

# What's on the Agenda? Automated Text Classification for Political Text

Monday 17:15 18:45 weekly 28.09.2020 - 07.12.2020 SOWI-Zoomroom 6

<https://uni-mannheim.zoom.us/j/5677743231?pwd=YzVDN0FqaWdYVlhNWXBwcVg1Zm1FQT09>

Marius Sältzer

msaeltze@mail.uni-mannheim.de

A 339

Office Hours: TBD

+49(0)621-181-2078

Syllabus will be adapted over the course of the semester.

## Course Description:

Politics is often about the importance of particular issues, from migration over climate change to COVID-19. To understand political processes from political communication to legislation, measuring what matters to political actors, newspaper or voters is of great importance. In recent times, the growing availability of text from social media to parliamentary procedures has increased the demand for automated content analysis. This methods course will introduce a number of basic clustering and machine learning approaches to automatically identify the topics of political text. It offers a hands-on introduction to big data with examples from newspapers, manifestos, speeches and social media data. Part of this course is an introduction to the statistical programming language R. It will be taught in the seminar and previous experience is not required, but the openness to learning a new language is key. Basic knowledge of statistics and respective software like STATA is required.

1. Data Science for Comparative Political Science
2. Get a first look at R
3. Find out what's out there in terms of data
4. Develop your own research project

Formalia:

- This is a methods course and I want you to get access to the tools to collect your own data.
- Your Grade depends on two Elements:
  - **Studienleistung:**
    - \* *Problem Sets*: You will hand in three problem sets (pass/fail) of which 2 have to be passed.
  - **Prüfungsleistung:**
    - A *Term Paper* that is somehow connected to salience, issue emphasis or classifying text. It can be empirical but does not have to be.
- Software: R
  - Freeware that is used everywhere around the world
  - Easy to use on your home computer without license etc.
  - Introduction to programming and great job opportunities
  - Steep learning curve, but worth it!
- Helpful Textbooks:
  - Technical
    - \* concerning R
      - Andrie de Vries and Joris Meys. *R for dummies: Learn to: use R for data analysis and processing : write functions and scripts for repeatable analysis : create high-quality charts and graphics : perform statistical analysis and build models*. Wiley, Hoboken, NJ, 2. ed. edition, 2015. ISBN 9781119055839 (RFD)
      - Additional for the really curious: Hadley Wickham. *Advanced R*. Taylor and Francis, 2015 (A)
  - \* Text Analysis
    - Daniel Jurafsky and James H. Martin. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Prentice Hall series in artificial intelligence. Prentice Hall, Upper Saddle River, NJ, 2000. ISBN 978-0130950697
  - \* Substantive
    - Ian Budge, Hans-Dieter Klingeman, Andrea Volkens, Judith Bara, and Michael D. McDonald. *Mapping policy preferences II: Estimates for Parties, Electors, and Governments in Eastern Europe, European Union, and OECD 1990-2003*. Oxford University Press, Oxford, 2001. ISBN 0199244006

## Part I: Theory:

The weekly coverage might change as it depends on the progress of the class. However, you must keep up with the reading assignments.

Week	Content
28.09.2020	<ul style="list-style-type: none"><li>• Introduction</li></ul>
05.10.2020 Salience and Agendas	<p>Why is it important what issues we talk about in politics?</p> <ul style="list-style-type: none"><li>• Peter Bachrach and Morton A. Baratz. Decisions and Nondecisions: An Analytical Framework. <i>American Political Science Review</i>, 57(3):632–642, 1963. ISSN 0003-0554</li><li>• Maxwell E. McCombs and Donald L. Shaw. The Agenda-Setting Function of Mass Media. <i>Public Opinion Quarterly</i>, 36(2):176, 1972. ISSN 0033362X. doi: 10.1086/267990</li></ul>
12.10.2020 Issues and Politics	<p>How do political parties engage with issues?</p> <p>Readings:</p> <ul style="list-style-type: none"><li>• Stephen Ansolabehere, James M. Snyder, and Jonathan Rodden. The Strength of Issues: Using Multiple Measures to Gauge Preference Stability, Ideological Constraint, and Issue Voting. <i>American Political Science Review</i>, 102(02):215–232, 2008. ISSN 0003-0554. doi: 10.1017/S0003055408080210</li><li>• Éric Bélanger and Bonnie M. Meguid. Issue salience, issue ownership, and issue-based vote choice. <i>Electoral Studies</i>, 27(3):477–491, 2008. ISSN 02613794. doi: 10.1016/j.electstud.2008.01.001</li></ul>
19.10.2020 Coding Political Text	<p>How can we find out what parties emphasize? We engage with classification schemes that are used to analyze political text.</p> <ul style="list-style-type: none"><li>• <a href="https://manifesto-project.wzb.eu/coding_schemes/mp_v5">https://manifesto-project.wzb.eu/coding_schemes/mp_v5</a></li><li>• <a href="https://www.comparativeagendas.net/pages/master-codebook">https://www.comparativeagendas.net/pages/master-codebook</a></li><li>• Reading:<ul style="list-style-type: none"><li>– Heike Klüver and Jae-Jae Spoon. Who Responds? Voters, Parties and Issue Attention. <i>British Journal of Political Science</i>, 46(03):633–654, 2016. ISSN 0007-1234. doi: 10.1017/S0007123414000313</li><li>– Thomas M. Meyer and Markus Wagner. It Sounds Like They are Moving: Understanding and Modeling Emphasis-Based Policy Change. <i>Political Science Research and Methods</i>, 7(04):757–774, 2019. ISSN 2049-8470. doi: 10.1017/psrm.2017.30</li></ul></li></ul>

