



Faculty of Health Sciences



Day 4: Analysis of Variance

Paul Blanche

Section of Biostatistics, University of Copenhagen

November 4, 2020



Outline

Simple pairwise comparisons

Analysis of Variance (ANOVA): one-way

Analysis of Variance (ANOVA): two-way



Case: Irritable Bowel Syndrome Dose Response

- ▶ Data from $n = 198$ women.
- ▶ Randomized (double-blind) to:
 - ▶ Placebo (Dose 0, $n = 50$)
 - ▶ Dose 1 ($n = 54$)
 - ▶ Dose 2 ($n = 49$)
 - ▶ Dose 3 ($n = 45$)

(Doses are blinded for confidentiality)

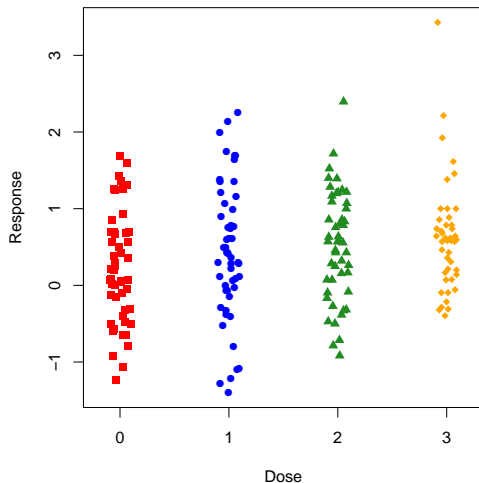


Outcome: baseline adjusted abdominal pain score at end of follow-up (12 weeks), approximately continuous variable, the larger the better.

Research questions:

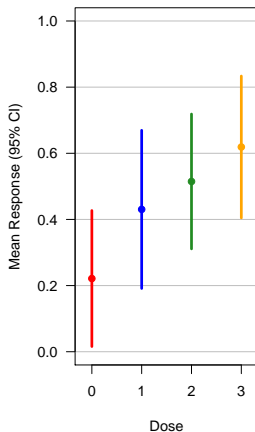
- ▶ Does the drug work?
- ▶ Are there differences between doses?

Outcome data



Dose	Mean	SD	n
0	0.22	0.72	50
1	0.43	0.88	54
2	0.51	0.71	49
3	0.62	0.71	45

Pairwise Welch's t-test: results



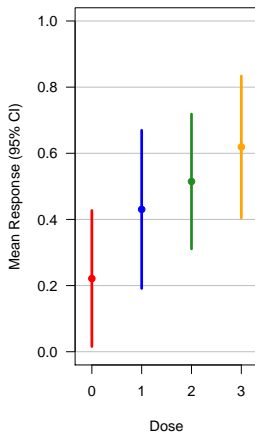
p-values from pairwise t-tests:

Dose	0	1	2
1	0.19		
2	0.04	0.59	
3	0.01	0.24	0.48

Note: the y-axis has changed!



Pairwise Welch's t-test: results



p-values from pairwise t-tests:

Dose	0	1	2
1	0.19		
2	0.04	0.59	
3	0.01	0.24	0.48

Have we not reported all relevant results? What is left to worry about?

Note: the y-axis has changed!



Interpretations (1/3)

Results include:

- ▶ Dose 0 not significantly different from dose 1.
- ▶ Dose 1 not significantly different from dose 2.
- ▶ But dose 0 significantly different from dose 1.

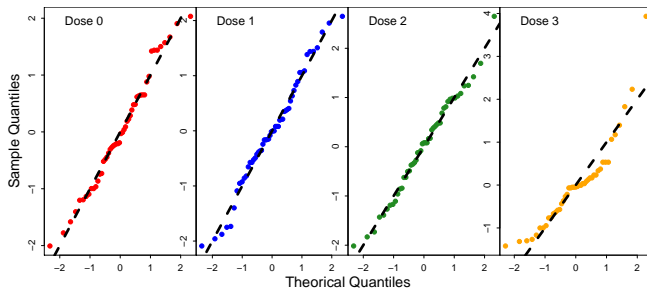
Are the results self-contradicting?

- ▶ No, this is just due to statistical uncertainty because of “small” sample sizes.
- ▶ Dose 1 may have a similar effect to either dose 0 or dose 2.



Interpretations (2/3)

What about the **assumptions**? Can we trust the results?



- ▶ QQplots for doses 0, 1, 2 look good but not so good for dose 3.
- ▶ However nothing “very” bad and decent sample size (≥ 45 per group), so **it seems fine**.

Interpretations (3/3)

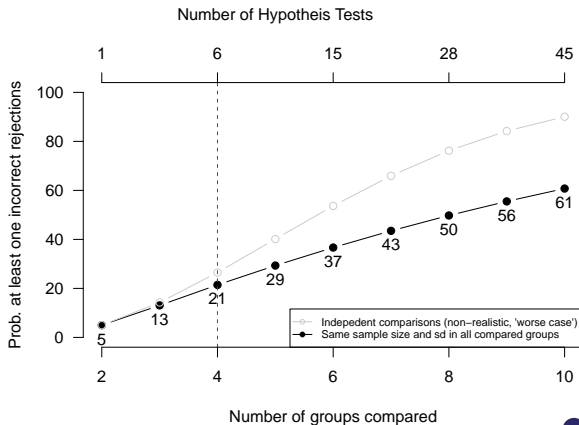
Beware of the **multiple testing issue!** We have computed 6 p-values, hence the risk of making at least one false “discovery” is $\geq 5\%$.



Interpretations (3/3)

Beware of the **multiple testing issue!** We have computed 6 p-values, hence the risk of making at least one false “discovery” is $\geq 5\%$.

How bad
can this be?



What statistical method should we use?

We want to control the FWER¹ at 5%.

Using Bonferroni? No.

It's a conservative (i.e. sub-optimal) approach which ignores the (strong) **correlation** between the comparisons.

¹Family-wise error rate (FWER): probability of making one or more false discoveries when performing multiple hypotheses tests (Lecture 2).

²That is why the method are still “new” and underused.



What statistical method should we use?

We want to control the FWER¹ at 5%.

Using Bonferroni? No.

It's a conservative (i.e. sub-optimal) approach which ignores the (strong) **correlation** between the comparisons.

More modern alternative? Yes.

Use specific method and software for multiple correction that do not make any additional assumptions. The details of the method and computation are more complicated² but not the **interpretation** and **user-friendly software** exist.

¹Family-wise error rate (FWER): probability of making one or more false discoveries when performing multiple hypotheses tests (Lecture 2).

²That is why the method are still “new” and underused.



