

Scaling Text

Kenneth Benoit

Quants 3: Quantitative Text Analysis

March 23th, 2018

Outline

- ▶ From classification to supervised scaling methods
 - ▶ Basics of supervised scaling methods
 - ▶ Wordscores
 - ▶ Practical aspects
 - ▶ From supervised learning to supervised scaling
 - ▶ Examples
- ▶ Unsupervised scaling of documents
 - ▶ Basics of supervised 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
- ▶ Unsupervised scaling of features

From Classification to Scaling

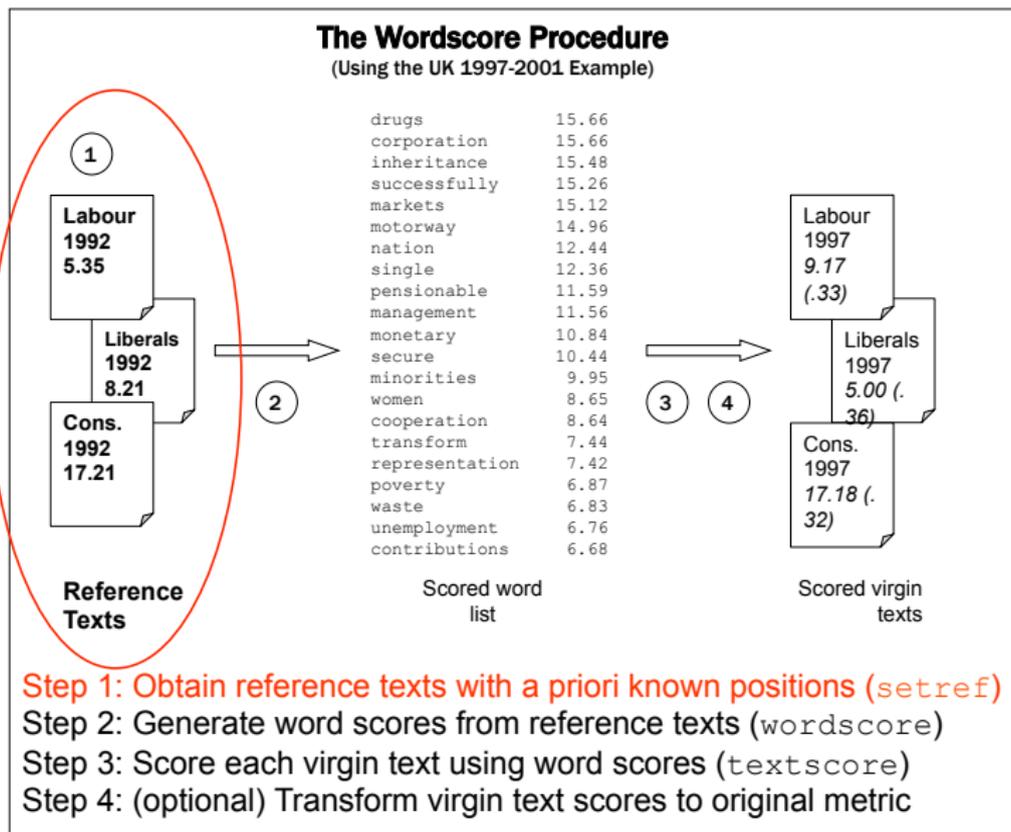
- ▶ Machine learning focuses on identifying classes (**classification**), while social science is typically interested in locating things on latent traits (**scaling**), e.g.:
 - ▶ Policy positions on economic vs social dimension
 - ▶ Inter- and intra-party differences
 - ▶ Soft news vs hard news
 - ▶ Petitioner vs respondent in legal briefs
 - ▶ ...and any other continuous scale
- ▶ But the two methods overlap and can be adapted – will demonstrate later using the Naive Bayes classifier
- ▶ In fact, the class predictions for a collection of words from NB can be adapted to scaling

Supervised scaling methods

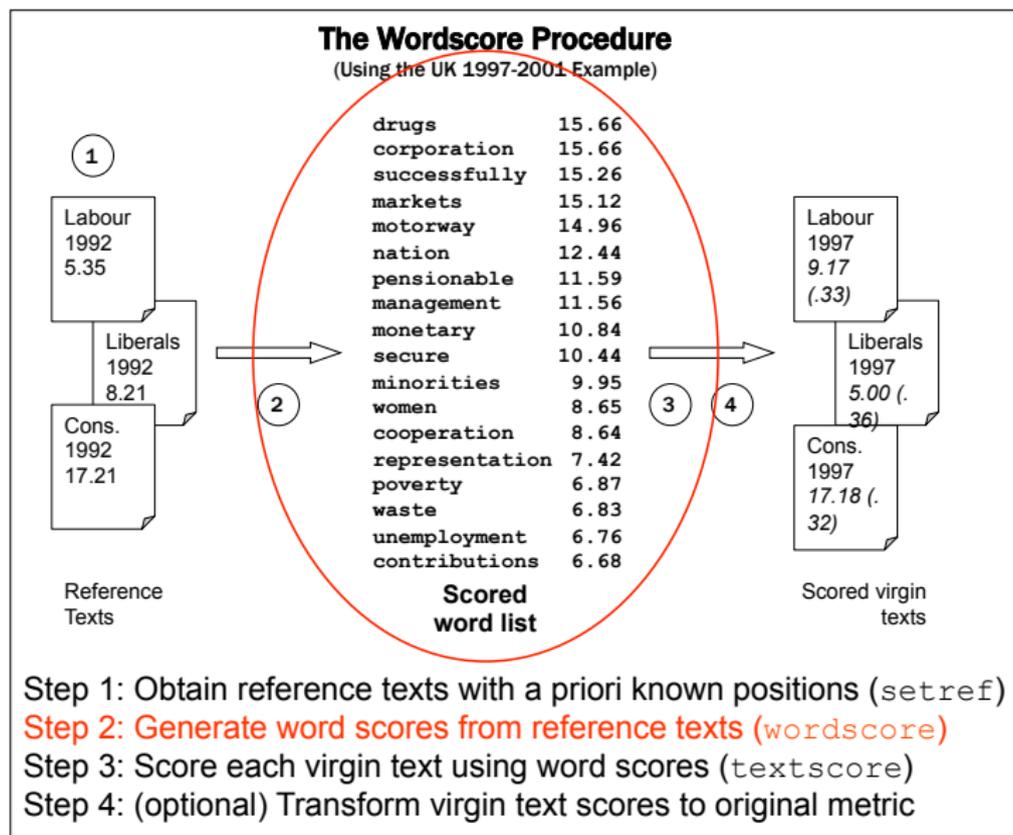
Wordscores method (Laver, Benoit & Garry, 2003):