## Part II: Programming in R:

Week	Content
26.10.2020 Introduction to R 1	<p>We start the introduction into R by installing R and R Studio, learning about basic operations, objects and functions. Please install</p> <ul style="list-style-type: none"><li>• R: <a href="https://cran.r-project.org/mirrors.html">https://cran.r-project.org/mirrors.html</a></li><li>• R Studio: <a href="https://rstudio.com/products/rstudio/download/download">https://rstudio.com/products/rstudio/download/download</a></li><li>• Suggested Reading:<ul style="list-style-type: none"><li>– Sessions 5 and 6 will cover RFD p.49-149</li></ul></li></ul>
02.11.2020 Introduction to R 2	<p>After we learned basic operations, we load data from the manifesto project! We learn how to import, display and manipulate data sets. We import the dataset of the Manifesto Project.</p> <ul style="list-style-type: none"><li>• Suggested Reading:<ul style="list-style-type: none"><li>– Sessions 5 and 6 will cover RFD p.49-149</li><li>– <a href="https://manifestoproject.wzb.eu/information/documents/manifestoR">https://manifestoproject.wzb.eu/information/documents/manifestoR</a></li></ul></li></ul>
09.11.2020 Text as Data in R	<p>After learning basic R, we directly enter the world of text as data. We learn about corpora, tokens, document-feature-matrices and word frequencies.</p> <ul style="list-style-type: none"><li>• Software: <a href="https://tutorials.quanteda.io/basic-operations/">https://tutorials.quanteda.io/basic-operations/</a></li><li>• Literature: Text as Data<ul style="list-style-type: none"><li>– Justin Grimmer and Brandon M. Stewart. Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. <i>Political Analysis</i>, 21(03):267–297, 2013. ISSN 1047-1987. doi: 10.1093/pan/mps028</li><li>– Daniel Jurafsky and James H. Martin. <i>Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition</i>. Prentice Hall series in artificial intelligence. Prentice Hall, Upper Saddle River, NJ, 2000. ISBN 978-0130950697: <b>pages 2-28</b></li></ul></li></ul>

### Part III: Text Classification:

Week	Content
16.11.2020 Supervised Learning	<p>A brief introduction to Machine Learning to start Working with Labeled Data: Supervised Learning.</p> <ul style="list-style-type: none"> <li>• Daniel Jurafsky and James H. Martin. <i>Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition</i>. Prentice Hall series in artificial intelligence. Prentice Hall, Upper Saddle River, NJ, 2000. ISBN 978-0130950697: <b>pages 56-74</b></li> </ul>
23.11.2020 Supervised Learning II	<p>Learn evaluating your model by measuring precision and accuracy</p> <ul style="list-style-type: none"> <li>• Pablo Barberá, Andrué Casas, Jonathan Nagler, Patrick J. Egan, Richard Bonneau, John T. Jost, and Joshua A. Tucker. Who Leads? Who Follows? Measuring Issue Attention and Agenda Setting by Legislators and the Mass Public Using Social Media Data. <i>American Political Science Review</i>, 113(4):883–901, 2019. ISSN 0003-0554. doi: 10.1017/S0003055419000352</li> </ul>
30.11.2020 Unsupervised Classification	<p>We start with unsupervised machine learning, letting the computer choose the categories. We introduce so called topic models.</p> <ul style="list-style-type: none"> <li>• Margaret E. Roberts, Brandon M. Stewart, Dustin Tingley, Christopher Lucas, Jetson Leder-Luis, Shana Kushner Gadarian, Bethany Albertson, and David G. Rand. Structural Topic Models for Open-Ended Survey Responses. <i>American Journal of Political Science</i>, 58(4):1064–1082, 2014. ISSN 00925853. doi: 10.1111/ajps.12103</li> <li>• David M. Blei. Probabilistic topic models. <i>Communications of the ACM</i>, 55(4):77–84, 2012. ISSN 0001-0782. doi: 10.1145/2133806.2133826</li> </ul>
07.12.2020 Validating Topic Models	<p>We look at ways to better interpret and optimize topic models.</p> <ul style="list-style-type: none"> <li>• Jonathan Chang, Sean Gerrish, Wang. Chong, Jordan Boyd-graber, and David M. Blei. Reading Tea Leaves: How Humans Interpret Topic Models. <i>Advances in neural information processing systems</i>, pages 288–296, 2009</li> <li>• Pablo Barberá, Andrué Casas, Jonathan Nagler, Patrick J. Egan, Richard Bonneau, John T. Jost, and Joshua A. Tucker. Who Leads? Who Follows? Measuring Issue Attention and Agenda Setting by Legislators and the Mass Public Using Social Media Data. <i>American Political Science Review</i>, 113(4):883–901, 2019. ISSN 0003-0554. doi: 10.1017/S0003055419000352</li> </ul>