Recommended analysis (see R-demo for code)

Statistical methods:

Comparisons between groups were made with a heteroscedastic ANOVA model (not assuming equal variances). P-values and 95% confidence intervals were adjusted for multiple testing using the min-P method as implemented in the multcomp-package [ref.³] of the statistical software R [ref.⁴] and described in [ref.⁵].

Results (adjusted for multiple testing):

Comparison	Est. Diff	95% CI	p-value
1 - 0	0.209	[-0.202; 0.620]	0.552
2 - 0	0.293	[-0.083; 0.670]	0.185
3 - 0	0.398	[0.011; 0.784]	0.041
2 - 1	0.084	[-0.325; 0.494]	0.951
3 - 1	0.189	[-0.230; 0.607]	0.647
3 - 2	0.104	[-0.281; 0.489]	0.896

Note: p-values ≤ 6 times the non-adjusted ones (Bonferroni).

³ Hothorn, Bretz & Westfall (2008). Simultaneous Inference in General Parametric Models. Biometrical Journal 50(3), 346–363.

⁴ R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

⁵ Herberich, Sikorski & Hothorn. "A robust procedure for comparing multiple means under heteroscedasticity in unbalanced designs." PLoS one 5.3 (2010): e9788.



What if we want only the comparisons to placebo?

Sometimes we only want all the comparisons to one (reference) group.

This is known as the **many-to-one** comparisons case ([Dunnett](#)), by contrast to the **all pairwise** comparisons case ([Tukey](#)).

The method for this case is similar and we can use the **same software**.

Case results (adjusted for multiple comparisons to placebo):

Comparison	Est. Diff	95% CI	p-value
1 - 0	0.209	[-0.168; 0.586]	0.418
2 - 0	0.293	[-0.052; 0.639]	0.115
3 - 0	0.398	[0.043; 0.752]	0.023

Note: p-values ≤ 3 times the non-adjusted ones (Bonferroni).



Pre-specification matters

The same comparison, e.g. Dose 3 versus Placebo (Dose 0), leads to the estimated mean difference 0.398, but different 95% confidence intervals (CI) and p-values (after adjusting for multiple testing) when we consider either:

- ▶ All-pairwise (6) comparisons: 95% CI=[0.011; 0.784], $p=0.041$.
- ▶ Many-to-one (3) comparisons: 95% CI=[0.043; 0.752], $p=0.023$.

Take home messages:

- ▶ The more comparisons the wider the 95% CI and the higher the p-values.
- ▶ Do not investigate more comparisons than interesting/possible (power).
- ▶ **The choice of investigating All-pairwise versus Many-to-one comparisons should be done before seeing the data, i.e. pre-specified.** Rigorously adjusting for multiple testing is not possible otherwise and how much we can trust the results without pre-specification is most unclear.



Digression: prespecified vs post hoc analyses

It is completely fine and often **useful** to performed **post hoc**⁶ analyses as long as:

- ▶ they are **reported as such in publications**⁷,
- ▶ **conclusions** based on them are **not too strong**.

“The main analyses should concentrate on the primary research questions to reduce the amount of testing of data-generated hypotheses. However, science would not proceed if analyses of questions not stated in the protocol were not allowed so, obviously, new ideas generated from the data can be pursued as long as conclusions based on such additional analyses are suitably calibrated.”⁸

⁶A post hoc analysis is an analysis specified after the data were seen.

⁷Otherwise this is “data fishing”, “data snooping” or “p-hacking” and this is considered as something in between “questionable research practice” and “scientific dishonesty and research misconduct”; see KU course “Responsible Conduct of Research”.

⁸Andersen & Skovgaard, *Regression with linear predictors*, page 473 (Springer, 2010).



Power and sample size calculation

When planning several comparisons, say K , with a FWER control at α , one can:

1. Define an adjusted type-I error $\alpha' = \alpha/K$.
2. Perform sample size and power calculation for each comparison as in the case of a unique comparison, using this adjusted type-I error α' as input of the formula instead of α .

Note: this is a “slightly” conservative approach⁹.



Outline

Simple pairwise comparisons

Analysis of Variance (ANOVA): one-way

Analysis of Variance (ANOVA): two-way



What is ANOVA about?

ANOVA stands for “**AN**alysis **Of** **VA**riance”, but this is a method to **compare means** (via the comparisons of **variances**).

Useful for answering **research questions** such as:

- ▶ Is this continuous outcome **associated** with this categorical variable?
- ▶ Is the **mean outcome** the same for all levels of this categorical outcome?

Examples:

- ▶ **Outcome:** weight, blood pressure, concentration, **pain score** ...
- ▶ **Categorical variable:** BMI group, age group, **dose level** ...

This is a very **commonly used**, well-known and “old” method.



ANOVA model (one-way)

The j -th observation from the i -th group is described as:

$$Y_{ij} = \mu_i + \varepsilon_{ij}$$

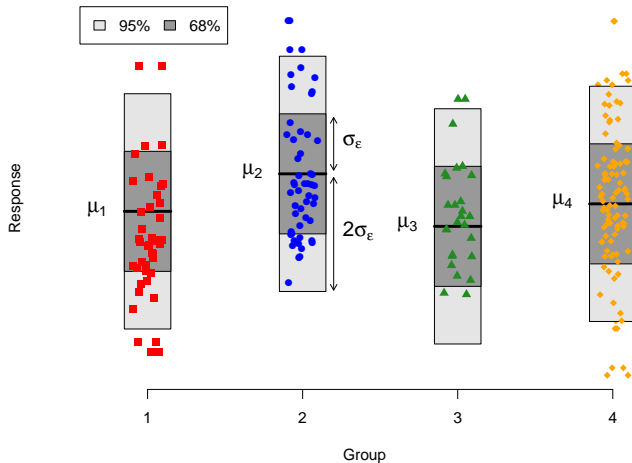
- ▶ μ_i is the (population) mean for group i .
- ▶ ε_{ij} 's are individual 'error' terms ("random/unexplained deviation from the mean") assumed normally distributed with zero mean and the same variance σ_ε^2 regardless of group.

Model assumptions (1-2 important, 3-4 not always):

1. Observations from different groups are independent.
2. Individual observations within each group are independent.
3. 'Error' terms are normally distributed.
4. The variance of 'error' terms is the same for all groups (homogeneity)



Visual interpretation ($Y_{ij} = \mu_i + \varepsilon_{ij}$, hypothetical data)



- σ_ε tells us how **vertically spread** are the points above and below each group mean μ_i , for each group.

ANOVA (one-way): why and how?

Why using a (traditional) ANOVA?

