

Quantitative text analysis: Unsupervised Scaling Methods

Blake Miller

MY 459: Quantitative Text Analysis

February 22, 2021

Course website: lse-my459.github.io

1. Overview and Fundamentals
2. Descriptive Statistical Methods for Text Analysis
3. Automated Dictionary Methods
4. Machine Learning for Texts
5. Supervised Scaling Models for Texts
6. *Reading Week*
7. Unsupervised Models for Scaling Texts
8. Similarity and Clustering Methods
9. Topic models
10. Word embeddings
11. Working with Social Media

Overview of text as data methods



Outline

- ▶ Basics of unsupervised scaling methods
- ▶ Parametric scaling models: Wordfish and Wordshoal
- ▶ Non-parametric scaling methods: correspondence analysis
- ▶ Practical aspects: computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Wordfish examples

Unsupervised methods scale distance

- ▶ Text gets converted into a quantitative matrix of features
 - ▶ words, typically
 - ▶ could be dictionary entries, or parts of speech
- ▶ Documents are scaled based on similarity or distance in feature use
- ▶ Fundamental problem: distance on which scale?
 - ▶ Ideally, something we care about, e.g. policy positions, ideology, preferences, sentiment
 - ▶ But often other dimensions (language, rhetoric style, authorship) are more predictive
- ▶ First dimension in unsupervised scaling will capture main source of variation, whatever that is
- ▶ Unlike supervised models, validation comes after estimating the model

Unsupervised scaling methods

Two main approaches

- ▶ **Parametric methods** model feature occurrence according to some stochastic distribution, typically in the form of a measurement model
 - ▶ for instance, model words as a multi-level Bernoulli distribution, or a Poisson distribution
 - ▶ word effects and “positional” effects are unobserved parameters to be estimated
 - ▶ e.g. Wordfish (Slapin and Proksch 2008) and Wordshoal (Lauderdale and Herzog 2016)
- ▶ **Non-parametric methods** typically based on the Singular Value Decomposition of a matrix
 - ▶ correspondence analysis
 - ▶ factor analysis
 - ▶ other (multi)dimensional scaling methods

Outline

- ▶ Basics of unsupervised scaling methods
- ▶ Parametric scaling models: Wordfish and Wordshoal
- ▶ Non-parametric scaling methods: correspondence analysis
- ▶ Practical aspects: computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Wordfish examples

Wordfish (Slapin and Proksch 2008)

- ▶ Goal: unsupervised scaling of ideological positions
- ▶ The frequency with which politician i uses word k is drawn from a **Poisson distribution**:

$$w_{ik} \sim \text{Poisson}(\lambda_{ik})$$

$$\lambda_{ik} = \exp(\alpha_i + \psi_k + \beta_k \times \theta_i)$$

- ▶ with **latent parameters**:
 - α_i is “loquaciousness” of politician i (document fixed effect, hence it's associated with the party or politician)
 - ψ_k is frequency of word k (word fixed effect)
 - β_k is discrimination parameter of word k
 - θ_i is the politician's ideological position
- ▶ **Key intuition**: controlling for document length and word frequency, words with negative β_k will tend to be used more often by politicians with negative θ_i (and vice versa)

Wordfish (Slapin and Proksch 2008)

Why Poisson?

- ▶ Poisson-distributed variables are bounded between $(0, \infty)$ and take on only discrete values $0, 1, 2, \dots, \infty$
- ▶ Exponential transformation: word counts are function of log document length and word frequency

$$w_{ik} \sim \text{Poisson}(\lambda_{ik})$$

$$\lambda_{ik} = \exp(\alpha_i + \psi_k + \beta_k \times \theta_i)$$

$$\log(\lambda_{ik}) = \alpha_i + \psi_k + \beta_k \times \theta_i$$

How to estimate this model

Conditional maximum likelihood estimation:

- ▶ If we knew ψ and β (the word parameters) then we have a Poisson regression model
- ▶ If we knew α and θ (the party / politician / document parameters) then we have a Poisson regression model too!
- ▶ So we alternate them and hope to converge to reasonable estimates for both
- ▶ Implemented in the `quanteda` package as `textmodel_wordfish`

An alternative is MCMC with a Bayesian formulation or variational inference using an Expectation-Maximization algorithm (Imai et al 2016)

