| V11 Multi-variate analysis | |
|---|---|
| Program for today: | |
| - what is a confounding variable? | |
| - Ways how to avoid confounding effects. | |
| - Staphylococcus aureus Africa project – analysis for confounding variables | |
| - Which slides are relevant for final oral exam? | |
| - Example questions. | |
| | |
| | |
| | |
| V11 Processing of Biological Data SS 2020 | 1 |

In today's lecture, we will discuss the meaning of the term "confounding effect", how these can be detected, and how these can be avoided by choosing an appropriate study design.



A confounding variable is closely related to both the independent and dependent variables in a study.

An independent variable represents the suppose cause, while the dependent variable is the supposed effect.

A confounding variable is a third variable that influences both the independent and dependent variables.

Confounding variables

Confounding variables can introduce biases.

Let's say you test 200 volunteers (100 men and 100 women) and you find that lack of exercise leads to weight gain.

One problem with this experiment is that it lacks any control variables, e.g. the use of placebos, or random assignment to groups.

One confounding variable could be how much people eat.

It is also possible that men eat more than women; this could also make **sex** a confounding variable.

Another possibility is **age**. E.g. if all of the women in the study were middle-aged, and all of the men were aged 16, **age** may have a direct effect on weight gain. That would make age a confounding variable.

V11 Processing of Biological Data SS 2020 https://www.statisticshowto.com/experimental-design/confounding-variable/

3

This is just an example ...

| Case-control studies There are various study designs that try to actively exclude or control confounding variables: |
|--|
| Case-control studies assign confounders equally to both groups, cases and controls. |
| E.g. if somebody wanted to study the cause of myocardial infarct and thinks that the age is a probable confounding variable, each 67-year-old infarct patient will be matched with a healthy 67-year-old "control" person. |
| In case-control studies, matched variables most often are the age and sex. |
| Drawback: Case-control studies are feasible only when it is easy to find controls. |
| |
| https://en.wikipedia.org/wiki/Confounding |
| V11 Processing of Biological Data SS 2020 4 |

Suppose a case-control study attempts to find the cause of a given disease in a person who is 1) 45 years old, 2) African-American, 3) from Alaska, 4) an avid football player, 5) vegetarian, and 6) working in education. A theoretically perfect control would be a person who, in addition to not having the disease being investigated, matches all these six characteristics and also has no other diseases that the patient also does not have. Finding such a control would be an enormous task.

| Cohort studies |
|--|
| A degree of matching is often realized by only admitting certain age groups or a certain sex into the study population. |
| This creates a cohort of people who share similar characteristics and thus all cohorts are comparable in regard to the possible confounding variable. |
| E.g. if age and sex are thought to be confounders, only 40 to 50 years old males would be involved in a cohort study that would assess the myocardial infarct risk in cohorts that either are physically active or inactive. |
| <u>Drawback</u> : In cohort studies, the over-exclusion of input data may lead researchers to define too narrowly the set of similarly situated persons for whom they claim the study to be useful. Then, other persons to whom the causal relationship does in fact apply may lose the opportunity to benefit from the study's recommendations. |
| Similarly, "over-stratification" of input data within a study may reduce the sample size in a given stratum to the point where generalizations drawn by observing the members of that stratum alone are not statistically significant. |
| V11 Processing of Biological Data SS 2020 |

https://en.wikipedia.org/wiki/Confounding

Wikipedia explains: In statistics, stratified sampling is a method of sampling from a population which can be partitioned into subpopulations.

| Double blinding |
|---|
| Double blinding conceals (hides) from the trial population and the observers the experiment group membership of the participants. |
| By preventing the participants from knowing if they are receiving treatment or not, the placebo effect should be the same for the control and treatment groups. |
| By preventing the observers from knowing of their membership, there should be no bias from researchers treating the groups differently or from interpreting the outcomes differently. |
| |
| https://en.wikipedia.org/wiki/Confounding |
| V11 Processing of Biological Data SS 2020 6 |

In a **double-blind** study (dt. Blindstudie), participants and experimenters do not know who is receiving a particular treatment.

"Double" specifies that both the participants AND the staff conducting the study (e.g. medical doctors) do not know to which group the participants belong.

Randomized controlled trial

This is a method where the **study population is divided randomly** in order to mitigate the chances of self-selection by participants or bias by the study designers.

