

TEXT ANALYSIS FOR SOCIAL SCIENCES IN PYTHON

(MSc in Economics, WiSo Universität Hamburg)

Winter Semester 2021 – 2022

Instructor : Huyen Nguyen

Website: <http://huyenttnguyen.com>

E-mail: huyen.nguyen-1@uni-hamburg.de

COURSE LOGISTICS & HOUSE RULES

- 1) **Lectures** are held in the computer room **WiWi 1005** on campus **every Wednesday** morning, from 13.10.2021 to 26.01.2022.

The *only exception* is Week 8 (i.e. Wed 08/12 session). During this week, there will be individual group consultation appointments via Zoom to discuss your mid-term project proposal assignment (more details below).

Every lecture is structured as follows: 45' of lecture → 15' break → 45' code exercise (self) → (5' break) → 45' code exercise (team) (from Week 4 onwards).

Classes start at **8.15 AM** sharp.

Students are expected to read the required readings and do assignments before class, as well as actively participating in all classes. Attendance on at least 11/14 sessions is required.

- 2) **Course materials** (slides, exercise links, data sources, tips) are uploaded **on GitHub** <https://github.com/htn21uhh/Text-Analysis-for-Social-Sciences-in-Python>
 - **Lecture slide** will be added to GitHub → Lecture Slides folder **every Thursday morning** of the same week, unless otherwise noticed. The list of required readings and exercise links are in the detailed syllabus overview below.
 - **Exercises** (on Jupyter notebook/Google Colab) that are used in class (and accessible after class for your own practice) will be updated in Github → Exercises folder.
 - **Data sources and tips** are updated regularly in the detailed Syllabus below as the course proceeds.
- 3) **Course forum** (introduction, teammate mix-and-match, exchange of ideas, code bugs) are **on Slack**: https://join.slack.com/t/textanalysisf-eia6533/shared_invite/zt-wkqywq13-3OQMM~S~gShngJ~VcTP2~Q

How to join in this Slack channel? Use your @uhh account to create an account, register and log into the channel (Desktop apps are also downloadable) using the above link.

NOTE (!) Do NOT share Slack channels to any non-course participants (!) Any requests must go through me first. Otherwise, I will personally kick them out and you will receive a minus point for doing so.

There are 4 channels, as follows:

- *#introduce-yourself*: [COMPULSORY, during 1st session] Introduce yourself, your research interest and your goals with the course
- *#group-matching* [COMPULSORY, deadline Wednesday 27.10]: For the project mid-term proposal and final match yourself to groups found on the Google Sheet in the top of the channel ([or use the same link here](#)). (!) **DO NOT CHANGE any details of other classmates.**
- *#qna-code-bugs*: The Q&A space for code issues, solutions, Stackoverflow materials and any useful exercises you found.
- *#research-projects*: The Q&A space across groups to discuss research project ideas, proposals, planning, and recording resources.

NOTE: Announcement and notices will be done mainly via STINE, as well as the #announcement channel in Slack forum.

COURSE SUMMARY

We are living in a rapidly digitized world, with an ever-increasing availability of large-scale textual corpora in law, politics and economics. This massive data development scene poses exciting challenges for social scientists, to understand the fabrics and functioning of our societies, beyond just numbers. Coupling the proliferation of legal and political corpora with the speedy growth of data science toolkits, we have at hand a powerful infrastructure to extract hidden novel insights about relevant institutional and human patterns in texts.

This course gives comprehensive introduction to the basic theory and hands-on applications of text analysis and machine learning for social science in Python. The course begins with quick introduction to Python languages and moves on to the challenge of representing texts as data. Next, it gives an overview of key techniques to clean texts, extract relevant information, and represent documents as vectors. These techniques include, but not limited to, for instance, measuring document similarity, clustering documents based on topics, as well as visualization methods such as word clouds and spatial relation plots between documents. Students are also provided with various sources to different text corpora, tips and techniques to query a programming issue online and self-study materials to deepen their understanding beyond the scope of the course.

Finally, we consider text-based prediction problems. For instance, given the evidence of a particular case, how will a judge decide on sentences? Given recorded speeches and transcripts of politicians, how ideological is a politician? Such predictions are then incorporated into social science analysis. Students will investigate and implement the relevant machine learning tools for making these types of predictions, including regression, classification, and deep learning models. If time permits, we will also touch upon causal inference methods using texts, either as treatment or outcome in a given data context.

Acknowledgement: The structure and materials of this course are drawn on a combination of pioneering text analysis courses of Elliot Ash (ETH Zurich), Brandon Stewart (Princeton), William Lowe (Hertie School) and Caroline Le Pennec-Caldichoury (HEC Montreal).