Conditional maximum likelihood for wordfish

Start by **guessing** the parameters (some guesses are better than others, e.g. SVD)

Algorithm:

1. Assume the current **legislator parameters** are correct and fit as a Poisson regression model
2. Assume the current **word parameters** are correct and fit as a Poisson regression model
3. **Normalize** θ s to mean 0 and variance 1

Iterate until convergence (change in values is below a certain threshold)

Identification

The *scale* and *direction* of θ is undetermined — like most models with latent variables

To **identify the model** in Wordfish

- ▶ Fix one α to zero to specify the left-right direction (Wordfish option 1)
- ▶ Fix the $\hat{\theta}$ s to mean 0 and variance 1 to specify the scale (Wordfish option 2)
- ▶ Fix two $\hat{\theta}$ s to specify the direction and scale (Wordfish option 3 and Wordscores)

Note: Fixing two reference scores does not specify the policy domain, it just identifies the model

“Features” of the parametric scaling approach

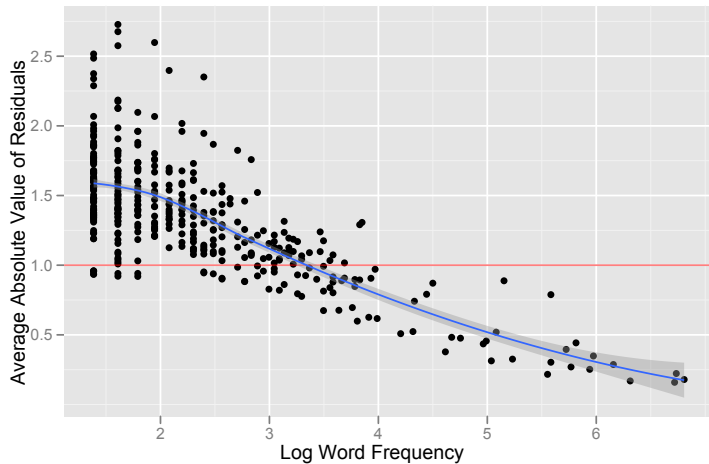
- ▶ Standard (statistical) **inference** about parameters
- ▶ **Uncertainty** accounting for parameters
- ▶ **Distributional assumptions** are made explicit (as part of the data generating process motivating the choice of stochastic distribution)
 - ▶ *conditional independence*
 - ▶ *stochastic process* (e.g. $E(Y_{ij}) = \text{Var}(Y_{ij}) = \lambda_{ij}$)
- ▶ Permits **hierarchical reparameterization** (to add covariates)
- ▶ Generative model: given the estimated parameters, we could **generate a document** for any specified length

Some reasons why this model is wrong

- ▶ Violations of conditional independence:
 - ▶ Words occur in sequence (serial correlation)
 - ▶ Words occur in combinations (e.g. as collocations)
“carbon tax” / “income tax” / “inheritance tax” / “capital gains tax” / “bank tax”
 - ▶ Legislative speech uses rhetoric that contains frequent synonyms and repetition for emphasis (e.g. “Yes we can!”)
- ▶ Heteroskedastic errors (variance not constant and equal to mean):
 - ▶ **over**dispersion when “informative” words tend to cluster together
 - ▶ **under**dispersion could (possibly) occur when words of high frequency are uninformative and have relatively low between-text variation (once length is considered)

Overdispersion in German manifesto data

(data taken from Slapin and Proksch 2008)



One solution to model overdispersion

Negative binomial model (Lo, Proksch, and Slapin 2014):

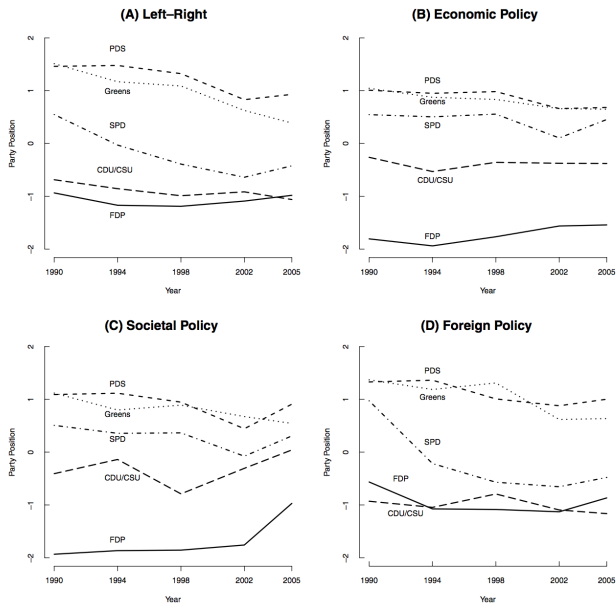
$$w_{ik} \sim \text{NB} \left(r, \frac{\lambda_{ik}}{\lambda_{ik} + r_i} \right)$$
$$\lambda_{ik} = \exp(\alpha_i + \psi_k + \beta_k \times \theta_i)$$

