

# Data Science and Big Data in Upper Secondary Schools: A Module to Build up First Components of Statistical Thinking in a Data Science Curriculum

Rolf Biehler, Daniel Frischemeier, Susanne Podworny, Thomas Wassong, Lea Budde, Birte Heinemann and Carsten Schulte

**Abstract** Within the framework of a design-based research project, computer science educators and statistics educators at Paderborn University designed a pilot course on the subject of data science and big data. It addresses upper secondary students and was realized by weekly sessions (three hours) over seven months. The whole course that is intended to introduce upper secondary school students to the field of data science consists of four modules. In module 1, the learners are introduced into the basics of statistics and big data and it aims at developing their data competence and data awareness. In the second module, learners are introduced to machine learning and programming based, among others, on examples from module 1. In the third and fourth module, learners can apply their knowledge gained in modules 1 and 2 and

---

Rolf Biehler · Daniel Frischemeier · Susanne Podworny · Thomas Wassong  
Institute of Mathematics, Paderborn University, Paderborn, Germany

Lea Budde · Birte Heinemann · Carsten Schulte  
Institute of Computer Science, Paderborn University, Paderborn, Germany

✉ biehler@math.upb.de  
✉ dafr@math.upb.de  
✉ podworny@math.upb.de  
✉ wassong@math.upb.de  
✉ lea.budde@uni-paderborn.de  
✉ heinemann@informatik.rwth-aachen.de  
✉ carsten.schulte@uni-paderborn.de

ARCHIVES OF DATA SCIENCE, SERIES A  
(ONLINE FIRST)  
KIT SCIENTIFIC PUBLISHING  
Vol. 5, No. 1, 2018

DOI: 10.5445/KSP/1000087327/28

ISSN 2363-9881



will work in small groups on real and meaningful data science projects. In this paper, we want to concentrate on the statistics components, especially of module 1, and we will present how we develop the data competence and data awareness of upper secondary school students to prepare them to work on data science projects in modules 3 and 4.

## 1 Introduction

In the era of big data and the omnipresence of data, data science has gained increased importance: Industrial and economic processes, marketing processes, and also monitoring processes in politics are based on statistics and massive amounts of data (Gould, 2017). That is why data science permeates many areas of life and, therefore, competent data handling is inevitable (Engel, 2017). To build up such a data competence, it is important to educate students and learners in this respect as early as possible. One way for example could be to implement and introduce elements of data science already in secondary school.

Our project “ProDaBi” (Project: data science and big data in secondary school, <https://www.prodabi.de>) which is a joint project of statistics educators and computer science educators at Paderborn University was initiated and is funded by the Deutsche Telekom Stiftung (<https://www.telekom-stiftung.de>) and has the aim to design and test an experimental data science curriculum for upper secondary school students. One major goal in the first year of the project is to design, implement, and evaluate a one-year-pilot course (3 hours per week) in grade 11–12 (ages 16–17) at an upper secondary school (Gymnasium) in Paderborn. The pilot course is a so-called non-obligatory “project course” that does not underly any curricular constraints from mathematics or computer science education. Thus, we have the freedom to experiment with new ideas.

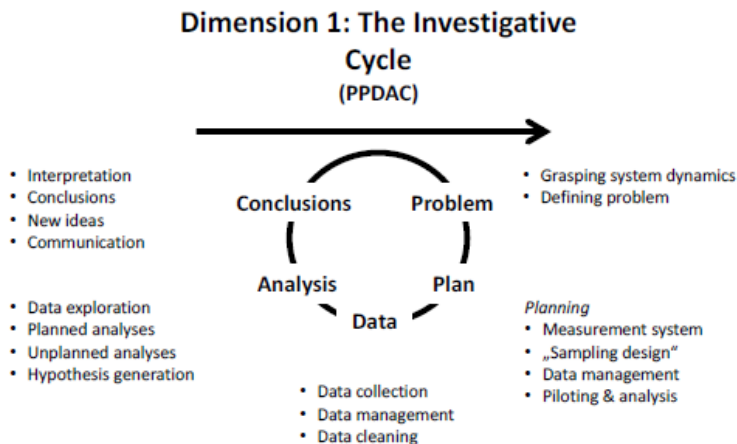
However, the domain of data science is a large field, combining the disciplines statistics, computer science, and also sociocultural issues (for an overview on several facets with regard to the different disciplines see Biehler et al., 2018). For implementing all these issues in a tight school curriculum, it has to be discussed which topics and which contents can and which have to be implemented. For instance, in the statistics education discussion fundamental ideas (see Burrill and Biehler, 2011) like data, representation, and variation have been identified, ideas which seem also to be necessary issues for a data science curriculum from a statistics education perspective. In this article we will present the design

of the module that is focussing on the basics of statistics for data science in secondary schools.

