

The ABC of Computational Text Analysis

#7 WORKING WITH (YOUR OWN) DATA

Alex Flückiger
Faculty of Humanities and Social Sciences
University of Lucerne

06 April 2024

Recap last Lecture

- describe text as pattern using RegEx
- extract + replace textual parts 
 - literal: `text a b c`
 - meta: `\w \s [^abc] *`
 - power of `.`^{*}
- questions regarding assignment 2

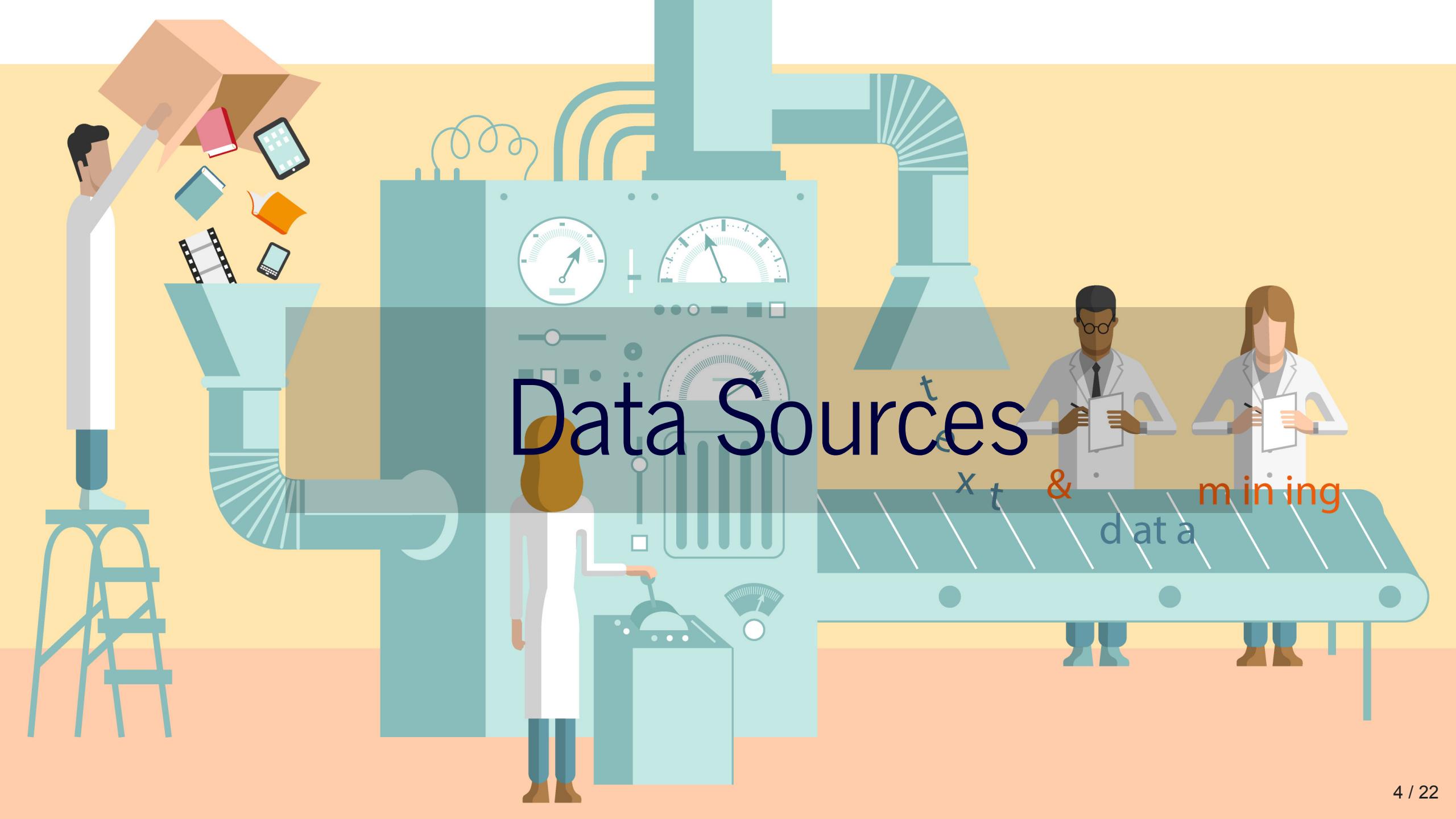
Outline

- learn about available data sources
- Build and justify (!) your own dataset

any text 

“any” format 

from anywhere 



An illustration depicting a data mining process. On the left, a man in a lab coat uses a shovel to move various items (books, a tablet, a film strip) into a large funnel. These items enter a teal-colored industrial machine. Inside the machine, a woman in a lab coat interacts with a control panel featuring several gauges and a bar chart. The machine has pipes and a fan at the top. On the right, two scientists in lab coats stand behind a conveyor belt. The conveyor belt has the words "Data Sources" written on it, along with "data mining" and "data". The background is a light yellow color.