where r_i is a variance inflation parameter that varies across documents.

It can have a substantive interpretation (**ideological ambiguity**), e.g. when a party emphasizes an issue but fails to mention key words associated with it that a party with similar ideology mentions.

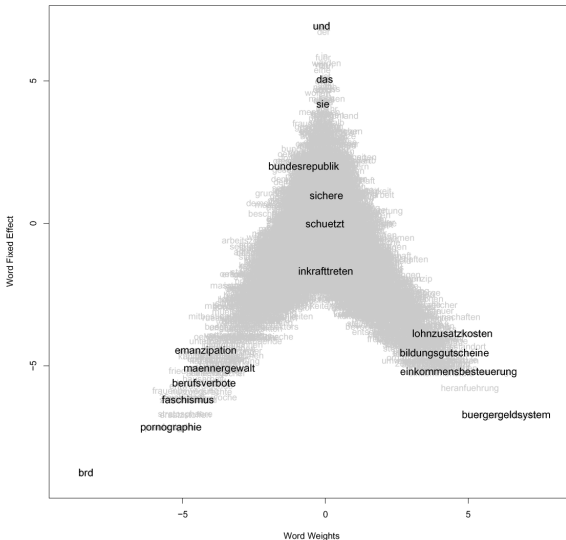
Example from Slapin and Proksch 2008

FIGURE 1 Estimated Party Positions in Germany, 1990–2005



Example from Slapin and Proksch 2008

FIGURE 2 Word Weights vs. Word Fixed Effects. Left-Right Dimension, Germany 1990–2005 (Translations given in text)



Example from Slapin and Proksch 2008

TABLE 1 Top 10 Words Placing Parties on the Left and Right

Dimension	Top 10 Words Placing Parties on the . .	
	Left	Right
Left-Right	Federal Republic of Germany (BRD) immediate (sofortiger) pornography (Pornographie) sexuality (Sexualität) substitute materials (Ersatzstoffen) stratosphere (Stratosphäre) women's movement (Frauenbewegung) fascism (Faschismus) Two thirds world (Zweidrittelwelt) established (etablierten)	general welfare payments (Bürgergeldsystem) introduction (Heranführung) income taxation (Einkommensbesteuerung) non-wage labor costs (Lohnzusatzkosten) business location (Wirtschaftsstandort) university of applied sciences (Fachhochschule) education vouchers (Bildungsgutscheine) mobility (Beweglichkeit) peace tasks (Friedensaufgaben) protection (Protektion)
Economic	Federal Republic of Germany (BRD) democratization (Demokratisierung) to prohibit (verbieten) destruction (Zerstörung) mothers (Mütter) debasing (entwürdigende) weeks (Wochen) quota (Quotierung) unprotected (ungeschützter) workers' participation (Mitbestimmungsmöglichkeiten)	to seek (anzustreben) general welfare payments (Bürgergeldsystem) inventors (Erfinder) mobility (Beweglichkeit) location (Standorts) negotiated wages (Tarif-Löhne) child-raising allowance (Erziehungsgeld) utilization (Verwertung) savings (Ersparnis) reliable (verlässlich)