## 2 Data Science and Statistics Education

A rich overview of teaching and research practices and pedagogical ideas in statistics education is provided by Garfield and Ben-Zvi (2008). For a stronger emphasis on the research in statistics education see Ben-Zvi et al. (2018). In the following, we will give an overview of trends in data analysis in statistics education and we try to derive some commonalities between statistics education and big data and on which stages the statistics education perspective has to be broadened for a data science perspective.

In his landmark article Ridgway (2016) – as its title says – discusses implications of the data revolution for statistics education. For instance, referring to the well-known PPDAC-cycle (Wild and Pfannkuch, 1999) of statistics education, which includes the steps problem (P), plan (P), data (D), analysis (A), and conclusions (C; in total “PPDAC”), some refinements have to be considered with regard to the investigative cycle (see Figure 1) from a data science perspective. This is done for example in the CRISP-DM cycle of Berthold et al. (2010, p. 9).



**Figure 1:** Dimension 1 (The Investigative Cycle) of the PPDAC-cycle, own diagram similar to Wild and Pfannkuch (1999, p. 226).

One issue is that modeling (Breiman, 2001) plays a fundamental role in data science and, therefore, should be implemented in such a cycle because data science uses new types of algorithmic models (Breiman, 2001).

The issue of model validation and the distinction between predictive and causal models and its implication on validation in data science is also missing in the PPDAC-cycle and associated with this issue prediction as a goal for modeling has to be emphasized in a data science cycle too.

Another “new” insight in the PPDAC cycle could be the issue of whether the phases of data and data collection are still up-to-date since data (especially in this era big data or open data) are already there and not necessarily have to be collected in a typical sense anymore (see e.g. Huber, 2012). So the shift from own data collection procedures to assessing the quality of already available data is very important in this sense. Furthermore the step “conclusions” of the PPDAC-cycle might be too narrow and should be extended for data science purposes, e.g. into “deployment” which can be also found in the cycle of Berthold et al. (2010, p. 9).

But – as mentioned before – not only the investigative process but also the view and types of data have to be expanded for data science. Whereas the datasets collected within the PPDAC-cycle are usually based on surveys, experiments, or observations, in data science new types of data and data collection methods may arise: Data collected by sensors, data collected by personal devices, and transactional data. Data can also be scraped from webpages or – with using GPS – one can also collect data with geographic information: All these types of data are multivariate (which is common and standard in statistics) but also messy and often ill-structured. This brings us to another issue: The selection of adequate digital tools for analyzing these kinds of data. An overview of statistical tools used in schools is provided by Biehler et al. (2013). The standard tools in (German) schools are, however, Microsoft Excel or GeoGebra or a graphics calculator. Educational software that builds bridges between the learning process of statistics and the process of doing data analysis are Fathom (Finzer, 2001) and TinkerPlots (Konold and Miller, 2011; for German versions see <https://www.stochastik-interaktiv.de>). Both tools are more or less only used in experimental classrooms but not in statistics education on a large scale. CODAP (Common Online Data Analysis Platform, <https://codap.concord.org/>) is a recent development by the Concord Consortium and has implemented many ideas of Fathom and TinkerPlots and

offers a web-based open source software for data analysis. In the discipline of statistics itself, there are professional tools such as SPSS, R, or Python, which allow sophisticated and deep data analyses but require support and considerable time for learners to learn the language and to build up their development environments. So at least with this short overview, we see the dilemma we face on the issue of tool selection in a data science course for schools. On the one hand, data science requires a powerful tool like R or Python to realize deep and sophisticated data analyses, on the other hand, one has to consider in which way and how an implementation of these tools in school are realistic and possible.

To sum this up, for the design of a data science course for upper secondary schools we have to expand the statistics education perspective, as we have mentioned in this paragraph: It is necessary to expand and refine the PPDAC-cycle of Wild and Pfannkuch (1999), to extend the statistical view of “data” and also to consider the issue of selecting adequate digital tools for the data analysis process.

### **3 The Pilot Course on Data Science and Big Data in Upper Secondary School 2018/2019**

Our idea is to bring elements of data science (and the broadened perspective on statistics) into the German school curriculum and we offered the first version of a pilot course on data science and big data in the school year 2018/2019. Within the frame of a design-based research project (Cobb et al., 2003), we designed a pilot course for data science and big data addressing upper secondary school. The course is offered in cooperation with two schools in Paderborn (Gymnasium Theodorianum & Gymnasium Reismann) and was realized in weekly sessions (each lasting three hours) over seven months in the school year 2018/2019. In practice, about 25 weeks were available. The whole course consists of four modules:

- In the first module, the learners are introduced to the basics of statistics and big data, following the metaphor of a data detective looking for patterns in data. This module aims at developing data competence and data awareness.

- The second module introduces learners to machine learning and programming based on examples from module 1.
- In the third module, learners can apply their knowledge gained in modules 1 and 2 and work in small groups on real and meaningful data science projects.
- In the fourth module, societal and cultural aspects are in the foreground based on insights from module 1 to 3.

