

Lecture 3: Simple linear regression

BIO144 Data Analysis in Biology

Stephanie Muff & Owen Petchey

University of Zurich

03 January, 2022

First an alert!

Downloading and opening data (CSV) files

Excel will try to be helpful, but in fact it can break things. Do not accept its help, unless you are sure.

2020_data

Home Insert Draw Page Layout Formulas Data Review View

Paste

Calibri (Body) 12

General

Conditional Formatting Format as Table Cell Styles

Insert Delete Format

Sort & Filter Find & Select

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format. Save As...

D3 X ✓ fx Echtes Johanniskraut

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
	Family	Species	FH_number	German_nar	Q1	Q2	Q3	Q4	Q5	Q6	Q7							
1	Equisetaceae	Equisetum a	13	Acker-Schachtel		1	0	1	1	0	0	0						
2	Hypericaceae	Hypericum p	497	Echtes Johar		1	0	0	1	1	1	0						
3	Rosaceae	Potentilla sp	938	Fingerkraut		1	0	0	1	1	1	1						
4	Apiaceae	Chaerophyllum	1409	Kälberkopf		0	0	0	0	0	0	0						
5	Gentianaceae	Gentiana pui	1511	Purpur-Enzian		0	0	1	0	0	0	0						
6	Scrophulariaceae	Euphrasia m	1834	Berg-Augent		1	1	1	0	1	1	1						
7	Scrophulariaceae	Rhinanthus s	1845	Zottiger Klapp		1	1	1	1	1	1	1						
8	Campanulaceae	Campanula t	1900	NA		0	0	0	0	0	0	0						
9	Campanulaceae	Campanula s	1907	Scheuchzeria		1	0	0	0	0	0	0						
10	Asteraceae	Solidago virg	2029	Echte Goldru		0	0	0	0	0	0	0						
11	Asteraceae	Arnica mont	2161	Arnika		0	0	0	0	0	0	0						
12	Asteraceae	Hypochaeris	2263	Einköpfiges H		1	0	0	0	0	0	0						
13	Cyperaceae	Carex sp.	2507	Sedg		0	0	0	1	0	1	0						
14	Poaceae	Festuca sp.	2619	Schwingel		1	0	0	0	0	0	0						
15	Poaceae	Briza media	2639	Mittleres Zit		0	1	0	0	1	0	0						
16	Poaceae	Nardus strict	2724	Borstgras		1	1	0	0	0	1	1						
17	Poaceae	Agrostis sp.	2757	Straussgras		0	0	1	0	0	1	0						
18	Liliaceae	Veratrum all	2838	Gemeiner Gi		0	0	0	0	0	0	0						

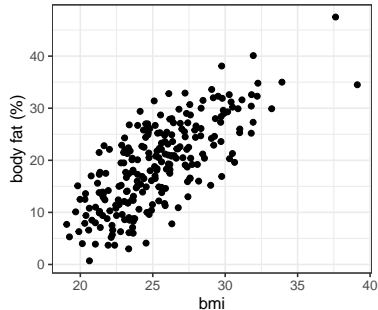
- ▶ Introduction of the linear regression model
- ▶ Parameter estimation
- ▶ Simple model checking
- ▶ Goodness of the model: Correlation and R^2
- ▶ Tests and confidence intervals
- ▶ Confidence and prediction ranges

The lecture material of today is based on the following literature:

- ▶ Chapter 2 of *Lineare Regression*, p.7-20 (Stahel script)

The body fat example

Remember: Aim is to find prognostic factors for body fat, without actually measuring it.
Even simpler question: How good is BMI as a predictor for body fat?



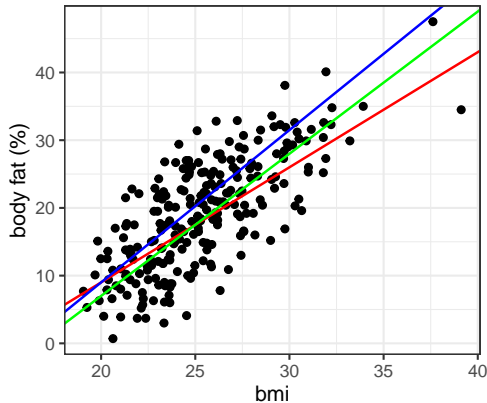
Linear relationship

- ▶ The most simple relationship between an *explanatory variable* (X) and a *target/outcome variable* (Y) is a linear relationship. All points (x_i, y_i) , $i = 1, \dots, n$, on a straight line follow the equation

$$y_i = \beta_0 + \beta_1 x_i .$$

- ▶ Here, β_0 is the **axis intercept** and β_1 the **slope** of the line. β_1 is also denoted as the regression coefficient of X .
- ▶ If $\beta_0 = 0$ the line goes through the origin $(x, y) = (0, 0)$.
- ▶ **Interpretation** of linear dependency: proportional increase in y with increase (decrease) in x .