Example from Slapin and Proksch 2008

TABLE 2 Cross-Validation: Correlations between German Party Position Estimates

	Poisson Scaling Model			
	Left-Right	Economic	Societal	Foreign
Hand-coding manifestos				
CMP: Left-Right (n = 15, 1990–1998)	−0.82			
CMP: Markeco (n = 15, 1990–1998)		0.81		
CMP: Welfare (n = 15, 1990–1998)			0.58	
CMP: Intpeace (n = 15, 1990–1998)				0.81
Expert Survey				
Benoit/Laver 2006: Left-Right (n = 5, 2002)	−0.91			
Benoit/Laver 2006: Taxes-Spending (n = 5, 2002)		0.86		
Wordscores				
Laver et al. 2003: Economic (n = 10, 1990–1994)		0.93		
Laver et al. 2003: Social (n = 10, 1990–1994)			−0.47	
Proksch/Slapin 2006: Economic (n = 5, 2005)		0.98		
Proksch/Slapin 2006: Social (n = 5, 2005)			−0.47	

Wordshoal (Lauderdale and Herzog 2016)

Two key **limitations** of wordfish applied to legislative text:

- ▶ Word discrimination parameters assumed to be **constant across debates** (unrealistic, think e.g. “debt”)
- ▶ May not capture left-right ideology but **topic variation**

Slapin and Proksch partially avoid these issues by scaling different types of debates separately.

But resulting estimates are confined to set of speakers who spoke on each topic.

Wordshoal solution: **aggregate debate-specific ideal points into a reduced number of scales.**

Wordshoal (Lauderdale and Herzog 2016)

- ▶ The frequency with which politician i uses word k in debate j is drawn from a **Poisson distribution**:

$$w_{ijk} \sim \text{Poisson}(\lambda_{ijk})$$

$$\lambda_{ijk} = \exp(\alpha_{ij} + \psi_{jk} + \beta_{jk} \times \theta_{ij})$$

$$\theta_{ij} \sim \mathcal{N}(\nu_j + \kappa_j \mu_i, \tau_i)$$

- ▶ with **latent parameters**:

α_{ij} is “loquaciousness” of politician i in debate j

ψ_{jk} is frequency of word k in debate j

β_{kj} is discrimination parameter of word k in debate j

θ_{ij} is the politician's ideological position in debate j

ν_j is baseline ideological position of debate j

κ_j is correlation of debate j with common dimension

μ_i is overall ideological position of politician i

- ▶ **Intuition**: debate-specific estimates are aggregated into a single position using dimensionality reduction

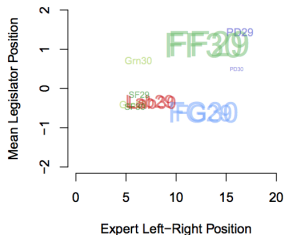
Wordshoal (Lauderdale and Herzog 2016)

New quantities of interest to estimate:

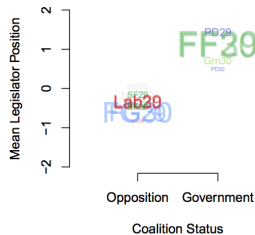
- ▶ Politicians' overall position vs debate-specific positions
- ▶ Strength of association between debate scales and general ideological scale
- ▶ Association of words with general scales, and stability of word discrimination parameters across debates

