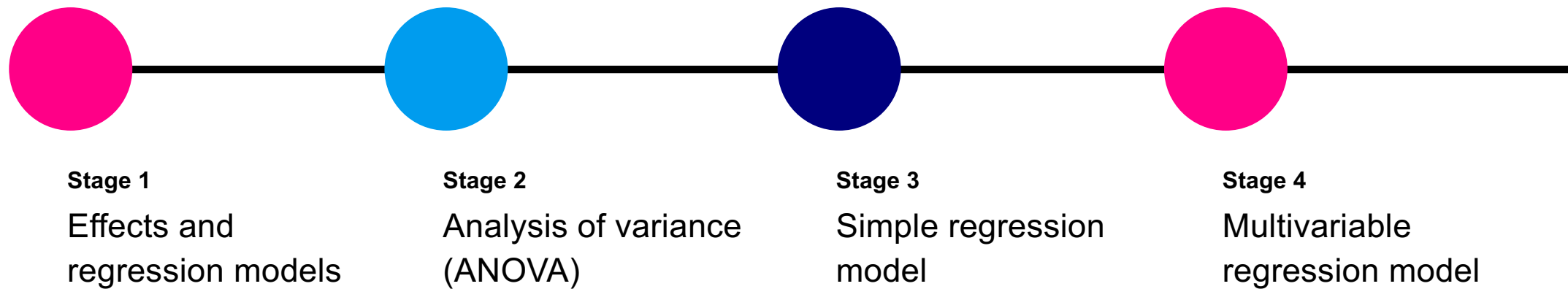
Luca Porcu & Chandra Chilamakuri (Bioinformatics core)

21st February 2025

Linear regression models

Fixed-effects models

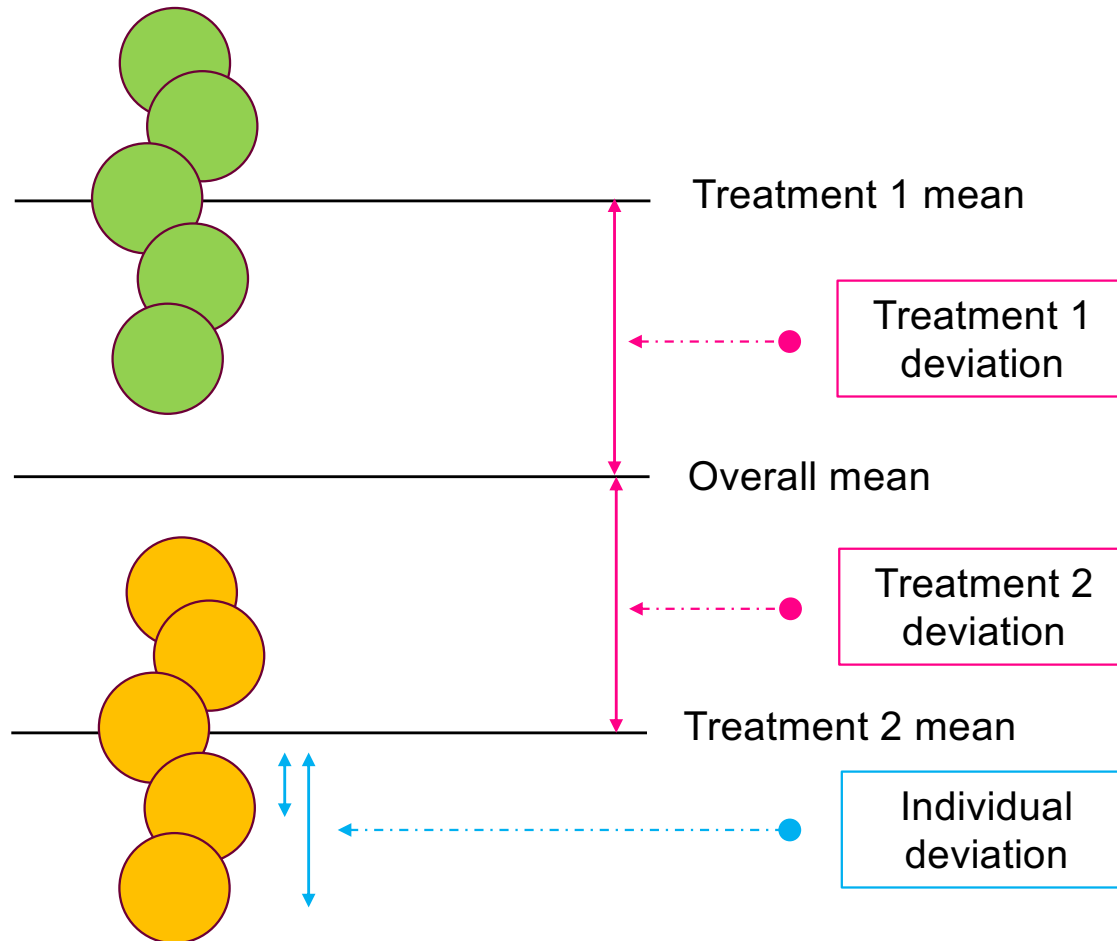
Process flow





CANCER
RESEARCH
UK

Cambridge
Institute



Analysis of variance (ANOVA)