To test the **global null hypothesis** “ H_0 : The mean of all (K) groups are all equal”, that is

$$H_0 : \mu_1 = \mu_2 = \cdots = \mu_K.$$

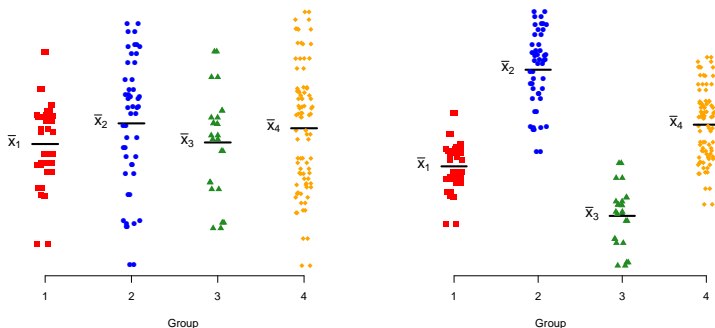
How does it work?

By using a **F-test** which compares the **between-group** variability to the **within-group** variability. If the between-group variability is large enough relative to the within-group variability, then we reject H_0 .

- ▶ Hence the name ANOVA: we analyze **variances**
- ▶ **Computation possible by hand**, hence the method became popular during the pre-computer age.



ANOVA: intuition of the F-test (hypothetical data)

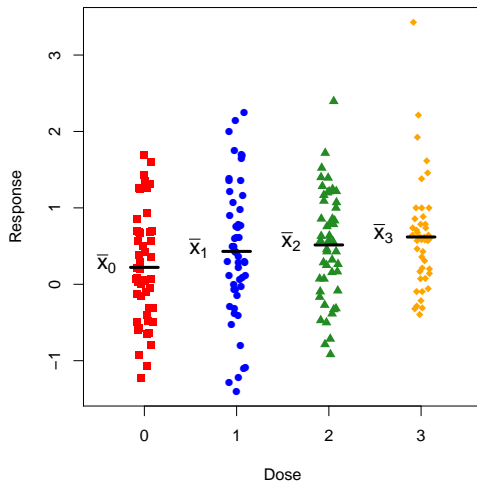


- ▶ **Left:** the **between-group** variance (i.e. the variance of sample means \bar{x}_i) is **small** relative to the **within-group** variability: **do not reject** H_0 .
- ▶ **right:** the **between-group** variance is **large** relative to the **within-group** variability: **reject** H_0 .

Note: of course “small”/”large” is also **relative to the sample size**.

Case: (traditional) ANOVA

- ▶ $\bar{x}_0 = 0.22$
- ▶ $\bar{x}_1 = 0.43$
- ▶ $\bar{x}_2 = 0.51$
- ▶ $\bar{x}_3 = 0.62$
- ▶ $\hat{\sigma}_\varepsilon = 0.76$
- ▶ p-value=0.07
- ▶ Do not reject!

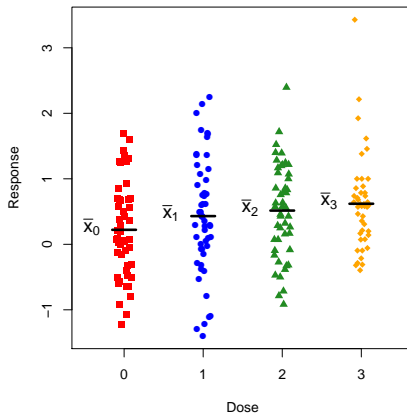


This is **not consistent** with the results
from the all-pairwise comparisons...

Case: (traditional) ANOVA (without assuming homogeneity)

(Flexible model relaxing assumption 4)

- ▶ $\bar{x}_1 = 0.22, \hat{\sigma}_1 = 0.72$
- ▶ $\bar{x}_2 = 0.43, \hat{\sigma}_2 = 0.88$
- ▶ $\bar{x}_3 = 0.51, \hat{\sigma}_3 = 0.71$
- ▶ $\bar{x}_4 = 0.62, \hat{\sigma}_4 = 0.71$
- ▶ p-value=0.055
- ▶ Do not reject!



Still **not consistent** with the results from the all-pairwise comparisons.
although the modeling assumptions are similar !

Power of F-test vs all-pairwise comparisons

We can test the global null hypothesis $H_0 : \mu_1 = \mu_2 = \dots = \mu_K$ using:

- ▶ F-test.
- ▶ **min-P test**: perform all-pairwise comparisons, compute the p-value for H_0 as the minimum of the (multiplicity adjusted) p-values of all the comparisons.

Which approach is the most powerful?

- ▶ F-test when the means of all groups are different although there is no particularly large difference between any two groups.
- ▶ min-P test: when there exists a particularly large difference between any two groups.



Critics of the F-test and recommendations (1/2)

- ▶ When the **F-test is significant** we can **conclude to differences** between the groups **but we do not know between which groups!** (**frustrating**....). Historically, people used to proceed in two steps: 1) test the global null $H_0 : \mu_1 = \mu_2 = \dots = \mu_K$, 2) if H_0 is rejected, proceed to make pairwise comparisons, but why not directly start with pairwise comparisons, especially because...
- ▶ **F-test and pairwise comparisons** are **inconsistent**: either may find a significant difference the other doesn't. When it happens it is **frustrating** and hard to explain.
- ▶ When the F-test is not significant it is difficult to know whether it is due to **lack of effect** or **lack of evidence** (no corresponding CIs)



Critics of the F-test and recommendations (2/2)

Recommendations:

- ▶ Unless you are not interested in the **pairwise comparisons** or have another specific reason in mind (e.g. power), prefer the all-pairwise comparisons and min-P approach to the F-test.
- ▶ If you are not interested in the pairwise comparisons but **only** in knowing whether a continuous outcome is **associated** with a categorical, and if you want to keep the analysis as “**simple and common**” as possible, then prefer to use the F-test to the more “modern” min-P approach.



Checking the ANOVA model assumptions

But the F-test and ANOVA are still very much used... so, what should we know about checking their modeling assumptions?

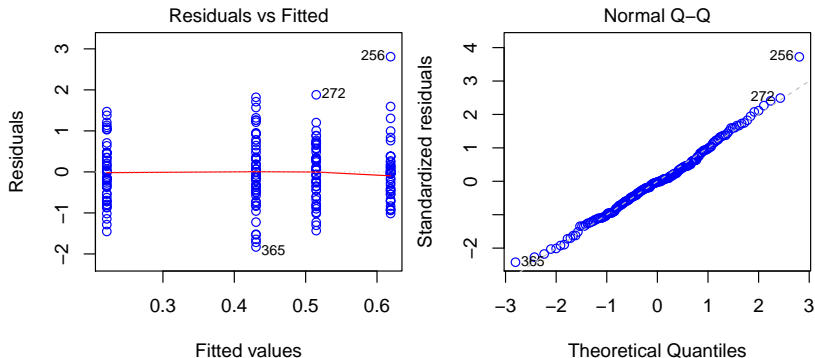
- ▶ Assumptions 1-2 (independence):
 - ▶ rely on the study design.

