TMA4267 Linear Statistical Models V2017 (L13)

Part 3: Hypothesis testing and analysis of variance Hypothesis testing: why, how and be aware Reproduciability

The universal F-test [F:3.3]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: March 3, 2017

1/36

Basal metabolic rate and the FTO-gene

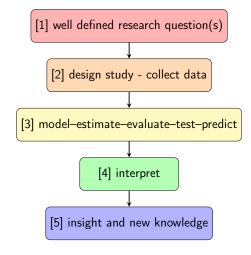
- ▶ The gene called FTO is known to be related to obesity
- ► The basal metabolic rate says how many calories you burn when you rest (hvilemetabolisme).
- ▶ Data has been collected for 101 patient from the obesity clinic at St. Olavs Hospital.
- ► Research question: is there an association between the variant of the FTO gene of the patient and the basal metabolic rate?
- ► Regression setting, other covariates include age, sex, weight, height, BMI, diet, exercise level, smoking, etc.

Today

- ► The scientific process.
- ▶ The basics of hypothesis testing and interpretation of *p*-value.
- ► The reproduciability "crisis".
- ▶ Properties of *p*-values.
- Linear hypotheses in regression vs. nested models.
- ► The universal F-test for linear hypotheses (nested models)

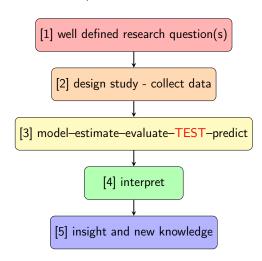
1/36

The scientific process



3/36

The scientific process



4/36

Hypothesis testing example (cont.)

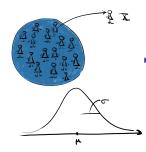


- We draw a random sample of size n = 100 from the blue population and measure systolic blood pressure: X_1, X_2, \dots, X_n .
- ▶ Test statistic: $\bar{X} \sim N(120, 1)$ when H_0 is true.
- ightharpoonup We find that $ar{x}=122$ mmHg.
- ▶ Data: n = 100, $\bar{x} = 122$, gives a p-verdi=0.02.

Questions:

- ► How have I calculated this *p*-value?
- ▶ Should I conclude that $\mu > 120$?

Hypothesis testing example



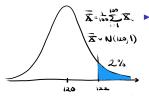
- It is known that in a population of women of age 20-29 years the systolic blood pressure is normally distributed with mean $\mu=120$ mmHg.
- ▶ We study a population of women of age 20-29 that have a specific disease (blue population), and also here we assume that the systolic blood pressure is normally distributed (with standard deviation 10 mmHg), but here we don't know the mean in the population.
- ▶ In addition to estimating this unknown mean we want to investigate if the mean blood pressure of the blue population is larger than 120 mmHg (because if it is, we need to start more investigations into the cause of this).
- ▶ $H_0: \mu = 120 \text{ vs. } H_1: \mu > 120.$

5/36

Q and A

- ► How have I calculated this *p*-value? $P(\bar{X} > 122 \mid H_0 \text{ true})$.
- Should I conclude that $\mu > 120$? Yes, if you choose significance level higher than 0.02. But, you should also report a (two-sided) confidence interval for μ : Here [120.04, 123.96].

Hypothesis testing example (end)



 $\overline{X} = \sum_{i=1}^{\infty} X$ The *p*-value is often based on a test statistic, and can be found in many ways (known distribution, enumerations, asymptotic).

- ► Significance level: highest probability of miscarriage of justice that we would tolerate.
- ▶ We reject the null hypothesis and say that we have a significant finding at significance level α if a/the *p*-value for the hypothesis test is below α .

8/36

What is a *p*-value

A more correct definition so that:

the p-value is the probability of your result or a more extreme result, given that H_0 is true.

or

the probability of your result or a more extreme result, given that it occurred randomly.

This is different from: the probability of your result occurring randomly.

Slide reconstructed from talk by Kristoffer H. Hellton, NR

10/36

What is a p-value

From The research handbook of Carlsen & Staff (2014)

... the p-value, the probability that the result could have occurred randomly, p=probability.

This is common, but not the correct definition of the p-value. What is wrong? Discuss!

Slide reconstructed from talk by Kristoffer H. Hellton, NR

9 / 36

A simple example

- ► Null hypothesis: It is sunny outside.
- ▶ Data: I enter the room soaking wet.
- ▶ Wrong *p*-value: the probability that it is sunny outside.
- ► Impossible to calculate.
- ▶ Right *p*-value: the probability that I'm wet, given that it is sunny.
- ▶ Should be small.

Important! From Bayes theorem:

 $P(observation \mid hypothesis) \neq P(hypothesis \mid observation)$

The probability of observing a result given some hypothesis is true not equivalent to the probability that the hypothesis is true given that some result has be observed.

To be able to calculate the right hand side, we need P(hypothesis), the probability of the hypothesis. This is exactly what is introduced in Bayesian statistics through the so-called prior, and some see the Bayes factor as the replacement for p-values.

Slide reconstructed from talk by Kristoffer H. Hellton, NR

Statistical significance and *p*-values

On March 7, 2016, the American Statistical Association posted a statement on statistical significance and p-values - "clarifying several widely agreed upon principles underlying the proper use and interpretation of the p-value".

12 / 36

Statement on proper use and interpretation of the p-value

Why is this needed: (2)

 $\label{eq:hack your way to scientific glory} Hack your way to scientific glory$

loannidis (2005): How many nonsignificant results have been studied before one research group has published its first significant finding?

Statement on proper use and interpretation of the p-value

Why is this needed: (1)

American Statistical Association discussion forum, 2014.

- ightharpoonup Q: Why do so many colleges and grad schools teach p = 0.05?
- ► A: Because that's still what the scientific community and journal editors use.
- ightharpoonup Q: Why do so many people still use p = 0.05?
- ► A: Because that's what they were taught in college or grad school.

Problem?

Urban knowledge: Unless an hypothesis test results in a p-value below 0.05 there is no finding. So, in some journals a researcher will not be able to publish his paper unless the test performed has a p-value below 0.05.

13/36

Statement on proper use and interpretation of the p-value

Why is this needed: (3)

The journal *Basic and Applied Social Psychology* (editors Trafimow and Marks, 2015) put a *ban* on null hypothesis significance testing.

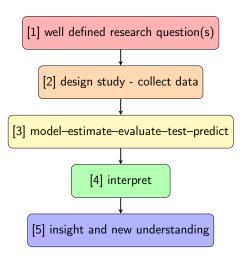
ASA Statement on Statistical Significance and *P*-values, March 2016

The ASA's statement on p-values: context, process, and purpose, Ronald L. Wasserstein & Nicole A. Lazar, The American Statistician, DOI:10.1080/00031305.2016.1154108.

- ▶ While the *p*-value can be a useful statistical measure, it is commonly misused and misinterpreted.
- ▶ Informally, a *p*-value is the probability under a specified statistical model that a statistical summary of the data would be equal to or more extreme than its observed value.
- ▶ P1: *P*-values can indicate how incompatible the data are with a specified statistical model.
- ▶ P2: *P*-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.
- ▶ P3: Scientific conclusions and business or policy decisions should not be based only on whether at *p*-value passes a specific threshold.

16 / 36

The scientific process



ASA Statement on Statistical Significance and P-values

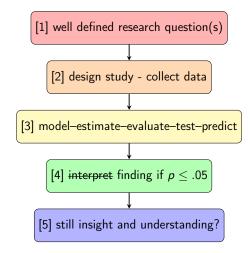
- ▶ P4: Proper inference requires full reporting and transparency.
- ▶ P5: A *p*-value, or statistical significance, does not measure the size of an effect or the importance of a result.
- ▶ P6: By itself, a *p*-value does not provide a good measure of evidence regarding a model or hypothesis.

Take home message: the *p*-value is a very risky tool ...

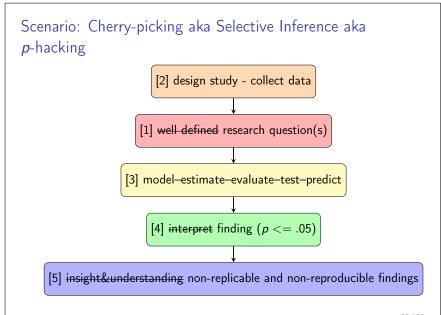
(Benjamini, 2016): but, replacing the *p*-value with other tools may lead to many of the same indeficiencies - so it would be better to instead focus on the appropriate use of statistical tools for addressing the crisis of reproducibility and replicability in science.

17 / 36

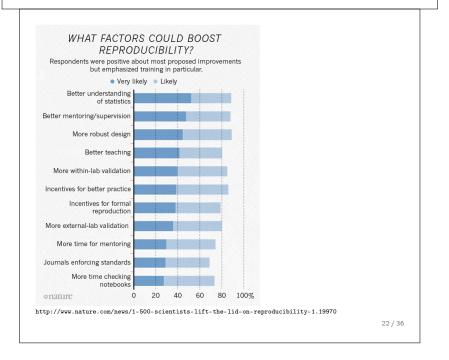
Scenario: finding only for $p \le 0.05$

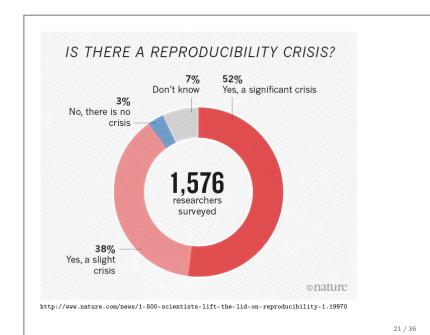


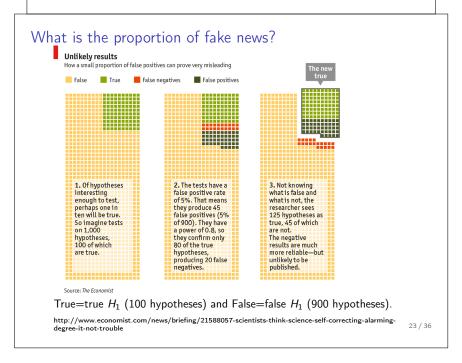
18 / 36











What is the proportion of fake news?

Color-coding for the far left figure:

- ▶ Yellow: all the hypotheses where H_0 is true (and H_1 is false), and H_0 is not rejected. All is good here, but this interesting(?) findings are very seldom published.
- ▶ Light green: all the hypotheses where *H*₀ is false (and *H*₁ is true) and the research reject the *H*₀ and make a correct discovery. This are our true news!
- ▶ Dark green: all the hypothesis where H₀ are true (and H₁ are false) but the researcher wrongly reject H₀. These are our fake news!
- ▶ Red: all the hypotheses where H_0 are false (and H_1 is true) but where the researcher fail to reject H_0 let guilty criminal go free. These are called false negatives and are usually not reported (unless the researcher is report a negative finding).

So, not 5% of published results are false positives (fake news), but rather at substantially larger number - 40-90% has be hinted to in different publications.

24 / 36

So far

- ► We (statisticians and other scientists) must focus on sound scientific process and step away from cherry-picking and the "finding=p-value < 0.05" urban truth.
- ▶ We must always report effect size.
- ▶ We must be aware that these two effects (selective inference and practical vs. statistical significance) are especially important for large than small data sets (both many samples and variables).
- Now, we move to hypothesis testing in linear regression and look at one unifying F-test can be used for all linear hypotheses.

Single hypothesis testing set-up

	H_0 true	H_0 false
Not reject H ₀	Correct	Type II error
Reject H_0	Type I error	Correct

Two types of errors:

- False positives = type I error =miscarriage of justice.
 These are our fake news.
- ► False negatives = type II error= guilty criminal go free.

The significance level of the test is α .

We say that : Type I error is "controlled" at significance level α .

The probability of miscarriage of justice (Type I error) does not exceed α

25 / 36

Happiness (n = 39)

Are love and work the important factors determining happiness?

