



CANCER
RESEARCH
UK

Cambridge
Institute

Together we are
beating cancer

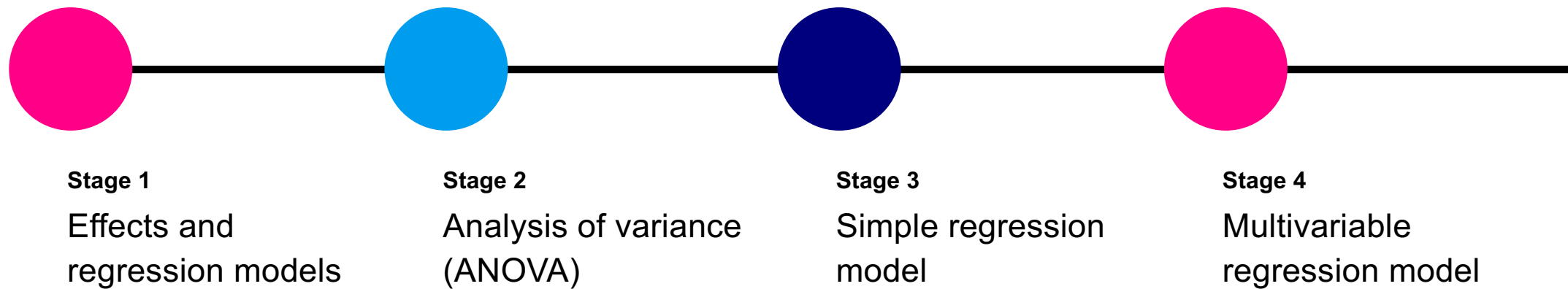
Luca Porcu & Chandra Chilamakuri (Bioinformatics core)

