

BA Seminar: The ABC of Computational Text Analysis

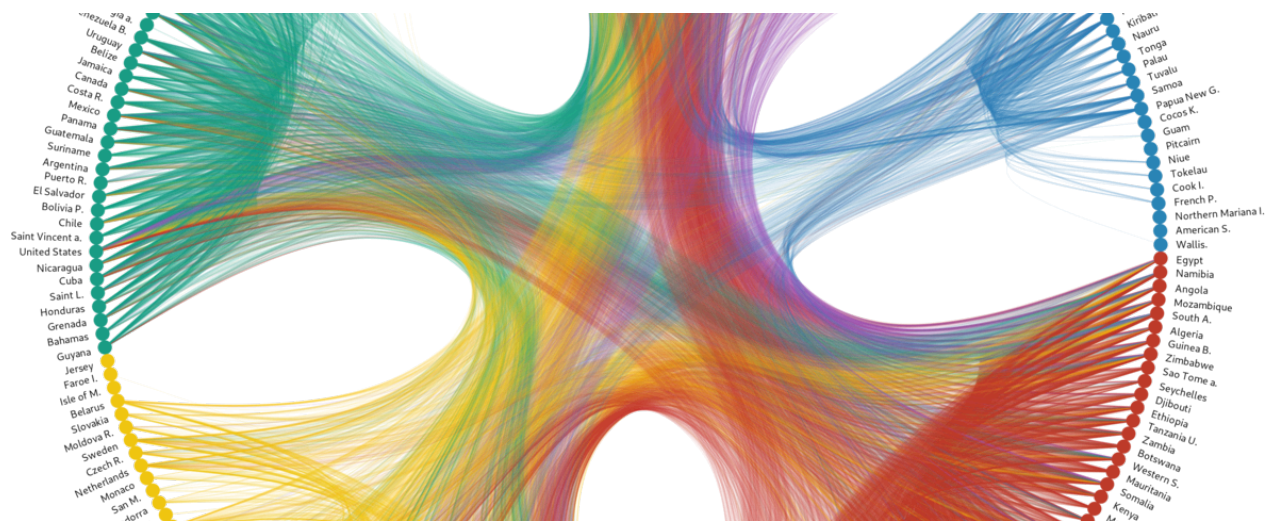
BA Seminar, Spring 2024, University of Lucerne

Alex Flückiger

16 December 2024

Contents

1	Schedule	2
2	Lectures	2
3	Assignments	5
4	Mini-Project	5
5	Optional Seminar Paper	6



In this hands-on seminar, bachelor students of social and cultural sciences learn the basics of programming, among other essential technical skills. Building on a modern technology stack, it aims to prepare students to conduct data-driven text analysis and to make everyday life easier by fostering technological fundamentals. While learning about the importance of computation in solving problems, we also discuss the current developments in information technology. In short, the course promotes digital literacy on a practical and theoretical level.

This seminar focuses on the computational processing of digital and digitized texts using Python and the command-line. For any empirical research, the systematic preparation and aggregation of data and the swift retrieval of information are critical. These tasks require the handling of various data forms, including data that is not yet structured in a tabular format. The seminar covers the complete workflow, from gathering textual data to analyzing an entire text collection to producing interactive visualizations. Sounds cool? It certainly is.

Along the way, we deal with questions like these:

- How can extensive collections of texts be evaluated quantitatively in order to strengthen content analysis?
- How has the discourse on a topic changed over time and how does communication between various actors differ?
- What are regular expressions and why are they so useful for text analytical applications?

1 Schedule

We have 12 seminar sessions together.

The plan below is provisional. I am happy to adapt the topics, as well as the schedule, to the needs and interests of the students. Likely, we will change some topics and orderings as we go.

Two lectures will be held via Zoom as I will not be in Lucerne.

Date	Topic
22 February 2024	Introduction + Where is the digital revolution?
29 February 2024	Text as Data
07 March 2024	Setting up your Development Environment
14 March 2024	Introduction to the Command-line
21 March 2024	Basic NLP with Command-line
28 March 2024	Introduction to Python in VS Code
04 April 2024	<i>no lecture (Osterpause)</i>
11 April 2024	Working with (your own) Data
18 April 2024	Data Analysis of Swiss Media
25 April 2024 (Zoom)	Ethics and the Evolution of NLP
02 May 2024 (Zoom)	NLP with Python
09 May 2024	<i>no lecture (Christi Himmelfahrt)</i>
16 May 2024	NLP with Python II + Working Session
23 May 2024	Mini-Project Presentations + Discussion
30 May 2024	<i>no lecture (Fronleichnam)</i>

2 Lectures

Below you find a brief description of all the lectures. I make the slides available before the lecture starts.