In the following, we present details on module 1 of the pilot course, with the title “Data and Data Detectives”. For details on module 2 of the pilot course, see Heinemann et al. (2018).

## **4 The Data Exploration Module of the Pilot Course on Data Science**

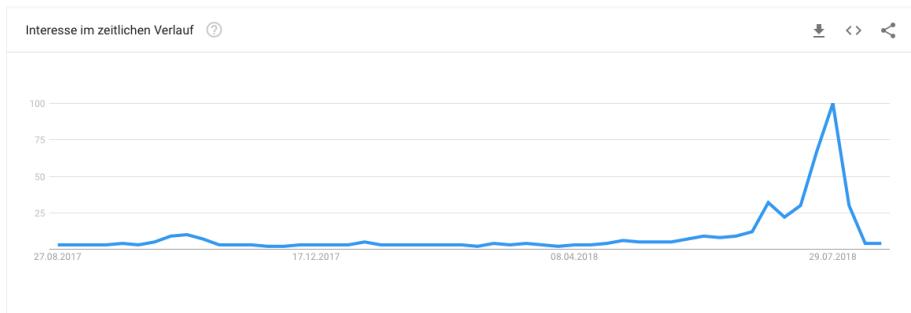
The module with a focus on the basics of statistics, module 1 (lasting from week 1 to week 5), has the major goals to enhance statistical thinking and to introduce into the use of digital tools for analyzing data, to develop data competence, and to sharpen an up-to-date data concept appropriate to data science. In the first unit (week 1–2) the students take on the role of data detectives and are directly confronted with concrete data science problems. The metaphor of data detectives searching for evidence in data is used throughout the module. In the second unit (week 3–5), the students are supposed to extend their role as data detectives and explore the multivariate JIM dataset on the use of information, technology, and social and other media by students aged 10–19 (JIM = Jugend, Information, Medien; translated: Youth, Information, Media) with CODAP (Common Online Data Analysis Platform). The students get to know fundamental statistical concepts for univariate and bivariate data analysis, such as different percentages in contingency tables, group comparisons, statistical associations and dependencies, and statistical representations. In the following, we will point out further details such as learning goals, activities and design ideas on both units.

## 4.1 Unit 1: Introduction & Data Detectives in Paderborn – On the Trail of Noise

As mentioned before, unit 1 has the aim to introduce the students to the module 1 named “Data and Data Detectives” and to confront them with real large data sets. More precisely we want our students to learn ...

- ... about basic concepts of data and data cleaning,
- ... to interpret, process, and visualize data (e.g., about noise data in a factual context and thus to obtain information),
- ... about basic terms of descriptive statistics (variables, types of variables, operationalization of variables),
- ... about basic knowledge in programming (use of Jupyter Notebooks),
- ... to analyze data with the help of a pre-built script.

As an appetizer and first activity, we confront our students with the concrete problem of interpreting charts taken from Google trends. In the example in Figure 2 we see the distribution of google searches of the term “Libori” (“Libori” is a famous funfair in Paderborn which takes place annually in the summer (end of July/begin of August) in the city of Paderborn).



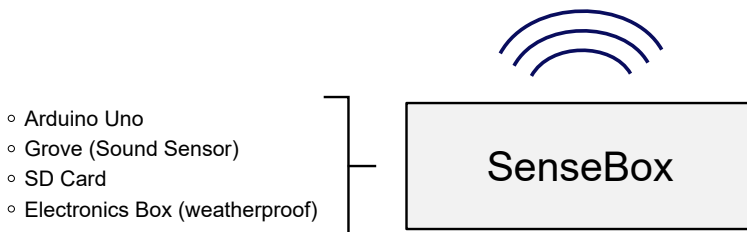
**Figure 2:** Google trends chart for the term “Libori” in Paderborn (Data source: Google trends (<https://www.google.com/trends>)).

Confronted with this kind of statistical display we designed several tasks to lead the students to an understanding of such a display: Our students are supposed to

1. describe what can be seen in the chart as a first open approach to read and interpret the chart,
2. understand the value of 100 on the y-axis (Google only shows relative data),
3. name both axes (x- and y-axis) to get an understanding of the meaning of the axes, and
4. finally sum up the insights of (1)–(3) for interpreting the chart.

Students learn how to relate patterns in data to context and vice versa. They are encouraged to do their searches with Google trends and find interpretations of the patterns found. Creating opportunities to explore questions of their own interest is a feature of the whole course.