Definition and classification

10.15 -10.40 am

Together we are
beating cancer

Fisher's one-way ANOVA

4

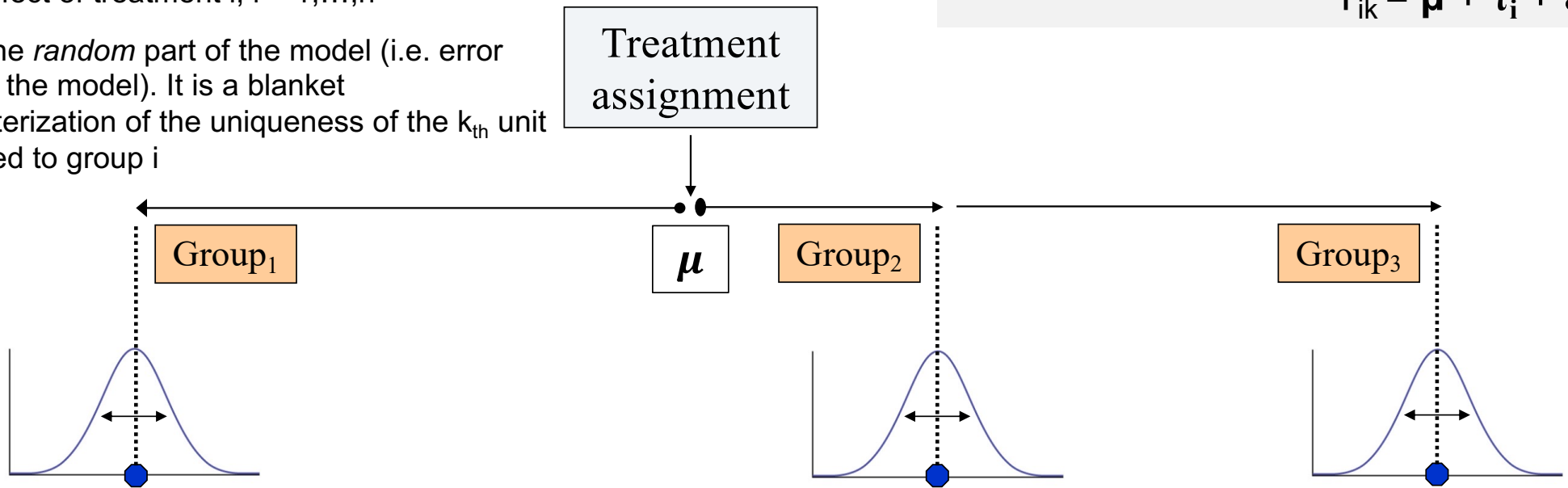
1. The unit k (e.g. mouse), $k = 1, \dots, u_i$; $N = \sum_i u_i$

2. τ_i : effect of treatment i , $i = 1, \dots, n$

3. ε_{ik} : the *random* part of the model (i.e. error term of the model). It is a blanket characterization of the uniqueness of the k_{th} unit assigned to group i

Equation of the statistical model:

$$Y_{ik} = \mu + \tau_i + \varepsilon_{ik}$$



Assumptions of ANOVA (ANalysis Of VAriance) models are the following:

- The effect of each factor is additive on μ (i.e. population mean) parameter
- ε_{ik} is assumed to be independent of one another and normally distributed with mean = 0 and common standard deviation = σ

Fisher's one-way ANOVA

5

Hypothesis to test: $\tau_1 = \dots = \tau_n$

Test statistic:

Source of variation	Sum of Squares	Degrees of freedom	Mean Squares	F _{df1,df2}	P-value
Treatment	SSB = $\sum_i u_i (m_i - M)^2$	df1 = n -1	MSB = SSB / (n - 1)	MSB / MSE	0.023
Residuals	SSE = $\sum_i \sum_k (x_{ik} - m_i)^2$	df2 = N - n	MSE = SSE / (N - n)		
Total	SST = SSB + SSE	df _{TOT} = N - 1			

Legend: m_i is the sample mean of group i

Note: the ANOVA divides the total variation in the response into parts.

R implementation		
Step	Aim	R function
1	We should fit our data to the ANOVA model	<code>fitModel = lm(Response ~ Treatment, data=dSet)</code>
2	We can get R to produce an ANOVA table	<code>anova(fitModel)</code>