21st February 2025

Linear regression models

Fixed-effects models

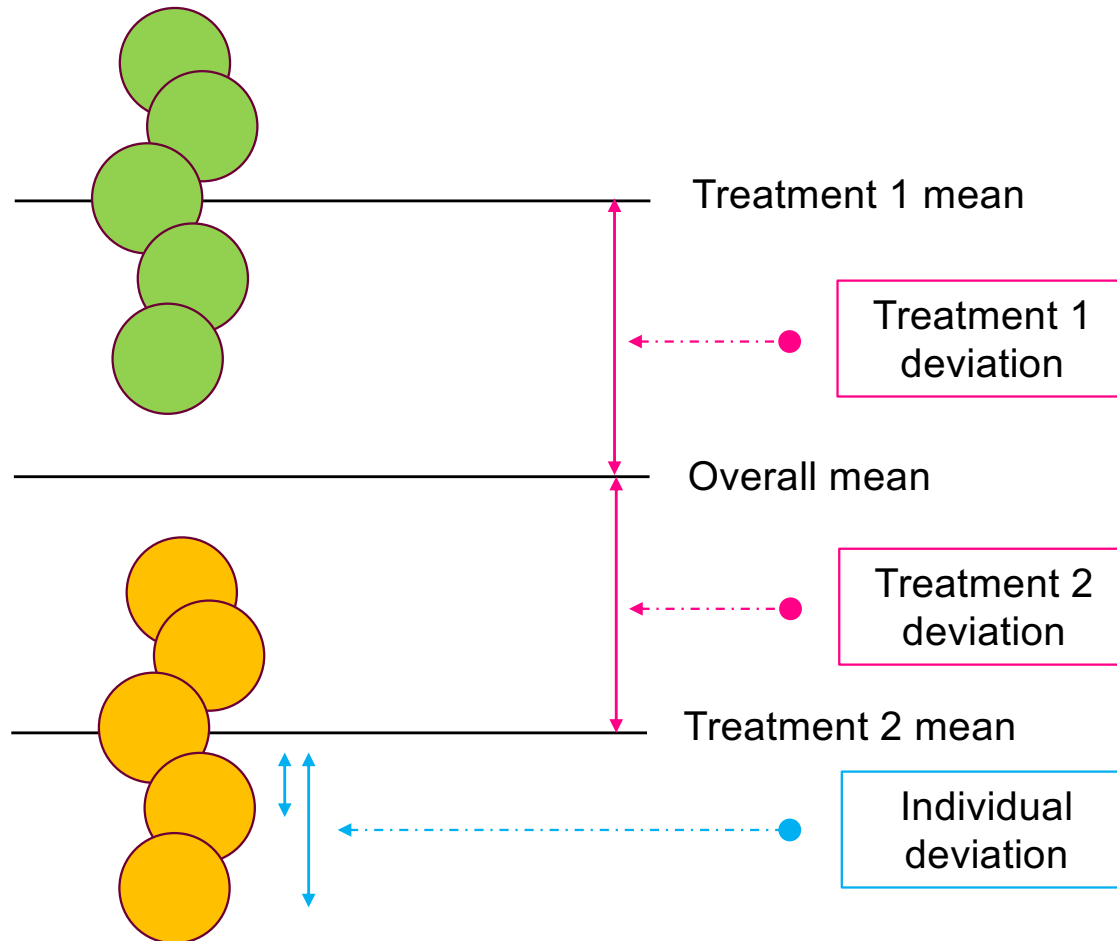
Process flow





CANCER
RESEARCH
UK

Cambridge
Institute



Analysis of variance (ANOVA)

Definition and classification

10.00 -10.20 am

Together we are
beating cancer

Fisher's one-way ANOVA

4

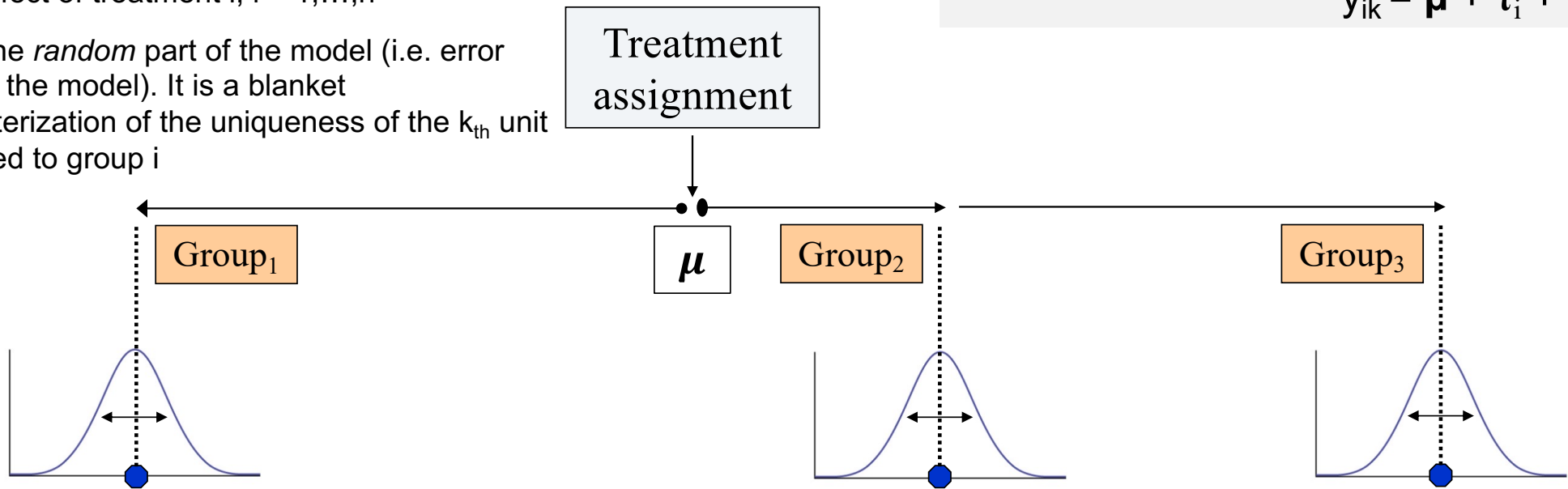
1. The unit k (e.g. mouse), $k = 1, \dots, u_i$; $N = \sum_i u_i$

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 treatment level i 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: $\tau_1 = \dots = \tau_n = 0$

Test statistic:

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

Legend: m_i is the sample mean of group i . M is the overall mean response

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>

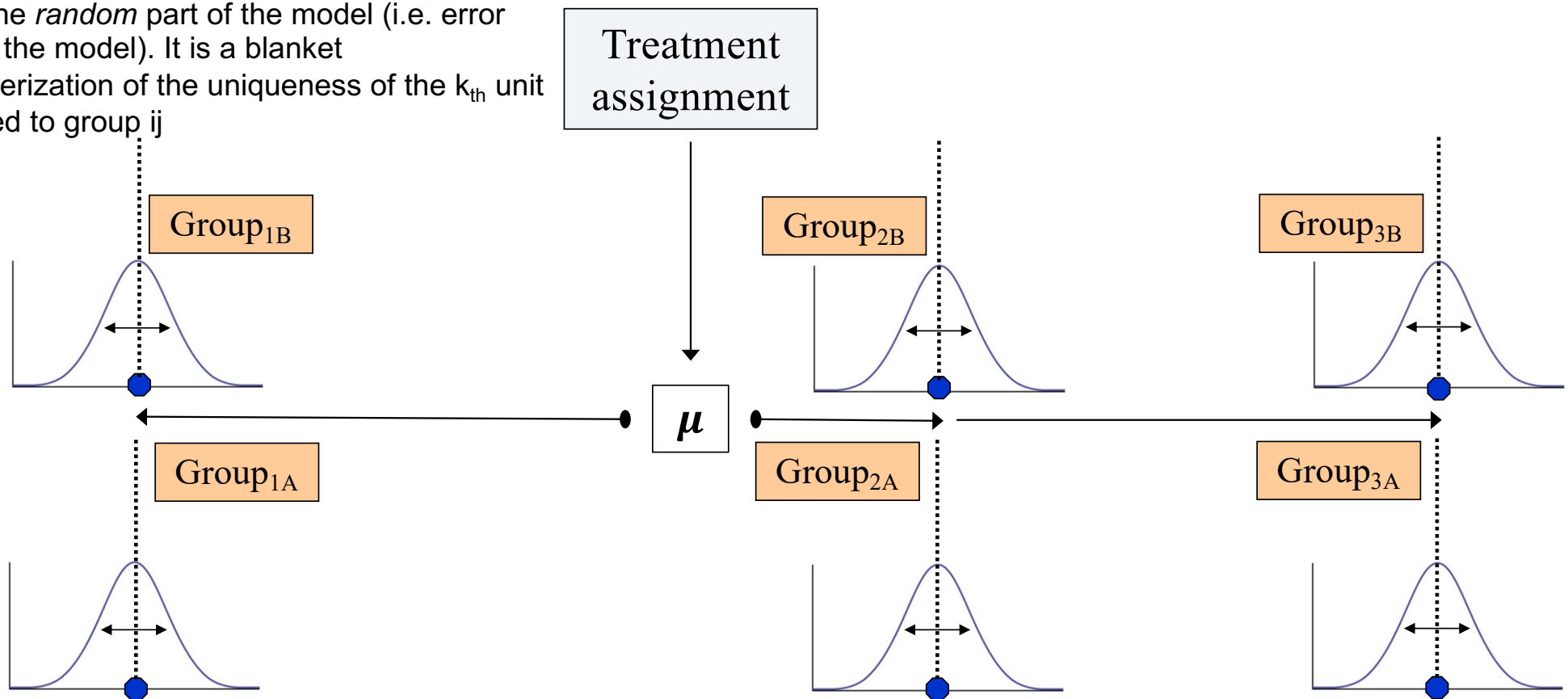
Fisher's two-way ANOVA

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: $\tau_1 = \dots = \tau_n = 0$

Hypothesis to test n.2: $\eta_1 = \dots = \eta_r = 0$

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
Residuals	$SSE = \sum_{ij} \sum_k (y_{ijk} - m_{ij})^2$	$df_2 = N - (n \cdot r)$	$MSE = SSE / df_2$		
Total	$SST = SSB_{\tau} + SSB_{\eta} + SSE$				