- ▶ Assumptions 3 (homogeneity of variances):
 - ▶ check with residual plots or compute sd in each group (best).
 - ▶ can be relaxed if needed (see R-demo for code).
 - ▶ log-transforming the data might help to obtain homogeneity of the variances.

- ▶ Assumptions 4 (normality):
 - ▶ check with qqplot.
 - ▶ not needed with large sample sizes in each group.



Case: model checking (default) plots



Note: these are similar plots to those of the linear models.

ANOVA: usual software parametrization (1/4)

R code for ANOVA

```
fitlm <- lm(resp~dosefact, data=d)
summary(fitlm)
```

which returns (among other things)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.2213	0.1079	2.051	0.0416 *
dosefact1	0.2091	0.1498	1.396	0.1643
dosefact2	0.2935	0.1534	1.913	0.0572 .
dosefact3	0.3977	0.1568	2.537	0.0120 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.7631 on 194 degrees of freedom

Multiple R-squared: 0.03513, Adjusted R-squared: 0.02021

F-statistic: 2.354 on 3 and 194 DF, p-value: 0.07335



ANOVA: usual software parametrization (2/4)

- ▶ (Intercept): est. mean in the **reference** group: $\bar{x}_0 = 0.2213$.
- ▶ dosefact1: est. mean **difference**: $\bar{x}_1 - \bar{x}_0 = 0.2091$.
- ▶ dosefact2: est. mean **difference**: $\bar{x}_2 - \bar{x}_0 = 0.2935$.
- ▶ dosefact3: est. mean **difference**: $\bar{x}_3 - \bar{x}_0 = 0.3977$.

- ▶ F-statistic: provides **F-test p-value**: 0.07335.
- ▶ Residual standard error: estimate of σ_ε : 0.7631



ANOVA: usual software parametrization (2/4)

- ▶ (Intercept): est. mean in the **reference** group: $\bar{x}_0 = 0.2213$.
- ▶ dosefact1: est. mean **difference**: $\bar{x}_1 - \bar{x}_0 = 0.2091$.
- ▶ dosefact2: est. mean **difference**: $\bar{x}_2 - \bar{x}_0 = 0.2935$.
- ▶ dosefact3: est. mean **difference**: $\bar{x}_3 - \bar{x}_0 = 0.3977$.

- ▶ F-statistic: provides **F-test p-value**: 0.07335.
- ▶ Residual standard error: estimate of σ_ε : 0.7631

- ▶ p-values for the mean differences are not adjusted for multiple testing.



ANOVA: usual software parametrization (2/4)

- ▶ (Intercept): est. mean in the **reference** group: $\bar{x}_0 = 0.2213$.
- ▶ dosefact1: est. mean **difference**: $\bar{x}_1 - \bar{x}_0 = 0.2091$.
- ▶ dosefact2: est. mean **difference**: $\bar{x}_2 - \bar{x}_0 = 0.2935$.
- ▶ dosefact3: est. mean **difference**: $\bar{x}_3 - \bar{x}_0 = 0.3977$.

- ▶ F-statistic: provides **F-test p-value**: 0.07335.
- ▶ Residual standard error: estimate of σ_ε : 0.7631

- ▶ p-values for the mean differences are not adjusted for multiple testing.

- ▶ “default” summary presents only **comparisons between the reference group and others** (3 out of 6 possible). This is **arbitrary**!
- ▶ Note that if **Dose 1** had been chosen as **the reference group**, among the 3 differences shown in the output **none would be significant**.



ANOVA: usual software parametrization (3/4)

R code for ANOVA when the reference Dose is now Dose 1.

```
d$dosefact <- relevel(d$dosefact,ref="1")
fitlm <- lm(resp-dosefact, data=d)
summary(fitlm)
```

which returns (among other things)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.4304	0.1038	4.145	5.08e-05 ***
dosefact0	-0.2091	0.1498	-1.396	0.164
dosefact2	0.0844	0.1505	0.561	0.576
dosefact3	0.1887	0.1540	1.225	0.222

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.7631 on 194 degrees of freedom

Multiple R-squared: 0.03513, Adjusted R-squared: 0.02021

F-statistic: 2.354 on 3 and 194 DF, p-value: 0.07335



ANOVA: usual software parametrization (4/4)

The ANOVA model is actually a specific kind of linear model.¹⁰ The mean of each group is described by the regression formula:

$$\mu_i = \alpha + \beta_1 \cdot I(\text{group}_i = \text{Dose 1}) + \beta_2 \cdot I(\text{group}_i = \text{Dose 2}) + \beta_3 \cdot I(\text{group}_i = \text{Dose 3})$$

where $I()$ is the **indicator function**:

$$I(\text{condition}) = \begin{cases} 1, & \text{if condition is true} \\ 0, & \text{if condition is false} \end{cases}$$

Group	Dose 0	Dose 1	Dose 2	Dose 3
Mean	α	$\alpha + \beta_1$	$\alpha + \beta_2$	$\alpha + \beta_3$
Estimate	0.2213	0.2213 + 0.2091	0.2213 + 0.2935	0.2213 + 0.3977



Outline

Simple pairwise comparisons

Analysis of Variance (ANOVA): one-way

Analysis of Variance (ANOVA): two-way



Two-way ANOVA

What is it about?

Analysis the mean of a continuous outcome depending on **two categorical variables**.

Why and when is it useful?

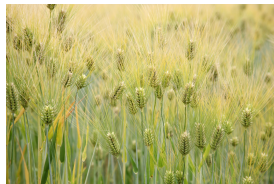
1. to increase **power** and precision of the estimates.
2. to correct/adjust for differences between the group that we primarily aim to compare (e.g. to adjust for baseline differences; to get closer to “causal” conclusions ¹¹).
3. to (sometimes) better handle missing data.

Note: points 2 and 3 are closely related.



Case: barley yield

- ▶ Data from $n = 30$ fields.
- ▶ Five varieties of barley.
- ▶ Six locations.
- ▶ Factorial experiment.



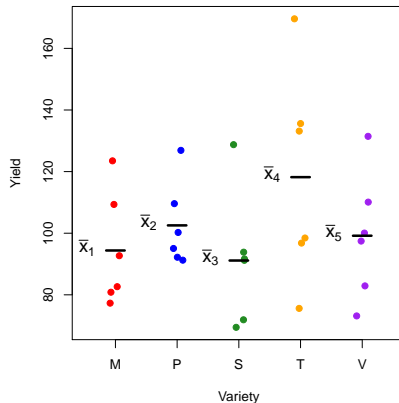
Outcome: average yield over two years (bushels per acre)

Research question:

Do some varieties of barley give better yields than others?