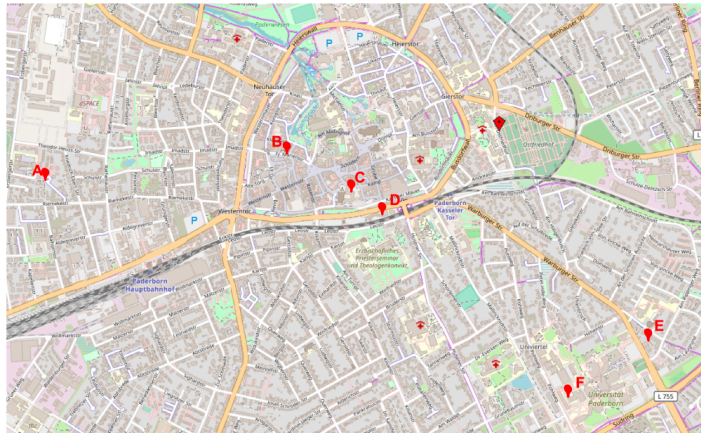
After this introduction, we introduce our students to the interplay between statistics, computer science, and real data collection. For this purpose, we designed a short project for the course called “noise project”. In this project, the students are supposed to investigate the noise situation at different initially unknown locations in Paderborn, to collect noise-data with so-called sense boxes (<https://sensebox.de/en/>) and to analyze the noise data using pre-prepared Jupyter Notebooks (<https://jupyter.org/>). A sense box is a tool which contains several devices like an Arduino Uno, a sound sensor and a SD memory card (see Figure 3).



**Figure 3:** Devices of a sense box: Arduino Uno, Grove-Sound Sensor, SD-card, and a weatherproof outdoor electronics box.



We use the Grove-Sound Sensor to collect noise data on different locations in Paderborn for one week. The data is stored as raw data on the SD-card containing two variables: time stamp and noise (collected as measurements in the unit of voltage). So the students are provided with large data sets from different locations in the city of Paderborn (see Figure 4).

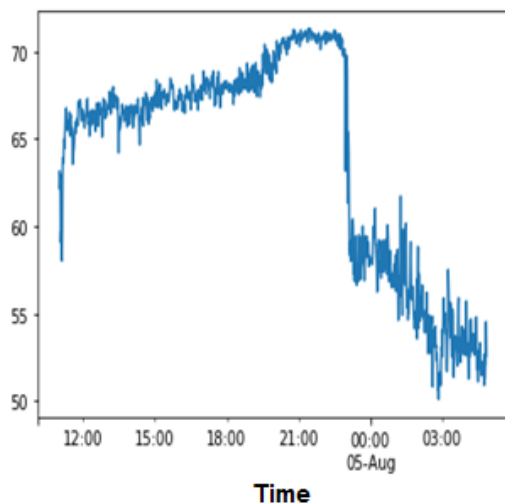


**Figure 4:** Map of the city of Paderborn with marked locations for noise collection.

To make it into an investigative task for data detectives, we do not tell the students which data set belongs to which location. The task for the students is to allocate a data set to a location, because this information was “lost”. In other words, we use the graphs as “mystery graphs”, which is a well-known task type in statistics education. This set-up engages students in explorations of the data sets and the context, where hypotheses of the noise development have to be generated. The students are supposed to use Python to (1) import the data, (2) clean the data, and (3) visualize the data on time and noise to allocate each noise profile to its location in the city of Paderborn correctly. In Figure 5, we see a possible data visualization of the noise profile at a certain location in the city of Paderborn. A pre-built Python notebook with explanations and pre-built code cells is given to students to support this process.

For example, the data visualization of Figure 5 can be allocated to the place “Liboriberg” (labeled D in Figure 4) because the noise will be louder from 12

o'clock until 23 o'clock, which corresponds to the opening hours of the funfair. We expect the students to notice the two drastic change points at 12 and 23 o'clock and, based on this, to analyze which of the five locations fits the profile. Since this noise profile occurs only once, and only the Liboriberg with its funfair on August 5 can explain exactly this “noisy” time slot, it should be possible to discover the correct location from looking at the data.



**Figure 5:** Snapshot of the data visualization.

Finally, the students are supposed to present their allocations and the corresponding hypotheses leading to these allocations to their classmates.

## **4.2 Unit 2: Data Detectives Using the JIM Data Set on Adolescents Media Use – Explore Multivariate Data With Codap and Deepen Statistical Thinking**

Unit 2 (week 3–5) builds on unit 1 in the sense that the students are supposed to analyze multivariate data (instead of noise data which contained only the two variables time and noise) with the digital tools CODAP and Python.