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>

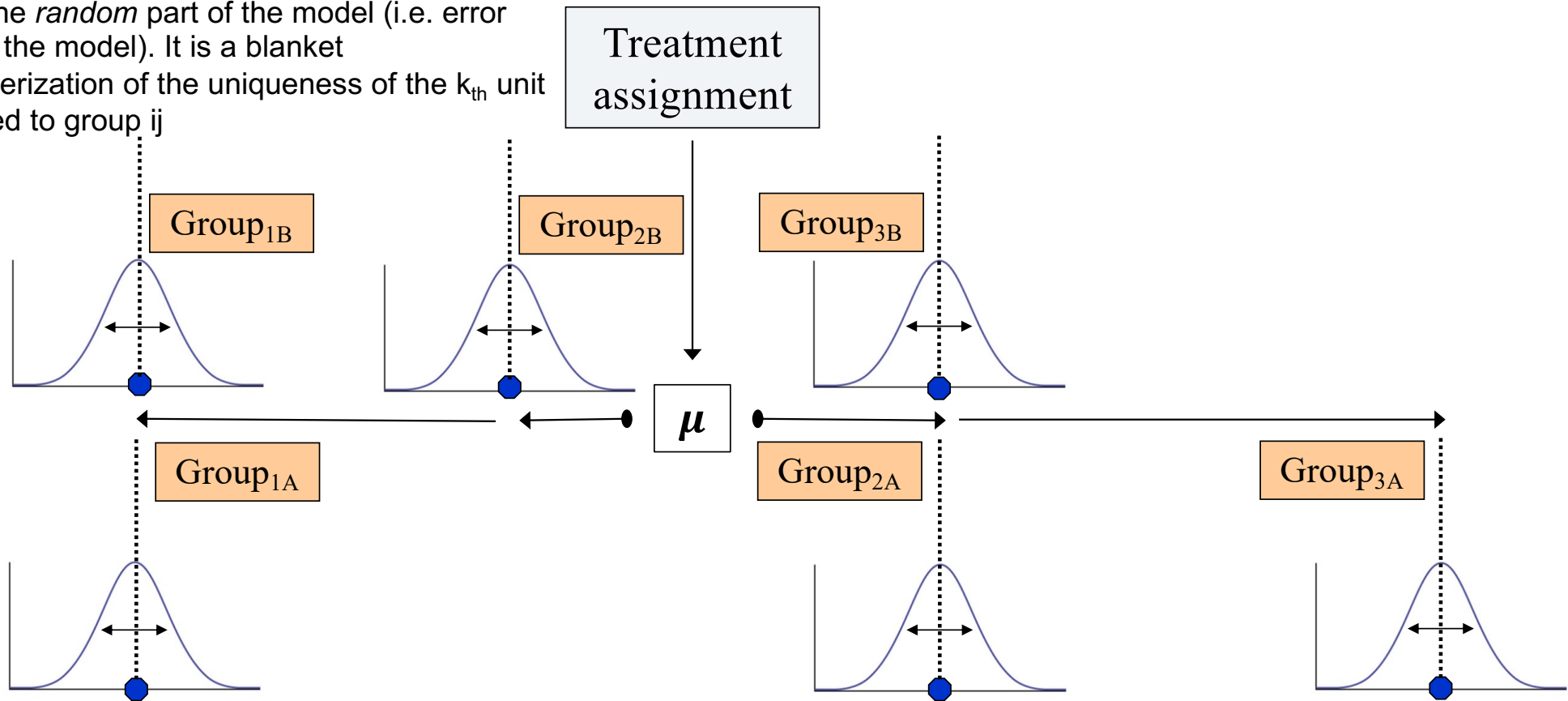
Fisher's two-way ANOVA with interaction

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: $\tau_1 = \dots = \tau_n = 0$

Hypothesis to test n.2: $\eta_1 = \dots = \eta_r = 0$

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

Test statistic:

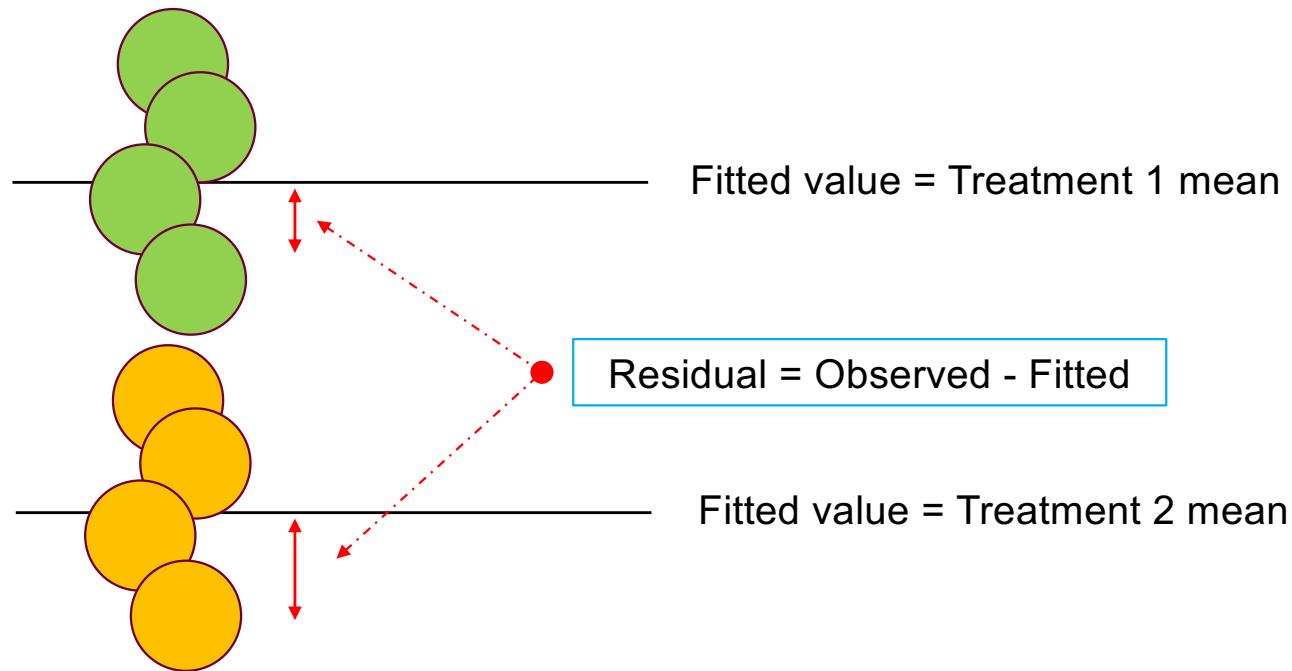
Source of variation	Sum of Squares	Degrees of freedom	Mean Squares	F _{df1,df2}	p-value
Treatment τ	$SSB_{\tau} = \sum_i u_i (m_i - M)^2$	$df1_{\tau} = n - 1$	$MSB_{\tau} = SSB_{\tau} / df1_{\tau}$	MSB_{τ} / MSE	0.023
Treatment η	$SSB_{\eta} = \sum_j u_i (m_j - M)^2$	$df1_{\eta} = r - 1$	$MSB_{\eta} = SSB_{\eta} / df1_{\eta}$	MSB_{η} / MSE	0.150
Interaction $\tau:\eta$	$SSB_{\tau:\eta} = \sum_{ij} u_{ij} (m_{ij} - m_j - m_i + M)^2$	$df1_{\tau:\eta} = (n - 1) \cdot (r - 1)$	$MSB_{\tau:\eta} = SSB_{\tau:\eta} / df1_{\tau:\eta}$	$MSB_{\tau:\eta} / MSE$	0.401
Residuals	$SSE = \sum_{ij} \sum_k (y_{ijk} - m_{ij})^2$	$df_2 = N - (n \cdot r)$	$MSE = SSE / df_2$		
Total	$SST = SSB_{\tau} + SSB_{\eta} + SSB_{\tau:\eta} + SSE$				

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 ~ <i>Treat_τ</i> * <i>Treat_η</i>, data=<i>dSet</i>)</code>
2	We can get R to produce an ANOVA table	<code>anova(fitModel)</code>