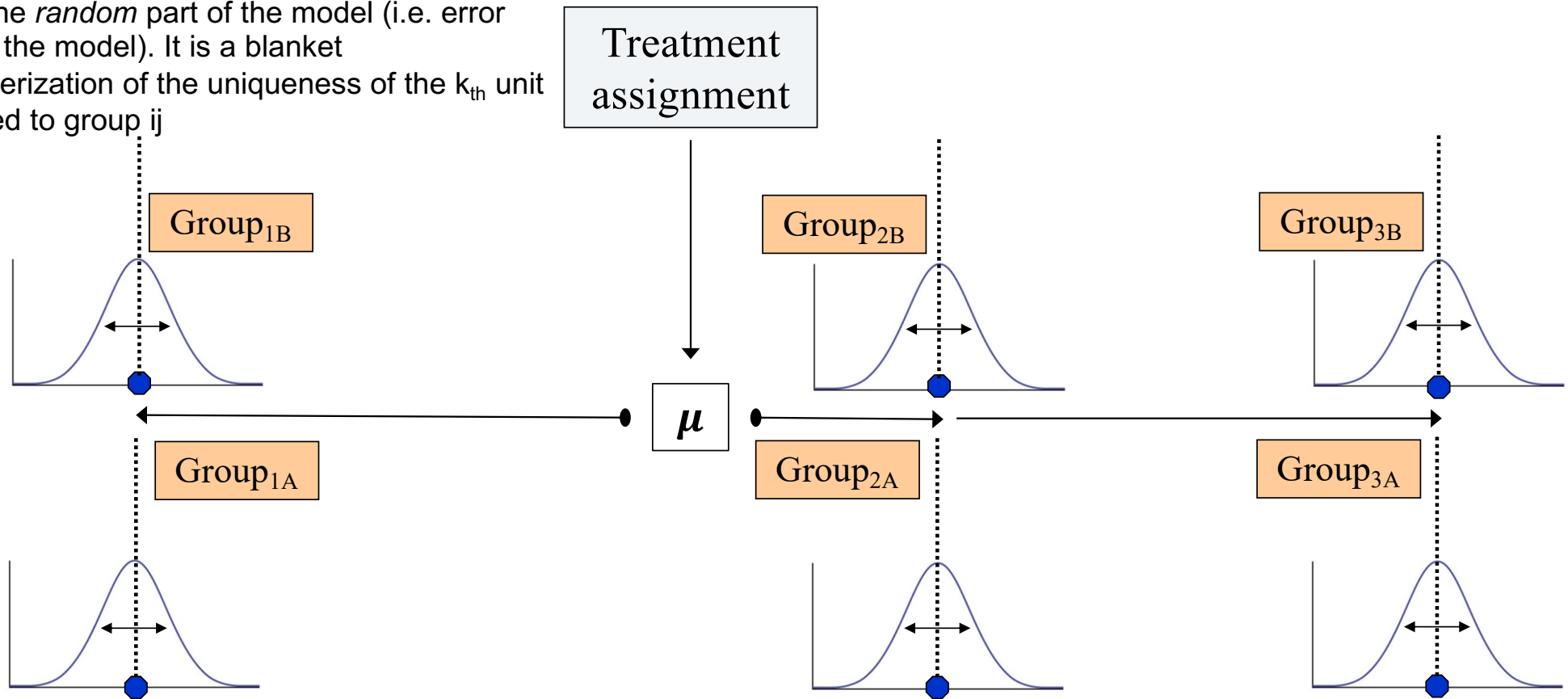
Fisher's two-way ANOVA

6

1. The unit k (e.g. mouse), $k = 1, \dots, u_{ij}$; $N = \sum_{ij} u_{ij}$
2. τ_i : effect of treatment i , $i = 1, \dots, n$; η_j : effect of treatment j , $j = 1, \dots, r$
3. ε_{ijk} : the *random* part of the model (i.e. error term of the model). It is a blanket characterization of the uniqueness of the k_{th} unit assigned to group ij

Equation of the statistical model:

$$Y_{ijk} = \mu + \tau_i + \eta_j + \varepsilon_{ijk}$$



Fisher's two-way ANOVA

7

Hypothesis to test n.1: $\tau_1 = \dots = \tau_n$

Hypothesis to test n.2: $\eta_1 = \dots = \eta_r$

Test statistic:

Source of variation	Sum of Squares	Degrees of freedom	Mean Squares	F _{df1,df2}	P-value
Treatment τ	$SSB_{\tau} = \sum_i u_i (m_i - M)^2$	$df1_{\tau} = n - 1$	$MSB_{\tau} = SSB_{\tau} / (n - 1)$	MSB_{τ} / MSE	0.023
Treatment η	$SSB_{\eta} = \sum_j u_i (m_j - M)^2$	$df1_{\eta} = r - 1$	$MSB_{\eta} = SSB_{\eta} / (r - 1)$	MSB_{η} / MSE	0.150
Residuals	$SSE = \sum_i \sum_k (x_{ijk} - m_{ij})^2$	$df2 = N - (n \cdot r)$	$MSE = SSE / [N - (n \cdot r)]$		
Total	$SST = SSB_{\tau} + SSB_{\eta} + SSE$	$df_{TOT} = N - 1$			

Note: the ANOVA divides the total variation in the response into parts.

R implementation		
Step	Aim	R function
1	We should fit our data to the ANOVA model	<code>fitModel = lm(Response ~ Treat$_{\tau}$ + Treat$_{\eta}$, data=dSet)</code>
2	We can get R to produce an ANOVA table	<code>anova(fitModel)</code>