- ▶ Two sets of texts
 - ▶ **Reference texts**: texts about which we know something (a scalar dimensional score)
 - ▶ **Virgin texts**: texts about which we know nothing (but whose dimensional score we'd like to know)
- ▶ These are analogous to a “training set” and a “test set” in classification
- ▶ Basic procedure:
 1. Analyze reference texts to obtain word scores
 2. Use word scores to score virgin texts

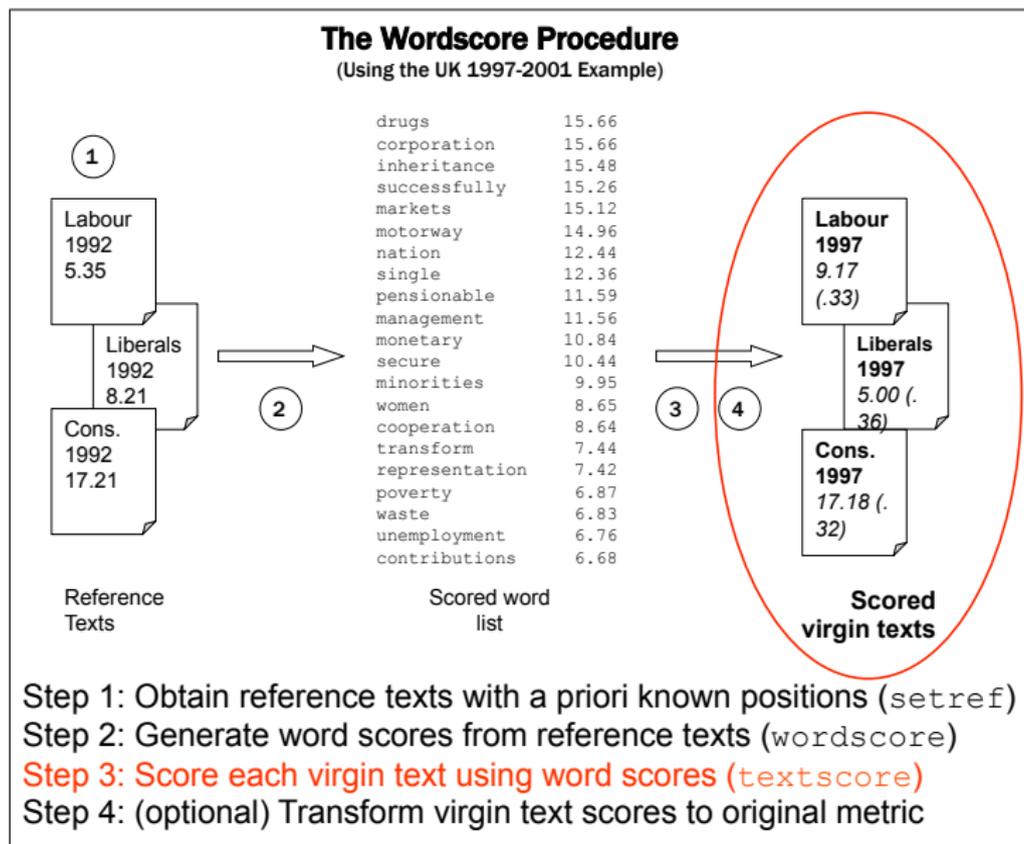
Wordscores Procedure



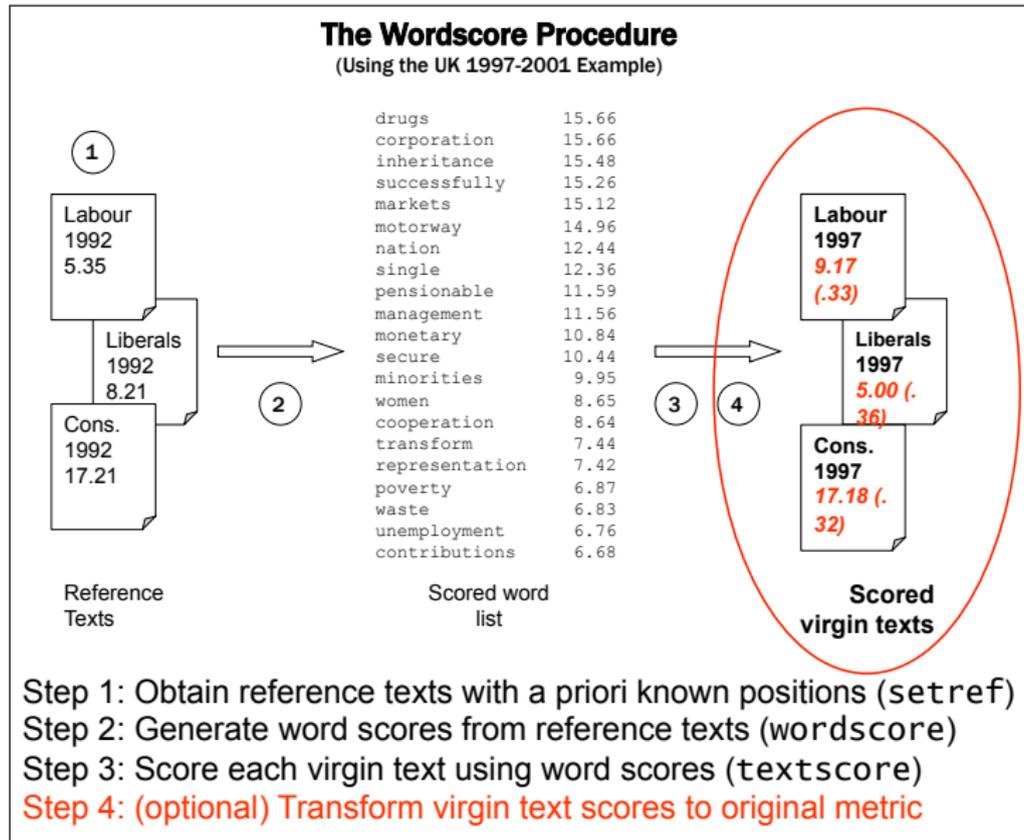
Wordscores Procedure



Wordscores Procedure



Wordscores Procedure



Wordscores mathematically: Reference texts

- ▶ Start with a set of I reference texts, represented by an $I \times J$ document-feature matrix C_{ij} , where i indexes the document and j indexes the J total word types
- ▶ Each text will have an associated “score” a_i , which is a single number locating this text on a single dimension of difference
 - ▶ This can be on a scale metric, such as 1–20
 - ▶ Can use arbitrary endpoints, such as -1, 1
- ▶ We *normalize* the document-feature matrix within each document by converting C_{ij} into a *relative* document-feature matrix (within document), by dividing C_{ij} by its word total marginals:

$$F_{ij} = \frac{C_{ij}}{C_{i.}} \quad (1)$$

where $C_{i.} = \sum_{j=1}^J C_{ij}$

Wordscores mathematically: Word scores

- ▶ Compute an $I \times J$ matrix of relative document probabilities P_{ij} for each word in each reference text, as

$$P_{ij} = \frac{F_{ij}}{\sum_{i=1}^I F_{ij}} \quad (2)$$

- ▶ This tells us the probability that given the observation of a specific word j , that we are reading a text of a certain reference document i

Wordscores mathematically: Word scores (example)

- ▶ Assume we have two reference texts, A and B
- ▶ The word “choice” is used 10 times per 1,000 words in Text A and 30 times per 1,000 words in Text B
- ▶ So F_i “choice” = $\{.010, .030\}$
- ▶ If we know only that we are reading the word choice in one of the two reference texts, then probability is 0.25 that we are reading Text A, and 0.75 that we are reading Text B