Diagnostics: residuals

10



The (raw) 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>fittedModel = lm(<i>Response</i> ~ <i>Predictor</i>, data=<i>dSet</i>)</code>
2	We want to obtain the <i>residuals</i> of the model	<code><i>dSet</i>\$resid = resid(fittedModel)</code>

Diagnostics: residuals

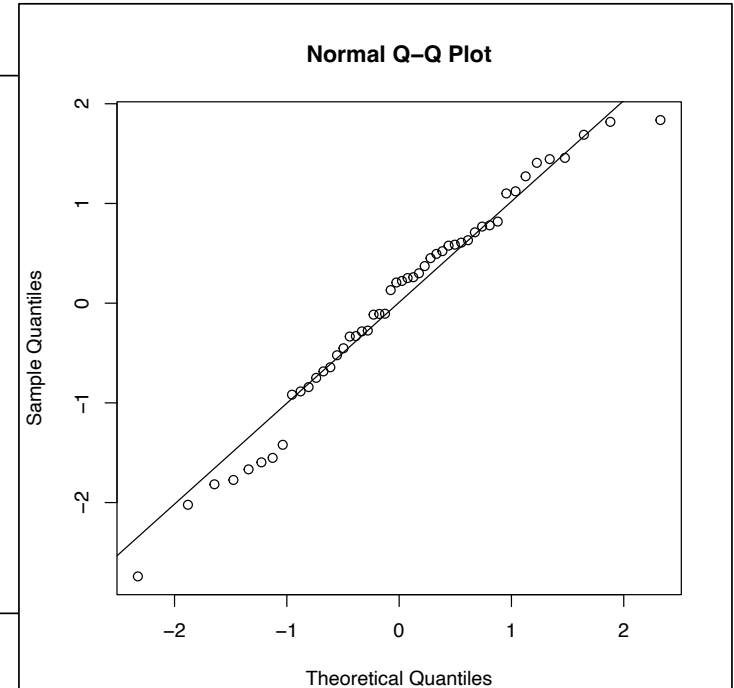
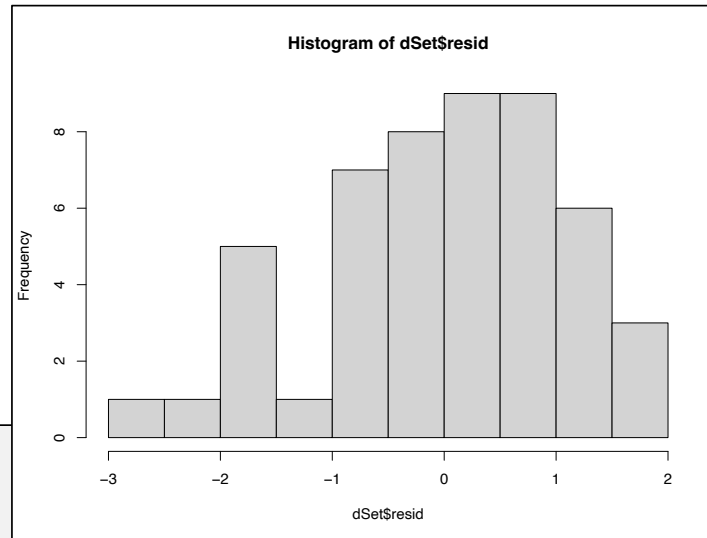
11

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	<code>hist(dSet\$resid)</code> <code>qqnorm(dSet\$resid); qqline(dSet\$resid)</code>
2	We could test the assumptions	Shapiro-Wilk <i>normality</i> test Bartlett's <i>homoscedasticity</i> test	<code>shapiro.test(dSet\$resid)</code> <code>bartlett.test(resid ~ Predictor, data = dSet)</code>

Diagnostics: residuals

12

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}). Pooled residuals should be examined in the diagnostic analysis.