One issue already mentioned is that students are now supposed to analyze multivariate data instead of just uni- or bivariate data. Another issue typical for unit 2 is that our students are now supposed to analyze data that have been collected via surveys. These are data typical for many statistical studies, but our students will not have seen such data sets in their normal school career. To confront our students with a motivating topic and to engage them as data detectives in survey data sets we decided to let the students themselves collect data of their schoolmates, using a questionnaire similar to the questionnaire that was used in the so-called JIM-study (<https://www.mpfs.de/>). The study deals with media use of German school students at the age of 12–19. The JIM-study survey consists of many questions with Likert scales like e.g.,

- *“How often do you use Whatsapp?”*,

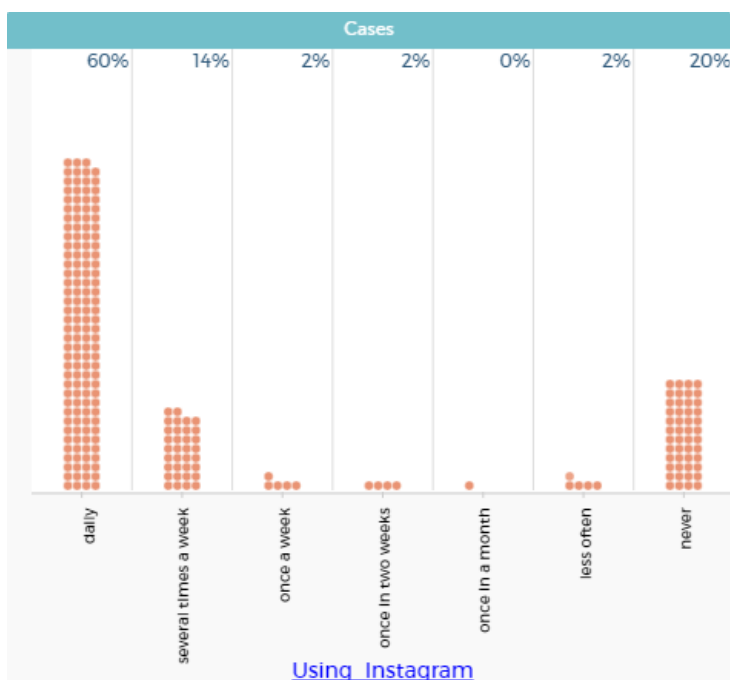
and answering options such as “daily”, “several times a week”, “once a week”, “two times a month”, “once a month”, “less often”, or “never”. To make the JIM-study, which is a nationwide study, even more interesting to our students, we decided to modify it and to include questions which tackle numerical variables, such as

- *“How many apps do you have on your smartphone?”*, or
- *“How many online accounts do you have?”*.

The published JIM report with its already analyzed, aggregated, and visualized data, is sometimes used as a ready-made report in media education at school level. However, we want our students to become involved in these data analysis processes themselves as data explorers and not only as consumers of data analysis reports. We want the students in unit 2 to ...

- ... regard the JIM data and the JIM questionnaire having a typical structure for multivariate survey data with different types of variables,
- ... learn to develop adequate statistical questions,
- ... explore a multivariate data set,
- ... use basic terms of descriptive statistics,

- ... analyze the JIM data with regard to selected questions and expand their knowledge about media use of classmates,
- ... apply/use statistical concepts (percentages, visualizations, group comparisons),
- ... use digital tools for data exploration, and
- ... reflect on the JIM data critically (What does it mean for us? Developing data awareness).



**Figure 6:** Distribution of the variable “Using\_Instagram” produced with CODAP for answering the question: How often do students use Instagram?.

At first, we introduce the students to the digital tool CODAP. CODAP is a free online data analysis platform and allows us to import and analyze multivariate data sets. CODAP stores the data in tables and allows learners to use graphs

(each case (row) of the table is represented by a dot in a graph, see Figure 6) to explore the data. Important and useful features of CODAP are handling different hierarchies in the data and getting easy access to a data set and the intermediate state of an exploration via an URL. A task for students is to pose their investigative questions and work as data detectives to answer their questions.

One first exploration in CODAP, the visualization of the distribution of the variable “Using\_Instagram” in form of a preliminary stage of a bar chart, can be seen in Figure 6 and is demonstrated to our students as a first exploration in CODAP. Figure 6 displays a U-shaped distribution, which we can interpret as follows: There are two main groups of students with regard to Instagram use: The ones who use Instagram frequently (approx. 3/4 use Instagram more than once a week or daily, 60 % of all students in the sample use Instagram even daily) and the ones who never use Instagram (20 %).

After this introduction into CODAP, the students are supposed to start a project consisting of their exploration of the JIM data set according to their questions of interest for the data in collaborative learning settings (in groups of 3 or 4 students). Our idea is that the students are enabled to generate statistical questions on their own to explore the JIM data from a perspective that is meaningful for them. Furthermore, we want to deepen several statistical concepts (percentages, comparing groups, representations) when working with the JIM data. In weeks 3 and 4, the students explore the JIM data and document their findings in form of a PowerPoint Presentation. In week 5, the students are supposed to present their findings to their classmates and to compare their findings with the findings of the official JIM-study. We emphasize the need for interpretations of graphs and tables to counteract the tendency to just collect representations and give superficial and not very thoughtful comments.