- y, happiness. 10-point scale, with 1 representing a suicidal state, 5 representing a feeling of «just muddling along», and 10 representing a euphoric state.
- ▶ x₁, money. Annual family income in thousands of dollars.
- $ightharpoonup x_2$, sex. Sex was measured as the values 0 or 1, with 1 indicating a satisfactory level of sexual activity.
- x₃, love. 3-point scale, with 1 representing loneliness and isolation, 2 representing a set of secure relationships, and 3 representing a deep feeling of belonging and caring in the context of some family or community.
- x4, work. 5-point scale, with 1 indicating that an individual is seeking other employment, 3 indicating the job is OK, and 5 indicating that the job is enjoyable.

Data taken from library faraway, data set happy.

27 / 36

What is \boldsymbol{C} and \boldsymbol{d} ?

Use the happiness data, with the four covariates x1=money, x2=sex, x3=love, x4=work, to construct the \boldsymbol{C} and \boldsymbol{d} to test $H_0: \boldsymbol{C}\boldsymbol{\beta} = \boldsymbol{d}$.

There is a linear effect in money? $H_0: \beta_1 = 0$ $\mathbf{C} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \end{bmatrix}, \mathbf{d} = 0$

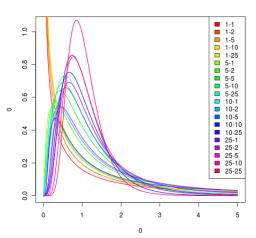
Is the regression significant? $H_0: \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$

$$m{C} = \left[egin{array}{cccc} 0 & 1 & 0 & 0 & 0 \ 0 & 0 & 1 & 0 & 0 \ 0 & 0 & 0 & 1 & 0 \ 0 & 0 & 0 & 0 & 1 \end{array}
ight], m{d} = \left[egin{array}{c} 0 \ 0 \ 0 \ 0 \end{array}
ight]$$

Is there a linear effect of money and/or sex? $H_0: \beta_1 = \beta_2 = 0$

$$oldsymbol{C} = \left[egin{array}{cccc} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{array}
ight], oldsymbol{d} = \left[egin{array}{c} 0 \\ 0 \end{array}
ight]$$

28 / 36



The Fisher distribution with different degrees of freedom ν_1 and ν_2 (given in the legend).

The Fisher distribution [F: B.1 Def 8.14], Exercise 2 Problem 5

"Tabeller og formeler i statistikk": If Z_1 and Z_2 are independent and χ^2 -distributed with ν_1 and ν_2 degrees of freedom, then

$$F = \frac{Z_1/\nu_1}{Z_2/\nu_2}$$

is F(isher)-distributed with ν_1 and ν_2 degrees of freedom.

- ▶ The expected value of F is $E(F) = \frac{\nu_2}{\nu_2 2}$.
- ► The mode is at $\frac{\nu_1 2}{\nu_1} \frac{\nu_2}{\nu_2 + 2}$.
- ► Identity:

$$f_{1-lpha,
u_1,
u_2} = rac{1}{f_{lpha,
u_2,
u_1}}$$

29 / 36

Unrestricted (Model A): all 4 covariates present

```
fitA <- lm(happy~.,data=happy)
summary(fitA)</pre>
```

Coefficients:

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

Residual standard error: 1.058 on 34 degrees of freedom Multiple R-squared: 0.7102, Adjusted R-squared: 0.6761 F-statistic: 20.83 on 4 and 34 DF, p-value: 9.364e-09

Restricted (Model B): only love and work

The estimate $\hat{\beta}_3$ (love) is 1.919 for model A and 1.959 for model B. Explain why these two estimates differ.

```
fitB <- lm(happy~love+work,data=happy)</pre>
summary(fitB)
```

Coefficients:

Estimate Std. Error t value Pr(>|t|) 0.7757 0.265 0.79241 (Intercept) 0.2057 love 1.9592 0.2954 6.633 9.99e-08 *** work 0.5106 0.1874 2.725 0.00987 **

Residual standard error: 1.08 on 36 degrees of freedom

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

Multiple R-squared: 0.6808, Adjusted R-squared: 0.6631 F-statistic: 38.39 on 2 and 36 DF, p-value: 1.182e-09

32 / 36

3.13 Testing Linear Hypotheses

Hypotheses

General linear hypothesis:

$$H_0: C\beta = d$$
 against $H_0: C\beta \neq d$

where C is a $r \times p$ -matrix with $\text{rk}(C) = r \leq p$ (r linear independent restrictions).

Test of significance (t-test):

$$H_0: \beta_i = 0$$
 against $H_1: \beta_i \neq 0$

3. Composite test of a subvector:

$$H_0: \boldsymbol{\beta}_1 = \mathbf{0}$$
 against $H_1: \boldsymbol{\beta}_1 \neq \mathbf{0}$

4. Test for significance of regression:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$
 against
 $H_1: \beta_j \neq 0$ for at least one $j \in \{1, \dots, k\}$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.135)

34 / 36

Model A vs model B

> anova(fitA.fitB) Analysis of Variance Table

33 / 36

Test Statistics

Assuming normal errors we obtain under H_0 :

1.
$$F = 1/r (C\hat{\beta} - d)' (\hat{\sigma}^2 C(X'X)^{-1}C')^{-1} (C\hat{\beta} - d) \sim F_{r,n-p}$$

$$2. t_j = \frac{\beta_j}{\text{se}_i} \sim t_{n-p}$$

3.
$$F = \frac{1}{r} (\hat{\beta}_1)' \widehat{\text{Cov}(\hat{\beta}_1)^{-1}} (\hat{\beta}_1) \sim F_{r,n-p}$$

4.
$$F = \frac{n-p}{k} \frac{R^2}{1-R^2} \sim F_{k,n-p}$$

Critical Values

Reject H_0 in the case of:

1.
$$F > F_{r,n-r}(1-\alpha)$$

1.
$$F > F_{r,n-p}(1-\alpha)$$

2. $|t| > t_{n-p}(1-\alpha/2)$
3. $F > F_{r,n-p}(1-\alpha)$
4. $F > F_{k,n-p}(1-\alpha)$

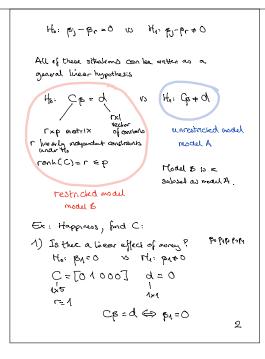
The tests are relatively robust against moderate departures from normality. In addition, the tests can be applied for large sample size, even with nonnormal errors.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.135)

Today

- ▶ Reproduciable research and the scientific method.
- ▶ Hypothesis testing and *p*-values in general.
- ► Type I errors=false positives=fake news.
- ▶ Linear hypotheses, and the F_{obs} test statistic.

36 / 36



PART S: HYPOTHESIS
TESTING AND ANALYSIS OF
VARIANCE (ANOVA)

21.1 FORPANT F102.80.50

A lecture + 1 RecEx + 1 Comulsony Ex

Hypothesis techniq in linear regression [f.3.3]

Y= XB+E, E~ Nn (0, 5° I)

So for we have looked at two types of hypotheses:

- 1) Test for significance of one Bj (Ex: Brown)
 Ho: Bj = 0 vs. th: Bj = 0

 => summer(Lin. radd) outermetricity added.
- 2) Is the regression significant?

 Ho: \$p_1 = p_2 = = pu = 0 US Ho: at least one
- in addition we might went to
- 3) Top of equality (Ex Hunch rent: top w good book on)

1

2) Is the regression agrificant?

It is $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ us the at least one $C = \begin{cases}
0.1 & 0.00 & 0.00 \\
0.0 & 0.00 & 0.00 \\
0.0 & 0.00 & 0.00
\end{cases}$ $C = \begin{cases}
0.1 & 0.00 & 0.00 \\
0.0 & 0.00 & 0.00 \\
0.0 & 0.00 & 0.00
\end{cases}$ $C = \begin{cases}
0.1 & 0.00 & 0.00 \\
0.00 & 0.00 & 0.00 \\
0.00 & 0.00 & 0.00
\end{cases}$

3) Is then a linear effect of money and/or sex?

to: \$1= 82=0 \$\text{ St. at leen one }\$\tag{25}\$

C= 2

3

Procedure (for technollinear hypotheses)

Urestricted nodel is Restricted model
Y= Xp+E, c~Na(0,0) > CR=d

Ex: "p1=0" money exemple
Unablickd: fit all eventio: x1, x2, x2, x4
Restricted: ft oats: x2, x3, x4

- i) Fit the unrestricted model (A) and compute SSE = ETE. Assume pregs. perem. Inted.
- ii) Fit the restricted model (3) and compute $85E_{R_0}=\hat{\mathcal{E}}_{R_0}^{+}\hat{\mathcal{E}}_{R_0}$

NB: the restricted award needs to be rested within the unrestricted.

A: Fill model: x, x, x, x, xy

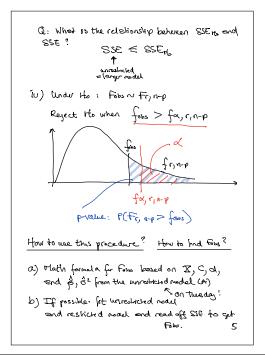
85E = 32. (n-p) = (1.056)284 = 88.087

B: Restricted model: X2 X4 3SEtt 2 (1.018. 26 = 41.952

p-value = P(Fz, 84 > 1.752) = 0.1934

⇒ do not reject the: we prefur the smaller model "x5 x xy".

H: p1= p1=0 H1: at least on + 0



TMA4267 Linear Statistical Models V2017 (L14)

Part 3: Hypothesis testing and analysis of variance The universal F-test [F:3.3] One-way ANOVA [H:8.1.1]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: March 7, 2017

Today

- Linear hypotheses in regression vs. nested models.
- ▶ The universal F-test for linear hypotheses: two formulas.
- ► The two formulas: one easy to use, one easy for proving F-distribution.
- Special cases of the universal F-test.
- ► New problem: categorical covariate with effect coding (for interpretation)

1/12

3.13 Testing Linear Hypotheses

Hypotheses

1. General linear hypothesis:

$$H_0: C\beta = d$$
 against $H_0: C\beta \neq d$

where C is a $r \times p$ -matrix with $\operatorname{rk}(C) = r \leq p$ (r linear independent restrictions).

2. Test of significance (t-test):

$$H_0: \beta_i = 0$$
 against $H_1: \beta_i \neq 0$

Composite test of a subvector:

$$H_0: \boldsymbol{\beta}_1 = \mathbf{0}$$
 against $H_1: \boldsymbol{\beta}_1 \neq \mathbf{0}$

4. Test for significance of regression:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$
 against

$$H_1: \beta_i \neq 0$$
 for at least one $j \in \{1, \dots, k\}$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.135)

Happiness (n = 39)

Are love and work the important factors determining happiness?

- y, happiness. 10-point scale, with 1 representing a suicidal state, 5 representing a feeling of «just muddling along», and 10 representing a euphoric state.
- \triangleright x_1 , money. Annual family income in thousands of dollars.
- ▶ x₂, sex. Sex was measured as the values 0 or 1, with 1 indicating a satisfactory level of sexual activity.
- x₃, love. 3-point scale, with 1 representing loneliness and isolation, 2 representing a set of secure relationships, and 3 representing a deep feeling of belonging and caring in the context of some family or community.
- x₄, work. 5-point scale, with 1 indicating that an individual is seeking other employment, 3 indicating the job is OK, and 5 indicating that the job is enjoyable.

Data taken from library faraway, data set happy.

2/12

Constrained and unconstrained estimate

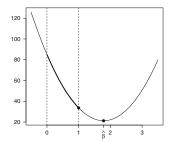


Fig. 3.15 Illustration of the difference in goodness of fit between the unconstrained least squares estimator and the estimator under the constraint $0 \le \beta \le 1$. The (unconstrained) least squares estimator is labeled as $\hat{\beta}$. For the constrained solution, we have $\hat{\beta}=1$

Figure 3.15 from our text book: Fahrmeir et al (2013): Regression. Springer. (p.1329)

3.13 Testing Linear Hypotheses

Hypotheses

1. General linear hypothesis:

$$H_0: C\beta = d$$
 against $H_0: C\beta \neq d$

where C is a $r \times p$ -matrix with $\text{rk}(C) = r \leq p$ (r linear independent restrictions).