2.1 Week 1: Introduction + Where is the digital revolution?

On the one hand, I present the goals and organization of the seminar. On the other hand, we look at recent advancements in the area of artificial intelligence (AI) and digital humanities (DH) that offer an impression of the fascinating prospects of modern computing.

2.2 Week 2: Text as Data

Computational text analysis comes with many challenges that are unique due to the fuzziness of natural language. In this session, we learn about its methodological foundation and conduct our first computational text analysis to understand how this translates into practice.

2.2.1 Required Reading

- Lazer, David, Alex Pentland, Lada Adamic, Sinan Aral, Albert-László Barabási, Devon Brewer, Nicholas Christakis, Noshir Contractor, James Fowler, Myron Gutmann, Tony Jebara, Gary King, Michael Macy, Deb Roy, and Marshall Van Alstyne. 2009. “Computational Social Science.” *Science* 323(5915):721–23. (on OLAT)

2.2.2 Optional Reading

- Graham, Shawn, Ian Milligan, and Scott Weingart. 2015. *Exploring Big Historical Data: The Historian's Macroscope*. Open Draft Version. Under contract with Imperial College Press. [online](#)

2.3 Week 3: Setting up your Development Environment

The title says it all. We are getting ready for the practical part of the course: Programming. As installing Python and non-standard command-line tools may be tricky, we dedicate time in class rather than doing it as homework only. Additionally, I will also introduce some jargon that guide your way in the programmer's brave new world.

2.3.1 Optional: pimp your workflow

- Healy, Kieran. 2019. "The Plain Person's Guide to Plain Text Social Science." [online](#).

2.4 Week 4: Introduction to the Command-line

The command-line is a powerful tool at your disposal. It is the working horse for many data wrangling tasks. In this session, you learn the basics of shells and how to perform many operations effectively by substituting clicks on the screen with commands. Admittedly, it is not overly exciting at this stage, yet it is essential for more sophisticated automation later on.

2.4.1 Recommended Resources

- [Cheatsheet](#) for this course
- [The Programming Historian](#)
- [DigitalOcean](#)

2.5 Week 5: Basic NLP with the Command-line

Counting words is the most basic method to look at texts from a computational perspective. The command-line provides tools to sift through a massive text collection with ease. Combined with pattern-based search, you get quick results regardless whether you want to get a quantitative overview across documents, count particular words or see their usage in context. Does it sound like the Swiss knife in your research toolkit that marks the starting point of almost any textual analysis? It certainly is.

2.5.1 Optional Reading

- Ben Schmidt. 2019. [Regular Expressions](#)
- Amit Chaudhary. 2021. [A Visual Guide to Regular Expression](#)

2.5.2 Online Regular Expression Editor

- [regex101](#) is a visual editor to check your regular expressions.
- [regexone](#) provides an interactive regex tutorial


2.6 Week 6: Introduction to Python in VS Code

As the folks say, Python is among the coolest programming languages, relatively easy to learn, and provides excellent NLP packages so that you do not have to implement everything yourself. In this session, we begin with an introduction to the basic syntax of Python. Starting with basics may be dry stuff; however, it will allow you to use third-party libraries and get a handle on more sophisticated NLP analyses later on.

2.7 Week 7: Working with Data

Up to this point, you have acquired the skills to cut a document into pieces and then count any text elements. Without interesting data, these tools are neat but of little use. Thus, we turn to relevant data sources for computational text analysis. When dealing with plain text, your tools will cut through the data like butter. For other formats, such as PDF, we learn some tricks to convert them into plain text. In particular, we use optical character recognition (OCR).


You can open the static code or run the code directly in your browser without any installation using Binder:

-  Run in your browser via Binder (OCR not working on Binder)

2.8 Week 8: Data Analysis of Swiss Media

Analyzing the media discourse from the past and the present has always been one of the most exciting approaches for social science. Only recently, however, Swissdox@LiRI has introduced a simple-to-use API to assemble and download comprehensive datasets of Swiss news media. In this session, we learn about the process of data curation and data analysis using the Pandas package.

You can open the static code or run the code directly in your browser without any installation using Binder.

-  Run in your browser via Binder


2.9 Week 9: Ethics and the Evolution of NLP

Ethics is not just an abstract topic of Philosophy. Modern NLP is more powerful than ever and, thus, pervasive in many aspects of life. Unfortunately, it also exhibits severe and poorly understood biases that may cause harm. With the recent *turn to data-driven deep learning*, NLP overcame many theoretical limitations, but this comes at a cost. It is our duty to better understand the workings and impact of this technology.

2.10 Week 10: NLP with Python

Python is the language of choice when it comes to advanced NLP. Have you ever wondered how the frequency of terms has evolved over the years? Or how the language differs between two groups whereby the groups may be formed by any metadata (people, organization, gender etc.)? In such an exploratory endeavor, using an interactive and visual mode is very effective and complements statistics. In short, we finally arrived at the serious stuff in our journey. To make sure you don't get lost in the forest of yet unknown concepts, you will also learn some more jargon of NLP.