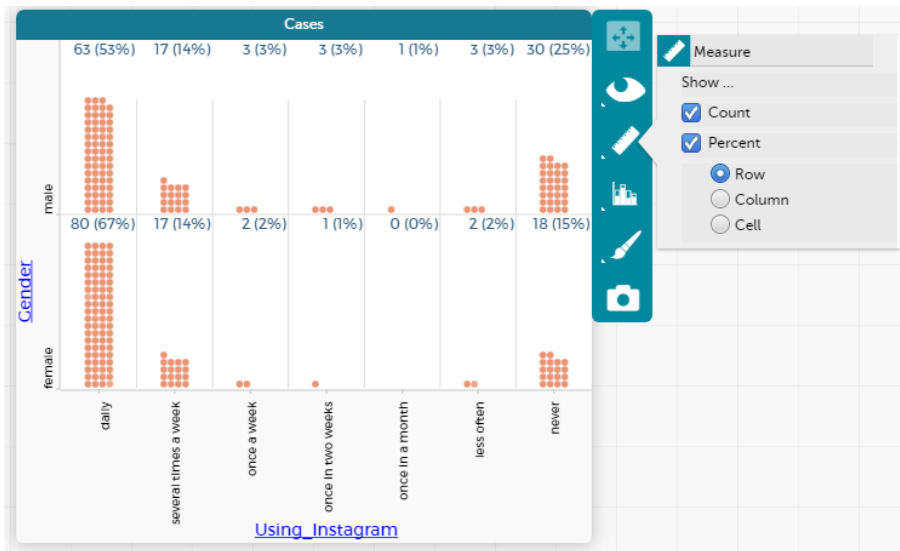
One fundamental idea is that the students are enabled to work on the JIM data project autonomously. At certain stages, we nevertheless prepare prompts to support their data exploration process. One prompt, for example, is given quite at the beginning of the project work when the groups are asked to generate statistical questions. According to Arnold (2013), the following categories of statistical questions exist: For example yes/no-questions like

- “Do boys use their smartphone more often than girls?”,

do not allow a deep data exploration because the question can be answered with yes/no. In contrast to yes/no-questions open and complex questions like

- “Which differences exist in the digital media use of 5th and 6th graders?”,

offer a good possibility of data explorations. In the study of Frischemeier and Biehler (2018) it was evident that learners, in this case preservice teachers, struggle when generating statistical questions.

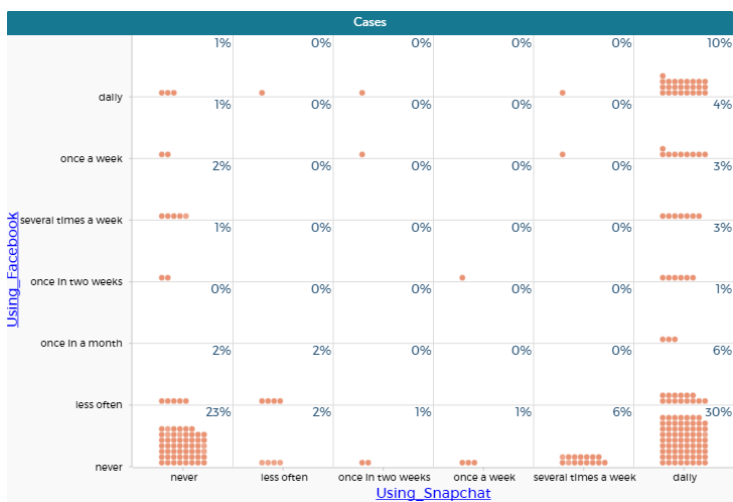


**Figure 7:** Screenshot of the comparing visualization for the variables “Using\_Instagram” and “Gender” with row percentages in CODAP.

For this reason, Frischemeier and Leavy (2020) have developed a think-pair-share setting with prompts to generate adequate statistical questions and to revise them during a think-pair-share process. In addition to the support in the problem phase, we also support our students in the data analysis phase, especially for the handling of different percentages (row, column, cell) in CODAP, they are given prompts when working on comparison questions like

- “In which way do male and female students differ in their use of Instagram?”.

In Figure 7, we exemplarily see a screenshot of CODAP when applying row percentages to tackle the question about whether male or female students tend to use Instagram more regularly. Looking at Figure 7, we see approx. 81 % of the female students but only approx. 2/3 of the male students in the sample use Instagram several times a week or even daily – so we can formulate that the female students tend to use Instagram more often than the male students. This statement can be further elaborated by comparing the percentage of non-use (25 % male vs. 15 % female). However, since we do not have any kind of a random stratified sample, we cannot make any further inferences about a larger population.

