

u^b

Text Mining with DSL

Sukanya Nath
Ahmad Alhineidi



Image generated by DALLE openai

Data Science Lab (DSL)

- DSL is an interfaculty core facility of the University of Bern.
- We provide university-wide support and training for researchers and research groups in data science, machine learning, artificial intelligence and research IT related matters.
- Our Services range from rapid advice by e-mail and code or algorithm reviews to long-term collaborations on infrastructure and research projects with co-analysis and co-authoring.
- We Consolidate several research support activities (Digital Humanities, Science IT Support, Computer Vision, Text Mining,...)
- UniBe analogue to ETHZ SIS, UZH S3IT, UniBas SciCORE, EPFL, SDSC, ...
- Mandate from 2023-01-01 to 2026-12-31 (prolongation upon positive evaluation)

u^b

Interdisciplinary Team

(domain experts, technical experts and data scientists)

Sigve Haug
Coordinator / All Rounder



Guillaume Witz
Computer Vision
Support / Training



NLP Support /
ChatGPT & Co



Alexander Kashev
Research IT
Support



Mykhailo Vladymyrov
Machine Learning
Support



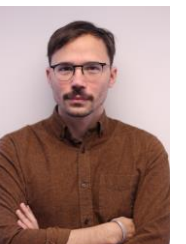
Ana Stojiljkovic
Computer Vision
Support



Ninoska Friedli
Administrative
Coordinator



Stephen Hart
Digital Humanities/
History



Peter Dängeli
Digital Humanities
Support



Marco Indermühle
HPC System
Administrator



Ahmad Alhineidi
NLP and HPC
Support



Sebastian Flick
Frontend developer,
Digital Hum. support



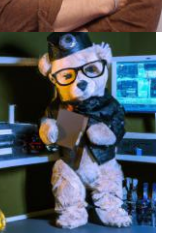
Sebastian Borkowski
Digital Humanities
Support