Barley: one-way ANOVA results

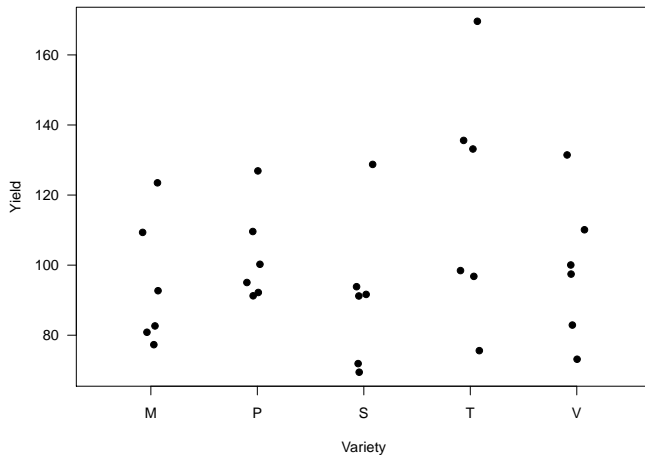
- ▶ $\bar{x}_1 = 94.4$
- ▶ $\bar{x}_2 = 102.5$
- ▶ $\bar{x}_3 = 91.1$
- ▶ $\bar{x}_4 = 118.2$
- ▶ $\bar{x}_5 = 99.2$
- ▶ p-value of F-test=0.30
- ▶ p-value of min-P test=0.26



Note: here we assume equal standard deviation for all varieties (for simplicity and because of the small sample size).



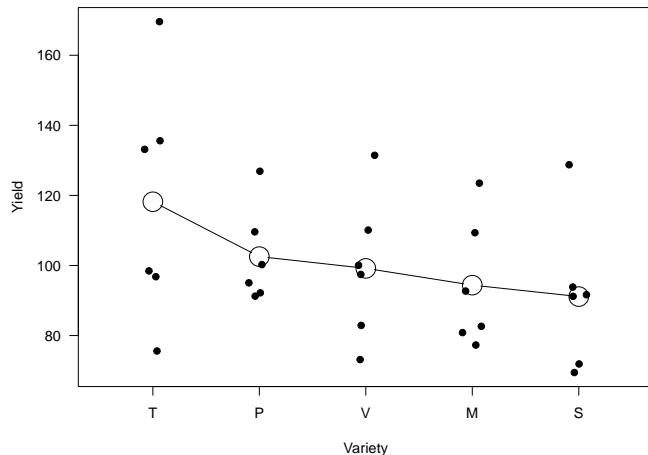
Two-way ANOVA: intuition for power gain



► Raw data.



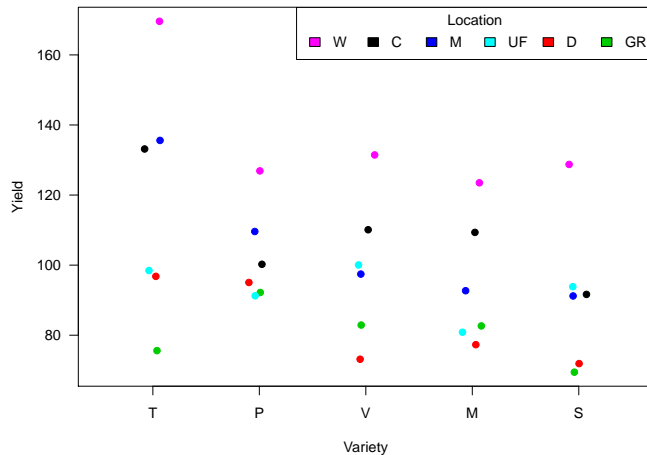
Two-way ANOVA: intuition for power gain



- Reorder the varieties by decreasing mean.

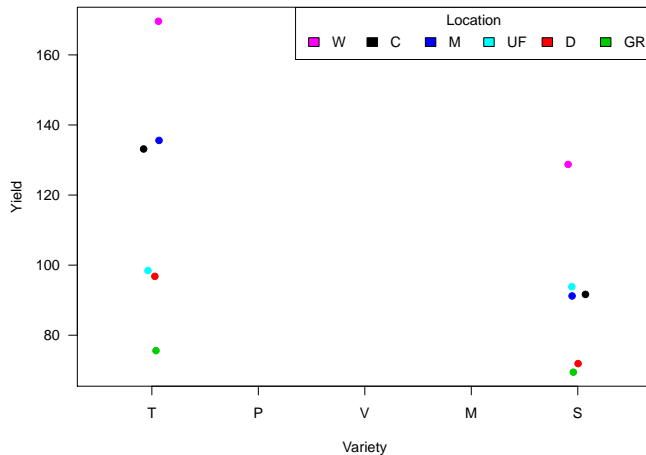


Two-way ANOVA: intuition for power gain



- ▶ The location of the field seems to matter (soil, sun, rain).
- ▶ The location explains some of the variability.

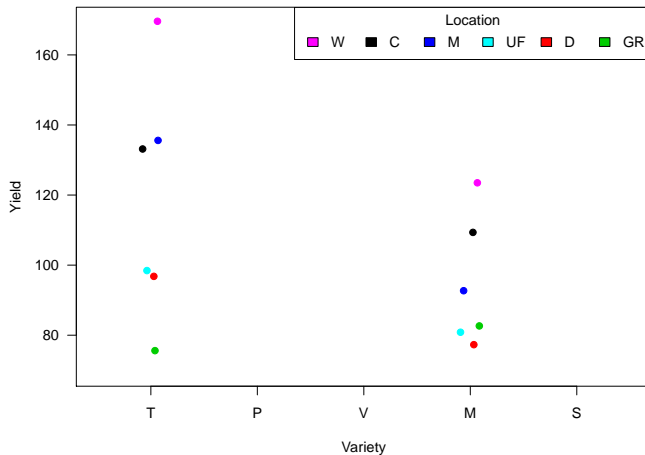
Two-way ANOVA: intuition for power gain



- ▶ T “seems” always better than S, when we compare **per location** (6/6).
- ▶ Whereas **overall**, T was not always better than S.