Test of significance (t-test):

$$H_0: \beta_i = 0$$
 against $H_1: \beta_i \neq 0$

3. Composite test of a subvector

$$H_0: \boldsymbol{\beta}_1 = \mathbf{0}$$
 against $H_1: \boldsymbol{\beta}_1 \neq \mathbf{0}$

4. Test for significance of regression:

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_k = 0$$
 against

$$H_1: \beta_j \neq 0$$
 for at least one $j \in \{1, \dots, k\}$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.135)

5/12

3.14 Confidence Regions and Prediction Intervals

Provided that we have (at least approximately) normally distributed errors or a large sample size, we obtain the following confidence intervals or regions and prediction intervals:

Confidence Interval for B

A confidence interval for β_i with level $1 - \alpha$ is given by

$$[\hat{\beta}_i - t_{n-p}(1-\alpha/2) \cdot se_i, \hat{\beta}_i + t_{n-p}(1-\alpha/2) \cdot se_i].$$

Confidence Ellipsoid for Subvector β_1

A confidence ellipsoid for $\beta_1 = (\beta_1, ..., \beta_r)'$ with level $1 - \alpha$ is given by

$$\left\{\boldsymbol{\beta}_1: \frac{1}{r}(\widehat{\boldsymbol{\beta}}_1-\boldsymbol{\beta}_1)'\widehat{\operatorname{Cov}(\widehat{\boldsymbol{\beta}}_1)^{-1}}(\widehat{\boldsymbol{\beta}}_1-\boldsymbol{\beta}_1) \leq F_{r,n-p}(1-\alpha)\right\}.$$

A confidence interval for $\mu_0 = E(y_0)$ of a future observation y_0 at location x_0 with level $1 - \alpha$ is given by

$$x_0'\hat{\beta} \pm t_{n-p}(1-\alpha/2)\hat{\sigma}(x_0'(X'X)^{-1}x_0)^{1/2}$$
.

Prediction Interval

A prediction interval for a future observation y_0 at location x_0 with level

$$x'_0 \hat{\beta} \pm t_{n-p} (1 - \alpha/2) \hat{\sigma} (1 + x'_0 (X'X)^{-1} x_0)^{1/2}$$
.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.137)

7 / 12

Test Statistics

Assuming normal errors we obtain under H_0 :

1.
$$F = 1/r (C\hat{\beta} - d)' (\hat{\sigma}^2 C(X'X)^{-1}C')^{-1} (C\hat{\beta} - d) \sim F_{r,n-p}$$

$$2. t_j = \frac{\beta_j}{se_j} \sim t_{n-p}$$

3.
$$F = \frac{1}{2}(\hat{\beta}_1)'\widehat{\text{Cov}(\hat{\beta}_1)^{-1}}(\hat{\beta}_1) \sim F_{r,n-1}$$

3.
$$F = \frac{1}{r}(\hat{\boldsymbol{\beta}}_{1})^{2} \operatorname{Cov}(\hat{\boldsymbol{\beta}}_{1})^{-1}(\hat{\boldsymbol{\beta}}_{1}) \sim \operatorname{F}_{r,n-p}$$
4. $F = \frac{n-p}{k} \frac{R^{2}}{1-R^{2}} \sim \operatorname{F}_{k,n-p}$

Critical Values

Reject H_0 in the case of:

1.
$$F > F_{r,n-p}(1-\alpha)$$

2. $|t| > t_{n-p}(1-\alpha/2)$
3. $F > F_{r,n-p}(1-\alpha)$
4. $F > F_{k,n-p}(1-\alpha)$

$$3 F > F_{max} \cdot (1 - \alpha)$$

2.
$$|t| > t_{n-p}(1-\alpha/2)$$

$$4 \quad F > F_{t-1} \quad (1 - \alpha)$$

The tests are relatively robust against moderate departures from normality. In addition, the tests can be applied for large sample size, even with nonnormal errors.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.135)

6/12

Concrete aggregates data

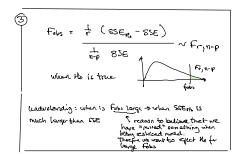
Aggregate:	1	2	3	4	5	
	551	595	639	417	563	
	457	580	615	449	631	
	450	508	511	517	522	
	731	583	573	438	613	
	499	633	648	415	656	
	632	517	677	555	679	
Total	3320	3416	3663	2791	3664	16,854
Mean	553.33	569.33	610.50	465.17	610.67	561.80

Table 13.1 of Walepole, Myers, Myers, Ye: Statistics for Engineers and Scientists - our textbook from the introductory TMA4240/TMA4245 Statistics course.

Today

- Linear hypotheses in regression vs. nested models.
- ▶ The universal F-test for linear hypotheses: two formulas.
- ▶ The two formulas: one easy to use, one easy for proving F-distribution.
- Special cases of the universal F-test.
- ▶ Next time: categorical covariate with effect coding (for interpretation)

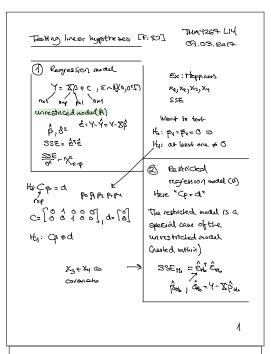
9 / 12



Now: - Why is Fobs ~ Fr, n-p under the - 15 it possible to write Fobs using X, C, d, &, &??

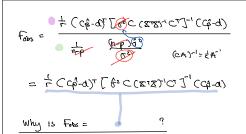
We start with a new version of fors:

2



Why is this e fromp - distribution?

- 3) 0^2 is unknown, but $\frac{2}{3} = \frac{3}{N-p}$ and we more that $\frac{85E}{\sigma^2} = \frac{(n-p)}{\sigma^2} \stackrel{d}{\sigma}^2 \sim \chi^2_{n-p}$
- 4) And of and SSE are independent wown from Part 2.



- 1) For the unrestricted snowl \$=(XTX)'XTY found by minimizing LS(p)=(Y-Xp)T(Y-Xp).
- 2) For the restricted model we minimize LS(P) subject to "CP=d", bey minimize LS(P) + 2 ht (CP-d) = LSR(P) $\hat{\beta}^2 = \hat{\beta} (\text{XTX})^{-1} \text{ CT} \left[c (\text{XTX})^{-1} \text{CT} \right] \text{CC}(P-d)$ see F: P 172-173

4

How can we use fobs for hypothesis bealing?

First: Is the regression significant ? Then may be conjucted for end reduced world (Conjucted PI) or text each $\mu_s=0$?

Solution 1 :

Happines

- a) Fit unrestricted model and get SSE. (Ex: X1+X1-X5+X1)
- b) Fit restricted model and get ESEH (Fix x3+x4)
- C) Fobs = + ASSE SSE N-P
- d) Celculate p-value $_{j}$ $P(Fr_{j}e_{-p}>fow)$, and reject or not to .

6

8)
$$\triangle SSE = \hat{\varepsilon}_{\alpha} \nabla \hat{\varepsilon}_{\alpha} - \hat{c}^{\dagger} \hat{\varepsilon}$$

$$= (Y - X \hat{c}^{\alpha})^{T} (Y - X \hat{c}^{\alpha}) - (Y - X \hat{c}^{\beta})^{T} (Y - X \hat{c}^{\beta})$$

$$= \dots = (C \hat{c}^{\dagger} - d)^{T} (C (Z \nabla \nabla^{\dagger} - C \nabla^{\dagger})^{T} (C \hat{c}^{\dagger} - d)$$

$$f \in PRS-RY$$

- 4) BSE = \$2 (n-p)
- € | + (C4-a) T = 6° C (818) C1] (C4 d)

5

Solution 2

- a) Fit unrestricted award > ph, 32
- b) What is Cend d
- C) Fobs = (calculate this
- d) Calculate the produc.
- => Hands-on: CompEx8. Problem1.

Q: We had to: \$j=0 vs Hi: \$j = 0 and used a t-tent Tj = \$j-0 ~ t.

and not an Fitch. Is it of the two same test as using fibe?

A: to: \$2 = 0 vs tt., \$2 = 0

where C= \(\text{C} = \text{C} \cdots \text{C} = \text{C} \text{C} = \text{C} \text{C} \text{C} = \text{C} \text{C} \text{C} = \text{C} \text{C} = \text{C} \text{C} = \text{C} \text{C} = \text{C} \text{C} \text{C} = \text{C} = \text{C} \text{C} = \text{C} \text{C} = \text{C} \text{C} = \text{C} = \text{C} \text{C} = \text{C} \text{C} = \text{C} \tex

7

$$C_{\beta}^{\beta} = \beta_{3}^{\beta}$$

$$[C(x^{\alpha}x^{\beta})^{-1}C^{\dagger}]^{-1} = G_{jj}^{\beta}$$

$$C_{\beta}^{\alpha} = A^{\beta} T \left[\hat{G}^{2} C(x^{\alpha}x^{\beta})^{-1}C^{\dagger} \right]^{-1} \left(c_{\beta}^{\alpha} - a \right)$$

$$= \frac{(\beta_{3}^{\beta} - o)^{2}}{\hat{G}^{2}} = T_{3}^{2} \in F_{A_{3}} n - \rho$$

$$(t_{\alpha, \beta}^{\alpha})^{2} = \frac{\chi^{2}}{\chi^{2}} \int_{0}^{\chi_{\alpha}^{\alpha}} v F_{A_{3}} v$$
From part A:
$$(F_{j})^{\alpha} = \left(\frac{Z}{\chi^{2}}\right)^{2} = \frac{\chi^{2}}{\chi^{2}} \int_{0}^{\chi_{\alpha}^{\alpha}} v F_{A_{3}} v$$

$$\Rightarrow \text{ add on , thus a an } F^{-1} \text{ text.}$$

Today: Analysis of variance (ANOVA) and analysis of covariance (ANCOVA)

- ► Good news: really nothing new, just linear regression where we have one or more categorical covariates.
- ► Bad news: a bit technical with respect to coding the covariates in the design matrix.
- Bad or good news: also tell the story of ANOVA without linear regression since that is the classical way to do things - so you will be able to recognize that this is a problem that you can solve.
- ► Good news: we are taking one step toward the last topic Part 4: Design of experiments.

TMA4267 Linear Statistical Models V2017 (L15)

Part 3: Hypothesis testing and analysis of variance One- and two-way ANOVA [H:8.1.1]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: March 10, 2017

1 / 47

Rothamsted Experimental Station

- founded in 1843 by John Bennet Lawes on his inherited 16t century estate, Rothamsted Manor,
 - wanted to investigate the impact of inorganic and organic fertilizers on crop yield
 - had founded a fertilizer manufacturing company in 1842
- Lawes appointed the chemist Joseph Henry Gilbert to the directorship of the chemical laboratory
- the two began a series of field experiments to examine the effects of inorganic fertilizers and organic manures on the nutrition and yield of a number of important crops

NAMES OF THE PARTY OF THE PARTY

http://www.stats.uwo.ca/faculty/bellhouse/stat499lecture13.pdf

1 / 47

The Broadbalk Field Trial at Rothamsted

- this was the first field trial started by Lawes and Gilbert
- began in 1843
- purpose was to investigate the relative importance of different plant nutrients (N, P, K, Na, Mg) on grain yield of winter wheat
- · weeds were controlled by hand hoeing and fallowing
 - now some herbicides are used
- · The experiment continues to this day

http://www.stats.uwo.ca/faculty/bellhouse/stat499lecture13.pdf



3 / 47

Concrete aggregates data

Aggregate:	1	2	3	4	5	
	551	595	639	417	563	
	457	580	615	449	631	
	450	508	511	517	522	
	731	583	573	438	613	
	499	633	648	415	656	
	632	517	677	555	679	
Total	3320	3416	3663	2791	3664	16,854
Mean	553.33	569.33	610.50	465.17	610.67	561.80

Table 13.1 of Walepole, Myers, Myers, Ye: Statistics for Engineers and Scientists - our textbook from the introductory TMA4240/TMA4245 Statistics course.