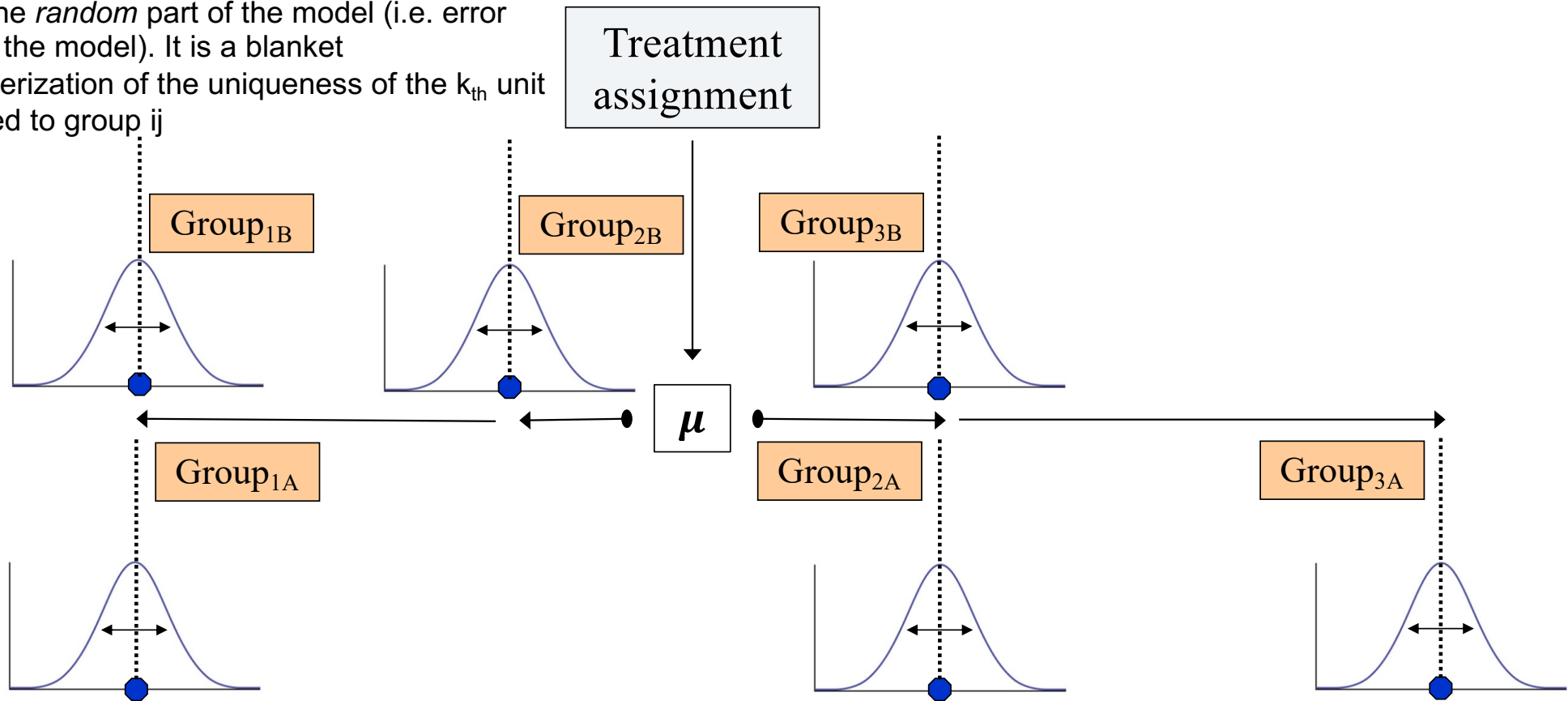
Fisher's two-way ANOVA with interaction

8

1. The unit k (e.g. mouse), $k = 1, \dots, u_{ij}$; $N = \sum_{ij} u_{ij}$
2. τ_i : effect of treatment i , $i = 1, \dots, n$; η_j : effect of treatment j , $j = 1, \dots, r$
3. ε_{ijk} : the *random* part of the model (i.e. error term of the model). It is a blanket characterization of the uniqueness of the k_{th} unit assigned to group ij

Equation of the statistical model:

$$Y_{ijk} = \mu + \tau_i + \eta_j + \tau_i:\eta_j + \varepsilon_{ijk}$$



Fisher's two-way ANOVA with interaction

9

Hypothesis to test n.1: $\tau_1 = \dots = \tau_n$

Hypothesis to test n.2: $\eta_1 = \dots = \eta_r$

Hypothesis to test n.3: $\tau:\eta = 0$

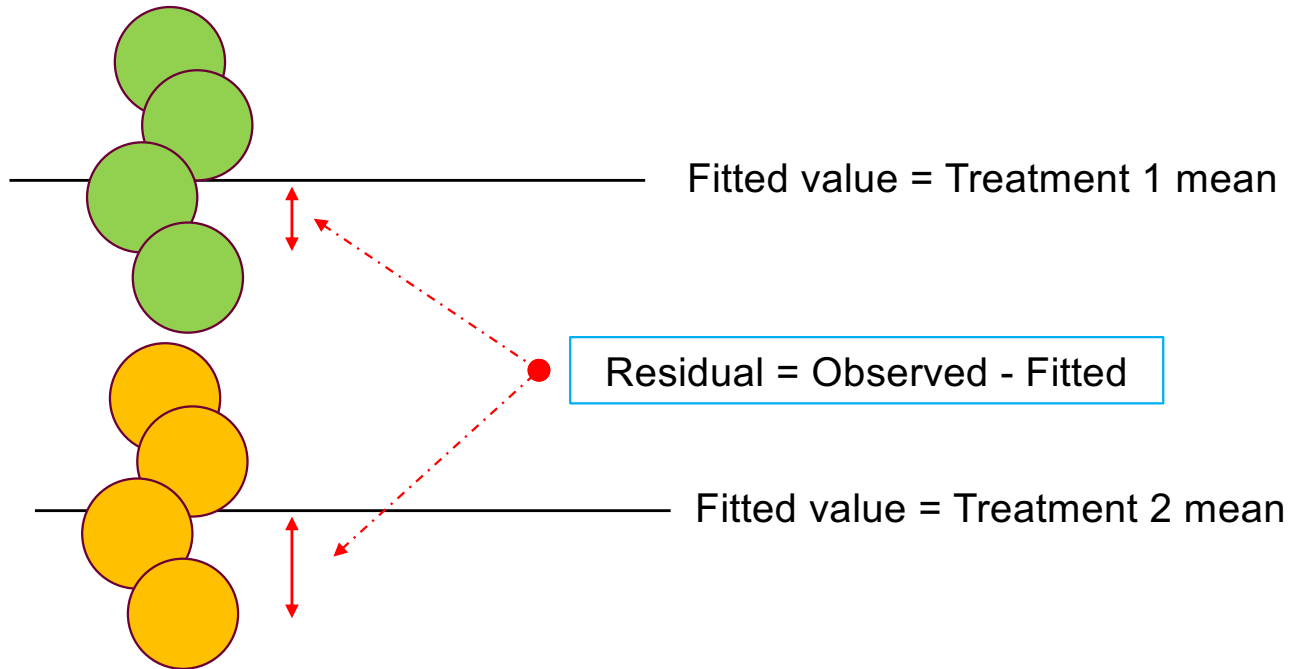
Test statistic:

Source of variation	Sum of Squares	Degrees of freedom	Mean Squares	F _{df1,df2}	P-value
Treatment τ	$SSB_{\tau} = \sum_i u_i (m_i - M)^2$	$df1_{\tau} = n - 1$	$MSB_{\tau} = SSB_{\tau} / df1_{\tau}$	MSB_{τ} / MSE	0.023
Treatment η	$SSB_{\eta} = \sum_j u_i (m_j - M)^2$	$df1_{\eta} = r - 1$	$MSB_{\eta} = SSB_{\eta} / df1_{\eta}$	MSB_{η} / MSE	0.150
Interaction $\tau:\eta$	$SSB_{\tau:\eta} = \sum_{jj} u_{ij} (m_{ij} - m_j - m_i + M)^2$	$df1_{\tau:\eta} = (n - 1) \cdot (r - 1)$	$MSB_{\tau:\eta} = SSB_{\tau:\eta} / df1_{\tau:\eta}$	$MSB_{\tau:\eta} / MSE$	0.401
Residuals	$SSE = \sum_i \sum_k (x_{ijk} - m_{ij})^2$	$df2 = N - (n \cdot r)$	$MSE = SSE / df2$		
Total	$SST = SSB_{\tau} + SSB_{\eta} + SSB_{\tau:\eta} + SSE$	$df_{TOT} = N - 1$			

Note: the ANOVA divides the total variation in the response into parts.

R implementation		
Step	Aim	R function
1	We should fit our data to the ANOVA model	<code>fitModel = lm(Response ~ $Treat_{\tau} * Treat_{\eta}$, data=dSet)</code>
2	We can get R to produce an ANOVA table	<code>anova(fitModel)</code>

Diagnostics: residuals



The residuals are equal to the difference between the observations and the corresponding fitted values.

R implementation

Step	Aim	R function
1	We should fit our data to the ANOVA model	<code>fitModel = lm(Response ~ Predictor, data=dSet)</code>
2	We want to obtain the <i>residuals</i> of the model	<code>dSet\$resid = resid(fitModel)</code>