Before the experiment begins, the testers will assign the members of the participant pool to their groups (control or intervention) using a randomization process such as the use of a random number generator.

E.g. in a study on the effects of exercise, the conclusions would be less valid if participants were given a choice if they wanted to belong to the control group which would not exercise or the intervention group which would be willing to take part in an exercise program.

The study would then capture other variables besides exercise, such as preexperiment health levels and motivation to adopt healthy activities.

From the observer's side, the experimenter may choose candidates who are more likely to show the results the study wants to see or may interpret subjective results (more energetic, positive attitude) in a way favorable to their desires.

```
V11 Processing of Biological Data SS 2020
https://en.wikipedia.org/wiki/Confounding
```

https://www.medicinenet.com/script/main/art.asp?articlekey=39532:

Randomized controlled trial: (RCT) A study in which people are allocated at random (by chance alone) to receive one of several clinical interventions. One of these interventions is the standard of comparison or control. The control may be a standard practice, a placebo ("sugar pill"), or no intervention at all. Someone who takes part in a randomized controlled trial (RCT) is called a participant or subject. RCTs seek to measure and compare the outcomes after the participants receive the interventions. Because the outcomes are measured, RCTs are quantitative studies.

In sum, RCTs are quantitative, comparative, controlled experiments in which investigators study two or more interventions in a series of individuals who receive them in random order. The RCT is one of the simplest and most powerful tools in clinical research.

7

| Stratification |
|--|
| As in the example just mentioned, physical activity is thought to be a behavior that protects from myocardial infarct; and age is assumed to be a possible confounder. |
| The data sampled is then stratified by age group – this means that the association between activity and infarct would be analyzed per each age group. |
| If the different age groups (or age strata) yield much different risk ratios, age must be viewed as a confounding variable. |
| There exist statistical tools that account for stratification of data sets. |
| |
| |
| |
| V11 Processing of Biological Data SS 2020 |
| https://en.wikipedia.org/wiki/Confounding |

https://asq.org/quality-resources/stratification

Stratification is defined as the act of sorting data, people, and objects into distinct groups or layers. It is a technique used in combination with other data analysis tools. When data from a variety of sources or categories have been lumped together, the meaning of the data can be difficult to see. This data collection and analysis technique separates the data so that patterns can be seen.

Here are examples of different sources that might require data to be stratified:

Equipment, Shifts, Departments, Materials, Suppliers, Day of the week, Time of day, Products

STRATIFICATION PROCEDURE

- Before collecting data, consider which information about the sources of the data might have an effect on the results. Set up the data collection so that you collect that information as well.

- When plotting or graphing the collected data on a scatter diagram, control chart, histogram, or other analysis tool, use different marks or colors to distinguish data from various sources. Data that are distinguished in this way are said to be "stratified."

- Analyze the subsets of stratified data separately.



Now, we will look back at the example that we discusses in the first 2 lectures, the collection of bacterial Staphylococcus aureus samples in three African countries and in three German university hospitals.



This is an overview to which clonal complexes the samples belong.

