

Applied Language Technology: NLP for the Humanities

Tuomo Hiippala

Department of Languages, University of Helsinki

P.O. Box 24, Unioninkatu 40 B.

00014 University of Helsinki, Finland

tuomo.hiippala@helsinki.fi

Abstract

This contribution describes a two-course module that seeks to provide humanities majors with a basic understanding of language technology and its applications using Python. The learning materials consist of interactive Jupyter Notebooks and accompanying YouTube videos, which are openly available with a Creative Commons licence.

1 Introduction

Language technology is increasingly applied in the humanities (Hinrichs et al., 2019). This contribution describes a two-course module named *Applied Language Technology*, which seeks to provide humanities majors with a basic understanding of language technology and practical skills needed to apply language technology using Python. The module is intended to *empower* the students by showing that language technology is both accessible and applicable to research in the humanities.

2 Pedagogical Approach

The learning materials seek to address two major pedagogical challenges. The first challenge concerns terminology: in my experience, the greatest hurdle in teaching language technology to humanities majors is not ‘technophobia’ (Öhman, 2019, 480), but the technical jargon that acts as a gatekeeper to knowledge in the field (cf. Maton, 2014). This issue is fundamental to teaching students with no previous experience of programming. To exemplify, beginners in my class occasionally interpret the term ‘code’ in phrases such as “Write your code here” as a numerical code needed to unlock an exercise, as opposed to a command written in a programming language. For this reason, the learning materials introduce concepts in Python and language technology in layperson terms and gradually build up the vocabulary needed to advance beyond the learning materials.

The second challenge involves the diversity of the humanities, which covers a broad range of disciplines with different epistemological and methodological standpoints. This results in considerable differences in previous knowledge among the students: linguistics majors may be more likely to be exposed to computational methods and tools than their counterparts majoring in philosophy or art history. Some students may have taken an introductory course in Java or Python, whereas others have never used a command line interface before. To address this issue, the learning materials are based on Jupyter Notebooks, which provide an environment familiar to most students – a web browser – for interactive programming. The command line is used for interacting with GitHub, which is used to distribute the learning materials and exercises.

The module also emphasises peer and collaborative learning: 20% of the course grade is awarded for activity on the course discussion forum hosted on GitHub. All activity – both asking and answering questions – counts positively towards the final grade. This allows the students with previous knowledge to help onboard newcomers. According to student feedback, this also fosters a sense of community. The discussion forum is also used to discuss weekly readings, which focus on ethics (e.g. Hovy and Spruit, 2016; Bird, 2020) and the relationship between language technology and humanities (e.g. Kuhn, 2019; Nguyen et al., 2020). These discussions are guided by questions that encourage the students to draw on their disciplinary backgrounds, which exposes them to a wide range of perspectives to language technology and the humanities.

3 Learning Materials

The learning materials cover two seven-week courses.

The first course starts by introducing rich, plain and structured text and character encodings, fol-

lowed by file input/output in Python, common data structures for manipulating textual data and regular expressions. The course then exemplifies basic NLP tasks, such as tokenisation, part-of-speech tagging, syntactic parsing and sentence segmentation by using the spaCy 3.0 natural language processing library (Honnibal et al., 2020) to process examples in the English language. This is followed by an introduction to basic metrics for evaluating the performance of language models. The course concludes with a brief tour of the pandas library for storing and manipulating data (McKinney, 2010).

The second course begins with an introduction to processing diverse languages using the Stanza library (Qi et al., 2020), shows how Stanza can be interfaced with spaCy, and how the resulting annotations can be searched for linguistic patterns using spaCy. The course then introduces word embeddings to provide the students with a general understanding of this technique and its role in modern NLP, which is also increasingly applied in research on the humanities. The course finishes with an exploration of discourse-level annotations in the Georgetown University Multilayer Corpus (Zeldes, 2017), which showcases the CoNLL-U annotation schema.

To what extent the students meet the learning objectives is measured in weekly exercises. The weekly assignments are distributed through GitHub Classroom and automatically graded using *nbgrader*¹, which allows generating feedback files with comments that are then pushed back to the student repositories on GitHub. The exercises are also revisited in weekly walkthrough sessions to allow the students to ask questions about the assignments. The students are also required to complete a final assignment for both courses: the first course concludes with a group project that involves preparing a set of data for further analysis, whereas the second course finishes with a longer individual assignment.