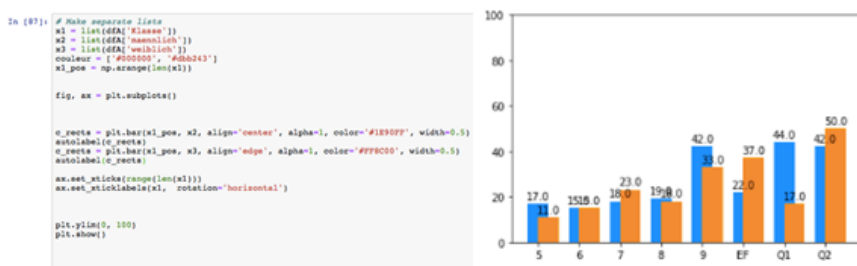


**Figure 8:** Screenshot of the comparing visualization for the variables “Using\_Snapchat” and “Using\_Facebook” with cell percentages in CODAP.

The investigation of a relationship between two categorical variables can also be made in a more complex way by comparing two categorical variables, which both have more than two values. For example, we could compare the variables “Using\_Snapchat” and “Using\_Facebook” (see Figure 8). Here, we expect our students to summarize specific cells and to use cell percentages to find relationships between Facebook use and Snapchat use. One possible summary might be to summarize the values “less often” and “never” to “rare use”, to summarize the values “once in two weeks” and “once a week” to “medium use”

and to summarize the values “several times a week” and “daily” to “high use”. With this categorization, the students then could conclude that there is a group who uses Facebook and Snapchat “rarely” (29 % use Facebook and Snapchat rarely). Furthermore, the students might identify that there is a group who uses Facebook rarely but shows frequent use in Snapchat (42 % use Facebook only rarely but show frequent use of Snapchat). Finally, our students could notice that there is a small group (13 %) who shows a frequent use of Facebook and Snapchat. Essentially this is a visual classification task based on two variables, which as such has not been part of the standard statistics school curriculum.

According to the prompts for comparing two categorical variables, we provide prompts for visualizing data and for comparing groups (numerical variable vs. categorical variable).



**Figure 9:** Python code for the modified visualization of Facebook use in classes 5–Q2, differentiated by Gender (left) and data visualization via Python in form of comparable bar graphs (right): The x-axis shows the distribution of different grades (5–Q2), the y-axis displays the relative frequency of Facebook use. Male students are colored blue, female students are colored orange. EF corresponds to grade 10, Q1 corresponds to grade 11 and Q2 corresponds to grade 12.

To conclude the project work in weeks 3–4, the students are confronted with a change of perspective when analyzing data with digital tools: Statistical evaluation tools (CODAP/Jupyter notebooks) are programmed and programmable. Therefore, the students are supposed to realize that programming in Jupyter Notebooks can help to avoid limitations of a pre-programmed tool such as CODAP and can support the data exploration process for certain questions. For example, if one wants to compare the use of Facebook along with the several grades and to compare in each grade the proportion of Facebook use by gender, one needs a display that offers all the necessary information on the



three variables. In this case, it can help to use Jupyter Notebooks to display the data like in Figure 9 (right) as comparable bar graphs which cannot be realized in CODAP (CODAP does not allow to generate grouped bar graphs as we can realize for example in Jupyter Notebooks; see Figure 9, right). This issue shows a strong connection between statistics (choosing an adequate display and reading and interpreting the display) and computer science (writing the program code to realize the creation of such a display).

Finally, in week 5 the groups are asked to present their findings in the form of statistical reports via PowerPoint and to compare their results with the results of the official JIM-study to conclude module 1 of the pilot course.

## 5 Outlook

The pilot course started in September 2018 and was taught in tandem teaching (one statistics educator, one computer science educator). The fundamental and overarching idea of module 1 is the exploration of real and multivariate data traffic (noise data, JIM data) with digital tools (CODAP, Jupyter notebooks) in the frame of project-based learning. Those project works include statistical activities like the generation of adequate statistical questions (e.g. within the JIM data), data management and data cleaning (noise data), data visualization (noise data and JIM data) and applying statistical concepts like percentages, visualizations, and group comparisons (JIM data).

According to the design-based research principles, we collected several datasets like video and audio recordings of each lesson, additional field notes, and the students' working sheets. Based on this collection, we analyzed the lessons retrospectively in our research team and we adapted and refined the lesson structure, the materials and the activities for the next cycle in 2019/2020. Details on the course and its realization in school year 2019/2020 can be found in Budde et al. (2020).

**Acknowledgements** Our program was made possible by the Deutsche Telekom Stiftung. The authors would like to thank the participants of "Perspectives for Data Science Education at School Level – Educational Contributions from Statistics, Computer Science, and Sociocultural Studies"-symposium for their support and for their comments, which helped to develop a curriculum draft.