But which is the “true” or “best” line?



Task: Estimate the regression parameters β_0 and β_1 (by “eye”) and write them down.

It is obvious that

- ▶ the linear relationship does not describe the data perfectly
- ▶ another realization of the data (other 243 males) would lead to a slightly different picture.

⇒ We need a **model** that describes the relationship between BMI and bodyfat.

The simple linear regression model

In the linear regression model the dependent variable Y is related to the independent variable x as

$$Y = \beta_0 + \beta_1 x + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

In this formulation Y is a random variable $Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$ where

$$Y = \underbrace{\text{expected value}}_{E(Y)=\beta_0+\beta_1 x} + \underbrace{\text{random error}}_{\epsilon}.$$

Note:

- ▶ The model for Y given x has **three parameters**: β_0 , β_1 and σ^2 .
- ▶ x is the **independent** / **explanatory** / **regressor** variable.
- ▶ Y is the **dependent** / **outcome** / **response** variable.

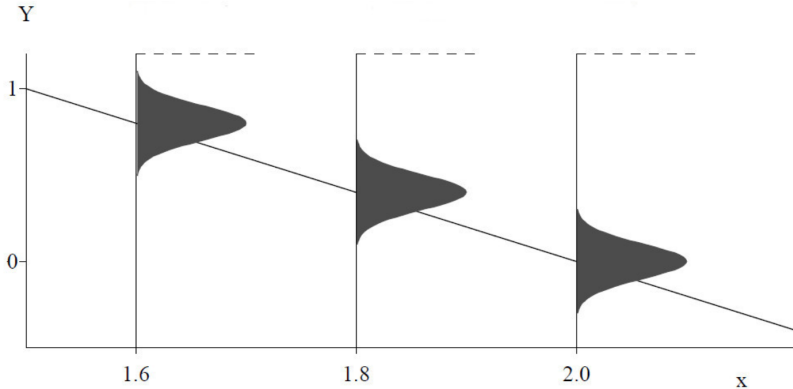
Note

- ▶ The linear model propagates the most simple relationship between two variables. When using it, please always think if such a relationship is meaningful/reasonable/plausible.
- ▶ Always look at the data **before** you start with model fitting.

Visualization of the regression assumptions

The assumptions about the linear regression model lie in the error term

$$\epsilon \sim N(0, \sigma^2) .$$



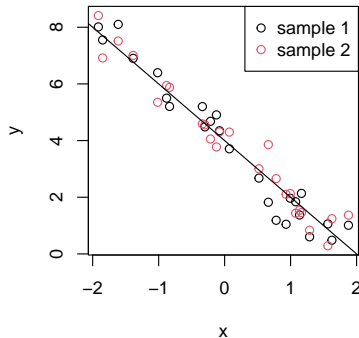
Note: The true regression line goes through $E(Y)$.

Insights from data simulation

(Simulation are *always* a great way to understand statistics!!)

Generate an independent (explanatory) variable **x** and **two** samples of a dependent variable **y** assuming that

$$y_i = 4 - 2x_i + \epsilon_i, \quad \epsilon_i \sim N(0, 0.5^2) .$$



→ Random variation is always present. This leads us to the next question.