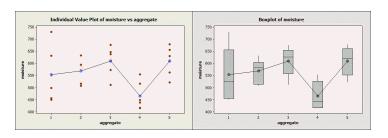
Concrete aggregates example



- Aggregates are inert granular materials such as sand, gravel, or crushed stone that, along with water and portland cement, are an essential ingredient in concrete.
- ▶ For a good concrete mix, aggregates need to be clean, hard, strong particles free of absorbed chemicals or coatings of clay and other fine materials that could cause the deterioration of concrete.
- ▶ We could like to examine 5 different aggregates, and measure the absorption of moisture after 48hrs exposure (to moisture).
- ▶ A total of 6 samples are tested for each aggregate.
- ▶ Research question: Is there a difference between the aggregates with respect to absorption of moisture?

4 / 47

Concrete aggregates example



One-way Analysis of Variance (ANOVA)

Model

$$Y_{ii} = \mu_i + \varepsilon_{ii}$$
 for $i = 1, 2, ..., p$ and $j = 1, 2, ..., n_i$

alternative parameterization

$$Y_{ii} = \mu + \alpha_i + \varepsilon_{ii}$$

The sample sizes for each group, n_i may vary. $\varepsilon_{ij} \sim N(0, \sigma^2)$. Let $n = \sum_{i=1}^p n_i$ be the total number of observations.

Aim: look at parameter estimates and test if there is any difference between the groups.

How can that be done using our linear regression model?

7 / 47

Concrete aggregates data

```
# the same with regression
> options(contrasts=c("contr.sum","contr.sum"))
> obj <-lm(moisture~as.factor(aggregate),data=ds)</pre>
> summary(obj)
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)
                                 12.859 43.688 < 2e-16 ***
as.factor(aggregate)1 -8.467
                                  25.719 -0.329 0.744743
                                 25.719 0.293 0.772005
as.factor(aggregate)2
                      7.533
as.factor(aggregate)3 48.700
                                 25.719 1.894 0.069910 .
as.factor(aggregate)4 -96.633
                                 25.719 -3.757 0.000921 ***
```

Concrete aggregates data

```
# means for each recipe
> means=
    aggregate(ds,by=list(ds$aggregate),FUN=mean)$moisture
> grandmean=mean(ds$moisture)
> grandmean
[1] 561.8
> alphas=means-grandmean
> alphas
[1] -8.466667 7.533333 48.700000 -96.633333 48.866667
```

8 / 47

Concrete aggregates data

Run R code from course lectures tab for model matrix.

9 / 47

Concrete aggregates data (1)

11 / 47

Concrete aggregates data (3)

Concrete aggregates data (2)

12 / 47

One factor: unequal sample sizes

Classical formulation with ANOVA decomposition

$$Y_{ij} - Y_{..} = (Y_{ij} - Y_{i.}) + (Y_{i.} - Y_{..})$$

$$\sum_{i=1}^{p} \sum_{j=1}^{n_i} (Y_{ij} - Y_{..})^2 = \sum_{i=1}^{p} \sum_{j=1}^{n_i} (Y_{ij} - Y_{i.})^2 + \sum_{i=1}^{p} \sum_{j=1}^{n_i} (Y_{i.} - Y_{..})^2$$

$$\sum_{i=1}^{p} \sum_{j=1}^{n_i} (Y_{ij} - Y_{..})^2 = \sum_{i=1}^{p} \sum_{j=1}^{n_i} (Y_{ij} - Y_{i.})^2 + \sum_{i=1}^{p} n_i (Y_{i.} - Y_{..})^2$$

$$SST = SSE + SSA$$

13 / 47

One factor: unequal sample sizes

ANOVA decomposition: what happened to the cross-term?

$$2\sum_{i=1}^{p}\sum_{j=1}^{n_{i}}(Y_{ij}-Y_{i.})(Y_{i.}-Y_{..})=2\sum_{i=1}^{p}(Y_{i.}-Y_{..})\sum_{j=1}^{n_{i}}(Y_{ij}-Y_{i.})=0$$

$$\sum_{j=1}^{n_{i}}(Y_{ij}-Y_{i.})=\sum_{j=1}^{n_{i}}Y_{ij}-\sum_{j=1}^{n_{i}}Y_{i.}=n_{i}Y_{i.}-n_{i}Y_{i.}=0$$

15 / 47

17 / 47

Machine example

- ▶ Response: time (s) spent to assemble a product.
- ► Factor: this is done by four different machines; M_1, M_2, M_3, M_4 .
- Question: Do the machines perform at the same mean rate of speed?

TABLE 13.12 Time, in Seconds, to Assemble Product

Machine	Operator:	1	2	3	4	5	6	Total
1		42.5	39.3	39.6	39.9	42.9	43.6	247.8
2		39.8	40.1	40.5	42.3	42.5	43.1	248.3
3		40.2	40.5	41.3	43.4	44.9	45.1	255.4
4		41.3	42.2	43.5	44.2	45.9	42.3	259.4
Total	1	163.8	162.1	164.9	169.8	176.2	174.1	1010.9

Data from Walepole, Myers, Myers, Ye: "Statistics for Engineers and Scientists", Example 13.6= our TMA4245/40 textbook.

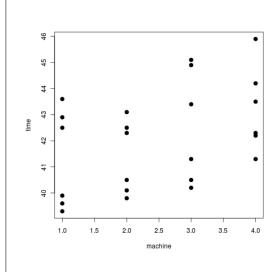
One factor: unequal sample sizes

 $H_0: \mu_1=\mu_2=\cdots=\mu_p=0$ vs. $H_1:$ At least one pair of μ_i different is then tested based on

$$F = \frac{\frac{\text{SSA}}{p-1}}{\frac{\text{SSE}}{n-p}}$$

Where H_0 is rejected if $f_{\text{obs}} > f_{\alpha}$, (p-1), (n-p).

16 / 47



One factor ANOVA

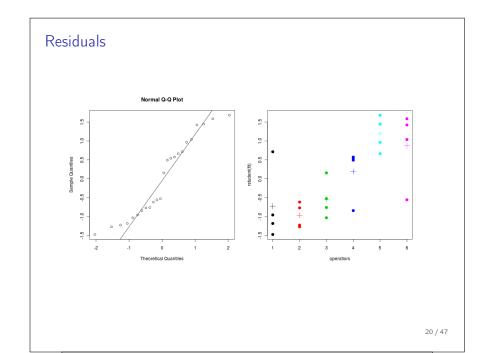
```
> options(contrasts=c("contr.sum","contr.sum"))
> fit <- lm(time~as.factor(machine),data=dsmat)</pre>
> summary(fit)
Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
                     42.1208
(Intercept)
                                0.3706 113.647
                                                  <2e-16 ***
as.factor(machine)1 -0.8208
                                0.6419 -1.279
                                                   0.216
as.factor(machine)2 -0.7375
                                0.6419 -1.149
                                                   0.264
as.factor(machine)3
                     0.4458
                                0.6419
                                         0.695
                                                  0.495
Residual standard error: 1.816 on 20 degrees of freedom
Multiple R-squared: 0.1945, Adjusted R-squared: 0.07372
F-statistic: 1.61 on 3 and 20 DF, p-value: 0.2186
> anova(fit)
Response: time
                   Df Sum Sq Mean Sq F value Pr(>F)
as.factor(machine) 3 15.925 5.3082 1.6101 0.2186
Residuals
                   20 65.935 3.2968
```

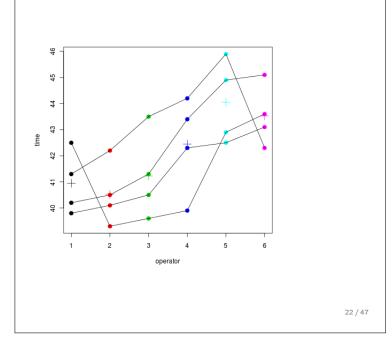
19 / 47

21 / 47

Machine example: operators

- ► The 6 repeated measurements for each machine was in fact made by 6 different operators.
- ► The operation of the machines requires physical dexterity and differences among the operators in the speed with which they operate the machines is anticipated.
- ▶ All of the 6 operators have operated all the 4 machines, and the machines were assigned in random order to the operators= randomized complete block design.
- ▶ By including a blocking factor called Operator, we will reduce the variation in the experiment that is du to random error. Thus, we reduce variation due to *anticipated factors*.
- ▶ By randomizing the order the machines were assigned to the operators we aim to reduce the variation due to *unanticipated* factors.





Model and Sums of squares

Model

$$Y_{ij} = \mu + \alpha_i + \gamma_j + \varepsilon_{ij}$$
 for $i = 1, 2, ..., r$ and $j = 1, 2, ..., s$

Sums of Squares Identity

$$Y_{ij} = Y_{..} + (Y_{i.} - Y_{..}) + (Y_{.j} - Y_{..}) + (Y_{ij} - Y_{i.} - Y_{.j} + Y_{..})$$

$$\sum_{i=1}^{r} \sum_{j=1}^{s} (Y_{ij} - Y_{..})^{2} = s \sum_{i=1}^{r} (Y_{i.} - Y_{..})^{2} + r \sum_{j=1}^{s} (Y_{.j} - Y_{..})^{2}$$

$$+ \sum_{i=1}^{r} \sum_{j=1}^{s} (Y_{ij} - Y_{i.} - Y_{.j} + Y_{..})^{2}$$

$$SST = SSA + SSB + SSE$$

$$r \cdot s - 1 = (r - 1) + (s - 1) + (r - 1)(s - 1)$$

23 / 47

RCBD ANOVA

```
> fit2 <- lm(time~as.factor(machine)+as.factor(operator),
data=dsmat)</pre>
```

> anova(fit2)

Df Sum Sq Mean Sq F value Pr(>F)
as.factor(machine) 3 15.925 5.3082 3.3388 0.047904 *
as.factor(operator) 5 42.087 8.4174 5.2944 0.005328 **
Residuals 15 23.848 1.5899

Effect of factor A:

 $H_0: \alpha_1 = \alpha_2 = \cdots = \alpha_r = 0$ vs. $H_1:$ At least one α_i different from 0

is then tested based on

$$F_1 = \frac{\frac{\text{SSA}}{r-1}}{\frac{\text{SSE}}{(r-1)(s-1)}}$$

Where H_0 is rejected if $f_1 > f_{\alpha}$, (r-1), (r-1)(s-1).

Block effect present?

 $H_0: \gamma_1 = \gamma_2 = \cdots = \gamma_s = 0$ vs. $H_1:$ At least one γ_i different from 0

is then tested based on

$$F_2 = \frac{\frac{\underline{SSB}}{s-1}}{\frac{\underline{SSE}}{(r-1)(s-1)}}$$

Where H_0 is rejected if $f_2 > f_{\alpha}$, (s-1), (r-1)(s-1).