All learning materials are openly available with a Creative Commons 4.0 CC-BY licence at the addresses provided in the following section. Access to the weekly exercises is available on request.

4 Technical Stack

The learning materials are based on Jupyter Notebooks (Kluyver et al., 2016) hosted in their own

¹<https://nbgrader.readthedocs.io>

GitHub repository.² This repository constitutes a submodule of a separate repository for the website, which is hosted on ReadTheDocs.³ The notebooks containing the learning materials are rendered into HTML using the Myst-NB parser from the Executable Books project.⁴ This allows keeping the learning materials synchronised, and enables the users to clone the notebooks without the source code for the website. Myst-NB also adds links to Binder (Project Jupyter et al., 2018) to each notebook on the ReadTheDocs website, which enables anyone to execute and explore the code.

The Jupyter Notebooks provide a familiar environment for interactively exploring Python and the various libraries used, whereas the ReadTheDocs website is meant to be used as a reference work. Both media embed videos from a YouTube channel associated with the courses.⁵ These short explanatory videos exploit the features of the underlying audiovisual medium, such as overlaid arrows, animations and other modes of presentation to explain the topics.

5 Conclusion

This contribution has introduced a two-course module that aims to teach humanities majors to apply language technology using Python. Targeted at a student population with diverse disciplinary backgrounds and levels of previous experience, the learning materials use multiple media and layperson terms to build up the vocabulary needed to engage with Python and language technology, complemented by the use of a familiar environment – a web browser – for interactive programming using Jupyter Notebooks.

References

- Steven Bird. 2020. *Decolonising speech and language technology*. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3504–3519.
- Erhard Hinrichs, Marie Hinrichs, Sandra Kübler, and Thorsten Trippel. 2019. Language technology for digital humanities: introduction to the special issue.

²<https://github.com/Applied-Language-Technology/notebooks>

³<https://applied-language-technology.readthedocs.io>

⁴<https://executablebooks.org>

⁵<https://www.youtube.com/c/AppliedLanguageTechnology>

- Language Resources and Evaluation*, 53(4):559–563.
- Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python. DOI: 10.5281/zenodo.1212303.
- Dirk Hovy and Shannon L. Spruit. 2016. [The social impact of natural language processing](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 591–598. Association for Computational Linguistics.
- Thomas Kluyver, Benjamin Ragan-Kelley, Fernando Pérez, Brian Granger, Matthias Bussonnier, Jonathan Frederic, Kyle Kelley, Jessica Hamrick, Jason Grout, Sylvain Corlay, Paul Ivanov, Damián Avila, Safia Abdalla, Carol Willing, and Jupyter development team. 2016. [Jupyter Notebooks – a publishing format for reproducible computational workflows](#). In *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, pages 87–90, Netherlands. IOS Press.
- Jonas Kuhn. 2019. Computational text analysis within the humanities: How to combine working practices from the contributing fields? *Language Resources and Evaluation*, 53(4):565–602.
- Karl Maton. 2014. *Knowledge and Knowers: Towards a Realist Sociology of Education*. Routledge, New York and London.
- Wes McKinney. 2010. Data structures for statistical computing in Python. In *Proceedings of the 9th Python in Science Conference*, pages 51–56.
- Dong Nguyen, Maria Liakata, Simon DeDeo, Jacob Eisenstein, David Mimno, Rebekah Tromble, and Jane Winters. 2020. [How we do things with words: Analyzing text as social and cultural data](#). *Frontiers in Artificial Intelligence*, 3.
- Emily Öhman. 2019. [Teaching computational methods to humanities students](#). In *Proceedings of the Digital Humanities in the Nordic Countries 4th Conference*, volume 2364 of *CEUR Workshop Proceedings*, pages 479–494.
- Project Jupyter, Matthias Bussonnier, Jessica Forde, Jeremy Freeman, Brian Granger, Tim Head, Chris Holdgraf, Kyle Kelley, Gladys Nalvarte, Andrew Osheroff, M Pacer, Yuvi Panda, Fernando Perez, Benjamin Ragan Kelley, and Carol Willing. 2018. [Binder 2.0 – Reproducible, interactive, sharable environments for science at scale](#). In *Proceedings of the 17th Python in Science Conference*, pages 113 – 120.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A Python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 101–108. Association for Computational Linguistics.
- Amir Zeldes. 2017. [The GUM corpus: Creating multilayer resources in the classroom](#). *Language Resources and Evaluation*, 51(3):581–612.