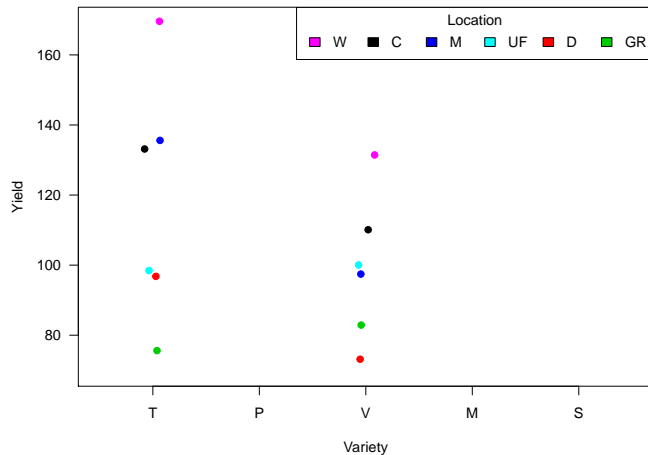
Two-way ANOVA: intuition for power gain



- ▶ T “seems” almost always better than M, when we compare **per location** (5/6).
- ▶ Whereas **overall**, T was not always better than S.



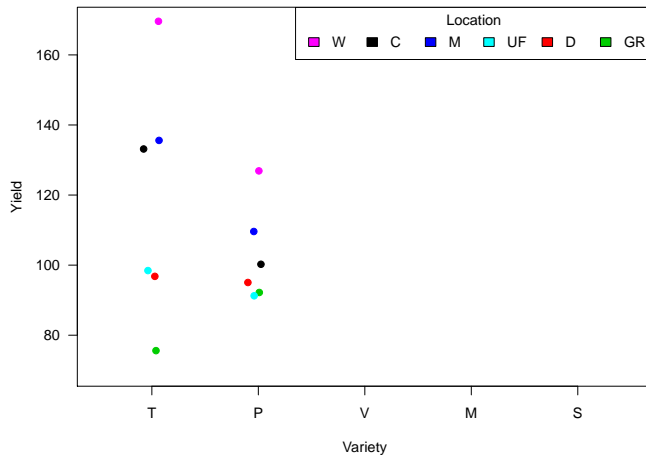
Two-way ANOVA: intuition for power gain



- ▶ T “seems” often better than V, when we compare **per location** (4/6).
- ▶ Whereas **overall**, T was not always better than v.



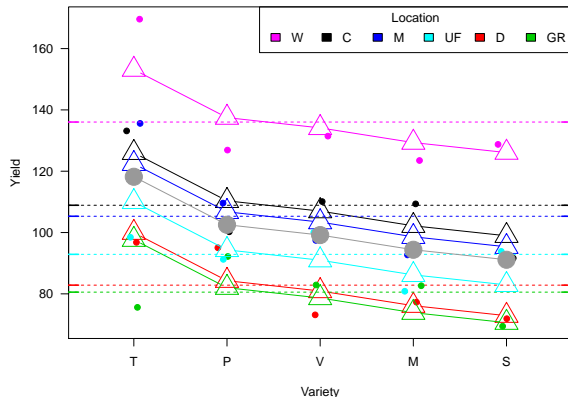
Two-way ANOVA: intuition for power gain



- ▶ T “seems” almost always better than P, when we compare **per location** (5/6).
- ▶ Whereas **overall**, T was not always better than P.



Two-way ANOVA: intuition for power gain



- ▶ With a two-way ANOVA (without interaction) we assume/model that the location systematically shifts the mean yield of each variety up or down.¹²

¹² **Note:** because of the factorial design, differences between varieties are estimated equal “overall” and within the same location: Grey line (overall) parallel to the others (location specific).



The two-way ANOVA model (without interaction)

The k -th observation from the (i, j) -th combination group (e.g. location i and variety j) is described as:

$$\begin{aligned} Y_{ijk} &= \mu_{ij} + \varepsilon_{ijk} \\ &= \gamma_i + \eta_j + \varepsilon_{ijk} \quad (\text{assuming no interaction}). \end{aligned}$$

- ▶ $\mu_{ij} = \gamma_i + \eta_j$ is the mean for the (i, j) -th combination group.
- ▶ ε_{ijk} 's are individual 'error' terms ("random/unexplained deviation from the mean") assumed normally distributed with zero mean and the same variance σ_ε^2 regardless of group.



Two-way ANOVA assumptions (without interaction)

Model assumptions (1-4 similar to that of the one-way ANOVA):

1. Observations from different groups are independent.
2. Individual observations within each group are independent.
3. 'Error' terms are normally distributed.
4. The variance of 'error' terms is the same for all groups (**homogeneity**).
5. There is no interaction (\rightarrow).

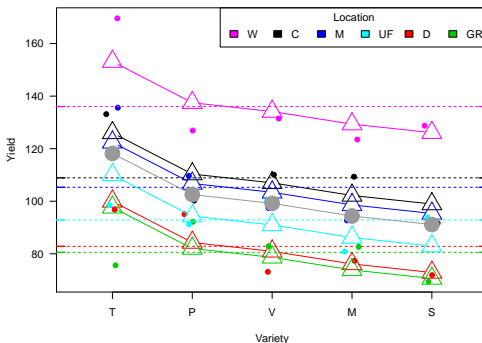
Note: 1-2 and 5 are important, 3-4 not always (as for one-way ANOVA).



The meaning of “no interaction”

No interaction models $\mu_{ij} = \gamma_i + \eta_j$ for the mean for the (i, j) -th combination group. In our example that means that:

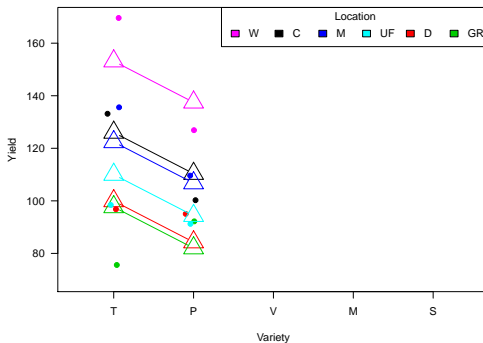
$\left\{ \begin{array}{c} \text{Location} \\ \text{Variety} \end{array} \right\}$ “shifts” (up or down) the mean of all $\left\{ \begin{array}{c} \text{Varieties} \\ \text{Locations} \end{array} \right\}$ in the same way.



The meaning of “no interaction”

No interaction models $\mu_{ij} = \gamma_i + \eta_j$ for the mean for the (i, j) -th combination group. In our example that means that:

$\left\{ \begin{array}{c} \text{Location} \\ \text{Variety} \end{array} \right\}$ “shifts” (up or down) the mean of all $\left\{ \begin{array}{c} \text{Varieties} \\ \text{Locations} \end{array} \right\}$ in the same way.



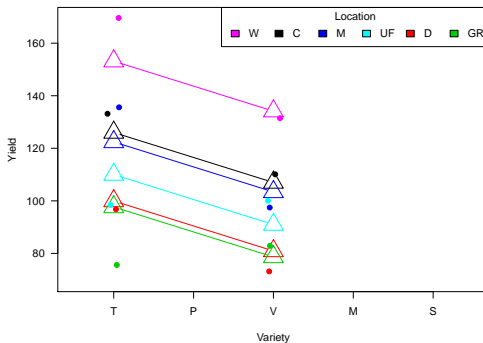
- For all locations, the mean yield difference between varieties T and P is the same (15.7).