| R | eview (V1) <i>:</i> A | ctiv | /iti | ity | of | in | di | vi | du | al | pr | ok | be | s f | or | C | C | 5 | |
|---------------------------|---|-----------|--------|-------------|------------|---------|---------|-------------|---------------|-------------|------------|------------|-----------|-------------|------------|---------------|----------|------------|-------|
| X | and the second se | - | _ | - | _ | - | S2 Exce | el Book - I | Microsoft | Excel | | | - | | | _ | - | - | S. |
| Datei Start Einfi | igen Seitenlayout Formein Daten Überpr | üfen Ansi | nt A | 8811 FineRe | ader 12 | Acrobat | | | | | | | | | | | | | |
| Ausschneiden | Calibri - 11 - A^ = = - | 20 | Tellen | umbruch | | Standar | d | | | 8 | 112 | Standa | ard | Gut | | | | 3 | Th |
| Kopieren * | | | | | | | - | | Cardio Cardio | | Taballa | Ale de | | Cable | - | | | | لتتور |
| * Format übertr | agen F K U * ⊡ * 💁 * 🗛 * 📰 🕾 🕸 | the she | Verbin | nden und ze | ntrieren * | | % 000 | 760 4,0 | Formatie | trung * for | matieren v | Neutra | 91 | Schie | cnt | | intugen | Loschen Po | * |
| Zwischenablage | G Schriftart G | Ausri | chtung | | 6 | | Zahl | - 6 | | | | Formatvor | rlägen | | | | _ | Zellen | |
| BSO | | | | | | | | | | | | | | | | | | | |
| A | 8 | С | D | E | F | G | н | T | J. | K | L | м | N | 0 | Р | Q | R | 5 | T |
| 1 | A | | | Address 1 | | | | | Gene | e distribut | ion in Afr | ican vs Ge | rman S. o | oureus isol | lates of t | he 10 pred | Jominant | CCs | 0 |
| 3 | Numbers (n) | 600 | 600 | 109 | 57 | 51 | 88 | 105 | 25 | 48 | 53 | 44 | 51 | 11 | 75 | African 83 | 2 | African G | 48 |
| 4 | Cional complex (CC) | all | CCs | CC1 | 15 | CC45 | | CC1 | 21 | CCE | в | CCS | | CC3 | 10 | CC1 | 52 | CC7 | 7 |
| 5 SPECIES MARKERS | rmD1Saureus. | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 6 | gapA | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 8 | coA | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 9 | nucl | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 10 | spa | 100% | 100% | 100% | 100% | 98% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 12 REGULATORY GENES | sbi | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 13 | saeS | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 14 | vraS | 100% | 100% | 100% | 100% | 100% | 100% | 99% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 15 | agritotal. | 35% | 55% | 0% | 0% | 41% | 99% | 0% | 0% | 100% | 100% | 0% | 0% | 0% | 0% | 100% | 100% | 100% | 10 |
| 10 | agr8.1 | 54% | 60% | 0% | 25 | 100% | 100% | 84% | 92% | 100% | 100% | 0% | 0% | 0% | 0% | 100% | 100% | 100% | - 10 |
| 18 | agr0.1 | 35% | 55% | ON | 0% | 41% | 99% | 0% | 0% | 100% | 100% | 0% | 0% | 0% | 0% | 100% | 100% | 100% | 10 |
| 19 | agril_total. | 27% | 25% | 99% | 100% | 0% | 0% | 0% | 0% | 0% | 0% | 100% | 100% | 0% | 0% | 0% | 0% | 0% | |
| 20 | agrB.II | 27% | 25% | 99% | 100% | 0% | 0% | 0% | 0% | 0% | 0% | 98% | 100% | 0% | 0% | 0% | 0% | 0% | |
| 21 | agrCII | 27% | 25% | 99% | 100% | 0% | 0% | 0% | 0% | 0% | 0% | 100% | 100% | 0% | 0% | 0% | 0% | 0% | |
| 23 | agrill_total | 16% | 14% | 05 | 05 | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 100% | 100% | 0% | 05 | 0% | |
| 24 | agr8.111 | 16% | 14% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 100% | 100% | 0% | 0% | 0% | |
| 25 | agrC.III | 15% | 14% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 91% | 97% | 0% | 0% | 0% | |
| 26 | agrD.III | 16% | 14% | 0% | | 50% | 0% | 100% | 100% | 6% | 2% | 0% | 0% | 100% | 100% | 100% | 100% | 0% | |
| 28 | agr8.IV | 53% | 41% | 0% | 0% | 59% | 1% | 100% | 100% | 96% | 98% | 0% | 0% | 0% | 0% | 100% | 100% | 100% | - |
| 29 | agrC.IV | 23% | 5% | 0% | 0% | 59% | 1% | 100% | 100% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | |
| 30 METHICILLIN RESISTANCE | hld | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 10 |
| 32 AND | delta mecR | 2% | 3% | 0% | 0% | 0% | 25 | 0% | 0% | 10% | 2% | 5% | 16% | 0% | 0% | 0% | 0% | 0% | |
| 33 SCCmec TYPING | UgpQ. | 3% | 4% | 0% | 0% | 2% | 2% | 0% | 0% | 13% | 2% | 5% | 16% | 0% | 0% | 0% | 0% | 0% | |
| 34 | ccrA.1 | 1% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 4% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | |
| 35 | ccr8.1 | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | |
| 37 | Q99868.dcs | 1% | 3% | ON | 05 | 0% | 2% | 0% | 0% | 10% | 0% | 5% | 12% | 0% | 0% | 0% | 0% | 0% | |
| 38 | ccrA.2 | 3% | 4% | 0% | 0% | 0% | 2% | 0% | 0% | 10% | 2% | 5% | 16% | 0% | 0% | 0% | 0% | 0% | |
| 39 | ccr8.2 | 3% | 4% | 0% | 0% | 0% | 2% | 0% | 0% | 10% | 2% | 5% | 16% | 0% | 0% | 0% | 0% | 0% | |
| 40 | kdpA kdpB | 1% | 15 | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 12% | 0% | 0% | 0% | 0% | 0% | |
| 42 | kdpC | 1% | 15 | ON | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 12% | 0% | 0% | 0% | 0% | 0% | |
| 43 V11 | kdpD.SCC | 1% | Proc | essing | g of₀B | iologic | al Da | ata SS | 6 2020 |) 0% | 0% | 0% | 12% | 0% | 0% | 0% | 0% | 0% | |
| 44 | kdpE.SCC | 1% | 1% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 0% | 12% | 0% | 0% | 0% | 0% | 0% | |
| | (mac) | 2004 | 100 | | | | - | | 016 | - | | 175 | 12% | 0.00 | 0% | | | - | |

