



Text Analysis Techniques



A Brief Overview



Text Analysis demystified:

It's just counting.

Sense and Sensibility		Moby Dick	
to	4063	the	13721
the	3861	of	6536
of	3565	and	6024
...			
secret	19	cold	30
discourse	19	jaws	30
hoped	19	won	30
age	19	harpoons	30

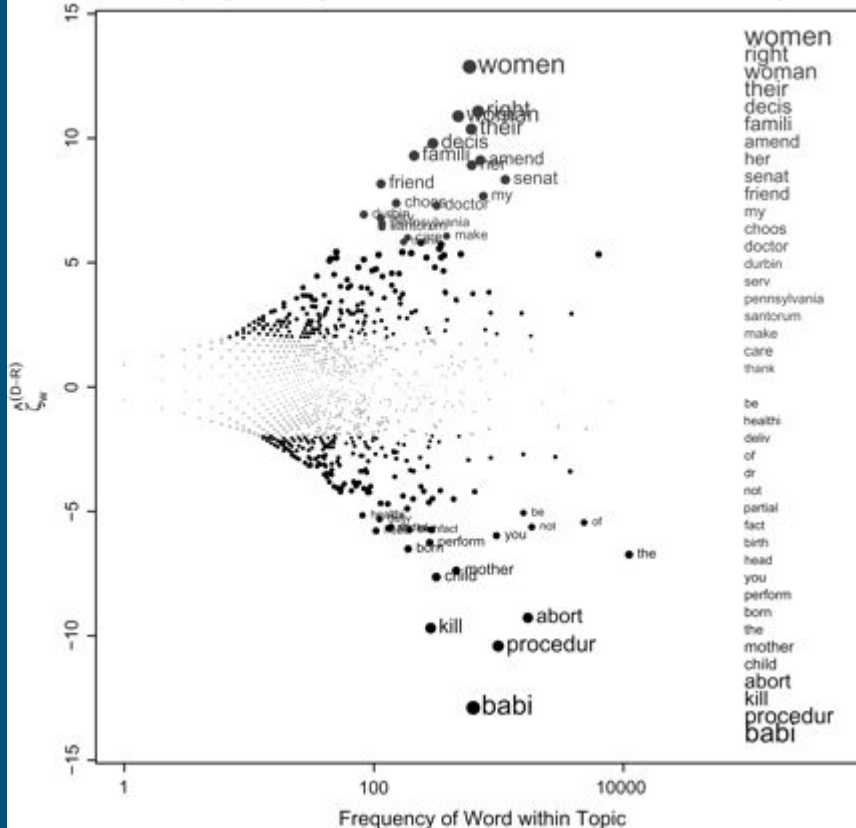
Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics

Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)

Partisan Words, 106th Congress, Abortion
(Weighted Log-Odds-Ratio, Uninformative Dirichlet Prior)



From: Fightin' Words: Lexical Feature Selection and
Evaluation for Identifying the Content of Political Conflict
Polit Anal. 2009;16(4):372-403.

doi:10.1093/pan/mpn018

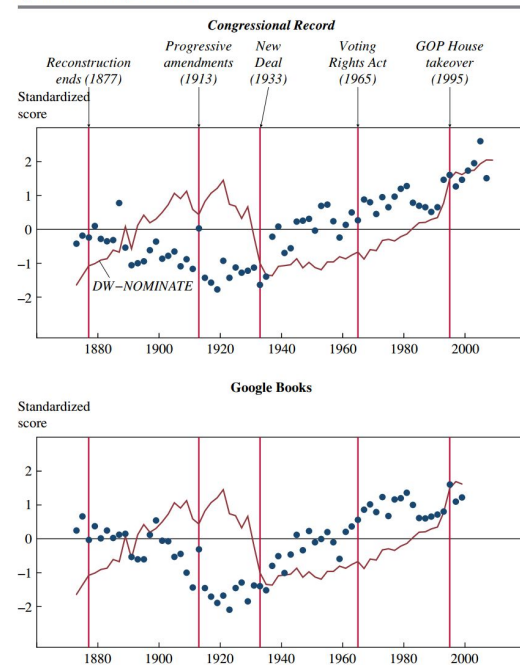
Polit Anal | © The Author 2009. Published by Oxford
University Press on behalf of the Society for Political
Methodology. All rights reserved. For Permissions,
please email: journals.permissions@oxfordjournals.org

Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)
- Dictionary Methods (themes)

Jensen, Jacob, Ethan Kaplan, Suresh Naidu, and Laurence Wilse-Samson. 2012. "Political Polarization and the Dynamics of Political Language: Evidence from 130 Years of Partisan Speech." *Brookings Papers on Economic Activity*.