Example from Lauderdale and Herzog 2016

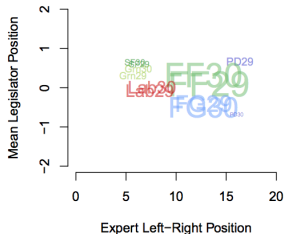
Wordshoal by Expert Location



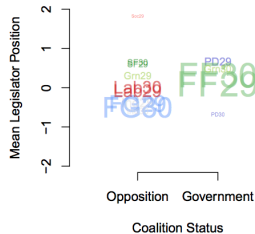
Wordshoal by Coalition Status



Wordfish by Expert Location



Wordfish by Coalition Status

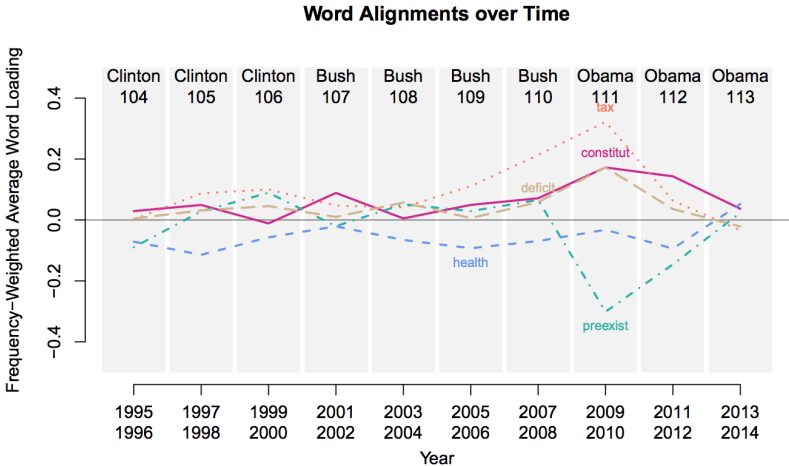


Example from Lauderdale and Herzog 2016

Table 2: The five debates with the highest and lowest loadings on the government versus opposition dimension, as measured by the absolute value of β_j ranging from 0 to 1.

<i>High government-opposition polarization</i>	Abs. β_j
Social Welfare and Pensions (No. 2) Bill 2009 (Second Stage)	0.942
Early Childhood Care and Education (Motion)	0.887
Private Members' Business – Vaccination Programme (Motion)	0.824
Capitation Grants (Motion)	0.819
Confidence in Government (Motion)	0.814
<i>Low government-opposition polarization</i>	
Cancer Services Reports (Motion)	0.003
Finance (No. 2) Bill 2007 (Committee and Remaining Stages)	0.002
Finance Bill 2011 (Report and Final Stages)	0.002
Private Members' Business – Mortgage Arrears (Motion)	0.002
Wildlife (Amendment) Bill 2010 (Committee and Remaining Stages)	0.001

Example from Lauderdale and Herzog 2016



Outline

- ▶ Basics of unsupervised scaling methods
- ▶ Parametric scaling models: Wordfish and Wordshoal
- ▶ Non-parametric scaling methods: correspondence analysis
- ▶ Practical aspects: computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Wordfish examples

Non-parametric methods

- ▶ Non-parametric methods are algorithmic, involving no “parameters” in the procedure that are estimated
- ▶ Hence there is no uncertainty accounting given distributional theory
- ▶ Advantage: don't have to make assumptions
- ▶ Disadvantages:
 - ▶ cannot leverage probability conclusions given distributional assumptions and statistical theory
 - ▶ results highly fit to the data
 - ▶ not really assumption-free, if we are honest

Correspondence Analysis

- ▶ CA is like principle component analysis for categorical data
- ▶ CA treats the document feature matrix as a **contingency table of counts** across two qualitative variables: documents and words.
- ▶ CA first preprocesses this contingency table using the logic of a χ^2 test of independence (is there a difference between expected cell frequencies under an **independence model** and the observed frequencies?)
- ▶ CA then performs **singular value decomposition** on this matrix to reduce the dimensionality of the document-feature matrix.
- ▶ This allows projection of the positioning of the words as well as the texts into d -dimensional space where $d \ll N$
- ▶ The number of dimensions can be chosen by zeroing out all but the principle d **singular values**.

Correspondence Analysis (details)

1. First normalize by document length and word frequency, compute matrix of standardized residuals, S :

$$Z = D_r^{1/2}(P - rc^T)D_c^{1/2}$$

where $P = X / \sum_i \sum_j x_{ij}$

r, c are row/column masses: e.g. $r_i = \sum_j p_{ij}$

$$D_r = \text{diag}(r), D_c = \text{diag}(c)$$

2. Calculate SVD of Z
3. Project rows and columns onto low-dimensional space:

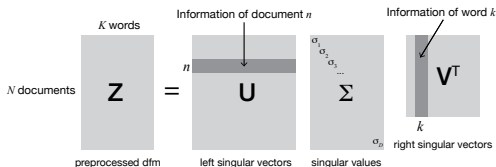
$$\theta = D_r^{1/2}U \text{ for rows (documents)}$$

$$\phi = D_c^{1/2}V \text{ for columns (words)}$$

Mathematically close to **log-linear poisson regression model**

(Lowe, 2008)

Singular Value Decomposition



- ▶ Matrix $Z_{n \times k}$ can be represented in a dimensionality equal to its rank d as:

$$Z_{n \times k} = U_{n \times d} \Sigma_{d \times d} V'^T_{d \times k}$$

- ▶ The U , Σ , and V matrixes “relocate” the elements of Z onto new coordinate vectors in d -dimensional Euclidean space
- ▶ Row variables of Z become points on the U column coordinates, and the column variables of Z become points on the V column coordinates
- ▶ The coordinate vectors are perpendicular (*orthogonal*) to each other and are normalized to unit length