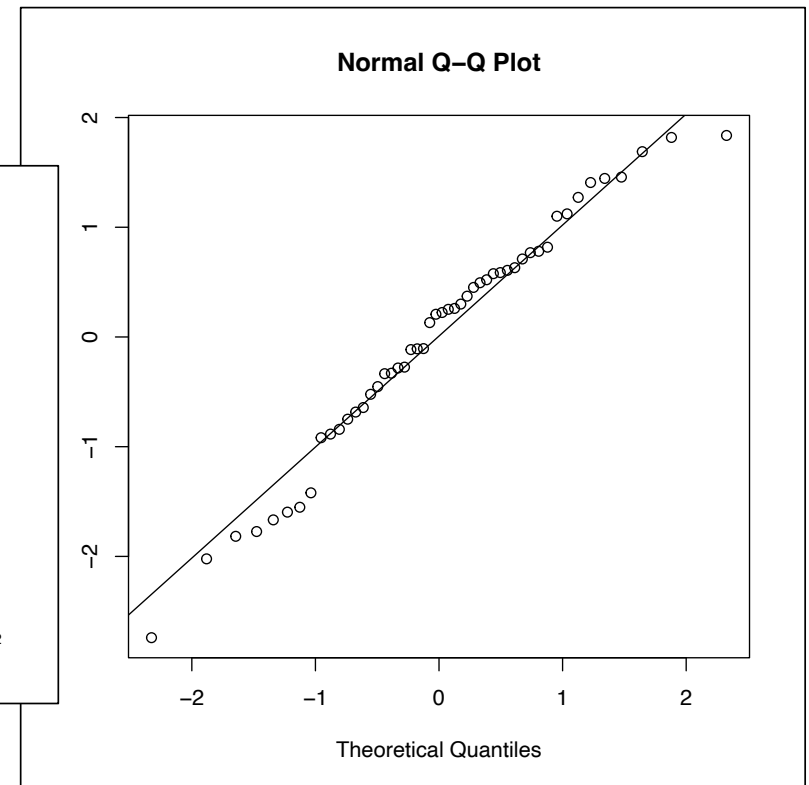
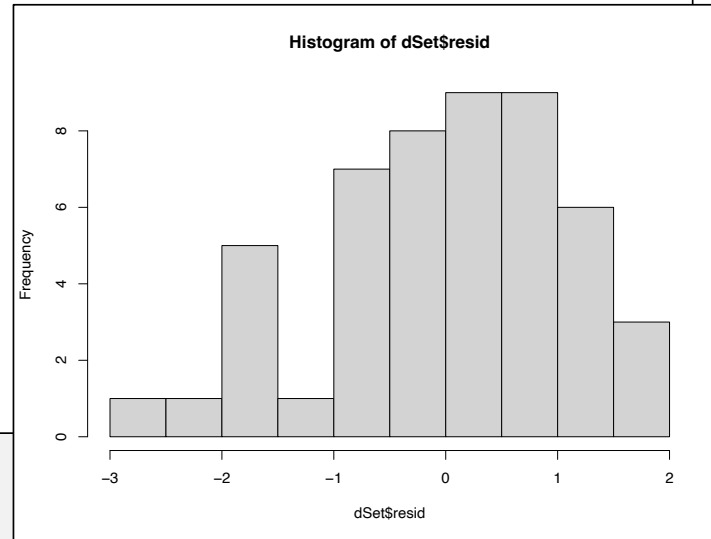
Diagnostics: residuals

Shapiro-Wilk normality test

data: dSet\$resid
W = 0.97324, p-value = 0.3119

Bartlett test of homogeneity of variances

data: resid by Predictor
Bartlett's K-squared = 1.5374, df = 1, p-value = 0.215



R implementation

Step	Aim	Tool	R function
1	We should plot the <i>residuals</i>	Histogram Q-Q plot	hist(dSet\$resid) qqnorm(dSet\$resid); qqline(dSet\$resid)
2	We could test the assumptions	Shapiro-Wilk <i>normality</i> test Bartlett's <i>homoscedasticity</i> test	shapiro.test(dSet\$resid) bartlett.test(resid ~ Predictor, data = dSet)

Diagnostics: residuals

Equation of the statistical model:

$$Y_{ijk} = \mu + \tau_i + \eta_j + \tau_i:\eta_j + \varepsilon_{ijk}$$

Assumptions of normality and homoscedasticity **must be satisfied** by residuals of single treatment group and **combined** treatment groups (e.g. **pooled** residuals of Group_{1A}, Group_{3A} and Group_{3B}).

Source of problems and possible solutions

Solution	Normality	Unequal variance	Outliers
Welch's one-way ANOVA		✓	
Weighting		✓	
Distribution-free methods [⊙]	✓	✓	✓
Data transformation	✓	✓	✓

[⊙] e.g. Kruskal-Wallis test

Welch's one-way ANOVA

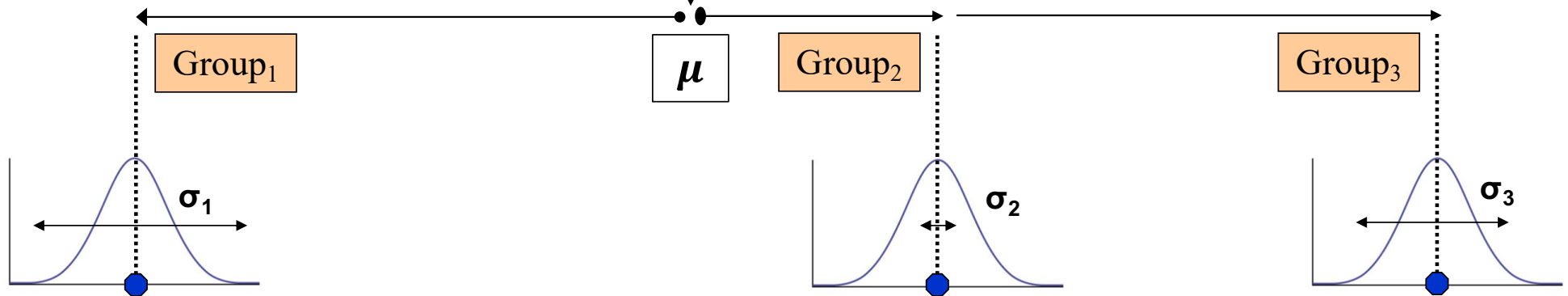
The Welch version of one-way ANOVA do not assume that all the groups are sampled from populations with equal variances.