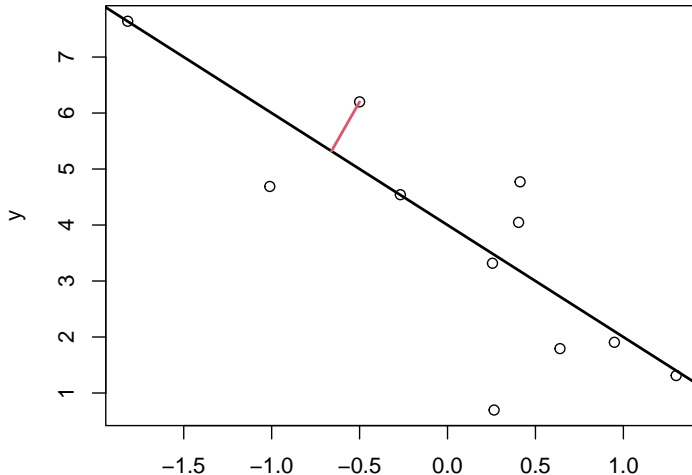
Parameter estimation

In a regression analysis, the task is to estimate the **regression coefficients** β_0, β_1 and the **residual variance** σ^2 for a given set of (x, y) data.

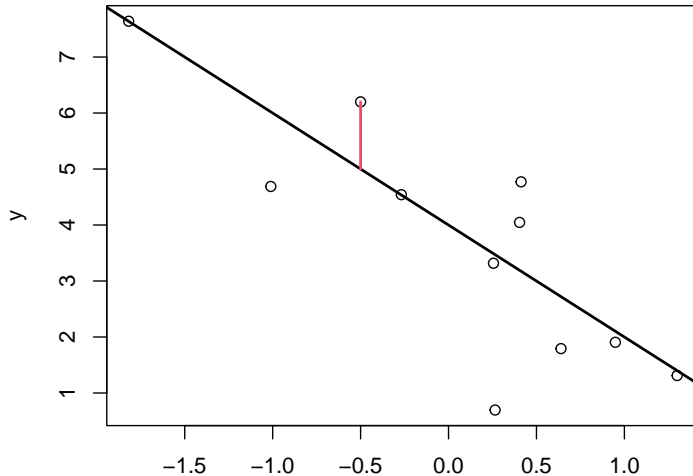
- ▶ **Problem:** For more than two points (x_i, y_i) , $i = 1, \dots, n$, there is generally no perfectly fitting line.
- ▶ **Aim:** We want to estimate the parameters (β_0, β_1) of the best fitting line $Y = \beta_0 + \beta_1 x$.
- ▶ **Idea:** Minimize the deviations between the data points (x_i, y_i) and the regression line.

But how?

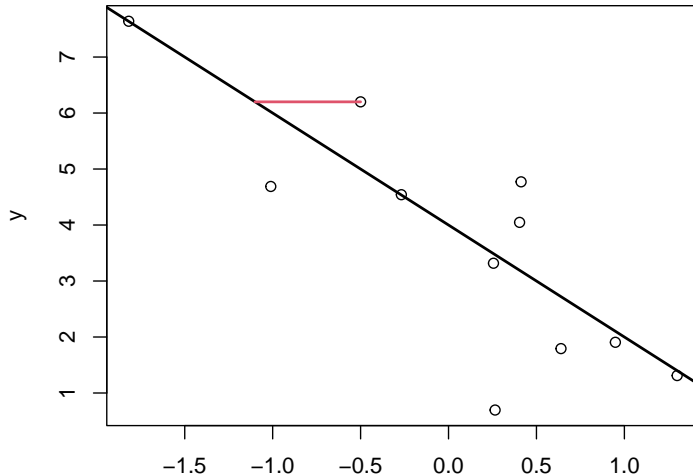
Should we minimize these distances...



Or these?



Or maybe even these?



Least squares

For multiple reasons (theoretical aspects and mathematical convenience), the parameters are estimated using the **least squares** approach. In this, yet something else is minimized:

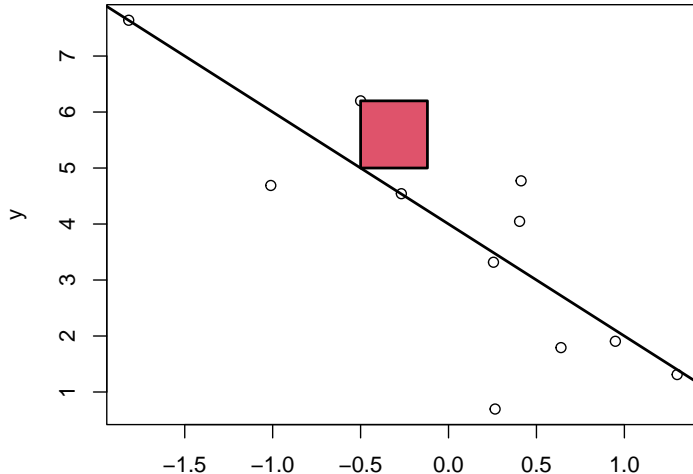
The parameters β_0 and β_1 are estimated such that the sum of **squared vertical distances** (sum of squared residuals)

$$SSE = \sum_{i=1}^n e_i^2, \quad \text{where} \quad e_i = y_i - \underbrace{(\beta_0 + \beta_1 x_i)}_{=\hat{y}_i}$$

is being minimized.