$$P_A \text{ "choice"} = \frac{.010}{(.010 + .030)} = 0.25 \quad (3)$$

$$P_B \text{ "choice"} = \frac{.030}{(.010 + .030)} = 0.75 \quad (4)$$

Wordscores mathematically: Word scores

- ▶ Compute a J -length “score” vector S for each word j as the average of each document i 's scores a_i , weighted by each word's P_{ij} :

$$S_j = \sum_{i=1}^I a_i P_{ij} \quad (5)$$

- ▶ In matrix algebra, $S = a \cdot P$
 $1 \times J \quad 1 \times I \quad I \times J$
- ▶ This procedure will yield a single “score” for every word that reflects the balance of the scores of the reference documents, weighted by the relative document frequency of its normalized term frequency

Wordscores mathematically: Word scores

- ▶ Continuing with our example:
 - ▶ We “know” (from independent sources) that Reference Text A has a position of -1.0 , and Reference Text B has a position of $+1.0$
 - ▶ The score of the word “choice” is then
$$0.25(-1.0) + 0.75(1.0) = -0.25 + 0.75 = +0.50$$

Wordscores mathematically: Scoring “virgin” texts

- ▶ Here the objective is to obtain a single score for any new text, relative to the reference texts
- ▶ We do this by taking the mean of the scores of its words, weighted by their term frequency
- ▶ So the score v_k of a virgin document k consisting of the j word types is:

$$v_k = \sum_j (F_{kj} \cdot s_j) \quad (6)$$

where $F_{kj} = \frac{C_{kj}}{C_k}$ as in the reference document relative word frequencies

- ▶ Note that **new words** outside of the set J may appear in the K virgin documents — these are simply ignored (because we have no information on their scores)
- ▶ Note also that nothing prohibits reference documents from also being scored as virgin documents

Wordscores mathematically: Rescaling raw text scores

- ▶ Because of overlapping or non-discriminating words, the raw text scores will be dragged to the interior of the reference scores (we will see this shortly in the results)
- ▶ Some procedures can be applied to rescale them, either to a unit normal metric or to a more “natural” metric
- ▶ Martin and Vanberg (2008) have proposed alternatives to the LBG (2003) rescaling