Hypothesis to test: $\tau_1 = \dots = \tau_n$

Treatment
assignment

Equation of the statistical model:

$$Y_{ik} = \mu + \tau_i + \varepsilon_{ik}$$



Assumptions of ANOVA (ANalysis Of VAriance) models are the following:

- The effect of each factor is additive on μ (i.e. population mean) parameter
- ε_{ik} is assumed to be independent of one another and normally distributed with mean = 0. **Standard deviation could be different between groups: $\sigma_i \neq \sigma_j, i \neq j$**

Weighted least square

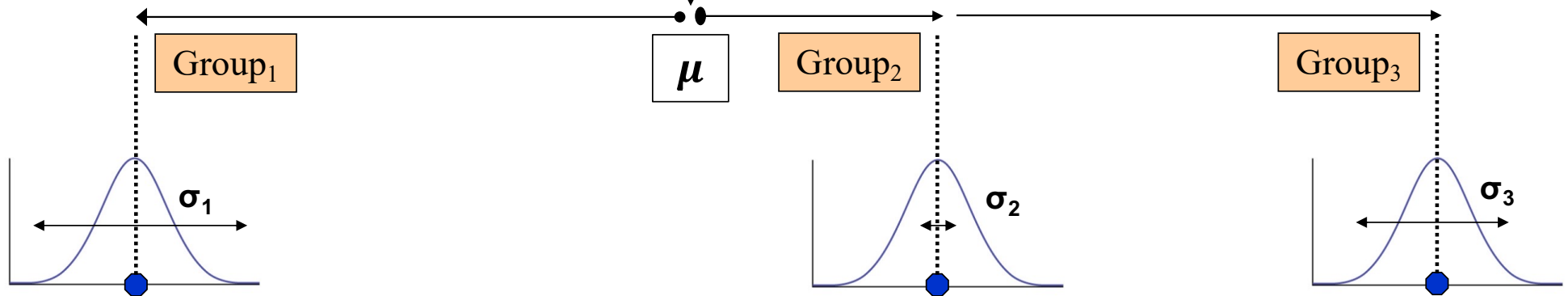
The *gls* function of the R package *nlme* using generalized least squares. The errors are allowed to be correlated and/or have unequal variances.

Hypothesis to test: $\tau_1 = \dots = \tau_n$

Treatment
assignment

Equation of the statistical model:

$$Y_{ik} = \mu + \tau_i + \varepsilon_{ik}$$

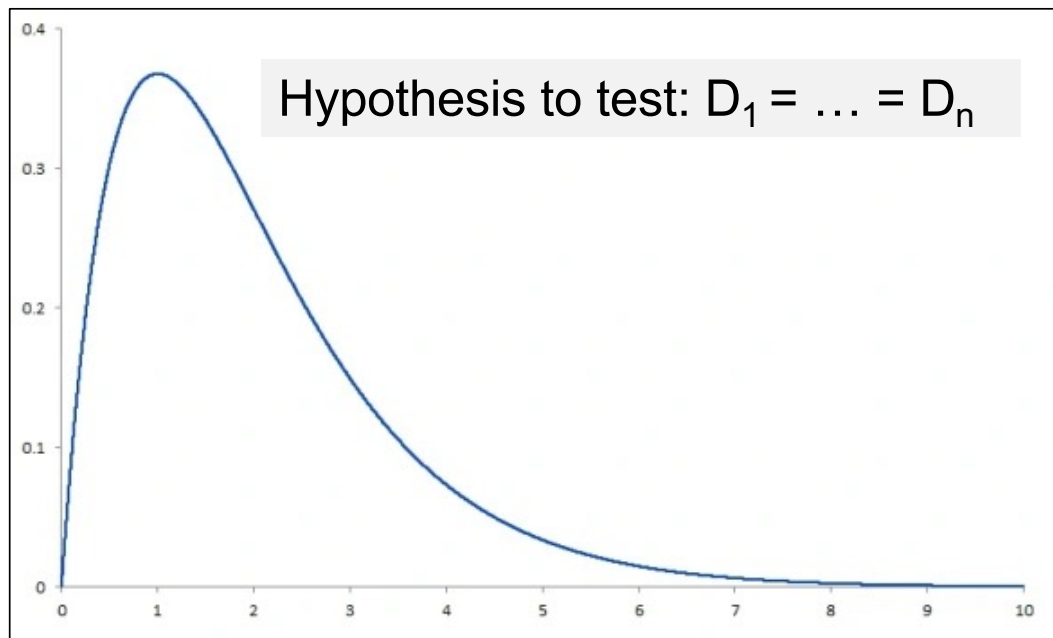


Assumptions of linear models fitted with the *gls* function are the following:

- The effect of each factor is additive on μ (i.e. population mean) parameter
- ε_{ik} could be **correlated**. They are normally distributed with mean = 0. **Standard deviation could be different between groups: $\sigma_i \neq \sigma_j, i \neq j$**

Kruskal-Wallis test

The Kruskal-Wallis test (i.e. one-way ANOVA on ranks) works on ranks. It tests whether samples originate from the same distribution.



10.2	24.7	33.2	..	99.7	99.9
------	------	------	----	------	------



Replacement of data
by their ranks

1	2	3	..	N-1	N
---	---	---	----	-----	---

Assumptions of Kruskal-Wallis test are the following:

- We only assume that the observations in the data set are independent of each other.