Figure B.2. Polarization Measured Using t -Statistic-Based Threshold and by DW-NOMINATE, 1873–2007^a



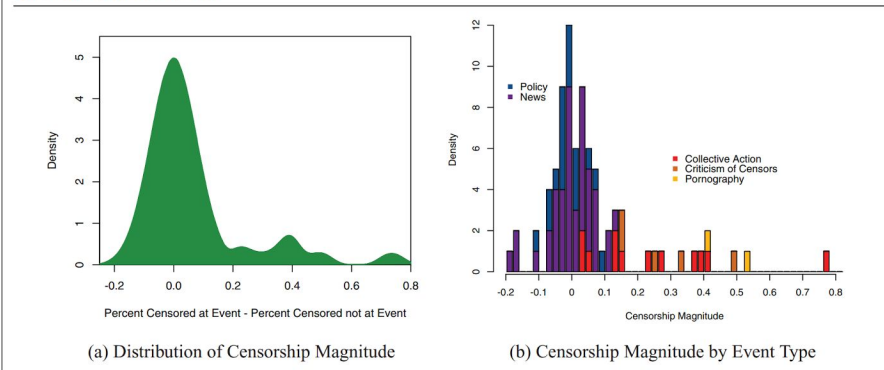
Sources: Authors' calculations using data from the digitized *Congressional Record*, Google Books, and the legislator estimates on voteview.com/dwnomin.htm.

a. All measures are standardized to have a mean of zero and a variance of 1.

Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)
- Dictionary Methods (themes)
- Supervised Machine Learning (classification)

Figure 3. “Censorship Magnitude,” The Percent of Posts Censored Inside a Volume Burst Minus Outside Volume Bursts.



King, Gary, Jennifer Pan, and Margaret E. Roberts. 2013. “How Censorship in China Allows Government Criticism but Silences Collective Expression.” *American Political Science Review* 107(2): 1-18.

Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)
- Dictionary Methods (themes)
- Supervised Machine Learning (classification)
- Vector Space Models
 - Clustering, document distance

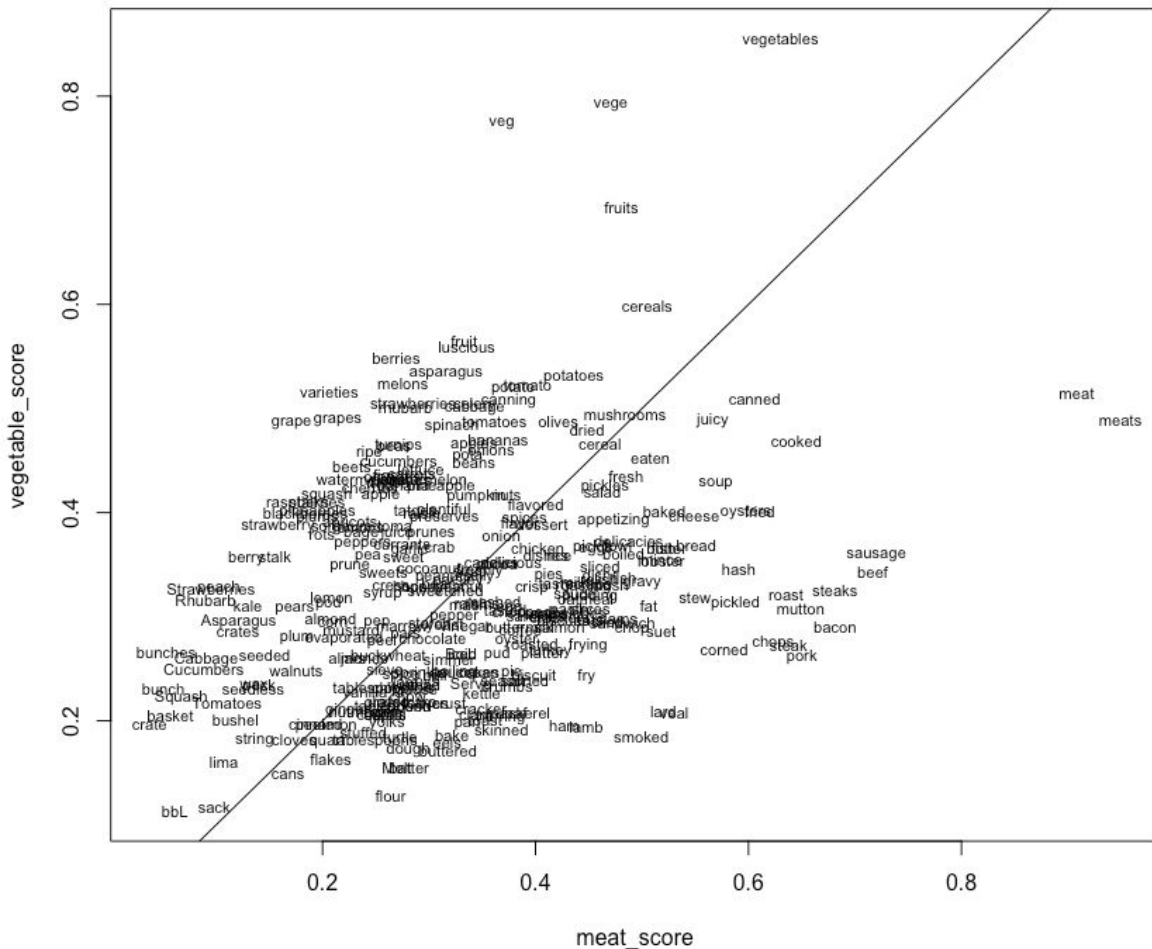
Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)
- Dictionary Methods (themes)
- Supervised Machine Learning (classification)
- Vector Space Models
 - Clustering, document distance
- Topic Modeling