This slide summarizes the results of the hybridization against the DNA microarray.

The columns denote what fraction of the African or German samples belonging to a particular clonal complex (third row) contain certain genes.

Red: practically all samples hybridize against this probe.

Dark Green: 0% show positive hybridization,

Light green : 1-2% show positive hybridization,

Yellow and orange: increasing fractions show positive hybridization.



This slide shows a principal component analysis of the microarray hybridization results (data in table on previous slide).

Circled clusters often contain only members of one clonal complex.

| Site | # of cases below 1 year | # of cases 1 to 5 years | # of cases 6 - 25 years | # of cases 26 – 65 years | # of cases above 66 years | | | |
|---|-------------------------------|-------------------------------|-------------------------------|--------------------------------|---------------------------------|--|--|--|
| Africa + Germany (clinical) | 88 | 109 | 90 | 225 | 88 | | | |
| Africa + Germany (commensal) | 19 | 34 | 363 | 175 | 9 | | | |
| Africa (clinical) | 86 | 106 | 53 | 54 | 1 | | | |
| Africa (commensal) | 17 | 34 | 156 | 89 | 4 | | | |
| Germany (clinical) | 2 | 3 | 37 | 171 | 87 | | | |
| Germany (commensal) | 2 | 0 | 207 | 86 | 5 | | | |
| Age distribution was heavily skewed : many small kids / babies in Africa – few seniors in Africa very few small kids / babies in Germany – many seniors in Germany Did this affect the analysis + interpretation? | | | | | | | | |
| | | | | | | | | |

This is something I did not mention in the lectures #1 and #2: the age distribution between African and German samples is quite different.

We were worried whether this would affect the outcomes of our study that focused on the analysis of bacterial samples, not on the human carriers.

Are elderly people preferably colonized by different bacterial strains than young people?

In principle, one could expect that the total numbers of samples in Africa vs. Germany would be biased by the higher life expectation in Germany than in Africa.

However, the fraction of commensal samples from elderly people is quite low in both continents (1 vs. 5).

In contrast, the ratios are very different for the clinical samples.

Africa had many clinical cases for infants and small children reflecting the problematic health situation in Africa.

Many African clinical cases were apparently due to traffic accidents, after which the victims had to be taken to far away hospitals, which could take days.

Germany had many clinical cases for elderly people reflecting their higher susceptibility toward infections.

Analyze whether age is a confounding variable

To test whether age is a **confounding variable**, one can compare the results from simple linear regression with those from multiple linear regression.

The principal difference between these two types of regression models is the number of explanatory variables:

(1) the simple linear regression (SLR) model uses only one dependent variable y and one explanatory variable x: $y = a + b \cdot x$

In our case, *y* stands for the binary output from the Alere-chip experiment for a particular gene. *y* therefore has values of 0 or 1.

With the binary variable x we could encode the sites Africa (x = 0) / Germany (x = 1). *a* and *b* are weights estimated by the model.

