

# Automated Text Analysis in Political Science

- POLS 5198 – Automated Text Analysis in Political Science
- Lecturer: Martijn Schoonvelde
- Email: [mschoonvelde@gmail.com](mailto:mschoonvelde@gmail.com)
- Credits: 2.0
- Program: 1 Year MA Political Science, 2 Year MA Political Science
- Spring Term 2017–2018
- Dates: 16–27 April 2018
- Course requirements: familiarity with R
- Office hours: upon appointment, either through email or in person

## Course introduction

Automated text analysis has become very popular across the social sciences over the last few years. With the massive availability of text data on the web, political scientists increasingly recognize automated text analysis as viable approach to analyzing social and political behavior. This course – in which we use R – introduces students to a variety of its methods and tools, ranging from dictionary methods and other supervised and unsupervised methods to learn about, for example, content, ideology and sentiment in text. The course – which combines lectures and practical sessions – will be hands-on, with an emphasis on dealing with practical issues in each step of the research process (ranging from collecting and pre-processing text data to validating and visualizing output of the analysis). Students who have finished this course are well-positioned to apply automated text analysis methods in their own work, and will be able to critically evaluate existing work.

**NB:** This course assumes familiarity with R. Students who have not used R before will need to get themselves up to speed before the start of the course, for example by working their way through a free online R resource, like <https://www.datacamp.com/courses/free-introduction-to-r> or the resources that are listed on <https://www.rstudio.com/online-learning/#R>. Students working on their own laptops will need to have R and RStudio installed.

## Learning Outcomes

This course introduces students to various approaches of automated text analysis in social science research, emphasizing hands on analysis of real (political) texts. Students will learn how to extract useful information from text, evaluate the outcomes and write up the results of an analysis that uses automated text analysis. Furthermore, students will be able to critically evaluate (social science) research that uses automated text analysis methods.

## Assessment

Students are assessed on the basis of 4 components:

1. **Attendance and participation in class (10% towards the final grade)**

- This course will involve quite some reading, some of which is technical. I expect that you come to class prepared, having read all required papers, and ready to discuss your questions, criticisms and thoughts. Furthermore, since this is a short class I expect you to attend all sessions. Missing class is only acceptable for urgent reasons and students will need to communicate this with me in advance.

## 2. Two coding assignments (15% each, 30% total towards the final grade)

- The coding assignments (one in week 1 and one in week 2) are designed to experience the workflow of a text research project. The first assignment will concern getting from text to data that can be analyzed. The second assignment will involve applying some of the methods we discuss in class to this data. Both assignments rely on the EUSpeech dataset.

## 3. Presentation of a research design (15% towards the final grade)

- On the last day of the course all students will give a brief presentation of a research design they have developed to address a topic they want to study using (one of) the methods discussed in class. This presentation should at least contain a research question, a discussion of the text sources, as well as the (expected) steps to address this question using the methods discussed in class. **NB:** Depending on time and enrollment, students will also act as discussant of the research design of another student with the goal of providing constructive comments to improve their work.

## 4. Research note (45% towards the final grade)

- All students will hand in a research note of about 1500 words (excluding references and appendices) in which they briefly but clearly write up the results of a small research project based on the design they presented in class. Students are free to collect their own data or use existing data (like the EUSpeech dataset or a replication file from a published research paper). Creativity is encouraged. This note should contain the following elements:
  - (a) Introduction & research question ( $\pm 300$  words): introduction to the topic.
  - (b) Data & methods ( $\pm 400$  words): description of the data sources as well as the methods employed.
  - (c) Analysis: ( $\pm 600$  words): a discussion (with figures and tables) of the results of the analysis.
  - (d) Conclusion ( $\pm 300$  words): a brief evaluation of the results and steps to push the research forward.

Since time is short this is not likely to be a very polished research project (nor is this expected). Rather the research note is a transparent write-up of the work the student put in to address a research question of their interest using text.

Table 1: Grade Breakdown

A	94–100
A-	87.00–93.99
B+	80.00–86.99
B	73.00–79.99
B-	66.00–72.99
C+	59.00–65.99
F	0–58.99

## Grading

Grading on a 100 point scale is reported in Table 1.

## Course outline

\* *This outline serves a general plan for the course; deviations (announced) may be necessary.*

### April 16: 09:00 - 10:40:

- Introduction to the course and to EUSpeech, a dataset which will use for running examples: <https://dataverse.harvard.edu/dataverse/euspeech>

#### – Required reading:

- \* Schumacher, G., Schoonvelde, M., Traber, D., Dahiya, T., & De Vries, E. (2016). EUSpeech: a New Dataset of EU Elite Speeches. In: *Proceedings of the International Conference on the Advances in Computational Analysis of Political Text*, 75–80.