Sukanya Nath
NLP Support



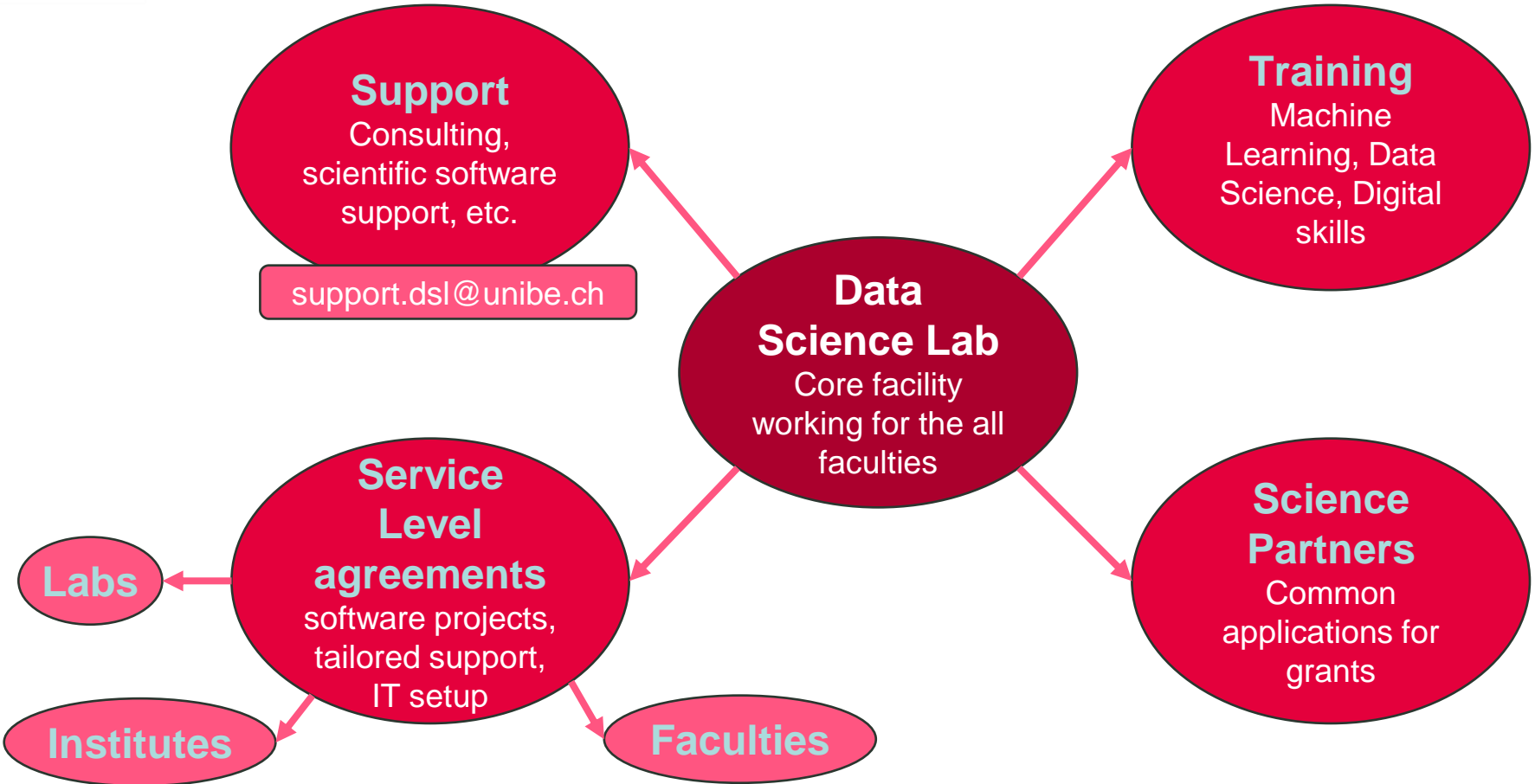
Michael Horn
Computer Vision
Support



Assistants

u^b

What is DSL?



u^b

DSL Walk-in Tuesdays

Do you have trouble running a data processing script?

Do you want to discuss how to use Machine Learning in your project?

Do you have any other question on Data Science or Research Software and IT?

Just come by the Data Science Lab and we will try to help you on the spot!

Data Science Lab Walk-In

Where: Main building (HS4), rooms 311 - 312

When: Tuesdays 13.00 - 15.00



DSL Trainings

DATA SCIENCE LAB

SUMMER TRAININGS 2024



www.dsl.unibe.ch

Mon, June 3 

Machine Learning with scikit-learn

A practical introduction to "classical" Machine Learning

Tue, June 4 

Introduction to Git and GitHub

Get started with version control.

Wed, June 5 

High Performance Computing (HPC) on UBELIX

Get started using the UniBe computing cluster

Thu, June 6 

Data Science and ML with MATLAB

An introduction to MATLAB for data science and ML

Fri, June 7

Organizing and Unions at the University (Mittelbau)

How can I organise myself at the university?

Mon, June 10 

Python packaging

Learn how to turn your scripts into an installable and testable package

Tue, June 11 

Web-scraping

Learn how to automate data collection from the web using Python

Wed-Fri, June 11.-12 

Data visualization with R and Quarto

Make beautiful visualizations and publish them

Max 25 people per course on a first come first served basis. **Free of charge** for all UniBe members.

Sign up now via ILIAS!

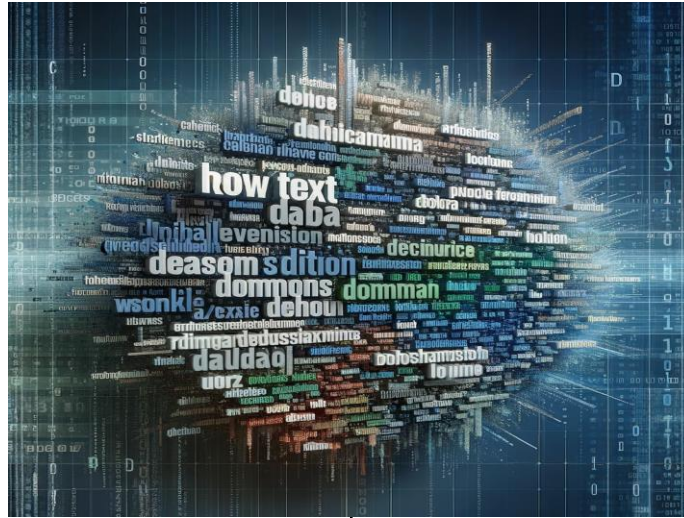


Winter schools

- Bern Winter School on Deep Learning
- Bern Winter School on Reinforcement Learning
- Bern Winter School on Natural Language Processing