Generally SLR tries to find such weights (values for \boldsymbol{a} and \boldsymbol{b}) so that the difference between the estimated \boldsymbol{y} and actual \boldsymbol{y} will be the smallest.

V11

Processing of Biological Data SS 2020

14

We will test whether age is a confounding variable by performing different types of linear regression analysis.

As a reminder, linear regression yields an optimal fit of a line y = a + b. x to the data.



We will compare the results of a linear fit to a multiple linear regression model against several variables.

| Steps of testing age categories for confounding (1) Estimate a linear regression model for the dependent variable and one or more explanatory variables. | | | | | | | | |
|---|--|----|--|--|--|--|--|--|
| (2) Repeat step 1 with age categories added as further explanatory variable. | | | | | | | | |
| (3) Compare the weights obtained in steps 1 and 2. | | | | | | | | |
| As a rule of thumb, if the weight (s) (regression coefficient(s)) from step 1 changes by more than 10%, then the additional variable (here: age) may be considered as a confounder . | | | | | | | | |
| By following these steps, one association) whether age is a | e can test for every significant finding (e.g. gene a confounder. | | | | | | | |
| Reasons for this could be e.g among samples. | g. a significant imbalance in the distribution of age | | | | | | | |
| V11 | Processing of Biological Data SS 2020 | 16 | | | | | | |

Obviously, the two different fit approaches will not give "the same" result.

What degree of difference should be considered as an "alarm" sign of a possible confounding effect?

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6447501/

states "A statistical approach to covariate selection … is what is sometimes called the "change-in-estimate" approach. In this approach covariate selection decisions are made based upon whether inclusion of a covariate changes the estimate of the causal effect for the exposure by more than some threshold, often 10%."

| Case study Case study: test whether age categories are a confounding va 2 genes lukS.PV and sdrCtotal | riable for the |
|--|--------------------|
| <u>Reason for selecting these genes</u> : these 2 genes have very different frequencies in African vs Germar in clinical vs commensal samples. | n sites as well as |
| Africa was encoded as $x_1 = 0$ and Germany as $x_1 = 1$. | |
| Clinical samples were encoded as $x_2 = 1$ and commensal with $x_2 =$ | 0. |
| Age categories were encoded from $x_3 = 1$ to 5. | |
| | |
| | |
| V11 Processing of Biological Data SS 2020 | 17 |

We performed this test for two genes lukS.PV and sdrC..total that showed very imbalanced frequencies.

Our reasoning was that if we could not find an age effect here, then we could argue that the age imbalance did not significantly affect the results we reported.

| Mul The Aler The site | tiple linear r e result for this go affiliation + clin/c | regressic ene for differ om values ar | on model fo ent samples is t re explanatory v | r the lukS.PV gen he dependent variable. ariables. | e | | | |
|---|---|--|--|---|----|--|--|--|
| The tabl | e lists the depend | dent(lukS.P\ | /) and explanate | ory (Africa_value, | | | | |
| clin_com | n_value) variables | s for 10 sam | oles out of 1200 | samples. | | | | |
| # 1 2 3 4 5 6 7 8 9 10 | samples FR-B001 FR-B003 FR-B004 FR-B005 FR-B007 FR-B008 FR-B009 FR-B010 FR-B011 FR-B012 | lukS.PV 0 0 0 0 0 0 0 0 0 0 0 | Africa_value 0 0 0 0 0 0 0 0 0 0 | clin_com_value 1 1 1 1 1 1 1 1 1 1 1 | | | | |
| Since all these samples are from a German site, the Africa_value = 0. Also, all samples are clinical (clin_com_value = 1). | | | | | | | | |
| V11 | | Processi | ng of Biological Data SS 202 | 0 | 18 | | | |

The table shows a small piece of the used data. In total, it contains 1200 rows.