COURSE OBJECTIVES

- Comprehensive overview of contemporary approaches in using Python to analyze text as data, and how it is applied in social science policy questions.
- Familiarity with statistical and practical issues around textual data.
- Ability to fit, interpret and apply the basic classes of text cleaning and analysis techniques to independent research projects.

PRE-REQUISITES

This course requires a basic understanding of Python, statistics, probability theories and applied econometric techniques used in social sciences. The first three sessions will quickly introduce you to Python basics and essential packages for text analysis, then we will quickly move onto text as data and relevant methods.

ALL students are required to check their math and stats background. you don't have a solid background in calculus, linear algebra, and probability, read [part 1 from this online book](#).

GETTING STARTED

1) Python

The examples in the course will mainly use Python. Additionally, when we move to regression exercises, Stata might be used in economic examples. You are strongly recommended to install and set up Python before the course, using the instructions in the following links:

[Python Setup Instructions](#)
[Codecademy Online Python Course](#)

2) Jupyter Notebook

Read the following step-by-step blog to get familiar with Jupyter Notebook.

<https://medium.com/codingthesmartway-com-blog/getting-started-with-jupyter-notebook-for-python-4e7082bd5d46>

For installation and trying out Jupyter Notebook, read this:

<https://jupyter.readthedocs.io/en/latest/install/notebook-classic.html>

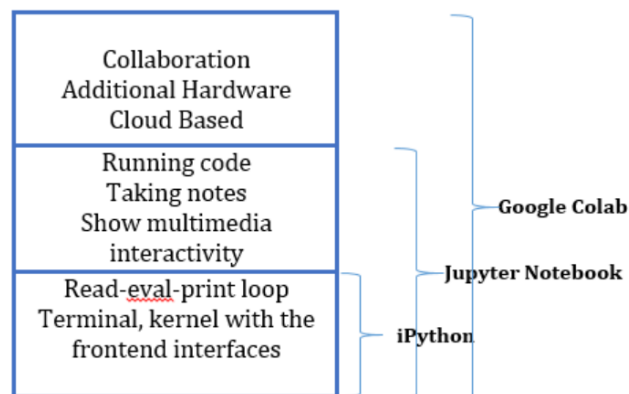
<https://satyaborg.com/posts/accessing-virtualenv-in-jupyter-notebooks>

3) Google Colab

If Jupyter Notebook on the web does not work for you during class exercise sessions, [use Google Colab](#).

Google Colab is a specialized version of Jupyter Notebook, which runs on the cloud and offers free computing resources. It does not require a setup, plus the notebooks that you will create can be simultaneously edited by your team members – in a similar manner you edit documents in Google Docs. The greatest advantage is that Colab supports most popular [machine learning libraries](#) which can be easily loaded in your notebook. In order to use Google Colab, you will need a Google Account. Your files are by default linked to your Google Drive. For more information and instructions, please check the following link: https://colab.research.google.com/?utm_source=scs-index

The following diagram summarizes the relationship between iPython, Google Colab and Jupyter Notebook:



4) StackOverflow <https://stackoverflow.com/>

Your go-to website (in addition to Google) for ANY programming bug, installation, package issues. Good programming and data wrangling skills = good Googling/StackOverflow skills.

NOTE! I will IGNORE any questions and e-mails related to programming/installation/package issues. The answers are highly likely available on StackOverflow or from your classmates (use Slack forum to ask, or ask them directly in class).

COURSE GRADING POLICY

The course will be graded on the usual grading scale with passing grades from 1.0 (very good) to 4.0 (sufficient), and with a failing grade 5.0 (insufficient). The grades are determined as follows:

30%: Midterm research pitch recording (5') + 1-page research proposal

(Deadline: 23.59 Wednesday 24.11)

50%: Final group research project oral presentation (25')

20%: Reproducible code package of your research

Additionally, all students are required to submit a confidential written contribution of him/her/themselves and what other group members of his/her/their own team have contributed. In general, free-riding and foul play are never an option in joint projects.

Submission rule: For every hour of late submission of the midterm pitch & proposal, everyone in the team will receive a 5% reduction of the score. Should you have verifiable extenuating circumstances (e.g. illness, personal loss, hardship, or caring duties), an extension can be granted. In such cases, please contact the course instructor as soon as possible before the deadline.

You can find more detailed instructions on the pitch, proposals, the consultation sessions and final oral presentations in a separate file in Github/Open Olat.