Note: $\hat{y}_i = a + bx_i$ are the **predicted values**.

So we minimize the sum of these areas!



Least squares estimates

For a given sample $(x_i, y_i), i = 1, \dots, n$, with mean values \bar{x} and \bar{y} , the least squares estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are computed as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{cov}(x, y)}{\text{var}(x)},$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Moreover,

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 \quad \text{with residuals } e_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

is an unbiased estimate of the residual variance σ^2 .

(The derivation of the parameters can be looked up in the Stahel script 2.A b. Idea: Minimization through derivating equations and setting them $=0$.)

Do-it-yourself “by hand”

Go to the Shiny gallery and try to “estimate” the correct parameters.

You can do this here:

https://gallery.shinyapps.io/simple_regression/

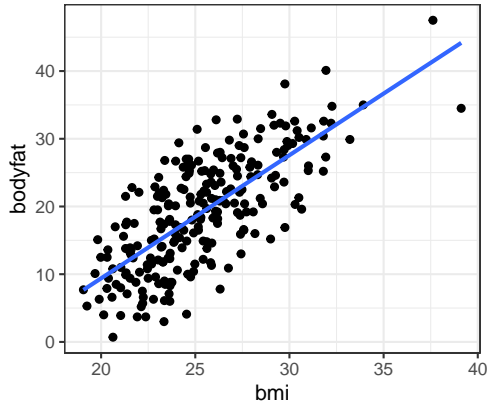
Estimation using R

Let's estimate the regression parameters from the bodyfat example

```
r.bodyfat <- lm(bodyfat ~ bmi, d.bodyfat)
summary(r.bodyfat)
```

```
##
## Call:
## lm(formula = bodyfat ~ bmi, data = d.bodyfat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.5485  -3.5583   0.0785   4.0384  12.7330
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -26.9844     2.7689  -9.746  <2e-16 ***
## bmi          1.8188     0.1083  16.788  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.573 on 241 degrees of freedom
## Multiple R-squared:  0.539, Adjusted R-squared:  0.5371
## F-statistic: 281.8 on 1 and 241 DF,  p-value: < 2.2e-16
```

The resulting line can be added to the scatterplot:



Interpretation: for an increase in the BMI by one index point, we roughly expect a 1.82% percentage increase in bodyfat.

Uncertainty in the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$

Important: $\hat{\beta}_0$ and $\hat{\beta}_1$ are themselves **random variables** and as such contain **uncertainty**!

Let us look again at the regression output, this time only for the coefficients. The second column shows the standard error of the estimate:

```
summary(r.bodyfat)$coef
```

```
##           Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) -26.984368  2.7689004 -9.745518 3.921511e-19
## bmi          1.818778  0.1083411 16.787522 2.063854e-42
```

→ The logical next question is: what is the distribution of the estimates?

Distribution of the estimators for $\hat{\beta}_0$ and $\hat{\beta}_1$

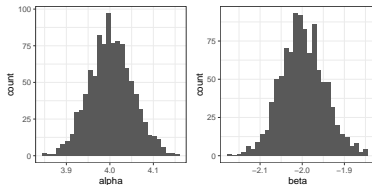
To obtain an idea, we generate data points according to model

$$y_i = 4 - 2x_i + \epsilon_i, \quad \epsilon_i \sim N(0, 0.5^2).$$

In each round, we estimate the parameters and store them:

```
niter <- 1000
pars <- matrix(NA, nrow=niter, ncol=2)
for (ii in 1:niter){
  x <- rnorm(100)
  y <- 4 - 2*x + rnorm(100, 0, sd=0.5)
  pars[ii,] <- lm(y~x)$coef
}
```

Doing it 1000 times, we obtain the following distributions for $\hat{\beta}_0$ and $\hat{\beta}_1$:



This looks suspiciously normal!