Our first regression model relates the value of the lukS.PV column (0 or 1) to the values of the Africa_value column (0 or 1) and the clin_com_value (0 or 1).

| IukS.PV Application of linear regression determines optimal weights w_1, w_2, w_3 . | | | | | | | | | | |
|--|--|---|--------------------------------------|---|--|--|--|--|--|--|
| For every sample we get | lukS.PV = w ₁ + | w₂ · Africa.value ⊣ | ⊦ w ₃ · clin c | om value . | | | | | | |
| For the first sample FR-B | For the first sample FR-B001, the formula would be $0 = w_1 + w_2 \cdot 0 + w_3 \cdot 1$. | | | | | | | | | |
| Results from multiple linear regression (coefficients marked in bold): | | | | | | | | | | |
| w_1 for intercept w_2 for Africa_value w_3 for clin_com_value | Estimate -0.07250 0.42833 0.19500 | Std. Error 0.01781 0.02057 0.02057 | t value -4.070 20.825 9.481 | Pr(> t) 5e-05 *** <2e-16 *** <2e-16 *** | | | | | | |
| In other words, the following model is estimated: lukS.PV = -0.07 + 0.42833 · Africa_value + 0.195 · clin_com_value | | | | | | | | | | |
| t value : equal to coefficient (estimate) divided by the standard error. Pr(> t) : p-value = probability of seeing a result as extreme in random data. | | | | | | | | | | |
| V11 | Processing of Bio | ological Data SS 2020 | | 19 | | | | | | |

The standard error of the regression (S), also known as the standard error of the estimate, represents the average distance that the observed values fall from the regression line.

Although the model is very simplistic, the std. error is quite small.

```
lukS.PV
We then added a further variable "age category" with weight w_4 to the model.
lukS.PV = w_1 + w_2 \cdot \text{Africa.value} + w_3 \cdot \text{clin com value} + w_4 \cdot \text{age}
                   Estimate
                                     Std. Error
                                                        t
                                                                  value Pr(>|t|)
(Intercept)
                   0.06211
                                     0.04559
                                                        1.362
                                                                  0.17333
                                                        16.556 < 2e-16 ***
Africa_value
                   0.39077
                                     0.02360
                                                        9.503
                                                                  < 2e-16 ***
clin_com_value 0.19470
                                     0.02049
                                                        -3.206 0.00138 **
age
                  -0.03618
                                     0.01129
lukS.PV =
0.06211 + 0.39077 · Africa value + 0.19470 · clin com value - 0.03618 · age
 V11
                                  Processing of Biological Data SS 2020
                                                                                     20
```

This is the same fit when we added age as a third variable.

Actually, we did not use the actual age of the cases, but grouped the cases into the 5 age categories listed on slide 13: below 1 year, 1 to 5 years, 6 - 25 years, 26 - 65 years, above 66 years

The categories were encoded as 0 to 4 or 1 to 5.

| lukS.PV | | | | | | | | |
|----------------------------------|---|-----|--|--|--|--|--|--|
| This result shows | tegery has a very small impact (its own weight is close to 0) and | | | | | | | |
| (b) the two other v | reights (for the site and clin/com) did not change much. | I | | | | | | |
| E.g. the weight of | the Africa_values changed in relative terms by : | | | | | | | |
| | $\frac{(0.42833 - 0.39077)}{0.42833} \cdot 100\% = 8.8\%$ | | | | | | | |
| The weight of clin | _com_value changed by only 0.15%. | | | | | | | |
| Both values are si | naller than 10% (rule of thumb). | | | | | | | |
| Conclusion: There is no (clea | r) statistical evidence that age acts as a confounding variab | le. | | | | | | |
| V11 | Processing of Biological Data SS 2020 | 21 | | | | | | |

The addition of an age variable had a very small effect on the linear regression of lukS.PV status.

| Same analysis for gene sdrC_total | | | | | |
|--|---|--|---|---|----|
| Before addi Coefficients (Intercept) Africa_value clin_com_va | ng age s: e alue | categories: Estimate Std. Error 1.02083 0.0125 -0.12833 0.0144 -0.05833 0.0144 | t value 81.711 -8.896 -4.044 | Pr(> t) < 2e-16 *** < 2e-16 *** 5.6e-05 *** | |
| After adding Coefficients (Intercept) Africa_value clin_com_va age-categor Weight of A | g age ca e alue ry frica_va | tegories: Estimate Std. Error 0.975445 0.0321 -0.115667 0.0166 -0.058232 0.0144 0.012198 0.0079 lue changed by 9.87 | t value 30.407 -6.964 -4.039 1.536 %, weight | Pr(> t) < 2e-16 *** 5.44e-12 *** 5.71e-05 *** 0.125 of clin_com_value changed by 0.17% | |
| V11 | | Proc | essing of Biologic | al Data SS 2020 | 22 |

Also for sdrC, the age category got a small weight and the probability (p-value) of a random-effect is actually quite high.