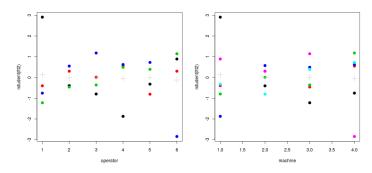
24 / 47

Effect of operator with linear hypotheses

```
fit2 <- lm(time~as.factor(machine)+as.factor(operator),</pre>
data=dsmat)
r=5
C=cbind(rep(0,5), rep(0,5), rep(0,5), rep(0,5), diag(5))
d=matrix(rep(0,r),ncol=1)
betahat=matrix(fit2$coefficients,ncol=1)
X=model.matrix(fit2)
sigma2hat=summary(fit2)$sigma^2
Fobs=(t(C%*\%betahat-d)%*\%solve(C%*\%solve(t(X)%*%X)%*%t(C))
%*%(C%*%betahat-d))/(r*sigma2hat)
> Fobs
         [,1]
[1,] 5.294435
> 1-pf(Fobs,r,n-dim(C)[2])
            [,1]
[1,] 0.005327541
```

25 / 47

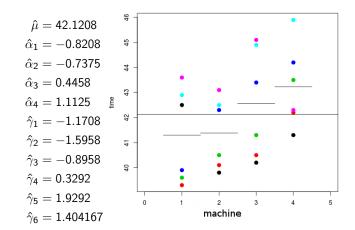
Residuals



27 / 47

29 / 47

Estimates



A second look at the RCBD: additive effects

Previously, randomized complete block design (RCBD) with the machine example:

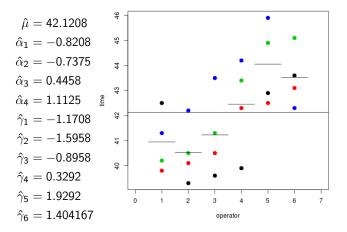
$$Y_{ij} = \mu + \alpha_i + \gamma_j + \varepsilon_{ij}$$

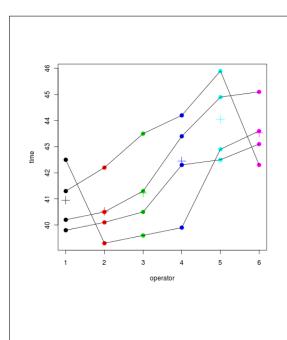
where $\sum_{i=1}^{r} \alpha_i = 0$ and $\sum_{j=0}^{s} \gamma_j = 0$. This is called *additive effects of treatment and blocks*.

- ▶ This means that if we compare two operators there is a constant difference in time to assemble the product,
- or, if we compare machines, these are ranked in the same order of (wrt time) for each operator.

28 / 47

Estimates





31 / 47

Interaction effect?

$$SSE = \sum_{i=1}^{r} \sum_{j=1}^{s} (Y_{ij} - Y_{.i} - Y_{j.} + Y_{..})^{2}$$
$$E(\frac{SSE}{(r-1)(s-1)}) = \sigma^{2} + \frac{\sum_{i=1}^{r} \sum_{j=1}^{s} (\alpha \gamma)_{ij}^{2}}{(s-1)(r-1)}$$

A large value of SSE will either mean that we have an interaction term present, or that σ^2 is large. We can not assess interaction in a RCBD. We need more than one observation for each observation to distinguish between $(\alpha \gamma)_{ii}$ and ε_{ii} .

Interaction effect?

But, it could be interactions present. What if one of the operators really could not manage one of the machines? Model with interaction between treatment and block:

$$Y_{ij} = \mu + \alpha_i + \gamma_j + (\alpha \gamma)_{ij} + \varepsilon_{ij}$$

where $\sum_{i=1}^{r} (\alpha \gamma)_{ij} = \sum_{j=1}^{s} (\alpha \gamma)_{ij} = 0$ (for all i and j) in addition to $\sum_{i=1}^{r} \alpha_i = 0$ and $\sum_{j=1}^{s} \gamma_j = 0$. But, since we only have one observation for each combination of i

and j, we can not separate $(\alpha \gamma)_{ii}$ and ε_{ii} .

32 / 47

Age and memory

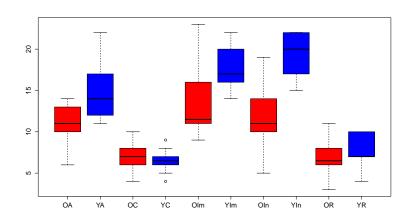
- ▶ Why do older people often seem not to remember things as well as younger people? Do they not pay attention? Do they just not process the material as thoroughly?
- ▶ One theory regarding memory is that verbal material is remembered as a function of the degree to which is was processed when it was initially presented.
- ▶ Eysenck (1974) randomly assigned 50 younger subjects and 50 older (between 55 and 65 years old) to one of five learning groups.
- ▶ After the subjects had gone through a list of 27 items three times they were asked to write down all the words they could remember.

Eysenck study of recall of older and younger subjects under conditions of differential processing, Eysenck (1974) and presented in Howell (1999).

The Age and Memory data set

- ▶ Number of words recalled: After the subjects had gone through the list of 27 items three times they were asked to write down all the words they could remember.
- ► Age: Younger (18-30) and Older (55-65).

35 / 47



Y=younger (blue), O=older (red), A=adjective, C=counting, Im=Imagery, In=intentional, R=rythming.

The Age and Memory data set: Process

- ➤ The Counting group was asked to read through a list of words and count the number of letters in each word. This involved the lowest level of processing.
- ► The Rhyming group was asked to read each word and think of a word that rhymed with it.
- ► The Adjective group was asked to give an adjective that could reasonably be used to modify each word in the list.
- ► The Imagery group was instructed to form vivid images of each word, and this was assumed to require the deepest level of processing.
 - None of these four groups was told they would later be asked to recall the items.
- ► Finally, the Intentional group was asked to memorize the words for later recall.

Data taken from: http://www.statsci.org/data/general/eysenck.html

36 / 47

Eysenck ANOVA

```
> res <- lm(Words~Age*Process)
> summary(res)
Call:
lm(formula = Words ~ Age * Process)
Residuals:
  Min 10 Median
                       30
  -7.0 -1.6 -0.5 2.0 9.6
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)
             11.6100
                         0.2833 40.982 < 2e-16 ***
              -1.5500
                         0.2833 -5.471 3.98e-07 ***
Age1
                         0.5666 2.277 0.025170 *
Process1
              1.2900
              -4.8600
                         0.5666
                                 -8.578 2.60e-13 ***
Process2
                         0.5666
Process3
                                 6.866 8.24e-10 ***
Process4
               4.0400
                         0.5666
                                 7.130 2.43e-10 ***
Age1:Process1 -0.3500
                         0.5666 -0.618 0.538312
Age1:Process2 1.8000
                         0.5666
                                3.177 0.002040 **
Age1:Process3 -0.5500
                         0.5666 -0.971 0.334288
Age1:Process4 -2.1000
                         0.5666 -3.706 0.000363 ***
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Residual standard error: 2.833 on 90 degrees of freedom
Multiple R-squared: 0.7293, Adjusted R-squared: 0.7022
F-statistic: 26.93 on 9 and 90 DF, p-value: < 2.2e-16
```

Eysenck model matrix

39 / 47

Two-way ANOVA questions

There are three main questions that we might ask in two-way \mbox{ANOVA} :

- ▶ Does the response variable depend on Factor A?
- ▶ Does the response variable depend on Factor B?
- ▶ Does the response variable depend on Factor A differently for different values of Factor B, and vice versa?

All of these questions can be answered using hypothesis tests, first we test the interaction.

Model and Sums of Squares

Model:

$$Y_{ijk} = \mu + \alpha_i + \gamma_j + (\alpha \gamma)_{ij} + \varepsilon_{ijk}$$

for $i = 1, 2, ..., r$ and $j = 1, 2, ..., s$ and $k = 1, ..., m$
 $\varepsilon_{ijk} \sim N(0, \sigma^2)$

40 / 47

Effect of interaction AB

$$H_0^A:(\alpha\gamma)_{11}=(\alpha\gamma)_{12}=\cdots=(\alpha\gamma)_{rs}=0$$
 vs. $H_1:$ At least one $(\alpha\gamma)_{ii}$ different from 0

is then tested based on

$$F_3 = \frac{\frac{SS(AB)}{(r-1)(s-1)}}{\frac{SSE}{rs(m-1)}}$$

Where H_0 is rejected if $f_3 > f_{\alpha}$, (r-1)(s-1), rs(m-1).

41 / 47

What do we do after testing for interaction?

- ▶ If the interaction is significant (we reject H_0^{AB}).
 - ▶ Then it is not recommended to test for main effects (that is, the marginal contributions of the two factors A and B separately). This is since the interpretation of the marginal "main effect" is unclear in the presence of interaction. How can we "separate out" the effect of A from the interaction?
 - ▶ Instead, it is usually preferable to examine contrasts in the treatment combinations.
- ▶ If the interaction is not found to be significant (do not reject
 - ▶ We are then interested in the main effects. These can now be tested within the complete model.

43 / 47

Eysenck ANOVA

```
> res <- lm(Words~Age*Process)
> summary(res)
Call:
lm(formula = Words ~ Age * Process)
Residuals:
  Min 10 Median
 -7.0 -1.6 -0.5 2.0 9.6
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 11.6100 0.2833 40.982 < 2e-16 ***
              -1.5500
                         0.2833 -5.471 3.98e-07 ***
Age1
                        0.5666 2.277 0.025170 *
Process1
                         0.5666 -8.578 2.60e-13 ***
Process2
Process3
                        0.5666 6.866 8.24e-10 ***
Process4
               4.0400
                                 7.130 2.43e-10 ***
Age1:Process1 -0.3500
                        0.5666 -0.618 0.538312
Age1:Process2 1.8000
                        0.5666 3.177 0.002040 **
Age1:Process3 -0.5500
                        0.5666 -0.971 0.334288
Age1:Process4 -2.1000
                        0.5666 -3.706 0.000363 ***
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Residual standard error: 2.833 on 90 degrees of freedom
Multiple R-squared: 0.7293, Adjusted R-squared: 0.7022
F-statistic: 26.93 on 9 and 90 DF, p-value: < 2.2e-16
```

45 / 47

Effect of factor A:

is then tested based on

 $H_0^A: \alpha_1 = \alpha_2 = \cdots = \alpha_r = 0$ vs. $H_1:$ At least one α_i different from 0

$$F_1 = \frac{\frac{\text{SSA}}{r-1}}{\frac{\text{SSE}}{rs(m-1)}}$$

Where H_0^A is rejected if $f_1 > f_{\alpha}$, (r-1), rs(m-1).

Effect of factor B:

 $H_0^B: \gamma_1=\gamma_2=\cdots=\gamma_s=0$ vs. $H_1:$ At least one γ_i different from 0

is then tested based on

$$F_2 = \frac{\frac{\text{SSB}}{s-1}}{\frac{\text{SSE}}{rs(m-1)}}$$

Where H_0^B is rejected if $f_2 > f_{\alpha}$, (s-1), sn(m-1).

44 / 47

Eysenck ANOVA

```
> res <- lm(Words~Age*Process)
```

> anova(res)

Analysis of Variance Table

Response: Words

```
Df Sum Sq Mean Sq F value
                                        Pr(>F)
            1 240.25 240.25 29.9356 3.981e-07 ***
Process
            4 1514.94 378.74 47.1911 < 2.2e-16 ***
                       47.58 5.9279 0.0002793 ***
Age:Process 4 190.30
Residuals 90 722.30
                        8.03
```

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

Next: maybe want to compare different combinations of age and process? Then, easiest to just combine the two factors into a new joint factor and skip the intercept.

Summing up

Topic today: the one-way and two-way ANOVA models.

- Classical formulation has focus on comparing sums of squares.
- ▶ We don't have to prove the classical results because we instead fit the ANOVA model using linear regression with effect coding of covariates.
- It is important to plot results and to understand when an interaction term is needed.
- ➤ To test ANOVA hypotheses we use linear hypotheses in the regression where we automatically have theoretical results for F-distributions.
- ▶ We will meet linear regression models with *k* factors with two levels each in Part 4: Design of Experiments (DOE).

47 / 47

Yes, the model can be filled as a linear regression with paremetrs (µ, x1, x2,...) dp).

Po P1 for ...

Previously: Clumary variable coding. Froblem: we went receive number of all because of all because of all occurrents

Now: effect coding.

Impose a restriction on the x's: sun-to-zeroconstraint \(\frac{1}{2-1} \pi_1 = 0 \), in practice only was

\(\text{X1,...}, \pi_{p1} \) and let \(\text{Mp} = -\frac{p^2}{1-1} \pi_1 \), it is given

effect-coding of design metrix

\(\text{Ex:} \quad \text{pose} \pi \text{ps=00} \quad \quad \text{ps=00} \quad \quad \text{ps=00} \quad \quad \text{ps=00} \quad \quad \quad \text{ps=00} \quad \quad \quad \text{ps=00} \quad \quad \quad \quad \text{ps=00} \quad \

Analysis of variance (ANOVA) 711 FORMANT 1708.80.01

Ex: Concrete recipes: 5 recipes to produce concrete, tooked on 6 samples each.

reaswed moisture a: 1sther a difference between the recipe wit moisture.

is the variability between the recipeo large compared vanability within.