## References

- Arnold PM (2013) *Statistical Investigative Questions – An Enquiry Into Posing and Answering Investigative Questions From Existing Data*. Thesis. URL: <https://researchspace.auckland.ac.nz/handle/2292/21305>.
- Ben-Zvi D, Makar K, Garfield J (2018) *International Handbook of Research in Statistics Education*. Springer International Handbooks of Education, Springer International Publishing, Cham (Switzerland). DOI: 10.1007/978-3-319-66195-7.
- Berthold MR, Borgelt C, Höppner F, Klawonn F (2010) *Guide to Intelligent Data Analysis: How to Intelligently Make Sense of Real Data*. Texts in Computer Science, Springer, London (United Kingdom). DOI: 10.1007/978-1-84882-260-3.
- Biehler R, Ben-Zvi D, Bakker A, Makar K (2013) Technology for Enhancing Statistical Reasoning at the School Level. In: *Third International Handbook of Mathematics Education*, Springer International Handbooks of Education, Vol. 27, chap. 21, pp. 643–689. Springer, New York (USA), Clements MA, Bishop AJ, Keitel-Kreidt C, Kilpatrick J, Leung FKS (eds.). DOI: 10.1007/978-1-4614-4684-2\_21.
- Biehler R, Budde L, Frischemeier D, Heinemann B, Podworny S, Schulte C, Wasong T (eds.) (2018) *Paderborn Symposium on Data Science Education at School Level 2017: The Collected Extended Abstracts*. Universitätsbibliothek Paderborn, Paderborn (Germany). DOI: 10.17619/UNIPB/1-374.
- Breiman L (2001) *Statistical Modeling: The Two Cultures (With Comments and a Rejoinder by the Author)*. *Statistical Science* 16(3):199–231. URL: [https://projecteuclid.org/download/pdf\\_1/euclid.ss/1009213726](https://projecteuclid.org/download/pdf_1/euclid.ss/1009213726).
- Budde L, Frischemeier D, Biehler R, Fleischer Y, Gerstenberger D, Podworny S, Schulte C (2020) *Data Science Education in Secondary School: How to Develop Statistical Reasoning When Exploring Data Using CODAP*. In: Arnold P (ed.), *IASE Roundtable 2020: New Skills in Statistics Education*, International Association for Statistical Education (IASE).
- Burrill G, Biehler R (2011) *Fundamental Statistical Ideas in the School Curriculum and in Training Teachers*. In: *Teaching Statistics in School Mathematics-Challenges for Teaching and Teacher Education*, New ICMI Study Series, Vol. 14, pp. 57–69. Springer, Dordrecht (The Netherlands), Batanero C, Burrill G, Reading C (eds.). DOI: 10.1007/978-94-007-1131-0\_10.
- Cobb P, Confrey J, diSessa A, Lehrer R, Schauble L (2003) *Design Experiments in Educational Research*. *Educational Researcher* 32(1):9–13.
- Engel J (2017) *Statistical Literacy for Active Citizenship: A Call for Data Science Education*. *Statistics Education Research Journal* 16(1):44–49.
- Finzer W (2001) *Fathom Dynamic Statistics (v1.0)*. Key Curriculum Press. Current version is 2.1.

- Frischemeier D, Biehler R (2018) Stepwise Development of Statistical Literacy and Thinking in a Statistics Course for Elementary Preservice Teachers. In: Dooley T, Gueudet G (eds.), 10th Congress of the European Society for Research in Mathematics Education, DCU Institute of Education and ERME, pp. 756–763.
- Frischemeier D, Leavy A (2020) The Challenges for Prospective Primary Teachers When Constructing Statistically Worthwhile Questions. In: Jankvist UT, Van den Heuvel-Panhuizen M, Veldhuis M (eds.), 11th Congress of the European Society for Research in Mathematics Education, Freudenthal Group & Freudenthal Institute, Utrecht University and ERME, Utrecht (The Netherlands), pp. 938–945.
- Garfield J, Ben-Zvi D (2008) *Developing Students Statistical Reasoning – Connecting Research and Teaching Practice*, 1st edn. Springer, Dordrecht (The Netherlands). DOI: 10.1007/978-1-4020-8383-9.
- Gould R (2017) Data Literacy is Statistical Literacy. *Statistics Education Research Journal* 16(1):22–25.
- Heinemann B, Budde L, Schulte C, Biehler R, Frischemeier D, Podworny S, Wassong T (2018) Data Science and Big Data in Upper Secondary Schools: What Should Be Discussed From a Perspective of Computing Education? *Archives of Data Science, Series A* 5(1):143–159. DOI: 10.5445/KSP/1000087327/26.
- Huber PJ (2012) *Data Analysis: What Can Be Learned From the Past 50 Years*, Wiley Series in Probability and Statistics, Vol. 874. John Wiley & Sons, Hoboken (USA). ISBN: 978-1-118010-64-8.
- Konold C, Miller C (2011) *TinkerPlots 2.0*. Key Curriculum Press.
- Ridgway J (2016) Implications of the Data Revolution for Statistics Education. *International Statistical Review* 84(3):528–549. DOI: 10.1111/insr.12110.
- Wild CJ, Pfannkuch M (1999) Statistical Thinking in Empirical Enquiry. *International Statistical Review* 67(3):223–248. DOI: 10.1111/j.1751-5823.1999.tb00442.x