COURSE SYLLABUS & READING MATERIALS

This detailed syllabus is subject to changes as the course proceeds. I will update it regularly on the Github repository, please check it on a weekly basis.

Required readings are to be read and analyzed thoroughly. Optional readings [OPT] are intended to broaden your knowledge in the respective area and it is recommended to at least get a gist of these materials. We will go through the step-by-step exercises together in class, and you are required to go through any remaining parts on your own before the next lecture.

Week	Date & location	Topic & Content	Readings	Exercises
1	Wed 13.10	Welcome & Installations Python basics – data types, variables, strings, list vs. dict, functions	A beginner's guide to Python https://wiki.python.org/moin/BeginnersGuide	W1_Python_basics.ipynb W1_Dictionaries and Collections; List; Variables, Expressions, Statements.pdf [OPT] Check this free course on Practical Python Programming
2	Wed 20.10	Python basics – loops, statements, logical operators, file handling, numpy	(same as above)	W2_Python_basics_cont.ipynb numpy tutorial [OPT] Hands-on exercises and

				solutions for Numpy Array
3	Wed 27.10 DEADLINE @ 23.59: group matching	Python basics – matplotlib, pandas, nltk, scikit-learn	Check the exercise materials and Python for Data Science Cheat Sheets (1-pager of the most useful commands for Python Basics, Jupyter Notebook, numpy, pandas, scikit-learn, matplotlib, seaborn and bokeh)	W3_Python(matplotlib, pandas, nltk, scikit-learn).ipynb Datasets: (pandas) titanic.csv (nltk) MobyDick.txt intro to scikit-learn overview pandas tutorial.
4	Wed 03.11	Text as data – state-of-the-art technique overview, NLP research design, pre-processing corpora (Part 1)	Gentzkow, Kelly, and Taddy, “ Text as Data. ” Grimmer and Stewart, “ Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. ” Denny and Spirling (2018), “ Text Preprocessing for Unsupervised Learning: Why It Matters, When It Misleads, and What to Do about It ” Ken Benoit (2020), “ Text as Data – an Overview ” Is your NLP project feasible? https://medium.com/@nqabell89/is-your-nlp-project-feasible-342c335d9ebe	W4_preprocessing_text_data.ipynb [HOME EXERCISE] W3_Python_refresher_6_ex_challenge.ipynb Check Chapter 4 to 6 of the NLTK book Text Analysis Starter Guide
5	Wed 10.11	Text as data – pre-processing corpora (Part 2)	Obtain data from the Internet – a guide to crawling for management	W5_feature_extraction_Amazon.ipynb

		& webscraping basics	<p>research (Claussen & Peukert)</p> <p>An example to webscraping with Selenium package https://towardsdatascience.com/how-to-use-selenium-to-web-scrape-with-example-80f9b23a843a</p> <p>Common pitfalls with pre-processing German language https://medium.com/ideal-o-tech-blog/common-pitfalls-with-the-preprocessing-of-german-text-for-nlp-3cfb8dc19ebe</p> <p>Text formatting with Regex https://medium.com/factory-mind/regex-tutorial-a-simple-cheatsheet-by-examples-649dc1c3f285</p> <p>https://medium.com/factory-mind/regex-cookbook-most-wanted-regex-aa721558c3c1</p> <p>Sebastian Raschka, "Turn your Twitter Timeline into a Word Cloud"</p>	<p>W5_webscraping.ipynb</p> <p>Step-by-step preprocessing example https://medium.com/hackerdawn/the-simple-nlp-text-preprocessing-tutorial-you-were-waiting-for-2e6b082962b2</p> <p>[OPT] Webscraping Twitter data https://towardsdatascience.com/hands-on-web-scraping-building-your-own-twitter-dataset-with-python-and-scrapy-8823fb7d0598</p>
6	Wed 17.11	Supervised methods – dimensionality reduction (PCA vs. PLS) , train-test-split data sets	<p>Maitra and Yan, PCA and PLS for Regression * Feature Selection in Scikit-Learn</p> <p>Training vs. validation vs. test sets https://tarangshah.com/blog/2017-12-03/train-validation-and-test-sets/</p>	<p>W6_dimreduction.ipynb</p> <p>[OPT] Sentiment analysis with scikit-learn https://www.twilio.com/blog/2017/12/sentiment-analysis-scikit-learn.html</p>