Computing confidence intervals

- ▶ The score v_k of any text represents a weighted mean
- ▶ LBG (2003) used this logic to develop a standard error of this mean using a *weighted variance* of the scores in the virgin text
- ▶ Given some assumptions about the scores being fixed (and the words being conditionally independent), this yields approximately normally distributed errors for each v_k
- ▶ An alternative would be to bootstrap the textual data prior to constructing C_{ij} and C_{kj} — see Lowe and Benoit (2012)

Pros and Cons of the Wordscores approach

- ▶ Fully automated technique with minimal human intervention or judgment calls – only with regard to reference text selection
- ▶ Language-blind: all we need to know are reference scores
- ▶ Could potentially work on texts like this:

ፍጥፍ ያደካዎቹ ርጅ ያደርጋሉጋህረግ ወደወ
ሃደላጎርጋሉጋ ጎርጎ ርጅ ሃገሪዎቹ ወደወጋ
ፊገግረጋህረግ ለሌላጎርጋሉጋ

(See <http://www.kli.org>)

Pros and Cons of the Wordscores approach

- ▶ Estimates unknown positions on a priori scales – hence no inductive scaling with a posteriori interpretation of unknown policy space
- ▶ Very dependent on correct identification of:
 - ▶ appropriate [reference texts](#)
 - ▶ appropriate [reference scores](#)

Suggestions for choosing reference texts

- ▶ Texts need to contain information representing a clearly dimensional position
- ▶ Dimension must be known a priori. Sources might include:
 - ▶ Survey scores or manifesto scores
 - ▶ Arbitrarily defined scales (e.g. -1.0 and 1.0)
- ▶ Should be as discriminating as possible: extreme texts on the dimension of interest, to provide reference anchors
- ▶ Need to be from the same lexical universe as virgin texts
- ▶ Should contain lots of words

Suggestions for choosing reference values

- ▶ Must be “known” through some trusted external source
- ▶ For any pair of reference values, all scores are simply linear rescalings, so might as well use $(-1, 1)$
- ▶ The “middle point” will not be the midpoint, however, since this will depend on the relative word frequency of the reference documents
- ▶ Reference texts if scored as virgin texts will have document scores more extreme than other virgin texts
- ▶ With three or more reference values, the mid-point is mapped onto a multi-dimensional simplex. The values now matter but only in relative terms (we are still investigating this fully)

Multinomial Bayes model of Class given a Word

Class posterior probabilities

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})}$$

- ▶ This represents the **posterior probability of membership in class k** for word j
- ▶ Under *certain conditions*, this is identical to what LBG (2003) called P_{wr}
- ▶ Under those conditions, the LBG “wordscore” is the **linear difference between $P(c_k|w_j)$ and $P(c_{\neg k}|w_j)$**

“Certain conditions”

- ▶ The LBG approach required the identification not only of texts for each training class, but also “reference” scores attached to each training class
- ▶ Consider two “reference” scores s_1 and s_2 attached to two classes $k = 1$ and $k = 2$. Taking P_1 as the posterior $P(k = 1|w = j)$ and P_2 as $P(k = 2|w = j)$, A generalised score s_j^* for the word j is then

$$\begin{aligned} s_j^* &= s_1 P_1 + s_2 P_2 \\ &= s_1 P_1 + s_2 (1 - P_1) \\ &= s_1 P_1 + s_2 - s_2 P_1 \\ &= P_1 (s_1 - s_2) + s_2 \end{aligned}$$

“Certain conditions”: More than two reference classes

- ▶ For more than two reference classes, if the reference scores are ordered such that $s_1 < s_2 < \dots < s_K$, then

$$\begin{aligned} s_j^* &= s_1 P_1 + s_2 P_2 + \dots + s_K P_K \\ &= s_1 P_1 + s_2 P_2 + \dots + s_K \left(1 - \sum_{k=1}^{K-1} P_k\right) \\ &= \sum_{k=1}^{K-1} P_k (s_k - s_K) + s_K \end{aligned}$$

A simpler formulation:

Use reference scores such that $s_1 = -1.0, s_K = 1.0$

- ▶ From above equations, it should be clear that any set of reference scores can be linearly rescaled to endpoints of $-1.0, 1.0$
- ▶ This simplifies the “simple word score”

$$s_j^* = (1 - 2P_1) + \sum_{k=2}^{K-1} P_k (s_k - 1)$$

- ▶ which simplifies with just two reference classes to:

$$s_j^* = 1 - 2P_1$$

Implications

- ▶ LBG's “word scores” come from a linear combination of class posterior probabilities from a Bayesian model of class conditional on words
- ▶ We might as well always anchor reference scores at $-1.0, 1.0$
- ▶ There is a special role for reference classes in between $-1.0, 1.0$, as they balance between “pure” classes — more in a moment
- ▶ There are alternative scaling models, such that used in Beauchamp's (2012) “Bayesscore”, which is simply the difference in logged class posteriors at the word level. For $s_1 = -1.0, s_2 = 1.0$,

$$\begin{aligned} s_j^B &= -\log P_1 + \log P_2 \\ &= \log \frac{1 - P_1}{P_1} \end{aligned}$$