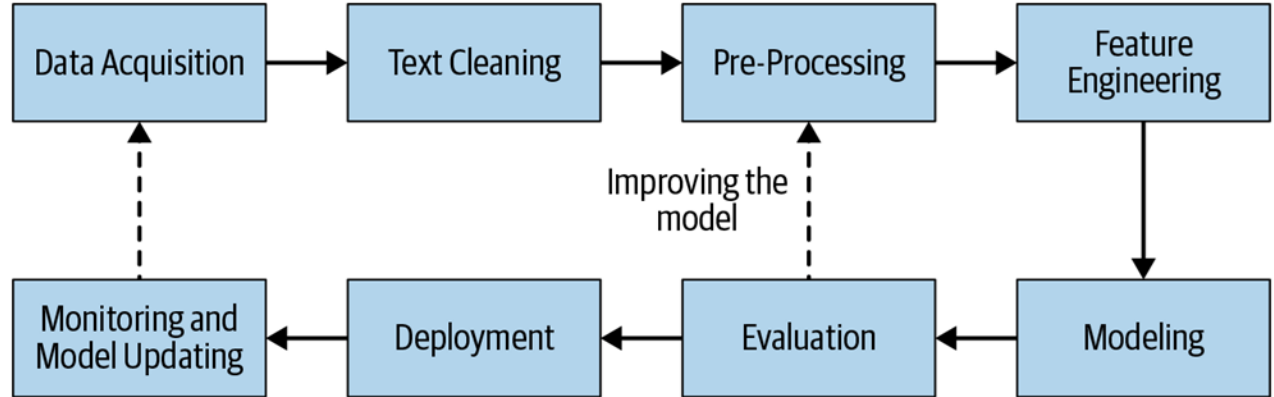
u^b Text Mining Overview



Knowledge!!

u^b Text Mining Overview

Generic
Text
Mining
Pipeline



Source: Vajjala et al. 2020

u^b Text Mining Overview

Text Mining problems

- Classification
- Summarization
- Information Extraction
- Named Entity Recognition
- Question Answering
- Topic Modelling
- Key Phrase Extraction
- Relation Extraction
- Syntactic/ Linguistic Feature Extraction
- Data Annotation/ Augmentation/ Anonymization

u^b

Text Mining Classification

- Binary or Multi-class classification
- Supervised ML

Example Usage:

- Sentiment analysis: Do people have a favourable opinion of the government?
- Hate speech detection: How does hate speech differ from free speech?
- Customized task: Is this text written by a man, woman or a bot?
- Spam detection
- Language identification

u^b

Text Mining Classification

- Token Classification (NER, POS)
- For Information Extraction
- Linguistics features and analysis

[Back to overview](#)

New Head of International Affairs Division

Bern, 15.4.2024 – From 1 August 2024, Barbara Schedler Fischer is to lead the International Affairs Division of the FOPH. Aged 48, she is a diplomat with extensive experience in bilateral and multilateral relations and is well acquainted with the Swiss political scene.



⚡ Inference API ⓘ

🔍 Token Classification

Examples ▾

Bern, 15.4.2024 – From 1 August 2024, Barbara Schedler Fischer is to lead the International Affairs Division of the FOPH. Aged 48, she is a diplomat with extensive experience in bilateral and multilateral relations and is well acquainted with the Swiss political scene.

Compute

Computation time on cpu: 0.041 s

Bern **LOC**, 15.4.2024 – From 1 August 2024, Barbara Schedler Fischer **PER** is to lead the International Affairs Division **ORG** of the **F ORG** **OPH ORG**. Aged 48, she is a diplomat with extensive experience in bilateral and multilateral relations and is well acquainted with the **Swiss MISC** political scene.

u^b

Text Mining Summarization

- Sequence-to-sequence models.
- Helpful for further steps such as text classification.
- Extractive or Abstractive.

Example usage:

- You want to know the change in stance of the UN with respect to a certain topic over time.
- You have access to the speeches made by UN secretary generals but these speeches are large, numerous and can contain topics irrelevant for your research.
- Summarization can help to condense the content of the speeches. It may also be possible to guide it in a certain direction.