1) One-way ANOVA model

Yij =
$$\mu$$
 + α i + ϵ ij

grand mean

 ρ iffere to grand

 ρ i = μ + α i $\Leftrightarrow \alpha$ i = μ i - μ

1

and
$$\beta = \begin{bmatrix} \mu \\ \alpha_{f-1} \\ \alpha_{f-1} \end{bmatrix}$$
, and $\hat{\beta} = (X^{T}X)^{-1}X^{T}Y$.

2) Hypothesis test

Ho: pue pe pe -- e pup vs Hi: at least one different

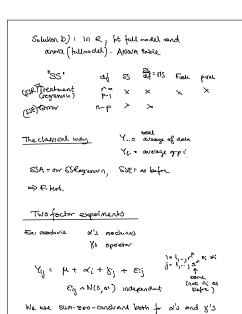
Ho: d1=d2: --- = dp=0 is Hi, at least

a: How can we do this with Unear hypotheses and Fobs from 114?

Solution a: Write as linear hypotheses; $\mu_{j,\alpha_1,\gamma_{i+1}}$ or $\mu_{j,\alpha_1,\gamma_{i+1}}$

from = 4.5, public =
$$P(F_{4}, e_{x} > 45)$$
 = 0.06 675
 \Rightarrow Reject to \Rightarrow difference between tections.

3



Additive effects end interedions

B

Ho: 04: 0e: 1=0 effect of mechanic Ex. smechanic 1=0 of 1=

TMA4267 Linear Statistical Models V2017 (L16)

Part 3: Hypothesis testing and analysis of variance Multiple testing [note]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: March 14, 2017

Today: Multiple testing

► Single hypothesis testing: H_0 and H_1 , test statistic and p-value.

1/33

Today: Multiple testing

- ▶ Single hypothesis testing: H_0 and H_1 , test statistic and p-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.
- ▶ Properties of *p*-values from true and false null hypotheses.

Today: Multiple testing

- ▶ Single hypothesis testing: H_0 and H_1 , test statistic and p-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.

1/33

Today: Multiple testing

- ▶ Single hypothesis testing: H_0 and H_1 , test statistic and p-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.
- ▶ Properties of *p*-values from true and false null hypotheses.
- ► Testing many hypotheses: why?

Today: Multiple testing

- ► Single hypothesis testing: *H*₀ and *H*₁, test statistic and *p*-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.
- ▶ Properties of *p*-values from true and false null hypotheses.
- ► Testing many hypotheses: why?
- Generalizing the type I error from single to multiple hypothesis testing: FWER and FDR.

1/33

1/33

Today: Multiple testing

- ► Single hypothesis testing: *H*₀ and *H*₁, test statistic and *p*-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.
- ▶ Properties of *p*-values from true and false null hypotheses.
- ► Testing many hypotheses: why?
- ► Generalizing the type I error from single to multiple hypothesis testing: FWER and FDR.
- ► Two methods (Bonferroni and Šidák) that control the FWER
- ► Summarizing Part 3 with a quiz.

Today: Multiple testing

- Single hypothesis testing: H₀ and H₁, test statistic and p-value.
- ► Controlling Type I error (false positive findings) by selecting a significance level.
- Properties of p-values from true and false null hypotheses.
- ► Testing many hypotheses: why?
- Generalizing the type I error from single to multiple hypothesis testing: FWER and FDR.
- ► Two methods (Bonferroni and Šidák) that control the FWER

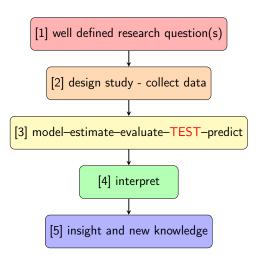
1/33

2/33

Basal metabolic rate and the FTO-gene

- ▶ The gene called FTO is known to be related to obesity
- ► The basal metabolic rate says how many calories you burn when you rest (hvilemetabolisme).
- ▶ Data has been collected for 101 patient from the obesity clinic at St. Olavs Hospital.
- ▶ Research question: is there an association between the variant of the FTO gene of the patient and the basal metabolic rate?
- ► Regression setting, other covariates include age, sex, weight, height, BMI, diet, exercise level, smoking, etc.

The scientific process



3 / 33

Hypothesis testing example (from L13)



- ▶ We draw a random sample of size n = 100 from the blue population and measure systolic blood pressure: $X_1, X_2, ..., X_n$.
- ▶ Test statistic: $\bar{X} \sim \textit{N}(120,1)$ when \textit{H}_0 is true.

Hypothesis testing example (from L13)



▶ We draw a random sample of size n = 100 from the blue population and measure systolic blood pressure: $X_1, X_2, ..., X_n$.

4/33

Hypothesis testing example (from L13)



- ▶ We draw a random sample of size n = 100 from the blue population and measure systolic blood pressure: $X_1, X_2, ..., X_n$.
- ▶ Test statistic: $\bar{X} \sim N(120,1)$ when H_0 is true.
- ▶ We find that $\bar{x} = 122$ mmHg.

4/33

Hypothesis testing example (from L13)



- ▶ We draw a random sample of size n = 100 from the blue population and measure systolic blood pressure: $X_1, X_2, ..., X_n$.
- ▶ Test statistic: $\bar{X} \sim N(120,1)$ when H_0 is true.
- We find that $\bar{x} = 122$ mmHg.
- ▶ Data: n = 100, $\bar{x} = 122$, gives a *p*-verdi=0.02.

4/33

Hypothesis testing example (from L13)

Questions:

- ► How have I calculated this *p*-value? $P(\bar{X} > 122 \mid H_0 \text{ true})$.
- ► How can I interpret this *p*-value? Informally, a *p*-value is the probability under a specified statistical model that a statistical summary of the data would be equal to or more extreme than its observed value.

Hypothesis testing example (from L13)

Questions:

- ► How have I calculated this *p*-value? $P(\bar{X} > 122 \mid H_0 \text{ true}).$
- ► How can I interpret this *p*-value? Informally, a *p*-value is the probability under a specified statistical model that a statistical summary of the data would be equal to or more extreme than its observed value.
- Should I conclude that $\mu > 120$? Yes, if you choose significance level higher than 0.02. But, you should also report a (two-sided) confidence interval for μ : Here [120.04, 123.96].

Single hypothesis testing set-up

	H_0 true	H_0 false
Not reject H ₀	Correct	Type II error
Reject H_0	Type I error	Correct

Two types of errors:

- ► False positives = type I error =miscarriage of justice.
- ► False negatives = type II error= guilty criminal go free.

The significance level of the test is α .

6 / 33

Single hypothesis testing set-up

	H_0 true	H_0 false
Not reject H_0	Correct	Type II error
Reject Ho	Type I error	Correct

Two types of errors:

- ► False positives = type I error =miscarriage of justice.
- ► False negatives = type II error= guilty criminal go free.

The significance level of the test is $\alpha.$

We reject the null hypothesis when the \emph{p} -value is $\emph{below}~\alpha.$

We say that : Type I error is "controlled" at significance level α .

Single hypothesis testing set-up

	H_0 true	H_0 false
Not reject H ₀	Correct	Type II error
Reject H_0	Type I error	Correct

Two types of errors:

- ► False positives = type I error =miscarriage of justice.
- ► False negatives = type II error= guilty criminal go free.

The significance level of the test is α .

We reject the null hypothesis when the *p*-value is *below* α .

6/33

Single hypothesis testing set-up

	H_0 true	H_0 false
Not reject H_0	Correct	Type II error
Reiect <i>H</i> ∩	Type I error	Correct

Two types of errors:

- ► False positives = type I error =miscarriage of justice.
- ► False negatives = type II error= guilty criminal go free.

The significance level of the test is α .

We reject the null hypothesis when the *p*-value is *below* α .

We say that : Type I error is "controlled" at significance level α .

The probability of miscarriage of justice (Type I error) does not exceed α .

Repeating the blood pressure experiment







 \bar{x} =120.9 p-value=0.18 $\bar{x} = 118.9$ *p*-value=0.86

 $\bar{x} = 121.2$ p-value=0.12

7 / 33

7 / 33

Repeating the blood pressure experiment

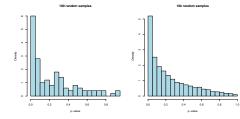




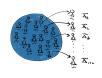


 \bar{x} =120.9 p-value=0.18 $\bar{x} = 118.9$ *p*-value=0.86

 $\bar{x} = 121.2$ p-value=0.12



Repeating the blood pressure experiment

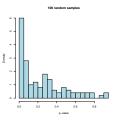






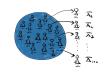
 \bar{x} =120.9 *p*-value=0.18 $\bar{x} = 118.9$ *p*-value=0.86

 $\bar{x} = 121.2$ p-value=0.12



7 / 33

Repeating the blood pressure experiment

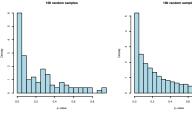


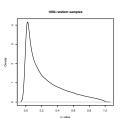




 \bar{x} =120.9 p-value=0.18 $\bar{x} = 118.9$ *p*-value=0.86

 $\bar{x} = 121.2$ p-value=0.12





Histogram - and smoothed histogram of p-values.

More about the *p*-value

► The *p*-value is just a function of the random sample and can be regarded as a random variable.

We had: $P(\bar{X} > \text{observed mean} \mid H_0 \text{ true})$.

8/33

More about the *p*-value

- ► The *p*-value is just a function of the random sample and can be regarded as a random variable.
 - We had: $P(\bar{X} > \text{observed mean} \mid H_0 \text{ true})$.
- ▶ But, isn't the *p*-value a probability? A number?
- ► A random variable (like the *p*-value) has a *probability* distribution.

More about the *p*-value

► The *p*-value is just a function of the random sample and can be regarded as a random variable.

We had: $P(\bar{X} > \text{observed mean} \mid H_0 \text{ true})$.

▶ But, isn't the *p*-value a probability? A number?

8 / 33

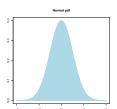
More about the *p*-value

- ► The *p*-value is just a function of the random sample and can be regarded as a random variable.
- We had: $P(\bar{X} > \text{observed mean} \mid H_0 \text{ true})$.
- ▶ But, isn't the *p*-value a probability? A number?
- ► A random variable (like the *p*-value) has a *probability* distribution.
- ▶ What is the distribution of a *p*-value?

8/33

Probability distribution for random variable Y

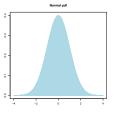
- ► Continuous random variable *Y* (could be the *p*-value).
- ▶ Probability distribution function (pdf): f(y).

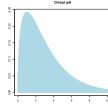


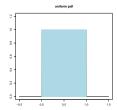
9 / 33

Probability distribution for random variable Y

- ► Continuous random variable *Y* (could be the *p*-value).
- ▶ Probability distribution function (pdf): f(y).

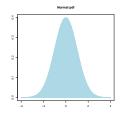


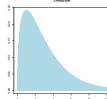




Probability distribution for random variable Y

- ► Continuous random variable *Y* (could be the *p*-value).
- ▶ Probability distribution function (pdf): f(y).





9/33

Distribution of *p*-values for false hypothesis?

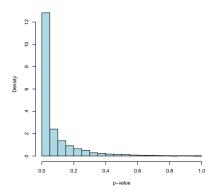
Blood pressure example:

Assume that $\mu=122$ so that H_0 is false, and that we collect a random sample of size 100. What is then the distribution of the p-value?

Distribution of *p*-values for false hypothesis?

Blood pressure example:

Assume that $\mu=122$ so that H_0 is false, and that we collect a random sample of size 100. What is then the distribution of the p-value?