Techniques

- Natural Language Processing (takes into account context)
 - Part of Speech, Named Entities, Concordances, Word Similarities, semantics
- Distinguishing Words
 - Difference of Proportions, TF-IDF, chi-squared, Dunning Log-Likelihood)
- Dictionary Methods (themes)
- Supervised Machine Learning (classification)
- Vector Space Models
 - Clustering, document distance
- Topic Modeling
- Word embeddings
 - Word2Vec, GloVe

Top 300 food words plotted by their similarity to meats (x axis) and vegetables (y axis).



Ben Schmidt, Word
Embeddings Models.

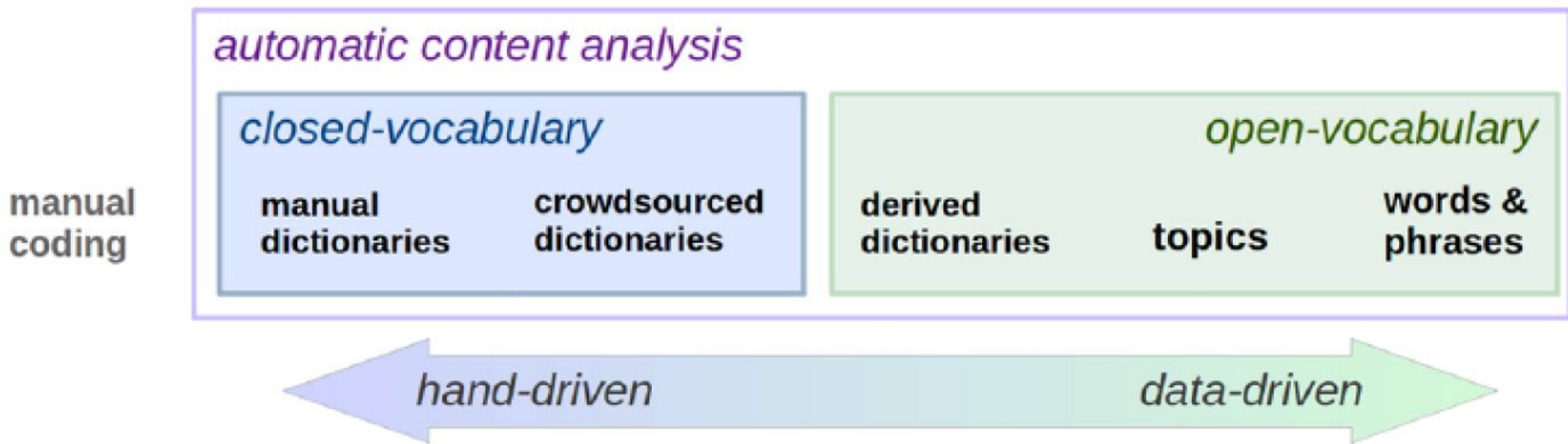


Figure 1 Categorization of Content Analysis Techniques

Published in: Dhavan V. Shah; Joseph N. Cappella; W. Russell Neuman;

Published in: H. Andrew Schwartz; Lyle H. Ungar; The ANNALS of the American Academy of Political and Social Science 659, 78-94.

Copyright © 2015 American Academy of Political & Social Science