Moving to the document level

- ▶ The “Naive” Bayes model of a joint document-level class posterior assumes conditional independence, to multiply the word likelihoods from a “test” document, to produce:

$$P(c|d) = P(c) \frac{\prod_j P(w_j|c)}{P(w_j)}$$

- ▶ So we *could* consider a document-level relative score, e.g. $1 - 2P(c_1|d)$ (for a two-class problem)
- ▶ But this turns out to be *useless*, since the predictions of class are **highly separated**

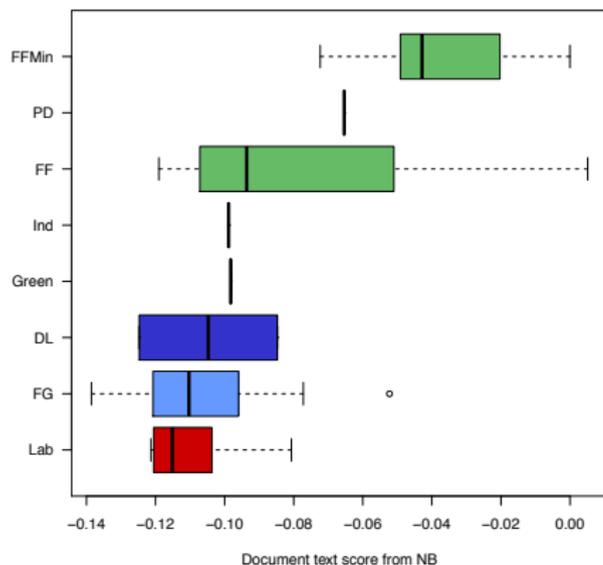
Moving to the document level

- ▶ A better solution is to score a test document as the **arithmetic mean** of the **scores of its words**
- ▶ This is exactly the solution proposed by LBG (2003)
- ▶ Beauchamp (2012) proposes a “Bayesscore” which is the arithmetic mean of the log difference word scores in a document – which yields extremely similar results

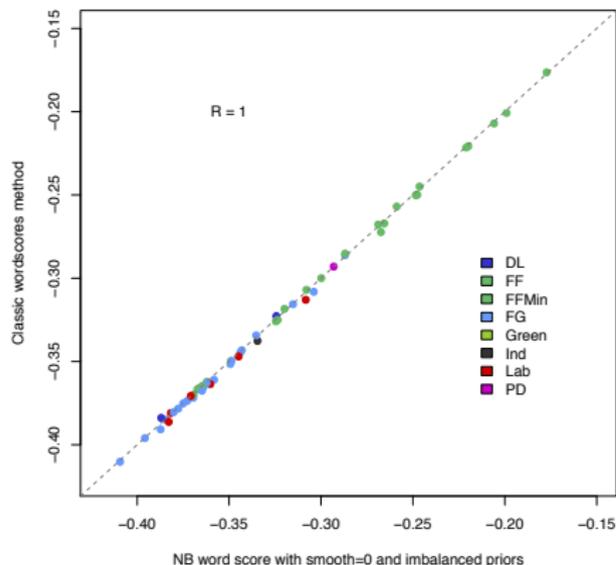
And now for some demonstrations with data...

Application 1: Daily speeches from LBG (2003)

(a) NB Speech scores by party, smooth=0, imbalanced priors



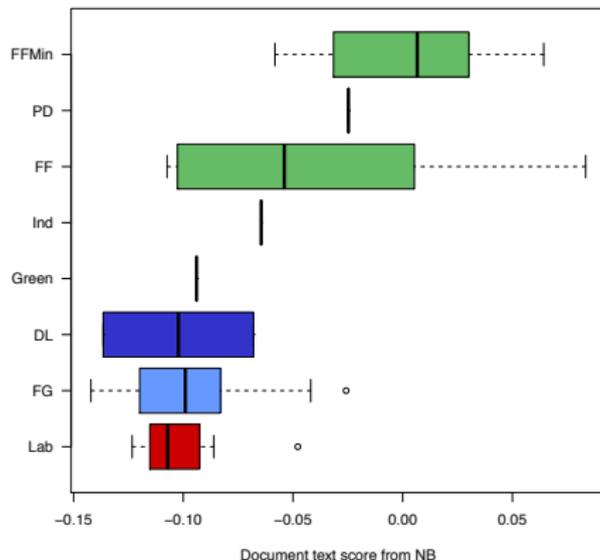
(b) Document scores from NB v. Classic Wordscores



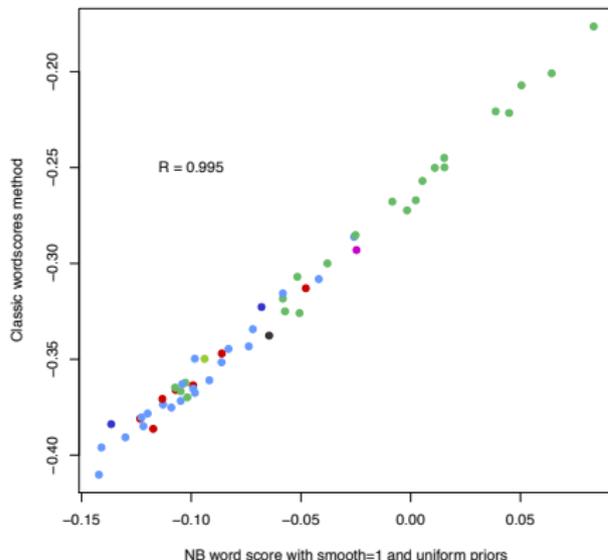
- ▶ three reference classes (Opposition, Opposition, Government) at $\{-1, -1, 1\}$
- ▶ no smoothing