Source of problems and possible solutions ¹³

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

[⊙] e.g. Kruskal-Wallis test (i.e. one-way ANOVA on ranks)

Welch's one-way ANOVA

14

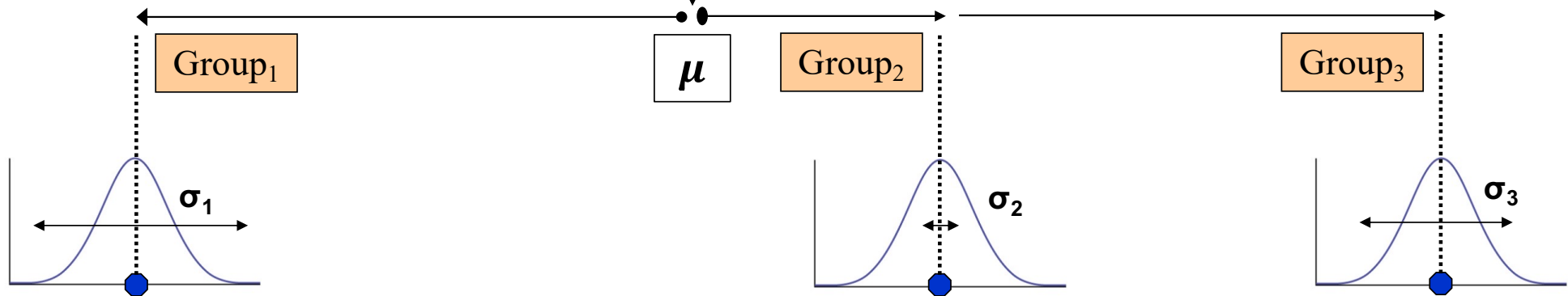
The Welch version of one-way ANOVA do not assume that all the groups are sampled from populations with equal variances.

Hypothesis to test: $\tau_1 = \dots = \tau_n = 0$

Treatment
assignment

Equation of the statistical model:

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



Assumptions of Welch's one-way ANOVA 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 squares

15

```
> library(nlme) # Load the nlme package
> fittedModel = gls(Response ~ Predictor, weight = varIdent(form= ~1 | Predictor) , data = dSet)
# Modeling variance as function of the predictor
> summary(fittedModel) # Output
```

Output

Variance function:

Structure: Different standard deviations per stratum

Formula: ~1 | gender

Parameter estimates:

	F	M
	1.0000000	0.7786179

Coefficients:

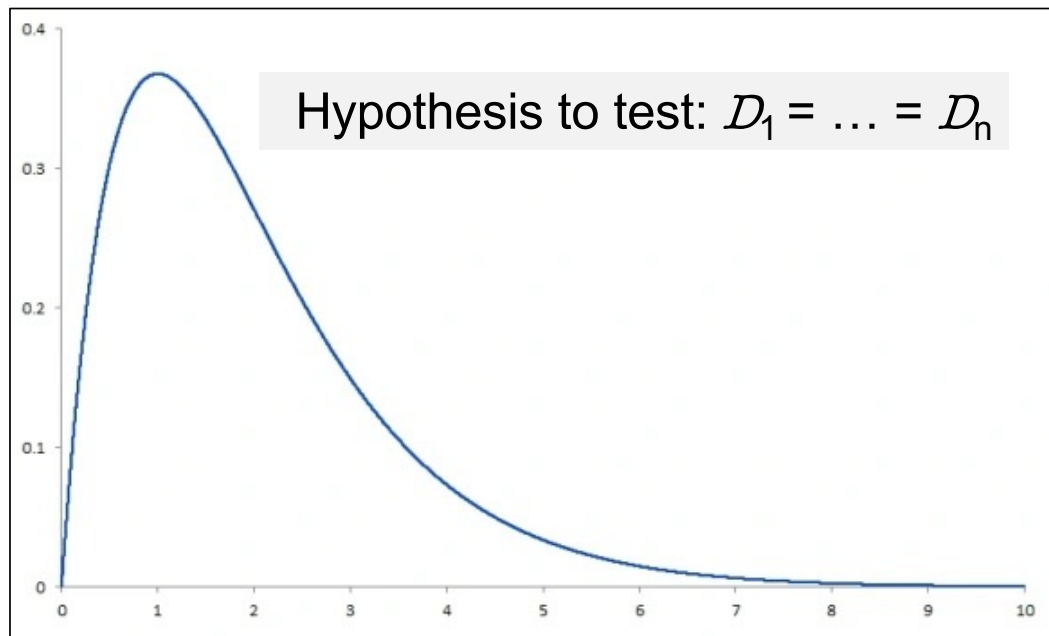
	Value	Std.Error	t-value	p-value
(Intercept)	147.4729	3.725061	39.58938	0.0000
gender M	4.9198	5.097216	0.96519	0.3368

Note: `gls()` function fits a linear model using generalized least squares.