			<p>Webscraping financial data and sentiment analysis https://www.linkedin.com/pulse/impact-news-stock-price-web-scraping-sentiment-2-part-hamid/</p> <p>Econometrics and sentiment analysis https://core.ac.uk/display/323232776?source=2 (Econometrics meets sentiment : an overview of methodology and applications 2020)</p>	
7	<p>Wed 24.11 DEADLINE @ 23.59: 5' pitch recording + 1-page research proposal</p>	<p>Supervised methods- (part 2) model evaluation methods, classification vs. regression ML techniques</p>	<p>[IMPORTANT] How to select ML algorithms for your analysis (with Cheatsheet)</p> <p>Metrics to evaluate your ML algorithms</p> <p>Leon Derczynski, "Collocations," http://www.derczynski.com/sheffield/teaching/inno/7c.pdf</p> <p>Yoav Goldberg and Omer Levy, "Word2Vec explained: Deriving Mikolov et al's Negative Sampling Word Embedding Method".</p> <p>Piero Molino, "Word embeddings: Past, present, and future".</p>	<p>[OPT] Stata to Python equivalent http://www.danielmsullivan.com/pages/tutorial_stata_to_pyhon.html</p>
8	<p>Wed 01.12 (consultation week)</p>	<p>(A big, step-by-step coding homework using all techniques so far) + 20' consultation session/group</p> <p>+ self-reading on space, cosine similarity, collocations, word embedding (W7 reading list)</p>		
9	<p>Wed 08.12</p>	<p>Unsupervised Methods – Word Embeddings</p>	<p>Matt Kusner, Yu Sun, Nicholas Kolkin, and Killian Weinberger, "From word</p>	<p>Andy Thomas, A Word2Vec Keras Tutorial</p>

			embeddings to document distances ". [OPT] Garg et al 2018, Word embeddings quantify 100 years of gender and ethnic stereotypes	
10	Wed 15.12	Unsupervised Methods - Topic Models (vs. Word Embeddings)	Brandon Rose, " Document clustering in python ." Blei, Ng and Jordan, 2003. " Latent Dirichlet Allocation ," <i>Journal of Machine Learning Research</i> . Topic Modelling: A deep dive into LDA, hybrid-LDA, and non-LDA approaches [CHRISTMAS READING] BERT word embeddings BERT embeddings tutorial	W10_DocEmbedding s_TopicModels.ipynb (+ package files: utils.py, txt_utils.py) Shivam Bansal, " Beginner's guide to topic modeling in Python ."
Christmas holiday (19.12.2021 – 02.01.2022)				
11	Wed 05.01.22	Machine Learning Refresher Challenge: Applications	Spiess and Mullainathan, Machine Learning: An Applied Econometrics Approach , <i>Journal of Economic Perspectives</i> . Geron's book Chapter 4	TBD
12	Wed 12.01.22	Causal inference – Text as Treatment	Text and Causal Inference: A Review of Using Text to Remove Confounding from Causal Estimates Katherine A. Keith, David Jensen, and Brendan O'Connor How to Make Causal Inferences Using Texts Naoki Egami, Christian J. Fong, Justin Grimmer,	Tutorial: Causal inference and Machine Learning in Practice with Industry cases (Microsoft, Tripadvisor and Uber)

			<p>Margaret E. Roberts, and Brandon M. Stewart</p> <p>[OPT] A comprehensive collection of research papers on causal inferences and language https://github.com/causaltext/causal-text-papers</p>	
13	Wed 19.01.22	Causal inference – Text as Outcome	<p>Causal Mediation Analysis for Interpreting Neural NLP: The Case of Gender Bias Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer and Stuart Shieber</p> <p>A Survey of Online Hate Speech through the Causal Lens Antigoni M. Founta, Lucia Specia</p>	<p>Tutorial: Causal inference and Machine Learning in Practice with Industry cases (Microsoft, Tripadvisor and Uber)</p>
14	Wed 26.01.22	Causal Inference – Text as Mediators	TBD	<p>Tutorial: Causal inference and Machine Learning in Practice with Industry cases (Microsoft, Tripadvisor and Uber)</p>
FINAL	TBA	Group project oral presentation (25') + Replicable code package		

SUPPLEMENTARY MATERIALS

The course will be mainly based on the weekly lecture slides and in-class exercises, but the following books and code exercises can be used as reference along with the slide content.

- Natural Language Processing in Python, Third Edition, available at nltk.org/book.
- Aurelien Geron, Hands-On Machine Learning with Scikit-Learn & TensorFlow, O'Reilly 2017 ([link](#))

- [Jupyter notebooks Github for Geron's book.](#)
- [Google Developers Text Classification Guide](#) (This guide contains some practical tips and code examples for using text data)
- For Python syntax programming, this book Fluent Python (O'Reilly 2015) serves as a good guideline <https://www.oreilly.com/library/view/fluent-python/9781491946237/>