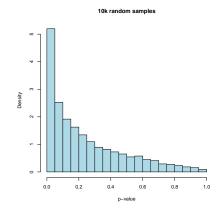


10 / 33

Distribution of *p*-values for false hypothesis?

Blood pressure example:

Assume that $\mu=121$ so that H_0 is false, and that we collect a random sample of size 100. What is then the distribution of the p-value?



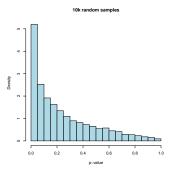
Distribution of *p*-values for false hypothesis?

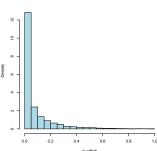
Blood pressure example:

Assume that $\mu=121$ so that H_0 is false, and that we collect a random sample of size 100. What is then the distribution of the p-value?

11/33

False null: $\mu=121$ left, and $\mu=122$ right, when H_0 : $\mu=120$





11 / 33

Distribution of *p*-values for true hypothesis?

Blood pressure example:

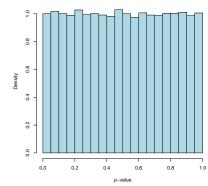
Assume that $\mu=120$ so that H_0 is true, and that we collect a random sample of size 100. What is then the distribution of the p-value?

13 / 33

Distribution of *p*-values for true hypothesis?

Blood pressure example:

Assume that $\mu=120$ so that H_0 is true, and that we collect a random sample of size 100. What is then the distribution of the p-value?



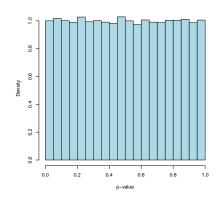
Urban myth: A *p*-value for a true null hypothesis is close to 1. No, all intervals of equal length are equally probable! =uniform distribution

13 / 33

Distribution of *p*-values for true hypothesis?

Blood pressure example:

Assume that $\mu=120$ so that H_0 is true, and that we collect a random sample of size 100. What is then the distribution of the p-value?



13/33

p-values from true null hypothesis is uniformly distributed

Why is this important:

- ▶ so you don't believe the urban myth, and
- ▶ it might be useful to understand plots (pdf or cdf) of p-values, and these are often used for quality control of statistical models.

p-values from true null hypothesis is uniformly distributed

Why is this important:

- ▶ so you don't believe the urban myth, and
- ▶ it might be useful to understand plots (pdf or cdf) of p-values, and these are often used for quality control of statistical models.

Assume that large values of the test statistic T leads to rejection of the null hypothesis, and that a value t of the test statistic T corresponds to a value w of the p-value W. This means that $P(T \ge t) = P(W \le w)$. On the other hand the p-value is $P(W \le w) = P(T \ge t) = w$ when H_0 is true.

This means that $P(W \le w) = w$ when H_0 is true. If W is a continuous random variable taking values from 0 to 1, the the p-value W must be uniformly distributed over the interval from 0 to 1.

This is true when the *p*-value is continuous and exact.

14 / 33

Valid *p*-value

A p-value $p(\mathbf{Y})$ is valid if

$$P(p(Y) \le \alpha) \le \alpha$$

for all α , $0 \le \alpha \le 1$, whenever H_0 is true, that is, if the *p*-value is valid, rejection on the basis of the *p*-value ensures that the probability of type I error does not exceed α .

Exact p-value

If $P(p(\mathbf{Y}) \le \alpha) = \alpha$ for all α , $0 \le \alpha \le 1$, the *p*-value is called an exact *p*-value.

15 / 33

From single to multiple hypothesis testing

In many situations we are not interested in testing only one hypothesis, but instead m hypotheses.

▶ In a regression setting m might be the number of covariates in the regression model, and we would test $H_0: \beta_j = 0$ vs $H_1: \beta_i \neq 0$ for all j = 1, ..., m.

From single to multiple hypothesis testing

In many situations we are not interested in testing only one hypothesis, but instead m hypotheses.

- ▶ In a regression setting m might be the number of covariates in the regression model, and we would test $H_0: \beta_j = 0$ vs $H_1: \beta_i \neq 0$ for all j = 1, ..., m.
- If we have a linear regression with one categorical covariate with k levels, called a one-way analysis of variance model, we might first want to test $H_0: \mu_1 = \mu_2 = \ldots = \mu_k$ against the alternative hypothesis, H_1 , that the means of at least two of the k levels are different from each other. If the null hypothesis is rejected we might want to continue to test which of all possible pairs of the means that are different giving $m = {k \choose 2}$ hypothesis tests, or compare the mean of all levels to a common reference level μ_1 , giving m = k 1 hypothesis tests.

17 / 33

From single to multiple hypothesis testing

In many situations we are not interested in testing only one hypothesis, but instead m hypotheses.

- In a regression setting m might be the number of covariates in the regression model, and we would test $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$ for all $j = 1, \ldots, m$.
- If we have a linear regression with one categorical covariate with k levels, called a one-way analysis of variance model, we might first want to test $H_0: \mu_1 = \mu_2 = \ldots = \mu_k$ against the alternative hypothesis, H_1 , that the means of at least two of the k levels are different from each other. If the null hypothesis is rejected we might want to continue to test which of all possible pairs of the means that are different giving $m = \binom{k}{2}$ hypothesis tests, or compare the mean of all levels to a common reference level μ_1 , giving m = k-1 hypothesis tests.

But, can't we still use cut-off α on the *p*-values to detect significant findings?

From single to multiple hypothesis testing

In many situations we are not interested in testing only one hypothesis, but instead m hypotheses.

- ▶ In a regression setting m might be the number of covariates in the regression model, and we would test $H_0: \beta_j = 0$ vs $H_1: \beta_i \neq 0$ for all j = 1, ..., m.
- If we have a linear regression with one categorical covariate with k levels, called a one-way analysis of variance model, we might first want to test $H_0: \mu_1 = \mu_2 = \ldots = \mu_k$ against the alternative hypothesis, H_1 , that the means of at least two of the k levels are different from each other. If the null hypothesis is rejected we might want to continue to test which of all possible pairs of the means that are different giving $m = {k \choose 2}$ hypothesis tests, or compare the mean of all levels to a common reference level μ_1 , giving m = k 1 hypothesis tests.

17 / 33

Westfall & Young (1993): Multicenter Oat Bran Study

- ▶ At each of ten study centers a control vs treated experiment is performed with 20 subjects per group.
- ▶ It is common to analyze the data for each center separately, as well as to combine over center.
- ► T-statistics are computed for each center as

$$\frac{\bar{y}_T - \bar{y}_C}{\sqrt{(s_T^2 + s_C^2)/20}}$$

with p-values calculated as lower tail probabilities from the t-distribution with 38 degrees of freedom.

FIRST Oat Bran Study

Table 1.2 First Multicenter Oat Bran Study Using Simulated Data

Center	Group	ÿ	5	t-Statistic	p-Vafue (Lower-Tailed)
1	Treated Control	219.1 218.3	7.0 9.8	.30	.616
2	Treated Control	212.6 218.5	11.3 9.8	-1.76	.043*
3	Treated Control	207.5 213.6	11.6 9.9	-1.79	.041*
4	Treated Control	212.5 209.6	10.4 13.5	.76	.774
5	Treated Control	211.9 206.6	8.5 9.1	1.90	.968
6	Treated Control	222.3 222.1	13.4 7.5	.06	.523
7	Treated Control	212.0 211.9	7.4 8.9	.04	.515
8	Treated Control	217.4 215.0	8.6 9.8	.82	.792
9	Treated Control	220.7 217.2	10.7 6.0	1.28	.895
10	Treated Control	222.9 224.4	9.1 11.6	45	.326

^{*} p-value less than .05.

19 / 33

SECOND Oat Bran Study

THE MULTIPLE TESTING PROBLEM

Table 1.3 Second Hypothetical Oat Bran Study

Center	Group	ÿ	S	t-Statistic	p-Value (Lower-Taile
1	Treated Control	214.6 209.3	9.2 8.4	1.90	.968
2	Treated Control	213.9 210.2	8.7 10.5	1.21	.884
3	Treated Control	217.6 216.0	7.6 9.5	.59	.720
4	Treated Control	215.5 211.7	6.2 8.7	1.59	.940
5	Treated Control	211.6 208.1	9.6 8.2	1.24	.889
6	Treated Control	220.1 219.9	8.7 9.6	.069	.527
7	Treated Control	210.3 215.0	5.9 8.7	-2.00	.026*
8	Treated Control	212.2 217.7	9.8 12.5	-1.55	.065
9	Treated Control	217.3 215.0	8.8 9.5	.79	.784
10	Treated Control	212.2 210.5	11.2 9.0	.53	.700

^{*} p-value less than .05.

FIRST Oat Bran Study

- ► Centres 2 and 3 show significant reduction in blood cholestreol for the treatment group.
- ► Centre 5 happens to show a significant increase, but that is not "noticed" since one-sided tests are performed.
- ▶ If the studies were run as uncoordinated trials, it is likely that the two significant studies would be reported and perhaps published in reputable journals.
- ► The eight nonsignificant studies would go to the file drawer and a "true, confirmed" effect would be established for the two sites where significance is found.
- ► The centres with insignificant results may decide to collect fresh data, and analyse only the new data.

20 / 33

Oat bran study: lessons to be learned

- ➤ These are SIMULATED data with equal means of the control and the treatment group, i.e. the truth is that there are no biological effects of the treatment.
- ▶ With simulated data: simple to point to the multiplicity issue as the *cause* for the small *p*-values for some centres.
- ► Real studies: not easy to determine if a seen effect is real or not.

Oat bran study: lessons to be learned

- Real studies: not easy to determine if a seen effect is real or not.
- ▶ At a particular centre showing significance: scientists would believe that the effect is real, because why should the existence of other centres in the study affect the outcome at the given centre?
- ► How should one verify that an unusual event is real or artificial?
- ► The possibility of false positive results is very real, and can lead to serious misinterpretation by analysts: it is human nature to rationalize any dramatic- statistically significant change.

23 / 33

Multiple hypothesis testing set-up

One hypothesis:

	Not reject H_0	Reject H_0
H_0 true	Correct	Type I error
H_0 false	Type II error	Correct

From single to multiple hypothesis testing

Set-up

- ▶ Let us assume that we perform *m* hypothesis tests,
- ▶ giving *m p*-values and then
- ▶ choose a cut-off on the *p*-values at some value α_{loc} (called a local significance level) to decide if we want to reject each null hypothesis.
- ▶ We then reject the null hypotheses where the *p*-value is smaller than α_{loc} , and this leads to rejection of *R* hypotheses.

24 / 33

Multiple hypothesis testing set-up

One hypothesis:

	Not reject H_0	Reject H_0
H ₀ true	Correct	Type I error
H_0 false	Type II error	Correct

m hypotheses:

	Not reject H_0	Reject H_0	Total
H_0 true	U	V	m_0
H_0 false	T	S	$m-m_0$
Total	m-R	R	m

- ► R rejected null hypotheses
- ► *V* false positives (type I errors)
- ► T false negatives (type II errors)

Only m and R are observed. What should we now control?

25 / 33

Overall Type I error control (1)

► In some situation one expects that just a few null hypothesis are false,

26 / 33

26 / 33

Overall Type I error control (1)

- ► In some situation one expects that just a few null hypothesis are false.
- ▶ therefore a *strict* criterion for controlling an overall version of the Type I error is chosen.
- ▶ Family-Wise Error Rate (FWER) is controlled at level α .

 $\mathsf{FWER} = P(V \ge 1) = P(\mathsf{the number of false positives is } \ge 1)$

(remark: V is not observed)

Overall Type I error control (1)

- ► In some situation one expects that just a few null hypothesis are false.
- ▶ therefore a *strict* criterion for controlling an overall version of the Type I error is chosen.

26 / 33

26 / 33

Overall Type I error control (1)

- ► In some situation one expects that just a few null hypothesis are false.
- ► therefore a *strict* criterion for controlling an overall version of the Type I error is chosen.
- ▶ Family-Wise Error Rate (FWER) is controlled at level α .