The meaning of “no interaction”

No interaction models $\mu_{ij} = \gamma_i + \eta_j$ for the mean for the (i, j) -th combination group. In our example that means that:

$\left\{ \begin{array}{c} \text{Location} \\ \text{Variety} \end{array} \right\}$ “shifts” (up or down) the mean of all $\left\{ \begin{array}{c} \text{Varieties} \\ \text{Locations} \end{array} \right\}$ in the same way.



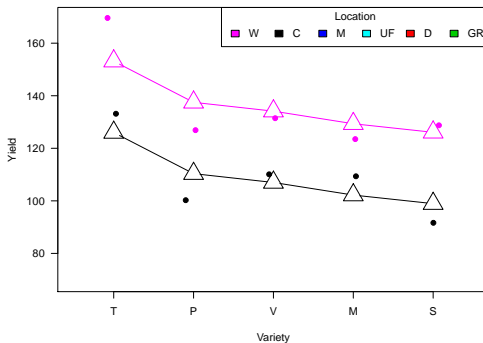
- For all locations, the mean yield difference between varieties T and V is the same (19.0).



The meaning of “no interaction”

No interaction models $\mu_{ij} = \gamma_i + \eta_j$ for the mean for the (i, j) -th combination group. In our example that means that:

$\left\{ \begin{array}{c} \text{Location} \\ \text{Variety} \end{array} \right\}$ “shifts” (up or down) the mean of all $\left\{ \begin{array}{c} \text{Varieties} \\ \text{Locations} \end{array} \right\}$ in the same way.



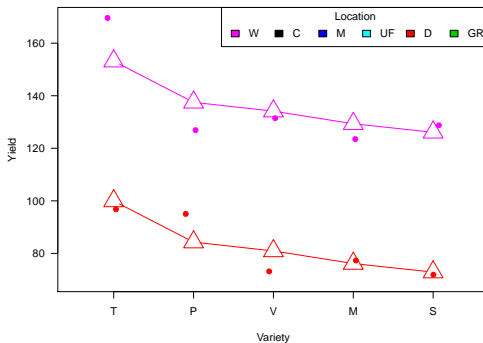
- For all varieties, the mean yield difference between locations W and C is the same (27.1).



The meaning of “no interaction”

No interaction models $\mu_{ij} = \gamma_i + \eta_j$ for the mean for the (i, j) -th combination group. In our example that means that:

$\left\{ \begin{array}{c} \text{Location} \\ \text{Variety} \end{array} \right\}$ “shifts” (up or down) the mean of all $\left\{ \begin{array}{c} \text{Varieties} \\ \text{Locations} \end{array} \right\}$ in the same way.



- For all varieties, the mean yield difference between locations W and D is the same (53.2).



Two-way ANOVA: usual software parametrization (1/4)

One-way ANOVA (for comparison)

```
OneWayRes <- lm(Y~ Var, data = d)
summary(OneWayRes)
```

which returns

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	94.392	9.247	10.207	2.12e-10	***
VarP	8.150	13.078	0.623	0.5388	
VarS	-3.258	13.078	-0.249	0.8053	
VarT	23.808	13.078	1.821	0.0807	.
VarV	4.792	13.078	0.366	0.7171	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 22.65 on 25 degrees of freedom

Multiple R-squared: 0.1715, Adjusted R-squared: 0.03893

F-statistic: 1.294 on 4 and 25 DF, p-value: 0.2993



Two-way ANOVA: usual software parametrization (2/4)

Two-way ANOVA

```
TwoWayRes <- lm(Y~ Var + Loc, data = d)
summary(TwoWayRes)
```

which returns

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	102.202	6.078	16.815	2.88e-13	***
VarP	8.150	6.078	1.341	0.194983	
VarS	-3.258	6.078	-0.536	0.597810	
VarT	23.808	6.078	3.917	0.000854	***
VarV	4.792	6.078	0.788	0.439728	
LocD	-26.060	6.658	-3.914	0.000860	***
LocGR	-28.340	6.658	-4.256	0.000386	***
LocM	-3.590	6.658	-0.539	0.595705	
LocUF	-16.010	6.658	-2.405	0.025996	*
LocW	27.140	6.658	4.076	0.000589	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 10.53 on 20 degrees of freedom

Multiple R-squared: 0.8568, Adjusted R-squared: 0.7924

F-statistic: 13.3 on 9 and 20 DF, p-value: 1.216e-06

Note: F-test p-value is not so interesting here (H_0 : “neither effect of variety nor of location”)



Two-way ANOVA: usual software parametrization (3/4)

- ▶ (Intercept): est. mean in the **reference** group (location C, variety M).
- ▶ VarP: est. mean **difference** between varieties P and M, grown **at the same location** (any).
- ▶ VarS: est. mean **difference** between varieties S and M, grown **at the same location** (any).
- ▶ LocD: est. mean **difference** between Location D and C, when growing **the same variety** (any).
- ▶ F-statistic: not so interesting (see previous slide).
- ▶ Residual standard error: estimate of σ_ε : 10.53



Two-way ANOVA: usual software parametrization (3/4)

- ▶ (Intercept): est. mean in the **reference** group (location C, variety M).
- ▶ VarP: est. mean **difference** between varieties P and M, grown **at the same location** (any).
- ▶ VarS: est. mean **difference** between varieties S and M, grown **at the same location** (any).
- ▶ LocD: est. mean **difference** between Location D and C, when growing **the same variety** (any).
- ▶ F-statistic: not so interesting (see previous slide).
- ▶ Residual standard error: estimate of σ_ε : 10.53
- ▶ p-values for the mean differences are not adjusted for multiple testing.



Two-way ANOVA: usual software parametrization (3/4)

- ▶ (Intercept): est. mean in the **reference** group (location C, variety M).
- ▶ VarP: est. mean **difference** between varieties P and M, grown **at the same location** (any).
- ▶ VarS: est. mean **difference** between varieties S and M, grown **at the same location** (any).
- ▶ LocD: est. mean **difference** between Location D and C, when growing **the same variety** (any).
- ▶ F-statistic: not so interesting (see previous slide).
- ▶ Residual standard error: estimate of σ_ε : 10.53
- ▶ p-values for the mean differences are not adjusted for multiple testing.
- ▶ “default” summary presents only **comparisons between the reference group and the others** (4 out of 10 possible for comparing varieties, 5 out of 15 for locations). This is **arbitrary**!