Kruskal-Wallis test

16

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	..	96.4	99.9
------	------	------	----	------	------

↓ Replacement of data by 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 one another.

R functions

17

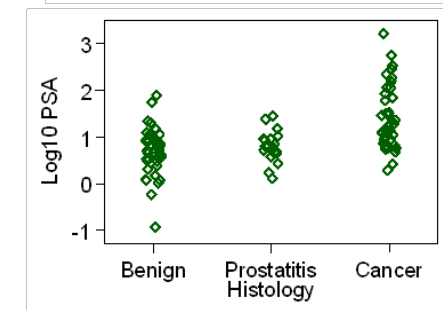
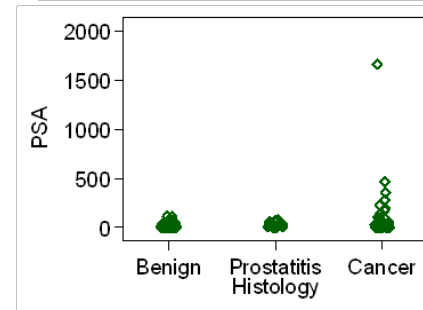
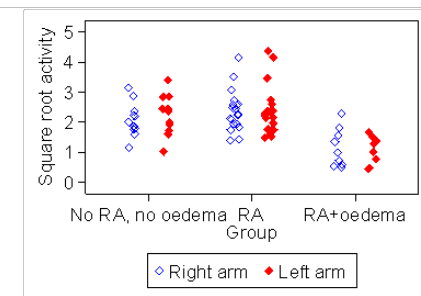
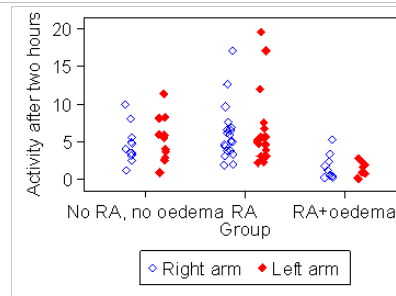
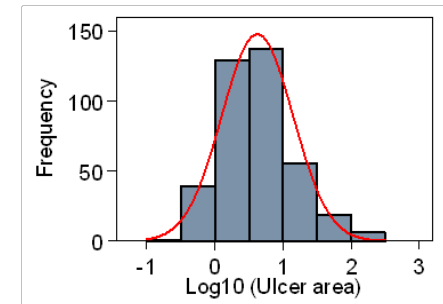
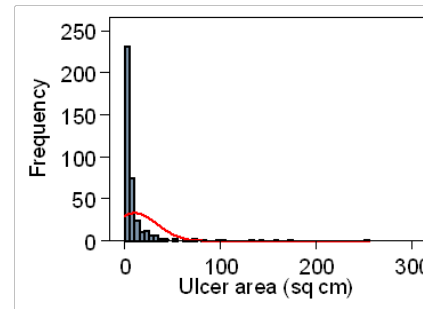
R implementation																
Test	R															
Welch's one-way ANOVA	Function	<ul style="list-style-type: none"><code>oneway.test(Response ~ Predictor, data = dSet, var.equal = FALSE)</code>														
	Output	<p>One-way analysis of means (not assuming equal variances)</p> <p>data: Response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

Data transformation

18

We can transform the data mathematically...

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



Data transformation

19

Common and useful transformations of the response variable:

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

[https://bioinformatics-core-shared-training.github.io/
Fixed-and-Mixed-effects-models/](https://bioinformatics-core-shared-training.github.io/Fixed-and-Mixed-effects-models/)



Hands on