Application 1: Daily speeches from LBG (2003)

(c) NB Speech scores by party, smooth=1, uniform class priors



(d) Document scores from NB v. Classic Wordscores

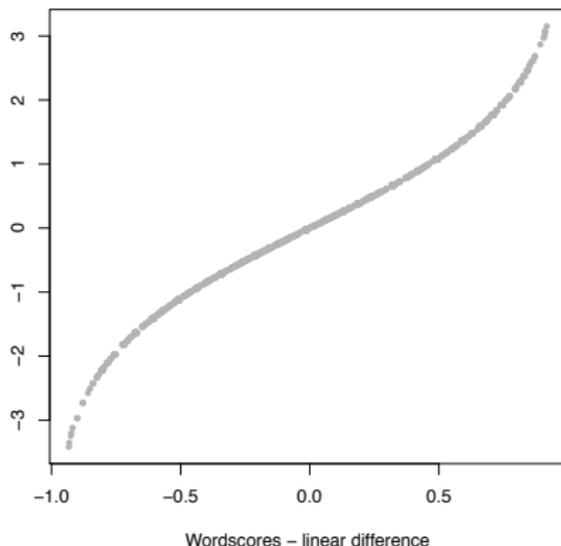


- ▶ two reference classes (Opposition+Opposition, Government) at $\{-1, 1\}$
- ▶ Laplace smoothing

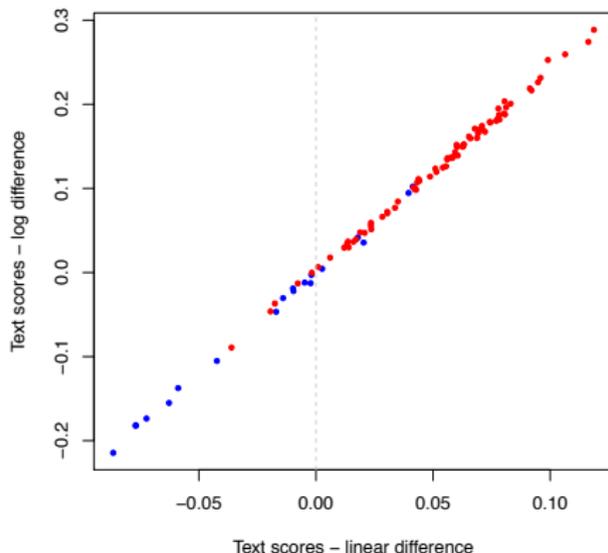
Application 2: Classifying legal briefs (Evans et al 2007)