Two-way ANOVA: usual software parametrization (4/4)

This ANOVA model is also a specific kind of linear model. The mean of each group is given by this regression formula:

$$\begin{aligned}\mu_{ij} = & \alpha + \beta_1 \cdot I(\text{variety}_i=P) + \beta_2 \cdot I(\text{variety}_i=S) + \beta_3 \cdot I(\text{variety}_i=T) \\ & + \beta_4 \cdot I(\text{variety}_i=V) + \beta_5 \cdot I(\text{location}_j=D) + \beta_6 \cdot I(\text{location}_j=GR) \\ & + \beta_7 \cdot I(\text{location}_j=M) + \beta_8 \cdot I(\text{location}_j=UF) + \beta_9 \cdot I(\text{location}_j=W)\end{aligned}$$

Examples (modeled means and estimates):

Location \ Variety	M	P	S
C	α 102.2	$\alpha + \beta_1$ 102.2 + 8.2	$\alpha + \beta_2$ 102.2 - 3.3
D	$\alpha + \beta_5$ 102.2 - 26.1	$\alpha + \beta_1 + \beta_5$ 102.2 + 8.2 - 26.1	$\alpha + \beta_2 + \beta_5$ 102.2 - 3.3 - 26.1
W	$\alpha + \beta_9$ 102.2 + 27.1	$\alpha + \beta_1 + \beta_9$ 102.2 + 8.2 + 27.1	$\alpha + \beta_2 + \beta_9$ 102.2 - 3.3 + 27.1



F-test with two-way ANOVA

An F-test can be performed for the null hypothesis

$$H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4$$

that is

H_0 : “All varieties give the same average yield,
when grown at the same location”.

this compares the mean yield of the varieties “adjusted” on location (i.e. within the same location).

- ▶ This is a very **commonly used**, well-known and “old” method.
- ▶ **Pros and cons**: similar to that of F-test for one-way ANOVA.



R code and results

F-test in two-way ANOVA

```
anova(lm(Y~ Loc, data = d),lm(Y~ Var + Loc, data = d))
```

which returns

Analysis of Variance Table

Model 1: Y ~ Loc

Model 2: Y ~ Var + Loc

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	24	4871.5				
2	20	2216.5	4	2655	5.9891	0.002453 **

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Comments:

- ▶ F-test p-value=0.002453 significant, whereas it was in the one-way ANOVA (p=0.299).
- ▶ To **avoid coding mistakes and misunderstandings** of R output do compare the two models: do not use “anova(TwoWayRes)”.



Recommended analysis (see R-demo for code)

Statistical methods:

Comparisons between varieties were made with a two-way ANOVA model (without interaction) to adjust for the location. P-values and 95% confidence intervals were adjusted for multiple testing using the min-P method as implemented in the multcomp-package [ref.¹³] of R [ref.¹⁴] and described in [ref.¹⁵].

Results (adjusted for multiple testing):

Comparison	Est. Diff	95% CI	p-value
P - M	8.1	[-10.0; 26.3]	0.670
S - M	-3.3	[-21.4; 14.9]	0.982
T - M	23.8	[5.6; 42.0]	0.007
V - M	4.8	[-13.4; 23.0]	0.931
S - P	-11.4	[-29.6; 6.8]	0.361
T - P	15.7	[-2.5; 33.8]	0.113
V - P	-3.4	[-21.5; 14.8]	0.980
T - S	27.1	[8.9; 45.3]	0.002
V - S	8.0	[-10.1; 26.2]	0.680
V - T	-19.0	[-37.2; -0.8]	0.038

Note: p-values ≤ 10 times the non-adjusted ones obtain from the default summary (Bonferroni).

¹³ Hothorn, Bretz & Westfall (2008). Simultaneous Inference in General Parametric Models. Biometrical Journal 50(3), 346–363.

¹⁴ R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

¹⁵ Bretz, Hothorn, & Westfall (2016). Multiple comparisons using R. CRC Press.

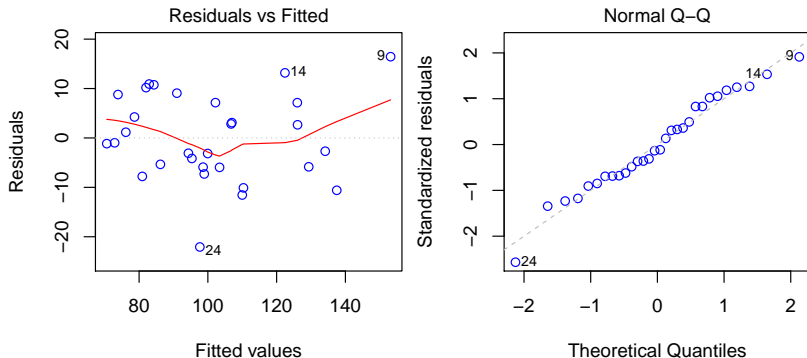


Assuming of no interactions

- ▶ This assumption can be important.
- ▶ It simplifies the interpretation of the results.
- ▶ It should be supported by **subject-matter knowledge**.
- ▶ This assumption can (most often) be checked.
- ▶ Usually, **the smaller the sample size the more assumptions we need to compensate**. This applies to the assumption of no interaction.
- ▶ More on interactions in Lecture 7.



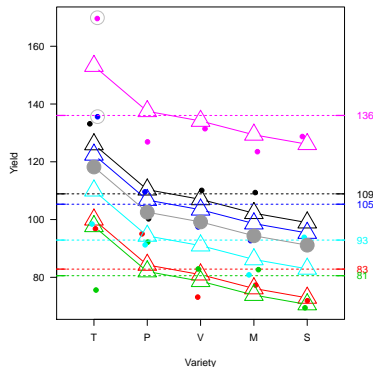
Model checking (default) plots



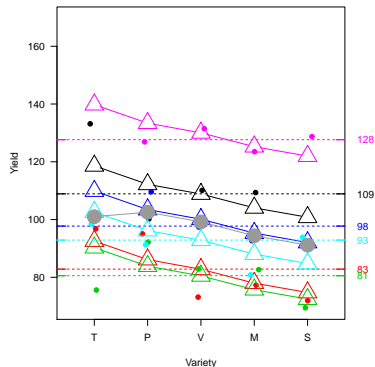
Note: these are similar plots to those of the linear models.

Missing data and unbalanced design

Complete data



Data with 2 observations missing



- ▶ Due to missing data, the average yield for variety “T” is no longer the largest, yet it is still estimated the largest when “adjusting” on location (right plot).
- ▶ Adjusting on location matters (more) for unbalanced data (e.g: “T” the best variety, is grown only at the 4 worst locations).
- ▶ Adjusting helps to “fairly” compare two varieties when those are not grown in similar locations (More on Lectures 6 & 7).