You can open the static code or run the code directly in your browser without any installation using Binder.

-  Run in your browser via Binder

2.11 Week 11: NLP with Python II + Working Session

We continue our deep dive into NLP with Python in today's session. It is the last piece of our puzzle. During this course, you have learned about the entire workflow, from assembling datasets of documents to analyzing their content and visualizing your findings. As soon as you have a structured text collection along with basic metadata (e.g., publication date), you can take numerous perspectives to look at your data.

You can open the static code or run the code directly in your browser without any installation using Binder.

-  Run in your browser via Binder

2.11.1 Explore interactively: 1 August Speeches by Swiss Federal Councilors

As a matter of tradition, Swiss Federal Councilors give an official speech on the Swiss National Day. Simon Schmid (Republik), in collaboration with Prof. Andreas Kley (Faculty of Law, UZH), collected many of these speeches and kindly shared the resulting dataset with me. The collection comprises 166 speeches, the newer ones are also publicly available [here](#). The original analysis of the dataset has been [published in the Republik](#).

The interactive visualization linked below shows how the language differs between speakers of *Social Democratic Party of Switzerland* (SP) and speakers of other parties. The top right corner shows terms that all parties have frequently used. In contrast, the top left and the lower right corner reveal words that have been used primarily by the members of the SP and correspondingly by the centre-right parties.

You can search for the terms of your interest. Moreover, you may click on the points in the plot to show the context of the corresponding words within speeches. These functions allow for a efficient investigation of the corpus along the dimensions of Swiss parties.

2.12 Week 12: Mini-Project Presentations + Discussion

In this session, it is your turn. Going beyond mere toy examples, you present what you have worked on and show off your first harvest of computational text analysis.

The seminar is coming to an end, yet it doesn't have to be a dead-end. You may have gotten more proficient in cursing your computer but you are also ready to fight your way through the jungle of technology. Keep going, cheers!

3 Assignments

You have to submit two assignments to complete the seminar successfully. The purpose of the assignments is not to make the course hard to pass but rather to foster your engagement with the topics covered.

#	Topic	Published	Deadline (by midnight)	Solution
1	Experimenting with LLMs	22 February 2024	1 March 2024	(individual feedback)
2	Word Counts using Python	11 April 2024	19 April 2024	Example solution

4 Mini-Project

You conduct a small computational text analysis and present the results in the final session. To give as many options as possible, you are free to choose your research question as well as the computational methods and data you will use. Working with data from your area of interest is certainly more fun.

This project doesn't aim to overwhelm students with requirements that are too ambitious. It should be the other way around. You will have as much freedom as you need to engage with your data creatively. The goal is that you realize that your knowledge is already good enough to perform powerful analyses. Complementing your claims with quantitative facts about the data is the only requirement.

More information will be provided later in the course.

4.1 Inspiring Student Projects

- [Doping in Swiss media](#), Aaron Steiner, Livia Köppel, Silvan Lachmuth, 2024
- [Crypto-currencies in Swiss media](#), Nick Rölly, Cedric Keiser, Nils Wolf, Valentin Casanova, 2023

- [Historical semantic of performance by University of Zurich](#), Denise Göddlin, Emma Notter, Melanie Meyer, 2023
- [Self-portrayal of Incels](#), Irfete Iseni, Asia Petrino, 2023
- [Communication of sustainability of Swiss supermarkets](#), Rebecca Ferraro & Khoana Hasanaj, 2023
- [Gender differences in 1 August speeches \[Code\]](#), Dario Haab, Valentina Meyer, Nils Brun, 2022

5 Optional Seminar Paper

You are welcomed to write an optional seminar paper (Hauptseminararbeit) for which you get additional credit points. As I am in the position of a guest lecturer, I will accept seminar papers in cooperation with [Prof. Sophie Mützel](#).

Due to the practical foundation of this seminar, you are well-prepared to subsequently apply computational text methods in a personal project. Although this is not a requirement, you may want to turn your mini-project into a seminar paper by deepening your empirical inquiry.

Students planning to write a seminar paper should send me an email with a short outline of their research idea until **15 May 2024** at the latest. When you would like to discuss your idea in person, feel free to do so any time after the seminar.

Requirements for the seminar paper (Hauptseminararbeit):

- Write your thesis in German or English.
- Use any computational methods to analyze your data.
- Your paper has a theoretical question guiding your methodical approach. In other words, methods are a means, not an end in themselves.
- Formal: 15 pages (A4), 12 pt Times New Roman, 3cm margin, 1.5 line spacing.
- Deadline for submitting the final paper: **31 August 2024**.