In this case, the weight of the Africa_value changed by almost 10%. We still considered this close enough to the simpler fit without age category.

| Conclusion | | | | |
|--|---|---------------------------|--|--|
| There is no evidence f sdrCtotal that age ac invasiveness and site | from our preliminary analysis for the genes lukS. Its as a confounder in the association of genes w affiliation. | PV and /ith | | |
| We wrote in our manuse "The discrepancy in pop potentially biases the 'tr the different geographic [but] application of a mu Panton-Valentine leucod confounding variable" | cript: oulation age between the German and African cohor rue' distribution of clones and genes between isolate c regions ultiple linear regression model for the detection rate o cidin genes failed to provide evidence that age acts a | t s from of as a | | |
| Ruffing et al. Sci. Rep. 7, 154 | (2017) | | | |
| V11 | Processing of Biological Data SS 2020 | 23 | | |

We are smart \bigcirc

In our manuscript, we first acknowledge openly that age may be confounding factor.

Then, we state that a multiple linear regression model did not provide evidence FOR it. But we also did not exclude that such an effect may exist.

The reviewers of our study were satisfied.



This is additional data and analysis that we did not report in our study.

We also tested whether the reported findings were affected in a possible imbalance in the occurrence of diabetes and HIV.

One convenient way to test this is Fisher's exact test.

| Analysis of HIV co-infection First, we will test the null hypothesis that "HIV is equally distributed in African and German samples". (a) For all African samples and all German samples we obtain the following | | | | |
|--|------------------------|---------------------------------------|---------|--|
| (you may s | say non-carriers) (HI\ | /-): | | |
| | | Africa | Germany | |
| | HIV+ | 41 | 0 | |
| | HIV- | 315 | 586 | |
| The p-value obtained for this table can be interpreted as the sum of evidence provided by the observed data—or any more extreme table—for the null hypothesis that "there is no difference in the proportions of HIV carriers among the African and German individuals tested in our study". | | | | |
| The smaller the value of p, the greater the evidence for rejecting the null hypothesis. | | | | |
| VII | | Processing of Biological Data SS 2020 | | |

Very sadly, Africa is strongly affected by HIV infections. This is also reflected in the study cohort.

In the German cohort, there was no HIV case included. This may either be coincidence or result from better monitoring of the population and exclusion of HIV infected cases.



The same Fisher's test was applied.

Indeed, the very small p-value is evidence that the HIV status differs significantly in African and German cases.



https://www.aerzteblatt.de/int/archive/article/175344/The-prevalence-and-incidence-of-diabetes-in-Germany-an-analysis-of-statutory-health-insurance-data-on-65-million-individuals-from-the-years-2009-and-2010

reports that the incidence of diabetes in Germany is roughly 10% of the population.

In Africa, the prevalence is estimated as about 4% of the population, see https://idf.org/our-network/regions-members/africa/welcome.html

So this difference would not explain the large imbalance observed in the confusion matrix.

The proper answer is provided by the prevalence in different age groups, see right table.

There is a steep rise from 1.6% for men between 40-49 years old to 26.3% for men between 80-89 years old.

Hence, the observed age imbalance strongly affects the prevelance of diabetes in African vs. German cases.

HIV/diabetes in individuals with selected CCs

Next, we tested the distribution of HIV/diabetes in individuals carrying *S. aureus* from selected clonal complexes (CC15, CC45, CC121, CC30 which showed significant imbalance in German/African samples).

These are the results (tables + p-values from Fisher's exact test)

RF_HIV

| CC15 | Africa | Germany |
|-------------|--------|--|
| hiv+ | 4 | 0 |
| hiv- | 65 | 57 |
| p-value | 0.126 | 0.25 (after correction for false discovery rate (FDR)) |
| CC45 | Africa | Germany |
| hiv+ | 1 | 0 |
| hiv- | 40 | 87 |
| p-value | 0.320 | 0.42 (FDR-corrected) |
| V11 | | Processing of Biological Data SS 2020 |

As shown before, there exists strong and statistically significant imbalances of HIV and diabetes in the study cohort.

But does this also affect the findings reported by us?

Our focus was not placed on the individuals and their infection status, but on the bacterial strains colonizing or infecting them.