### April 17: 09:00 - 10:40:

- A survey of automated text analysis in political science. Supervised and unsupervised methods. Validation, validation, validation. Text Analysis in R.

#### – Required reading:

- \* Grimmer, J., & Stewart, B. M. (2013). Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political Analysis*, 21(3), 267–297.
- \* Welbers, K., Van Atteveldt, W., & Benoit, K. (2017). Text Analysis in R. *Communication Methods and Measures*, 11(4), 245–265.

### April 18: 09:00 - 10:40 & 11:00 - 12:40:

- Pre-processing data. Going from text to data, including a few notes of caution. Discussion of the research design and research note.

#### – Required reading:

- \* Denny, M. J., & Spirling, A. (2017). Text preprocessing for unsupervised learning: why it matters, when it misleads, and what to do about it. *Working paper*.
- \* Greene, Z., Ceron, A., Schumacher, G., & Fazekas, Z. (2016). The Nuts and Bolts of Automated Text Analysis. Comparing Different Document Pre-Processing Techniques in Four Countries. *Working paper*

### April 19: 09:00 - 10:40:

- Systematically describing and comparing texts.

#### – Required reading: Chapters 3 and 4 of Silge, J., & Robinson, D. (2018). Text Mining with R: A Tidy Approach. O'Reilly Media, Inc. Available at <https://www.tidytextmining.com>

### April 20: 09:00 - 10:40:

- Using dictionaries to measure sentiment, happiness and other things we're interested in.

#### – Required reading:

- \* Pennebaker JW & King L (1999) Linguistic styles: language use as an individual difference. *Journal of Personality and Social Psychology*, 77(6), 1296–1312.

- \* Young, L., & Soroka, S. (2012). Affective news: The automated coding of sentiment in political texts. *Political Communication*, 29(2), 205–231.

– **Suggested reading:**

- \* Rooduijn, M., & Pauwels, T. (2011). Measuring populism: Comparing two methods of content analysis. *West European Politics*, 34(6), 1272–1283.
- \* Rheault, L., Beelen, K., Cochrane, C., & Hirst, G. (2016). Measuring Emotion in Parliamentary Debates with Automated Textual Analysis. *PLoS One*, 11(12).

– **17:00, Coding Assignment 1 due**

**April 23: 09:00 - 10:40:**

- Scaling methods locating text on an underlying (political) dimension. What do they mean? And how do they work?

– **Required reading:**

- \* Slapin JB & Proksch SO (2008) A Scaling Model for Estimating Time-Serial Positions from Texts. *American Journal of Political Science* 52, 705–722.
- \* Hjorth, F., Klemmensen, R., Hobolt, S., Hansen, M. E., & Kurrild-Klitgaard, P. (2015). Computers, coders, and voters: Comparing automated methods for estimating party positions. *Research & Politics*, 2(2).

– **Suggested reading:**

- \* Lo, J., Proksch, S. O., & Slapin, J. B. (2016). Ideological clarity in multiparty competition: A new measure and test using election manifestos. *British Journal of Political Science*, 46(3), 591–610.

**April 24: 09:00 - 10:40:**

- Topic models, unsupervised models for summarizing what a text is about.

– **Required reading:**

- \* Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4), 77–84.
- \* Roberts, M et al.. (2014). Structural Topic Models for Open-Ended Survey Responses. *American Journal of Political Science*, 58(4), 1064–1082.

– **Suggested reading:**

- \* Boumans JW & Trilling D (2016) Taking stock of the toolkit: An overview of relevant automated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism* 4(1): 8–23.
- \* <http://www.scottbot.net/HIAL/index.html?p=19113.html>

**April 25: 09:00 - 10:40 & 11:00 - 12:40:**

- New developments in automated text analysis: (i) crowd-sourcing and (ii) measurement of elite personality, (iii) measurement of semantic shifts.

– **Required reading:**

- \* Ramey, A. J., Klingler, J. D., & Hollibaugh, G. E. (2016). Measuring elite personality using speech. *Political Science Research and Methods*, 1–22.
- \* Benoit, K., Conway, D., Lauderdale, B. E., Laver, M., & Mikheylov, S. (2016). Crowd-sourced text analysis: Reproducible and agile production of political data. *American Political Science Review*, 110(2), 278–295.

– **Suggested reading:**

\* Azarbonyad, H., Deghani, M., Beelen, K., Arkut, A., Marx, M., & Kamps, J. (2017). Words are Malleable: Computing Semantic Shifts in Political and Media Discourse. *Proceedings of the 2017 ACM Conference on Information and Knowledge Management*, 1509–1518.

– **17:00 Coding Assignment 2 due**

**April 26: 09:00 - 10:40:**

- Loose ends, review, and general discussion of pros and cons of automated text analysis.

**April 27: 09:00 - 10:40:**

- Research design presentations.

**4 May, 17:00: Final Assignment Due: Research Note**