Data Sources

What Data Sources are there?

- broadly social
 - newspapers + magazines
 - websites + social media
 - reports by NGOs/GOs
- scientific articles
- economic
 - business plans/reports
 - contracts
 - patents

👉 basically, any textual documents...

Some great (historical) Corpora

Ready off-the-shelf, machine-readable

- 1 August speeches by Swiss Federal Councilors
 - provided via [course repo](#)
- Human Rights Reports by various NGOs
- United Nations General Debate Corpus
- Inaugural Speeches by US Presidents

...



There are still not many.

Swissdox: A game changer

Assemble corpus and export as .csv

- over 250 Swiss newspapers
- historical and updated daily
- needs registration

The screenshot shows the 'Corpus query' page of the Swissdox@LiRI website. At the top, there is a navigation bar with links: Swissdox@LiRI, Start, Projects, Corpus query (which is highlighted in blue), Retrieved datasets, and Manual. The main area is titled 'Corpus query'. It contains three dropdown menus: 'Languages *' (with 'Select languages'), 'Source *' (with 'Select sources'), and 'Date ranges' (set to '2023-04-01 ~ 2023-04-30'). Below these is a note: '* no filtering is applied if no option is selected'. At the bottom are two buttons: 'Reset filters' and 'Next'.

Scraping PDF from Websites

Use case: Swiss voting booklet

- wget to download any files from the internet

```
# get a single file
wget EXACT_URL

# get all linked pdf from a single webpage
wget --recursive --accept pdf -nH --cut-dirs=5 \
--ignore-case --wait 1 --level 1 --directory-prefix=data \
https://www.bk.admin.ch/bk/de/home/dokumentation/abstimmungsbuechle

# --accept FORMAT_OF_YOUR_INTEREST
# --directory-prefix YOUR_OUTPUT_DIRECTORY
```

::: notes

- bis hierher: Wie Daten in txt Format bringen

More Interesting Publishers

- Project Gutenberg and HathiTrust
 - massive collection of books
 - open, HathiTrust requires agreement
- Nexis Uni
 - international newspaper, business + legal reports
 - licensed by the university
- Constellate
 - scientific articles of JSTOR across disciplines
 - provides an easy dataset builder



check out other resources licensed by ZHB

::: notes - Nexis ähnlich Swissdox, aber international - Zugang tlw. über ezproxy -
Constellate - kurze Demo von Constellate - brandneue Platform - einfache
Zusammenstellung von JSTOR Artikeln - sehr gute Metadaten - auch gut für
... (Illegible text)

Search existing Datasets

- Harvard Dataverse
 - open scientific data repository
- Google Dataset Search
 - Google for datasets basically

👉 search for a topic followed by **corpus, text collection or text as data**

::: notes - fortschrittliche Wissenschaft veröffentlicht nicht nur Papers, sondern auch Daten und Code - computergestützte Textanalyse ist aber immer noch Nische - Suchmaschinen für Datensätze - allerlei Datensätze, primär aus Wissenschaft :::

Search Techniques



Make your web search more efficient by using dedicated tags. Examples:

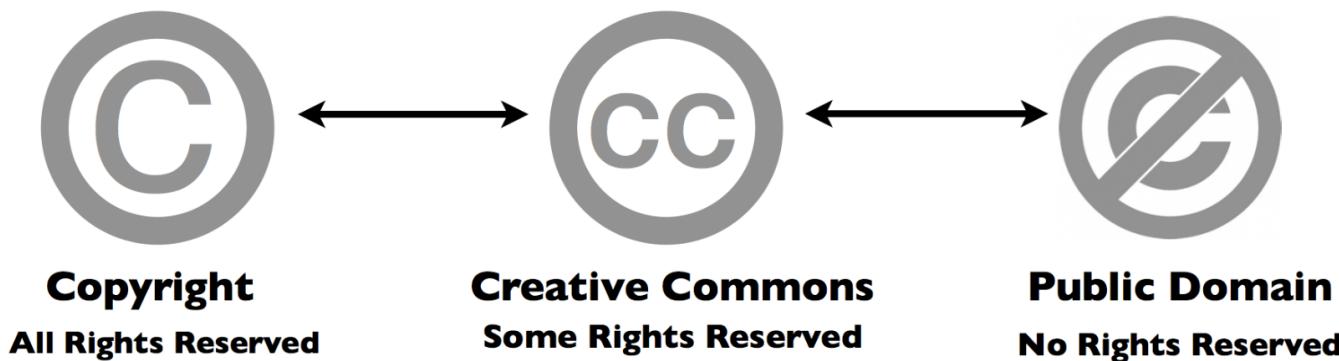
- "computational social science"
- site:nytimes.com
- nature OR environment

::: notes - Quotes für Wörter die zusammen gehören - Summe von Wörter haben evtl andere Bedeutung - Boolean Search - OR / AND :::

Data is Property

... and has rights too

- copyright may further limit access to high quality data
- check the rights before processing data



::: notes - Zugang zu Daten nicht immer einfach - open data unterschiedlich unterstützt - Datenbereitstellung oftmals Teil von Geschäftsmodell - dann restriktiv - oftmals ist Verwendung nicht geregelt - nutzt Graubereich :::

Interesting Resources

- **Party Programmes across Europe**
covers over 1000 parties from 1920 until today in over 50 countries
- **Swissvotes**
collection of resources on Swiss public votings
- **Swiss voting booklets**
from 1978 until today
 - 1 August speeches by Swiss Federal Councillors
 - Nestlé Annual Reports
 - University of Zurich Annual Reports
 - ... any organization of your interest 

What data are you interested in?

Think about the topic of your mini-project



| 6

Illustration of text analysis generated by [Image Creator from Microsoft Bing](#)

Preparing your own Data

.DOC

.JPG

.PNG

.PSD

.EPS

.CDR

.TXT

.GIF

.PPT

.MP3

.WAV

.AI

.MOV

.EXE

.DMG

.RAR

.ZIP

.PDF

**A world for humans ...
... and a jungle of file
formats.**

Common Conversions

news, press releases, reports from organizations



digital native documents
.pdf, .docx, .html



scans of (old) documents
.pdf, .jpg, .png



convert to .txt



Optical Character Recognition (OCR)

machine-readable A green square icon containing a white checkmark symbol.

Imperfect Data: A Tail of Bias

- **noise in text**
non-content (e.g. table of content), inconsistent spelling
- **archive holes**
lost or uncollected data
- **selective corpus curation**
supposition that key-word(s) captures topic
- **social bias**
view from somewhere, stereotypes



think about the data and mitigate issues



Questions?

In-class: Exercises I

1. Make sure that your local copy of the Github repository KED2024 is up-to-date with `git pull`. Check out the data samples in `materials/data` and the scripts to extract their text in `materials/code`.
2. Decide on one use-case that interests you most. Install the missing tool with the commands given on the respective slides (e.g., `pandoc`, `imagemagick`, `poppler`).
3. **Apply the commands to reproduce on the given data. Test them on your own data. Check the resources. Ask questions. Think about your mini-project.**

In-class: Exercises II

1. Use `wget` to download *cogito* and its predecessor *uniluAKTUELL* issues (PDF files) from the [UniLu website](#). Start with downloading one issue first and then try to automatize the process to download all the listed issued using arguments for the `wget` command.
2. Convert the *cogito* and *uniluAKTUELL* PDF files into TXT files using `pdftotext`. Try with a single issue first and then write a loop to batch process all of them.
3. What is the University of Lucerne talking about in its issues? Use the commands of the previous lectures to count the vocabulary.
4. Do the same as in 3.), yet analyze the vocabulary of *cogito* and *uniluAKTUELL* issues separately. Does the language and topics differ between the two magazines?

In-class: Exercises III

1. Use `wget` to download a book from Project Gutenberg and count some things (e.g., good/bad, joy/sad).
2. `wget` is a powerful tool. Have a look at its arguments and search for more examples in tutorials on the web. Have a nice Easter break! {data-background-image=../images/easter-eggs.jpg}