Summary of a Multiple Linear Regression Fit

```
> lm4 <- lm(log(wage) ~ education + I(experience) + I(experience^2) + gender, data = CPS1985)
                                            numerical variables
            lm(formula = log(wage) \sim education + I(experience) + I(experience^2) +
            gender, data = CPS1985)
Vanable Residuals:
                                            30
                         1Q Median
                                                   Max
            -2.24980 -0.29235 0.01609 0.29184 2.13816
                                                           (>|t|) P-values
from t tests
            Coefficients:
            (Intercept) 0.6007445 0.1194927 5.027 6.81e-07 ***
         education 0.0912936 0.0080049 11.405 < 2e-16 ***
I(experience) 0.0360522 0.0054352 6.633 8.14e-11 ***
I(experience^2) -0.0005412 0.0001197 -4.520 7.64e-06 ***
Cients | genderfemale -0.2570355 0.0387066 -6.641 7.77e-11 ***
            Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
            Residual standard error: 0.4442 on 529 degrees of freedom
            Multiple R-squared: 0.2968, Adjusted R-squared: 0.2915 
Rad;
            F-statistic: 55.81 on 4 and 529 DF, p-value: < 2.2e-16
                                                                  p-value from Ftest
              for this example: h-(p+1) = h-p-1
```

```
Explanations on the terms

R<sup>2</sup>, adjusted R<sup>2</sup> [Slides 13, 21-25]

coefficient of determination

in measures in percent

how good the estimated

regression function fits to the

data

100% perfect fit (unrealistic!)

all data points are fitted exactly

0% Y is not explained by the

explanatory variables at all

(also unrealistic!)

Problem: R<sup>2</sup> increases with p (no. of explanatory variables)
```

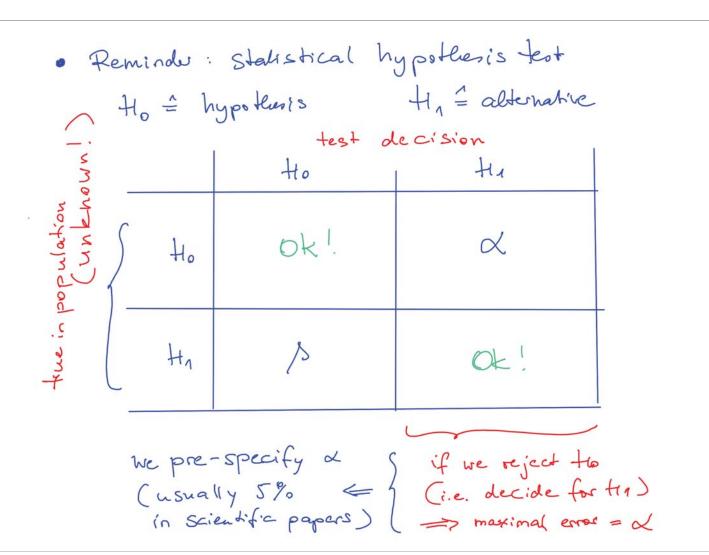
be assume $Y_i = \beta_0 + \beta_1 \times i_1 + ... + \beta_p \times i_p + \epsilon_i$ Use assume $Y_i = \beta_0 + \beta_1 \times i_1 + ... + \beta_p \times i_p + \epsilon_i$ $= \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{k=1}^{n} \sum_{k$

= δ is the estimate of $\delta = \sqrt{\delta^2}$

for the derivation of 3^2 \Rightarrow See propostion of RSS [Slides 33-34]

we have p+1 Coefficients \Rightarrow $6^2 = \frac{RSS}{n-p-1}$ residual degrees of freedom

by potheses tests und mormality
assumption for the E; [see Slides 36ff]



the tests for coefficients [Slides 40-41]

the is
$$\beta_j = 0$$
 vs. $\beta_j = 0$

the test is based on the lest statistic

 $\frac{\hat{\beta}_j}{\hat{\beta}_j} = \frac{\hat{\beta}_j}{\hat{\beta}_j} = \frac{\hat{\beta}_j}{\hat{\beta}_j} \sim t_{n-p-1}$

estimated

Standard deviation

of $\hat{\beta}_j$

the is rejected if

 $\frac{1}{2}t - value = t_{n-p-1, 1-\frac{n}{2}} = t_{n-p-1, 1-\frac{n}{2}}$

Why are we interested in this t test?

remember: $Y = \beta_0 + \beta_1 \times_1 + ... + \beta_p \times_p + \epsilon$ if a coefficient $\beta_j = 0$ then χ_j is not relevant

for the model

if we reject the hypothesis to,

then the variable is relevant

So we are interested, which of the coefficients are significantly (at a) tifterent from zero

10%, 5%, 1% ...

T test (in Summary) [Slides 42-45]

Ho: $m(X) = \beta_0 \leftarrow constant model$ VS. $H_1: m(X) = \beta_0 + \beta_1 X_1 + ... + \beta_p X_p$ regression

estimated model

function

Why are we interested in this?

=> if we reject to then at least one of the variables is relevant
(we do not know, which one)

More agneral: F test to compare 2 nested models
[slides 46-48]

- · definition / use of p values [Slide 41]
- · factor variables [slide 49]
 - => in R factors are categorical variables e.g. gender with values male / female
 - => our linear regression approach (up to now)
 only handles numerical explanatory
 variables
 - => we need to code (recode) factors into numerical Columns of the design matrix of
 - => most easy: dummy variables

for example for gender:

add a Column for females

(say Column k) where

Xik = { 0 otherwise

The R lm function does this automatically!

Why not a column for males in the same way?

* with the first column of I (all equal to 1)

and 2 dummy Columns for females / males:

=> the IC matrix would not have full

rank anymore as

dummy female + dummy male = 1

* so "male" is the reference contegory

* R typically uses the first level of the factor variable as reference

=> check in the data: levels (gender)

=> one can change this using the R function relevel

For factor variables with more than 2 levels;

** again , R uses the first level as reference

** for all other levels / categories then

dummy columns (variables) are added

=> check in the data: levels (occupation)