Start with a dfm

	fool	wit	soldier	enemy
As You Like It	36	20	2	5
Twelfth Night	58	15	2	5
Julius Caesar	1	10	10	10
Henry V	4	3	32	10

Calculate the marginals

	fool	wit	soldier	enemy	sum (doc. len.)
As You Like It	36	20	2	5	63
Twelfth Night	58	15	2	5	80
Julius Caesar	1	10	10	10	23
Henry V	4	3	32	10	49
sum (term freq.)	21	345	99	40	215

$$P = X / \sum_i \sum_j x_{ij}$$

- ▶ Calculate a table of observed proportions by dividing each row by 215.
- ▶ The margins are the **row masses** r and the **column masses** c

	fool	wit	soldier	enemy	row mass r_i
As You Like It	0.1676	0.0935	0.0092	0.023	0.460
Twelfth Night	0.2694	0.0694	0.0092	0.023	0.186
Julius Caesar	0.0046	0.0092	0.0462	0.046	0.213
Henry V	0.0185	0.0138	0.1480	0.046	0.139
column mass c_j	0.293	0.372	0.106	0.227	1

$$E = rc'$$

- rc' gives us the expected proportions under an independence model (rows and columns have no relationship)

	fool	wit	soldier	enemy
As You Like It	0.134	0.171	0.049	0.104
Twelfth Night	0.054	0.069	0.019	0.042
Julius Caesar	0.062	0.079	0.022	0.048
Henry V	0.040	0.051	0.014	0.031

$$R = P - E$$

- ▶ Subtract the expected proportions E under independence from the observed proportions P
- ▶ This gives us the residuals R
- ▶ Interpretation: **fool** is used more frequently in the comedies than it is in the dramas, and much more frequently in **Twelfth Night**.

	fool	wit	soldier	enemy
As You Like It	0.032	-0.07804	-0.039	-0.0819
Twelfth Night	0.215	0.00039	-0.010	-0.0198
Julius Caesar	-0.058	-0.07077	0.023	-0.0026
Henry V	-0.022	-0.03764	0.133	0.0145

$$I = R/E$$

- To normalize by document length, we divide by E

	fool	wit	soldier	enemy
As You Like It	0.240	-0.457	-0.811	-0.778
Twelfth Night	3.948	0.007	-0.532	-0.451
Julius Caesar	-0.925	-0.883	1.032	-0.046
Henry V	-0.544	-0.731	8.971	0.462

$$Z = I * \sqrt{E}$$

- ▶ Give cells with a higher expected value a higher weight
- ▶ This controls for higher sampling error in cells with lower counts.

	fool	wit	soldier	enemy
As You Like It	0.088	-0.189	-0.180	-0.252
Twelfth Night	0.921	0.002	-0.075	-0.092
Julius Caesar	-0.231	-0.249	0.156	-0.010
Henry V	-0.110	-0.1666	1.096	0.082