In fact, from theory it is known that

$$\hat{\beta}_1 \sim N(\beta_1, \sigma^{(\beta_1)^2}) \quad \text{and} \quad \hat{\beta}_0 \sim N(\beta_0, \sigma^{(\beta_0)^2})$$

For formulas of the standard deviations $\sigma^{(\beta_1)^2}$ and $\sigma^{(\beta_0)^2}$, please consult Stahel 2.2.h.

To remember:

- ▶ $\hat{\beta}_0$ and $\hat{\beta}_1$ are **unbiased estimators** of β_0 and β_1 .
- ▶ the parameters estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are **normally distributed**.
- ▶ the formulas for the variances depend on the residual variance σ^2 , the sample size n and the variability of X ($SSQ^{(X)(*)}$).

(*)

$$SSQ^{(X)} = \sum_{i=1}^n (x_i - \bar{x})^2$$

Are the modelling assumptions met?

In practice, it is advisable to check if all our **modelling assumptions are met**.

→ Otherwise we might draw invalid conclusions from the results.

Remember: Our assumption is that $\epsilon_i \sim N(0, \sigma^2)$. This implies

- a) The expected value of ϵ_i is 0: $E(\epsilon_i) = 0$.
- b) All ϵ_i have the same variance: $Var(\epsilon_i) = \sigma^2$.
- c) All ϵ_i are normally distributed.

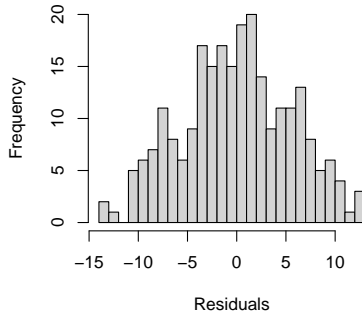
In addition, it is assumed that

- d) $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are independent.

Note: We do not actually observe ϵ_i , but only the residuals e_i . Let us introduce two simple graphical model checking tools for our residuals e_i .

Model checking tool II: Histogram of residuals

Look at the histogram of the residuals:



The normal distribution assumption (c) seems ok as well.

How good is the regression model?

This is, per se, a difficult question. . . .

One often considered index is the **coefficient of determination (Bestimmtheitsmass)** R^2 . Let us again look at the regression output from the bodyfat example:

```
summary(r.bodyfat)$r.squared
```

```
## [1] 0.5390391
```

Compare this to the squared correlation between the two variables:

```
cor(d.bodyfat$bodyfat,d.bodyfat$bmi)^2
```

```
## [1] 0.5390391
```

→ In simple linear regression, R^2 is the squared correlation between the independent and the dependent variable.

- ▶ R^2 indicates the proportion of variability of the response variable y that is **explained by the ensemble of all covariates**.
- ▶ Its value lies between 0 and 1.

The **larger** R^2

- ⇒ the **more** variability of y is captured (“explained”) by the covariate
- ⇒ the **"better"** is the model.

(However, it's a bit more complicated, see later in the course. . .)

Testing and Confidence Intervals

After the regression parameters and their uncertainties have been estimated, there are typically two fundamental questions:

1. **“Are the parameters compatible with some specific value?”**

Typically, the question is whether the slope β_1 might be 0 or not, that is: “Is there an effect of the covariate x or not?”

⇒ This leads to a **statistical test**.

2. **“Which values of the parameters are compatible with the data?”**

⇒ This leads us to determine **confidence intervals**.

Let's first go back to the output from the bodyfat example:

```
summary(r.bodyfat)$coef
```

```
##              Estimate Std. Error  t value    Pr(>|t|)
## (Intercept) -26.984368  2.7689004 -9.745518 3.921511e-19
## bmi          1.818778  0.1083411 16.787522 2.063854e-42
```

Besides the estimate and the standard error (which we discussed before), there is a **t value** and a probability **Pr(>|t|)** that we need to understand.

How do these things help us to answer the two questions above?

Testing the effect of a covariate

Remember: in a statistical test you first need to specify the *null hypothesis*. Here, typically, the null hypothesis is

$$H_0 : \beta_1 = 0 .$$

In words: H_0 = "no association"

(Included in H_0 is the assumption that the data follow the simple linear regression model!)