Wordscores v. Bayesscore

(a) Word level



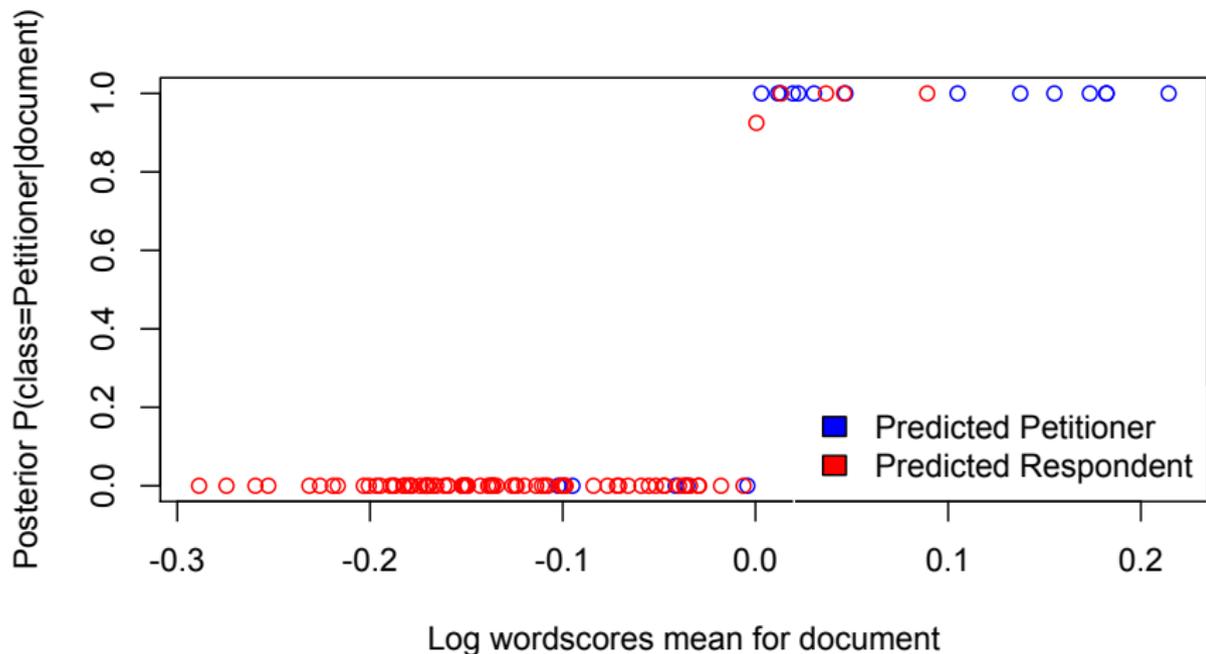
(b) Document level



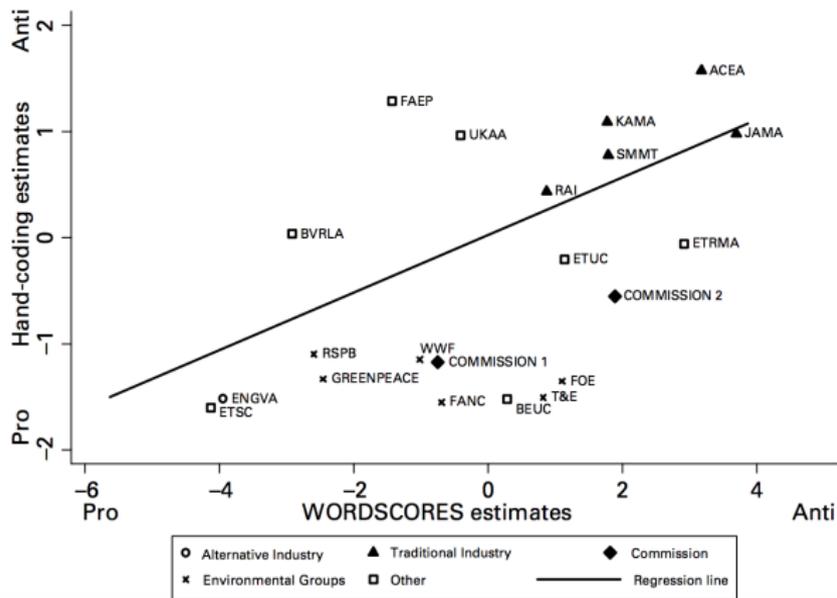
- ▶ Training set: **P**etitioner and **R**espondent litigant briefs from *Grutter/Gratz v. Bollinger* (a U.S. Supreme Court case)
- ▶ Test set: 98 amicus curiae briefs (whose **P** or **R** class is known)

Application 2: Classifying legal briefs (Evans et al 2007)

Posterior class prediction from NB versus log wordscores



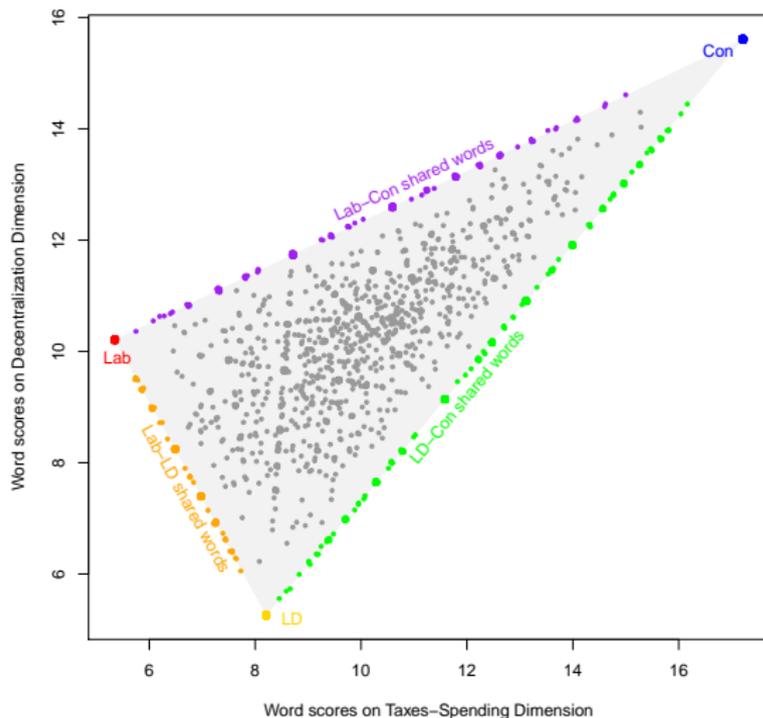
Application 3: Scaling environmental interest groups (Klüver 2009)



- ▶ Dataset: text of online consultation on EU environmental regulations
- ▶ Reference texts: most extreme pro- and anti-regulation groups

Application 4: LBG's British manifestos

More than two reference classes



- ▶ x-axis: Reference scores of $\{5.35, 8.21, 17.21\}$ for Lab, LD, Conservatives
- ▶ y-axis: Reference scores of $\{10.21, 5.26, 15.61\}$

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/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

- ▶ Unsupervised scaling of documents
 - ▶ Basics of supervised 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
- ▶ Unsupervised scaling of features
 - ▶ Word embeddings
 - ▶ Examples with word2vec

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
 - ψ_k is frequency of word k
 - 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