 $\mathsf{FWER} = P(V \ge 1) = P(\mathsf{the number of false positives is } \ge 1)$

(remark: V is not observed)

▶ The FWER can be controlled by defining a local significance level α_{LOC} for each test and reject the H_0 of that test if the p-value of the test is less than the α_{LOC} .

Basal metabolic rate and the FTO-gene: revisited

- ▶ The gene called FTO is known to be related to obesity
- ► The basal metabolic rate says how many calories you burn when you rest (hvilemetabolisme).
- ▶ Data has been collected for 101 patient from the obesity clinic at St. Olavs Hospital.
- Research question: is there an association between the variant of the FTO gene of the patient and the basal metabolic rate?
- ► Regression setting, other covariates include age, sex, weight, height, BMI, diet, exercise level, smoking, etc.

If we had not only collected data on this one gene, but instead for many (e.g. m=10000) genetic markers positioned along the chromosome, and then wanted to test m hypotheses, we would not expect to find many true associations. This strategy is called a genome-wide association analysis and for this purpose FWER is usually controlled.

27 / 33

Overall Type I error control for GWA data: FWER control

- ▶ GWAS often use $\alpha_{LOC} = 5 \cdot 10^{-8}$.
- ► The most popular method controlling the FWER is the Bonferroni method, which can always be used.

Overall Type I error control for GWA data: FWER control

▶ GWAS often use $\alpha_{LOC} = 5 \cdot 10^{-8}$.

28 / 33

Overall Type I error control for GWA data: FWER control

- ▶ GWAS often use $\alpha_{LOC} = 5 \cdot 10^{-8}$.
- ► The most popular method controlling the FWER is the Bonferroni method, which can always be used.
- ▶ The Bonferroni method might be slightly conservative (too low α_{LOC}), since it is constructed to control FWER for all types of dependency structures between the test statistics for the different hypotheses- including independence.

28 / 33

Overall Type I error control for GWA data: FWER control

- ▶ GWAS often use $\alpha_{LOC} = 5 \cdot 10^{-8}$.
- ► The most popular method controlling the FWER is the Bonferroni method, which can always be used.
- ▶ The Bonferroni method might be slightly conservative (too low α_{LOC}), since it is constructed to control FWER for all types of dependency structures between the test statistics for the different hypotheses- including independence.
- https://arxiv.org/abs/1603.05938: Efficient and powerful familywise error control in genome-wide association studies using generalized linear models, K. K. Halle, Ø. Bakke, S. Djurovic, A. Bye, E. Ryeng, U. Wisløff, O. A. Andreassen, M. Langaas.

28 / 33

Overall Type I error control (2)

- ► For other types of data one expects that many null hypotheses are false.
- ▶ and therefore a less strict criterion for controlling an overall version of the Type I error is chosen.

Overall Type I error control (2)

► For other types of data one expects that many null hypotheses are false

29 / 33

Overall Type I error control (2)

- ► For other types of data one expects that many null hypotheses are false.
- ▶ and therefore a less strict criterion for controlling an overall version of the Type I error is chosen.
- ▶ The False Discovery Rate (FDR) by Benjamini & Hochberg (1995) is controlled at level α .

29 / 33

Overall Type I error control (2)

- ► For other types of data one expects that many null hypotheses are false.
- ▶ and therefore a less strict criterion for controlling an overall version of the Type I error is chosen.
- ▶ The False Discovery Rate (FDR) by Benjamini & Hochberg (1995) is controlled at level α .
- ▶ Informally, the FDR is the expected proportion of Type I errors among the rejected hypotheses.

$$FDR = E(Q)$$
 where by definition

$$Q = \begin{cases} V/R & \text{if } R > 0, \text{ or} \\ 0 & \text{if } R = 0 \end{cases}$$

29 / 33

Overall Type I error control for gene expression data

► Popular algorithm for controlling the FDR: the Benjamini-Hochberg step-up procedure.

Hedenfalk et al (2001) gene expression dataset

Available from library(qvalue) from Bioconductor

- ► The data from the breast cancer gene expression study of Hedenfalk et al. (2001) were obtained and analyzed.
- ▶ A comparison was made between 3,226 genes of two mutation types, BRCA1 (7 arrays) and BRCA2 (8 arrays).
- ► The data included here are p-values, test-statistics, and permutation null test-statistics obtained from a two-sample t-test analysis on a set of 3170 genes, as described in Storey and Tibshirani (2003).

For such gene expression data researchers expect to find may genes that are differently expressed between conditions and therefore the false discovery rate (FDR) is usually controlled. Hedenfalk I et al. (2001). Gene expression profiles in hereditary breast cancer. New England Journal of Medicine, 344: 539-548. Storey JD and Tibshirani R. (2003). Statistical significance for genome-wide studies. Proceedings of the National Academy of Sciences, 100: 9440-9445. http://www.pnas.org/content/100/16/9440.full

30 / 33

Overall Type I error control for gene expression data

- ► Popular algorithm for controlling the FDR: the Benjamini-Hochberg step-up procedure.
- ► Focus on minimal interesting biological effect: is possible that you don't want to test *difference between treatments*=0, but instead ≥ minimal biological interesting effect.

31/33

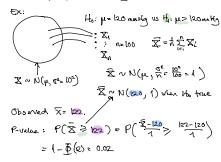
Multiple testing

- ▶ Note from course www-page.
- ► RecEx5.Problem 2.
- ► CompulsoryPart3 Problem 2.
- ▶ This topic is new on the reading list in 2017.
- ▶ It replaces the topics of regularization with the lasso and ridge regression, which will be covered in TMA4268 Statistical Learning.

32 / 33

Multiple hypothesis teolog (note available from Bb) LIG, TUMY267 1403.2017

First single hypothesis testing



Informally: the p-vidue is the probability that our test statistic (Ξ) is observed to be \bar{x} ="122 or a more extreme value" (that is $\bar{x} \ge 122$), when

the truth is that $\mu=120$ so that $\widehat{\mathbb{Z}}\sim \mathcal{N}(120,1)$.

4

Summarizing Part 3

with quiz in Kahoot!

33 / 33

Then we choose if we have enough evidence against the bay looking at the produce.

If the produce is small than what we have observed for more extreme class.) is not very probable when the 10 true. \Rightarrow so for small produce we believe that the must be false and reject the.

Smaller than chosen significence level or 10%, 5%, 1%

⇒ So, the p-value can be seen as a probability?

G: What happens if I collect dota on n=100New persons from the population. We dissure a new X, and will get a new produc.

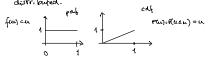
=> The product is a rendom veneble - and it has a probability distribution.

Ex: blood grassive. Easy to sample 100 from $N(\mu, \sigma^2 = 100)$, calculate $\overline{\times}$ and $\overline{\mu}$ value $\overline{\mu}$ hotegram.

h=121 N ← journalum h=121

2

When Ho is true the p-values are uniformly distributed.



⇒ see note for e-code & proof:

This is (usually) non-intuitive to people ... but rather useful to know ...

See note: Obline relid and exact p-value.

3

Lat Richreject Honri, re Pischoo?

Ri=dnot reject Monri, Pischoo?

assume all Hobra

$$P(V>0) = 1 - P(V=0) = 1 - P(\bar{K}_1 \cap \bar{K}_2 \dots \cap \bar{K}_m)$$

need the joint
distribution of the m

Took statistics Type, The

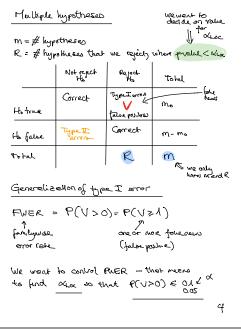
perform a multiple integral. Difficult
to solve. See note on detects.

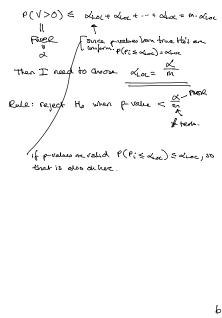
Bonferoni's nethod: Assume all the one true.

$$P(U>0) = P(R_1 \cup R_2 \cup R_3 \cup \dots \cup R_m)$$

 $\leq P(R_1) + P(R_2) + P(R_n)$

$$P(U>0) \in \frac{R_1}{P(\text{regionny Ho not})} + \dots + P(\text{reg. Howo m})}{P(P_i \leq \alpha_{LOC})}$$





TMA4267 Linear statistical models

Part 3: Hypothesis testing and ANOVA

March 14, 2017

Type I errors

What is a commonly used name for the type I errors?

- **B** false positives A true positives
- false negatives **D** true negatives

Happiness

Estimate Std. Error t value Pr(>|t|) (Intercept) -0.072081 0.852543 0.009578 0.005213 1.837 0.0749 money 0.7240 sex -0.149008 0.418525 - 0.356love 1.919279 0.295451 6.496 1.97e-07 0.476079 0.199389 2.388 0.0227 work

For which covariates would we reject the null hypothesis $\beta = 0$ at significance level 1%?

A money

B sex

C love

D work

Linear hypotheses

 H_0 : $C\beta = d$ in a regression model $Y = X\beta + \varepsilon$. *n*=number of observations,

p = number of estimated regression coefficients r=number of linear hypotheses (rank of \boldsymbol{C}).

What is the distribution of F_{obs} $= \frac{1}{r} (\boldsymbol{C} \hat{\boldsymbol{\beta}} - \boldsymbol{d})^T (\hat{\boldsymbol{\sigma}}^2 \boldsymbol{C} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{C}^T)^{-1} (\boldsymbol{C} \hat{\boldsymbol{\beta}} - \boldsymbol{d})?$

 \mathbf{A} $F_{r,n-p}$

- $\mathsf{B} \mathsf{F}_{\mathsf{p},\mathsf{n-r}}$
- \mathbf{C} $N(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1})$ \mathbf{D} $N(0, \sigma^2\boldsymbol{I})$

ANOVA

Which type of covariate coding is used in the oneway ANOVA model with design matrix given as:

- A Continuous
- **B** Effect coding
- C Dummy variable coding
- **D** Categorical

p-value from true null hypothesis

For a continuous test statistic that gives an exact p-value, what is the distribution the pvalue when the null hypothesis is true?

- Normal Chisquared
- Exponential Uniform

ANOVA

Is the interaction term significant at significance level 0.01?

- > res <- lm(Words~Age*Process)</pre>
- > anova(res)

```
Df Sum Sq Mean Sq F value
                                        Pr(>F)
            1 240.25 240.25 29.9356 3.981e-07 ***
Age
Process
            4 1514.94 378.74 47.1911 < 2.2e-16 ***
Age:Process 4 190.30
                       47.58 5.9279 0.0002793 ***
           90 722.30
Residuals
                        8.03
```

A Yes

Not enough information to decide

C No

FWER

V=number of false positives and R=number of rejections. The familywise error rate FWER is

 $\mathbf{A} \quad E(V/R)$

- \mathbf{B} E(V)
- P(V/R > 0.05) P(V > 0)

Bonferroni

 α =level for control of FWER.

 α_{loc} =cut-off on *p*-value

m =number of tests.

What is the Bonferroni rule?

$$\Delta \alpha_{loc} = m\alpha$$

B
$$\alpha_{loc} = \frac{\alpha}{m}$$

$$\alpha_{loc} = \alpha^m$$

D
$$\alpha_{loc} = (1 - \alpha)^{1/m}$$

Correct?

Are you sure you want to read the correct answers? Maybe try first? The answers are explained on the next two slides.

Answers

- 1. C: only love is significant on level 1%, since this is the only p-value below 0.01 (last column).
- 2. B: type I errors are called false positive findings
- 3. A: linear hypotheses with $F_{r,n-p}$ -distributed statistic.
- 4. B: Effect coding is used in ANOVA.

Answers

- 5. A: Interaction term has *p*-value below 0.01.
- 6. D: p-values from true nulls are uniform.
- 7. D: FWER is the probability of one or more false positives.
- 8. B: Bonferroni rule is α/m .