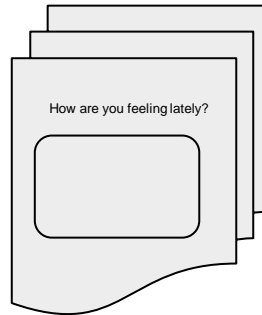
u^b

Text Mining Topic Modeling

- Unsupervised learning
- Based on clustering algorithms
- Extract knowledge from massive amount of documents

Example usage:

- You have performed an open ended survey of people suffering from anxiety.
- You want to know what topics were commonly mentioned.
- Topic Modelling can help to identify the topics and can also be guided towards a preferred direction.



u^b

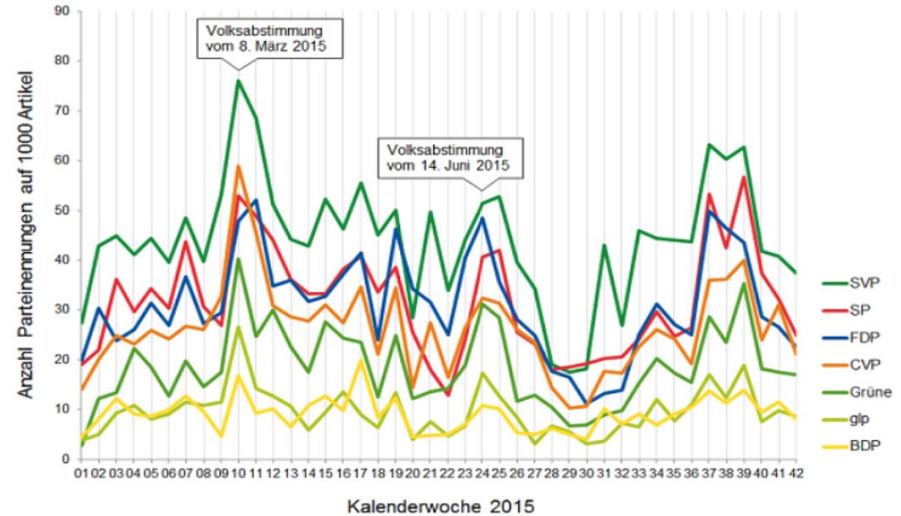
Data Annotation/ Augmentation/ Anonymization

- Data Annotation/ Labelling is the process of identifying and labelling the context of the text.
- Data Augmentation is the process of artificially generating newer data points from existing ones.
- Data Anonymization is the process of removing personally identifiable information from the data set.

u^b

Text Mining Output Example

- Number of mentions of Swiss Political parties
- How to achieve such text mining output?

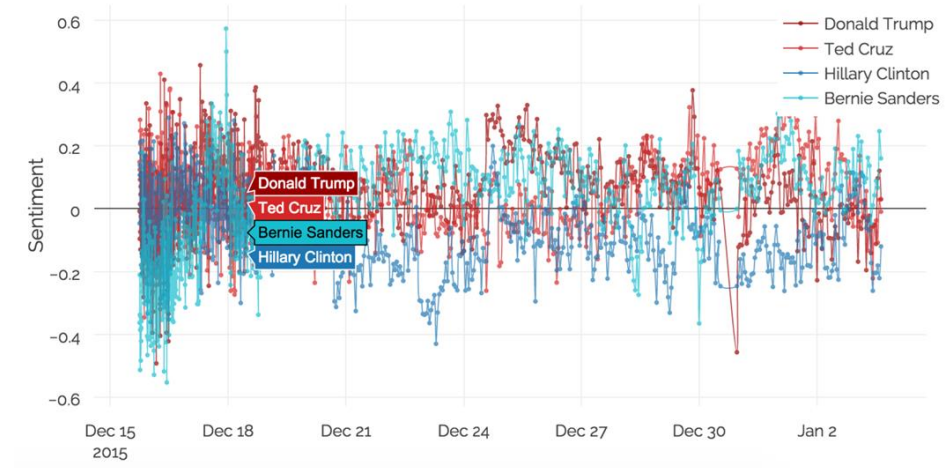


Source: Eurospider Newsletter

u^b

Text Mining Output Example

- Twitter data on sentiment of presidential candidates in 2015
- How to achieve such text mining output?



Source: Election Analysis and Data Science. [link](#)

u^b

Closing Thoughts

- Text Mining methods have improved massively over the past few years especially with the advent of generative AI.
- However
 - Models / Data may be prone to bias/ stereotypes.
 - Access to data must respect privacy laws.
 - Explaining and visualising the decisions by models is important but can be difficult.

u^b

THANK YOU!
QUESTIONS?