Here, the *alternative hypothesis* is given by

$$H_A : \beta_1 \neq 0$$

Remember: To carry out a statistical test, we need a *test statistic*.

What is a test statistic?

→ It is some type of **summary statistic** that follows a known distribution under H_0 .
For our purpose, we use the so-called **T -statistic**

$$T = \frac{\hat{\beta}_1 - C}{se(\beta_1)} . \quad (1)$$

Again: typically, $C = 0$, so the formula simplifies to $T = \frac{\hat{\beta}_1}{se(\beta_1)}$.

Under H_0 , T has a t -distribution with $n - 2$ degrees of freedom (n = number of data points).

(You should try to recall the t -distribution. Check Mat183, keyword: t -test.)

So let's again go back to the bodyfat regression output:

```
summary(r.bodyfat)$coef
```

##		Estimate	Std. Error	t value	Pr(> t)
##	(Intercept)	-26.984368	2.7689004	-9.745518	3.921511e-19
##	bmi	1.818778	0.1083411	16.787522	2.063854e-42

Task:

→ Please use equation (1) to find out how the first three columns (Estimate, Std. Error and t value) are related! Check by a calculation...

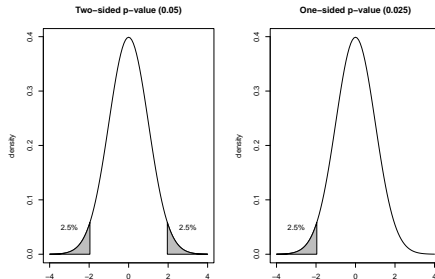
Note: The last column contains the **p-value** of the test $\beta_1 = 0$.

Recap: Formal definition of the p -value

The **formal definition of p -value** is the probability to observe a data summary (e.g., an average) that is at least as extreme as the one observed, given that the Null Hypothesis is correct.

Example (normal distribution): Assume the observed test-statistic leads to a z -value = -1.96

$$\Rightarrow \Pr(|z| \geq 1.96) = 0.05 \text{ and } \Pr(z \leq -1.96) = 0.025.$$



The regression output on slide 33 indicates that the p -value for BMI is very small ($p < 0.0001$).

Conclusion: there is **very strong evidence** that the BMI is associated with bodyfat, because p is extremely small (thus it is very unlikely that such a slope $\hat{\beta}_1$ would be seen if there was no effect of BMI on body fat).

This basically answers question 1 from slide 29.

A cautionary note on the use of p -values

Maybe you have seen that in statistical testing, often the criterion $p \leq 0.05$ is used to test whether H_0 should be rejected. This is often done in a black-or-white manner.

However, we will put a lot of attention to a more reasonable and cautionary interpretation of p -values in this course!

Confidence intervals of regression parameters

Question 2 from slide 29:

Which values of the parameters are compatible with the data?

To answer this question, we can determine the confidence intervals of the regression parameters.

Facts we know about $\hat{\beta}_1$

- ▶ $\hat{\beta}_1$ is estimated with a standard error of $\sigma^{(\beta_1)}$
- ▶ The distribution of $\hat{\beta}_1$ is normal, namely $\hat{\beta}_1 \sim N(\beta_1, \sigma^{(\beta_1)^2})$.
- ▶ However, since we need to estimate $\sigma^{(\beta_1)}$ from the data (the standard error), we have a t -distribution.

Doing some calculations (similar to those in chapter 8.2.2 of Mat183 script) leads us to the 95% confidence interval

$$[\hat{\beta}_1 - c \cdot \hat{\sigma}^{(\beta_1)}; \hat{\beta}_1 + c \cdot \hat{\sigma}^{(\beta_1)}] ,$$

where c is the 97.5% quantile of the t -distribution with $n - 2$ degrees of freedom.