Calculate θ and ϕ with $SVD(Z)$

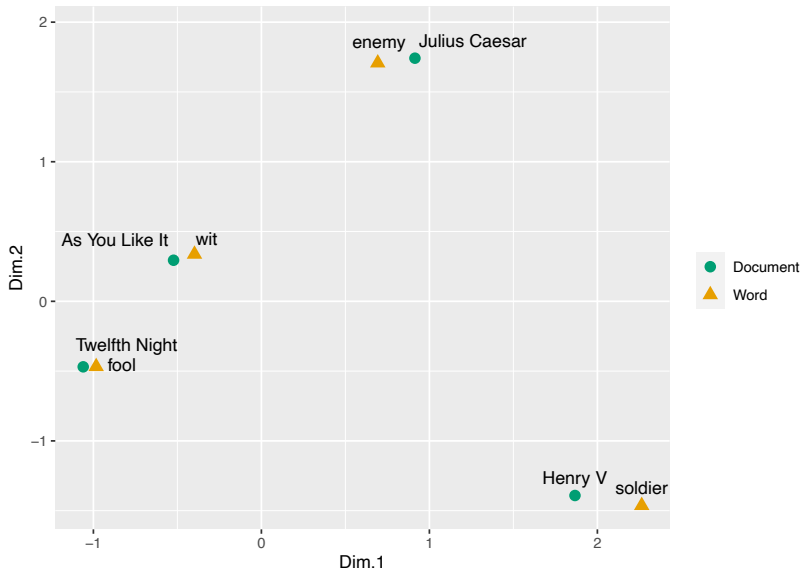
$$\theta = U/\sqrt{c}$$

	Dim.1	Dim.2	Dim.3	Dim.4
As You Like It	-0.522	0.294	1.084	0.797
Twelfth Night	-1.060	-0.469	-1.424	1.414
Julius Caesar	0.913	1.741	-0.552	0.707
Henry V	1.866	-1.391	0.337	1.278

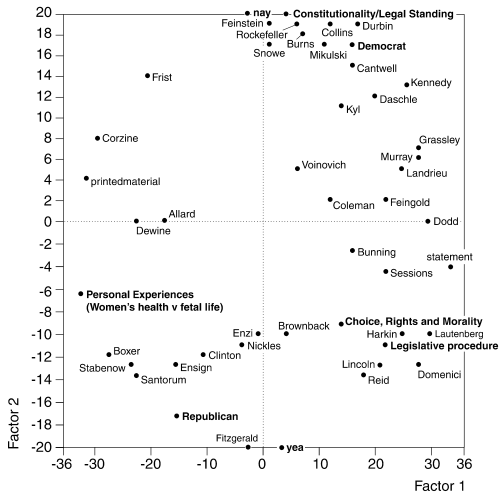
$$\phi = V/\sqrt{r}$$

	Dim.1	Dim.2	Dim.3	Dim.4
fool	-0.983	-0.467	-0.809	1.253
wit	-0.398	0.337	1.383	0.707
soldier	2.263	-1.463	0.285	1.414
enemy	0.692	1.707	-0.616	0.782

Plot of first two component values of θ and ϕ



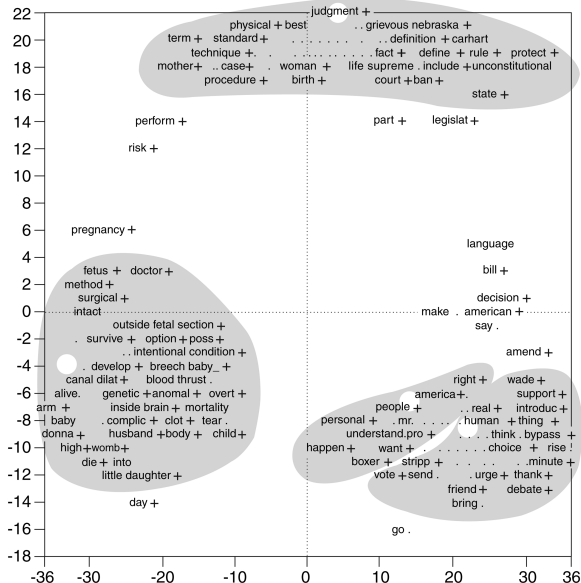
Example: Schonhardt-Bailey (2008) - speakers



	Eigenvalue	% Association	% Cumulative
Factor 1	0.30	44.4	44.4
Factor 2	0.22	32.9	77.3

Fig. 3. Correspondence analysis of classes and tags from Senate debates on Partial-Birth Abortion Ban Act

Example: Schonhardt-Bailey (2008) - words



Outline

- ▶ Parametric scaling models: Wordfish and Wordshoal
- ▶ Non-parametric scaling methods: correspondence analysis
- ▶ Practical aspects: computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Wordfish examples

Interpreting scaled dimensions

How can we validate that we are measuring a construct of interest?

1. Semantic validity

- ▶ Most discriminant words correspond to extremes of dimension of interest

2. Convergent/discriminant construct validity

- ▶ Estimated positions match other existing measures where they should match, and depart where they should depart

3. Predictive validity

- ▶ Variation in positions or word usage corresponds with expected events

4. Hypothesis validity

- ▶ Variation in positions or word usage can be used effectively to test substantive hypotheses

How to account for uncertainty in parametric models

► Option 1: Analytical derivatives

- Reformulating the Poisson model as a multinomial model, we can compute a Hessian for the log-likelihood function
- The standard errors on the θ_i parameters can be computed from the covariance matrix from the log-likelihood estimation (square roots of the diagonal)
- The covariance matrix is (asymptotically) the inverse of the negative of the Hessian
(where the negative Hessian is the observed Fisher information matrix, a.k.a. the second derivative of the log-likelihood evaluated at the maximum likelihood estimates)
- Problem: These are *too small*

How to account for uncertainty in parametric models

- ▶ Option 2: **Parametric bootstrapping** (Slapin and Proksch, Lewis and Poole)

Assume the distribution of the parameters, and generate data after drawing new parameters from these distributions.