## References

- Stephen Ansolabehere, James M. Synder, and Jonathan Rodden. The Strength of Issues: Using Multiple Measures to Gauge Preference Stability, Ideological Constraint, and Issue Voting. *American Political Science Review*, 102(02):215–232, 2008. ISSN 0003-0554. doi: 10.1017/S0003055408080210.
- Peter Bachrach and Morton A. Baratz. Decisions and Nondecisions: An Analytical Framework. *American Political Science Review*, 57(3):632–642, 1963. ISSN 0003-0554.
- Pablo Barberá, Andrue Casas, Jonathan Nagler, Patrick J. Egan, Richard Bonneau, John T. Jost, and Joshua A. Tucker. Who Leads? Who Follows? Measuring Issue Attention and Agenda Setting by Legislators and the Mass Public Using Social Media Data. *American Political Science Review*, 113(4):883–901, 2019. ISSN 0003-0554. doi: 10.1017/S0003055419000352.
- Éric Bélanger and Bonnie M. Meguid. Issue salience, issue ownership, and issue-based vote choice. *Electoral Studies*, 27(3):477–491, 2008. ISSN 02613794. doi: 10.1016/j.electstud.2008.01.001.
- David M. Blei. Probabilistic topic models. *Communications of the ACM*, 55(4):77–84, 2012. ISSN 0001-0782. doi: 10.1145/2133806.2133826.
- Ian Budge, Hans-Dieter Klingeman, Andrea Volkens, Judith Bara, and Michael D. McDonald. *Mapping policy preferences II: Estimates for Parties, Electors, and Governments in Eastern Europe, European Union, and OECD 1990-2003*. Oxford University Press, Oxford, 2001. ISBN 0199244006.
- Jonathan Chang, Sean Gerrish, Wang. Chong, Jordan Boyd-graber, and David M. Blei. Reading Tea Leaves: How Humans Interpret Topic Models. *Advances in neural information processing systems*, pages 288–296, 2009.
- Andrie de Vries and Joris Meys. *R for dummies: Learn to: use R for data analysis and processing : write functions and scripts for repeatable analysis : create high-quality charts and graphics : perform statistical analysis and build models*. Wiley, Hoboken, NJ, 2. ed. edition, 2015. ISBN 9781119055839.
- Justin Grimmer and Brandon M. Stewart. Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. *Political Analysis*, 21(03):267–297, 2013. ISSN 1047-1987. doi: 10.1093/pan/mps028.
- Daniel Jurafsky and James H. Martin. *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Prentice Hall series in artificial intelligence. Prentice Hall, Upper Saddle River, NJ, 2000. ISBN 978-0130950697.
- Heike Klüver and Jae-Jae Spoon. Who Responds? Voters, Parties and Issue Attention. *British Journal of Political Science*, 46(03):633–654, 2016. ISSN 0007-1234. doi: 10.1017/S0007123414000313.
- Maxwell E. McCombs and Donald L. Shaw. The Agenda-Setting Function of Mass Media. *Public Opinion Quarterly*, 36(2):176, 1972. ISSN 0033362X. doi: 10.1086/267990.
- Thomas M. Meyer and Markus Wagner. It Sounds Like They are Moving: Understanding and Modeling Emphasis-Based Policy Change. *Political Science Research and Methods*, 7(04):757–774, 2019. ISSN 2049-8470. doi: 10.1017/psrm.2017.30.

Margaret E. Roberts, Brandon M. Stewart, Dustin Tingley, Christopher Lucas, Jetson Leder-Luis, Shana Kushner Gadarian, Bethany Albertson, and David G. Rand. Structural Topic Models for Open-Ended Survey Responses. *American Journal of Political Science*, 58(4):1064–1082, 2014. ISSN 00925853. doi: 10.1111/ajps.12103.

Hadley Wickham. *Advanced R*. Taylor and Francis, 2015.