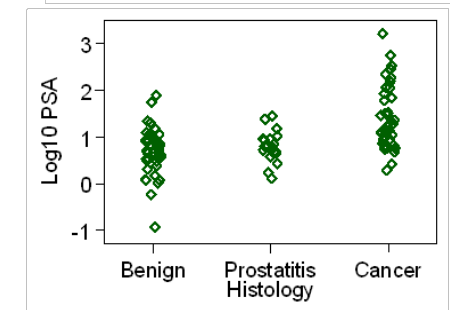
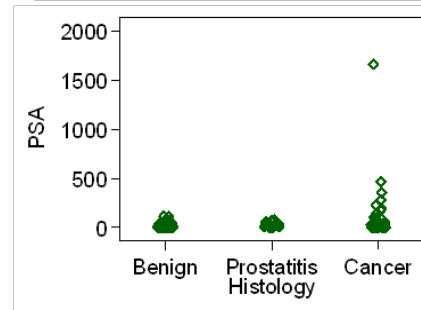
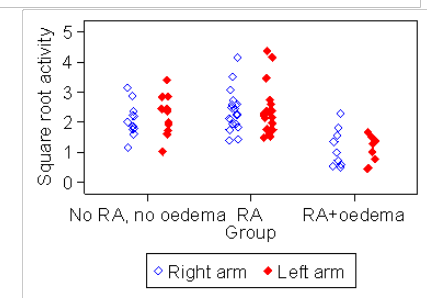
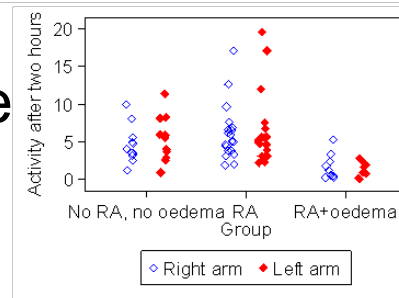
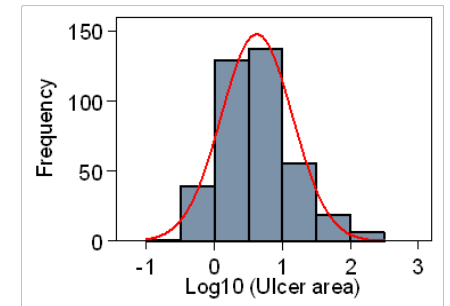
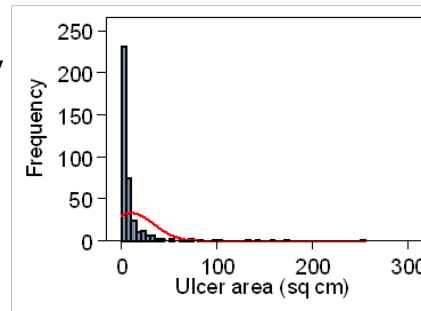
R functions

R implementation																
Test	R															
Welch's one-way ANOVA	Function	<ul style="list-style-type: none"><code>oneway.test(<i>Response</i> ~ <i>Predictor</i>, data = <i>dSet</i>, var.equal = FALSE)</code>														
	Output	<div>One-way analysis of means (not assuming equal variances)</div> <div>data: Response and Predictor</div> <div>F = 118.34, num df = 1.000, denom df = 45.143, p-value = 3.342e-14</div>														
Weighted least square	Function	<ul style="list-style-type: none"><code>fitModel <- gls(<i>Response</i> ~ <i>Predictor</i>, weights = varIdent(form= ~ 1 <i>Predictor</i>), data = <i>dSet</i>)</code><code>summary(fitModel)</code>														
	Output	<div>Variance function: Structure: Different standard deviations per stratum Formula: ~1 Predictor Parameter estimates: 1 2 1.000000 1.293192</div> <div>Coefficients:</div> <table><tr><td></td><td>Value</td><td>Std.Error</td><td>t-value</td><td>p-value</td></tr><tr><td>(Intercept)</td><td>-0.001177</td><td>0.1890228</td><td>-0.006228</td><td>0.9951</td></tr><tr><td>Predictor</td><td>3.361487</td><td>0.3090014</td><td>10.878548</td><td>0.0000</td></tr></table>		Value	Std.Error	t-value	p-value	(Intercept)	-0.001177	0.1890228	-0.006228	0.9951	Predictor	3.361487	0.3090014	10.878548
	Value	Std.Error	t-value	p-value												
(Intercept)	-0.001177	0.1890228	-0.006228	0.9951												
Predictor	3.361487	0.3090014	10.878548	0.0000												
Kruskal-Wallis	Function	<ul style="list-style-type: none"><code>kruskal.test(<i>Response</i> ~ <i>Predictor</i>, data = <i>dSet</i>)</code>														
	Output	<div>Kruskal-Wallis rank sum test</div> <div>data: Response by Predictor</div> <div>Kruskal-Wallis chi-squared = 34.222, df = 1, p-value = 4.917e-09</div>														

Data transformation

We can transform the data mathematically...

- to make them fit the normality more closely
- to obtain more similar variance
- to handle outliers



Data transformation

Common and useful transformations of the response variable:

1. the logarithm ($x_i > 0$, $i=1, \dots, n$)
2. the square root ($x_i \geq 0$, $i=1, \dots, n$)
3. the square power ($x_i \geq 0$, $i=1, \dots, n$)
4. the ranks (e.g. Welch's one-way ANOVA on ranks)

<http://bioinformatics-core-shared-training.github.io/IntroductionToStats/practical.html>



Hands on