Issues:

- ▶ slow
 - ▶ relies heavily (twice now) on parametric assumptions
 - ▶ requires some choices to be made with respect to data generation in simulations
- ▶ Option 3: **Non-parametric bootstrapping**
 - ▶ draw new versions of the texts, refit the model, save the parameters, average over the parameters
 - ▶ slow
 - ▶ not clear how the texts should be resampled
- ▶ (and yes of course) Posterior sampling from MCMC

How to account for uncertainty in non-parametric models

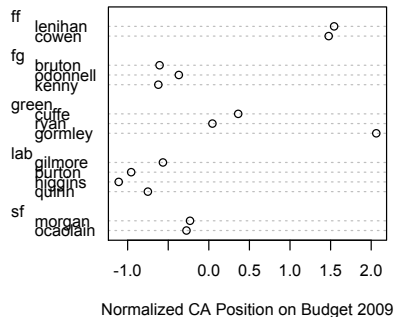
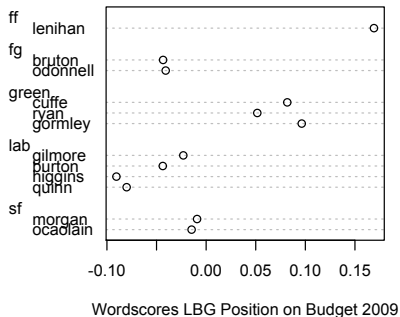
- ▶ There are problems with bootstrapping results from correspondence analysis (Milan and Whittaker 2004):
 - ▶ rotation of the principal components
 - ▶ inversion of singular values
 - ▶ reflection in an axis
- ▶ Ignore the problem and hope it will go away?
 - ▶ SVD-based methods (e.g. correspondence analysis) typically do not present errors
 - ▶ and traditionally, point estimates based on other methods have not either

Interpreting multiple dimensions

To get one dimension for each policy area, split up the document by hand and use the subparts as documents (the Slapin and Proksch method). There is currently *no* implementation of Wordscores or Wordfish that extracts two or more dimensions at once.

- ▶ But since Wordfish is a type of factor analysis model, there is no reason in principle why it could not
- ▶ Correspondence analysis by definition gives you multiple dimensions

What happens if we include irrelevant text?



What happens if we include irrelevant text?



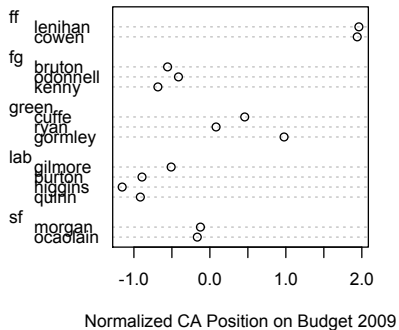
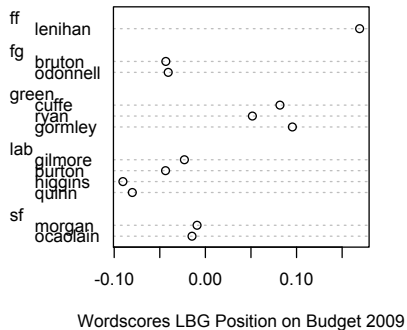
John Gormley: leader of the Green Party and Minister for the Environment, Heritage and Local Government

“As leader of the Green Party I want to take this opportunity to set out my party’s position on budget 2010...”

[772 words later]

“I will now comment on some specific aspects of my Department’s Estimate. I will concentrate on the principal sectors within the Department’s very broad remit ...”

Without irrelevant text



Outline

- ▶ Basics of unsupervised scaling methods
- ▶ Parametric scaling models: Wordfish and Wordshoal
- ▶ Non-parametric scaling methods: correspondence analysis
- ▶ Practical aspects: computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Wordfish examples