Doing this for the bodfat example “by hand” is not hard. We have 241 degrees of freedom:

```
coefs <- summary(r.bodyfat)$coef
beta <- coefs[2,1]
sdbeta <- coefs[2,2]
beta + c(-1,1) * qt(0.975,241) * sdbeta
```

```
## [1] 1.605362 2.032195
```

Even easier: directly ask R to give you the CIs.

```
confint(r.bodyfat, level=c(0.95))
```

```
##                2.5 %      97.5 %
## (Intercept) -32.438703 -21.530032
## bmi         1.605362   2.032195
```

In summary,

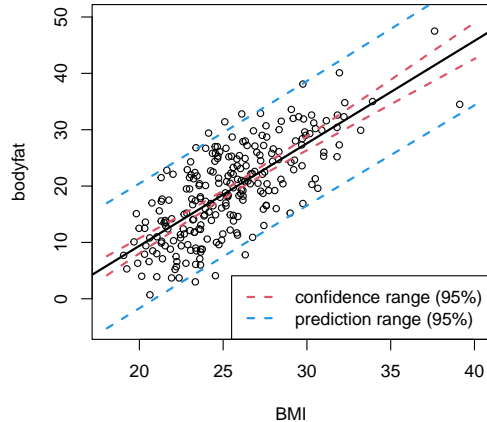
	Coefficient	95%-confidence interval	<i>p</i> -value
Intercept	-26.98	from -32.44 to -21.53	< 0.0001
bmi	1.82	from 1.61 to 2.03	< 0.0001

Interpretation: for an increase in the bmi by one index point, roughly 1.82% percentage points more bodyfat are expected, and all true values for β_1 between 1.61 and 2.03 are compatible with the observed data.

Confidence and Prediction Ranges

- ▶ Remember: When another sample from the same population was taken, the regression line would look slightly different.
- ▶ There are two questions to be asked:
 1. Which other regression lines are compatible with the observed data?
⇒ This leads to the **confidence range**.
 2. Where do future observations with a given x coordinate lie?
⇒ This leads to the **prediction range**.

Bodyfat example



Note: The prediction range is much broader than the confidence range.

Calculation of the confidence range

Given a fixed value of x , say x_0 . The question is:

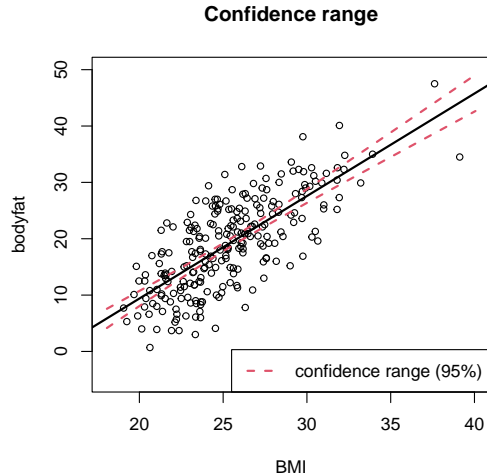
Where does $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ lie with a certain confidence (i.e., 95%)?

This question is not trivial, because both $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimates from the data and contain uncertainty.

The details of the calculation are given in Stahel 2.4b.

Plotting the confidence interval around all \hat{y}_0 values one obtains the **confidence range** or **confidence band for the expected values** of y .

Note: For the confidence range, only the uncertainty in the estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ matters.



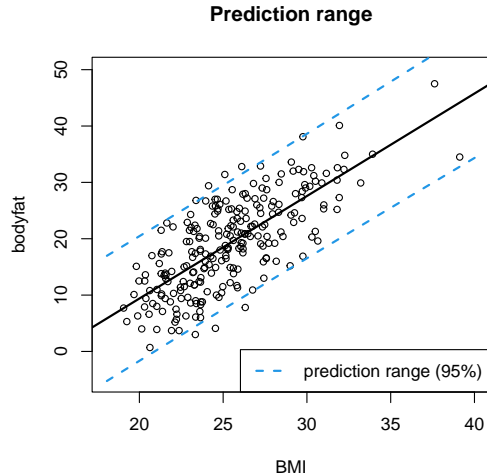
Calculations of the prediction range

Given a fixed value of x , say x_0 . The question is:

Where does a **future observation** lie with a certain confidence (i.e., 95%)?

To answer this question, we have to **consider not only the uncertainty in the predicted value** $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$, but also the **error in the equation** $\epsilon_i \sim N(0, \sigma^2)$.

This is the reason why the **{prediction range is always wider than the confidence range}**.



Tasks until the next practical (Thu/Fri)

The idea of the course is that as a preparation for the practical part you will do the following:

- ▶ Consolidate your understanding of today's lecture.
- ▶ Go to openedX and do all the “Before class (BC)” tasks.

→ **The same procedure applies to all course weeks.**