Hence, we repeated the same analysis shown before for the major subpopulations of clonal complexes.

Now it turns out that the imbalances were mostly statistically insignificant.

Because we performed multiple tests, we had to apply an FDR correction.

28



Colored red are those scenarios with FDR-corrected p-value less or equal to the significance threshold of 0.05.

HIV/diabetes in individuals with selected CCs

| RF_CCS | SI_Diab_ | mel |
|-----------------------------------|-----------------------|--|
| CC15 | Africa | Germany |
| diab+ | 0 | 1 |
| diab- | 88 | 56 |
| p-value | 0.393 | 0.52 (FDR-corrected) |
| CC45 | Africa | Germany |
| diab+ | 0 | 12 |
| diab- | 47 | 75 |
| p-value | 0.0081 | <mark>0.03</mark> (FDR-corrected) |
| CC121 | Africa | Germany |
| diab+ | 0 | 1 |
| diab- | 57 | 24 |
| p-value | 0.305 | 0.52 (FDR-corrected) |
| CC30 diab+ diab- p-value | Africa 0 9 1 | Germany 7 68 1 (FDR-corrected) Processing of Biological Data SS 2020 |

| | Interpretation | | | | |
|---|--|--|--|--|--|
| In most cases, there is no evidence based on our data to reject the null hypothesis of assuming a similar distribution of HIV and diabetes carriers among African and German samples belonging to particular clonal complexes . | | | | | |
| The only exceptions to this are and | CC45 (diabetes – p=0.008/q=0.03) CC121 (HIV – borderline p=0.013/q=0.05). | | | | |
| Therefore, we summarized in an internal report of the project team: we observed statistically significant imbalances in the frequencies of all these clonal complexes XXX, YYY between African and Germany. we tested based on Fisher's exact test that these imbalances were not due to an imbalance of HIV and diabetes carriers in both groups. The only exceptions to this are CC45 (diabetes) and CC121 (HIV) where such associations cannot be ruled out, but are only weakly supported by the data. | | | | | |
| V11 | Processing of Biological Data SS 2020 31 | | | | |

We did not even mention the outcome of these checks in the manuscript because we believe that the reported results are not effected by the HIV/diabetes status of the individuals.

| - | Summary It is important to think about the occurrence of possible confounding effects | | | |
|---|---|----|--|--|
| - | - Multi-variate vs. single-variate analysis reveals possible confounding effects | | | |
| - | - Choosing an appropriate study design reduces the likelihood that confounding effects may occur. | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | V11 Processing of Biological Data SS 2020 | 32 | | |

| Relevant slides for oral exam on July 27-31, 2020 | | | |
|---|---|--|--|
| Lecture | Slides | | |
| 1 | 14, 18, 19, 26-30, 35-40 | | |
| 2 | 4-7, 9, 13-15, 18, 22-23, 26-32 | | |
| 3 | 6-7, 13-16, 18-23, 32 | | |
| 4 | 11-17, 19-23, 42-43 | | |
| 5 | 4-6, 17, 19-22, 33-34 | | |
| 6 | 3 (only Hi-C), 7-11, 14-23, 26-37 | | |
| 7 | 16, 18-24, 35-39 | | |
| 8 | 4-5, 8 | | |
| 9 | 10-25, 39-43 | | |
| 10 | 5, 8-14, 25- 27, 33-35 | | |
| 11 | 2, 4-8, 13-31 | | |
| Assignments | Ass.#2 (full), Ass. #3 (full) | | |
| l will not ask ques | tions about the other slides & assignments. | | |
| V11 | Processing of Biological Data SS 2020 | | |

| Example questions | |
|---|--|
| (1) What types of methods exist to perform integrated analysis of multiple omics datasets? | |
| (Answer: multi-staged approaches and meta-dimensional approaches) | |
| To which category does the tool iCluster belong? | |
| (Answer: it is a meta-dimensional method – uses model-based integration) | |
| How does it work? | |
| (2) How does the watershed algorithm work? | |
| (3) What is the connection between PCA and singular value decomposition? | |
| Can one perform PCA without SVD? | |
| (Answer: yes, e.g. by geometric construction, see slide 29 in V1). | |
| V/41 Decessoring of Dislogical Data SS 2020 | |
| VII Processing of diological Data 55 2020 34 | |