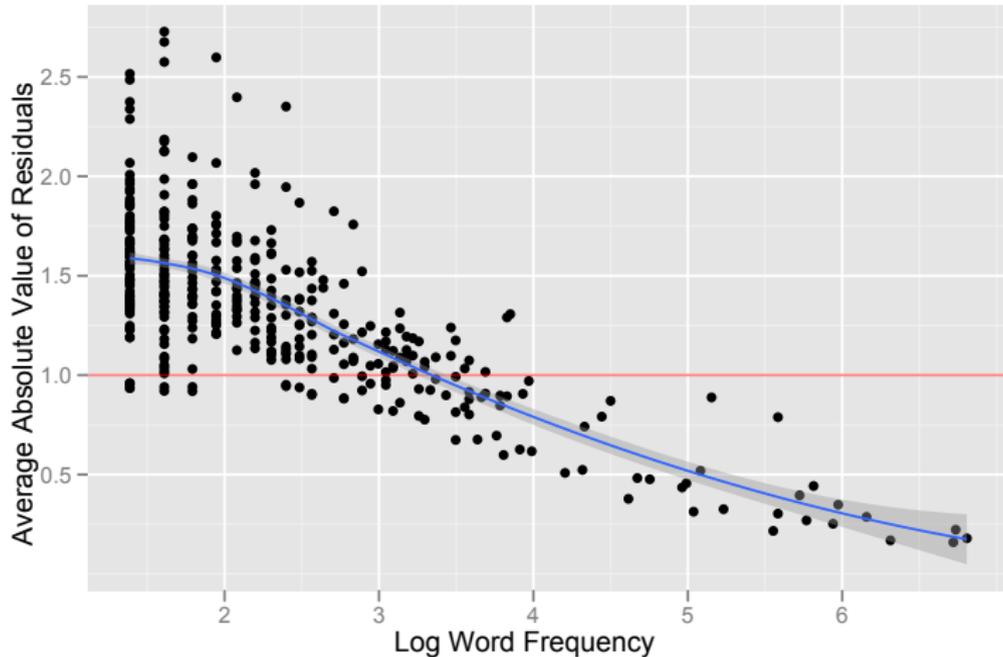
- ▶ Standard (statistical) **inference** about parameters
- ▶ **Uncertainty** accounting for parameters
- ▶ **Distributional assumptions** are made explicit (as part of the data generating process motivating the choice of stochastic distribution)
 - ▶ *conditional independence*
 - ▶ *stochastic process* (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

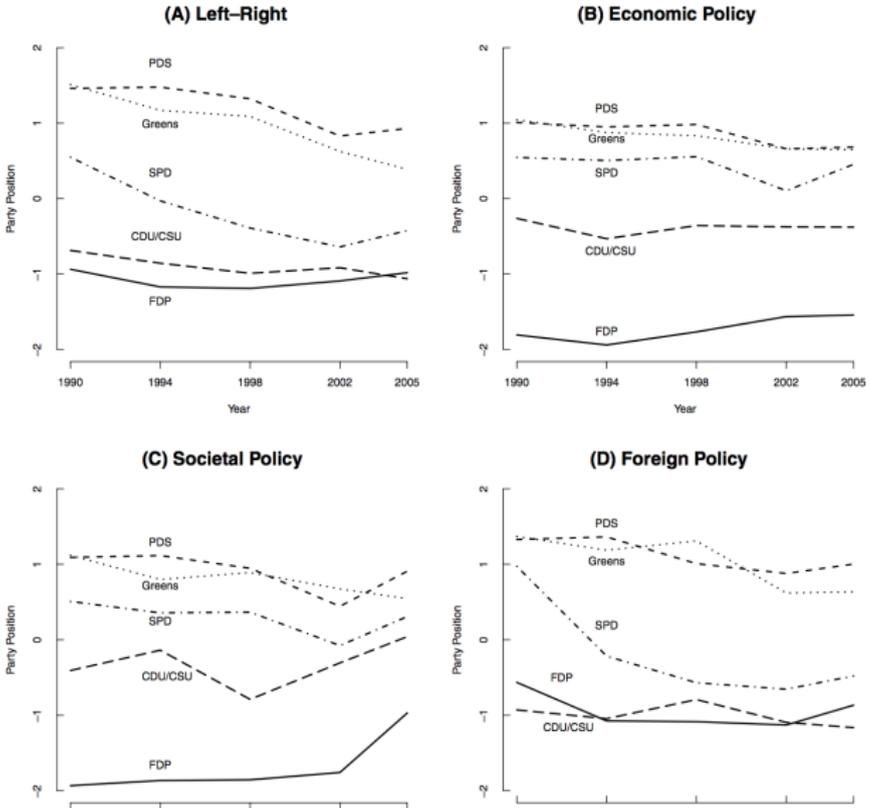
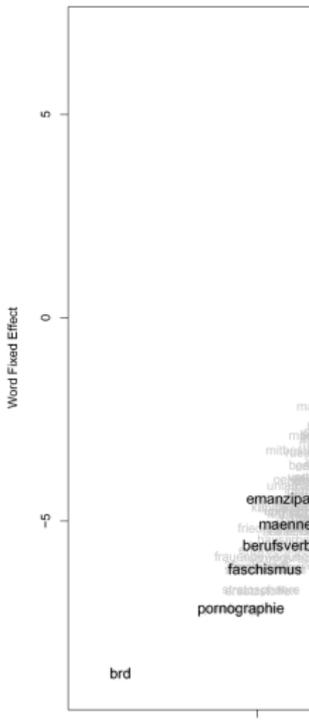


FIGURE 2 Word Weights vs. Word Fixed Effect 1990–2005 (Translation)



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

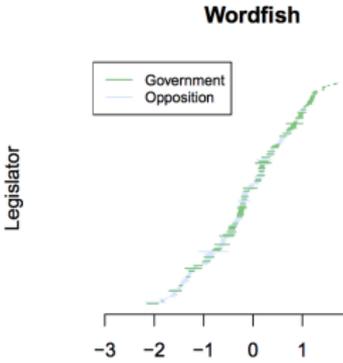
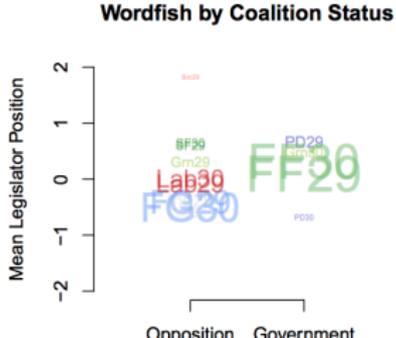
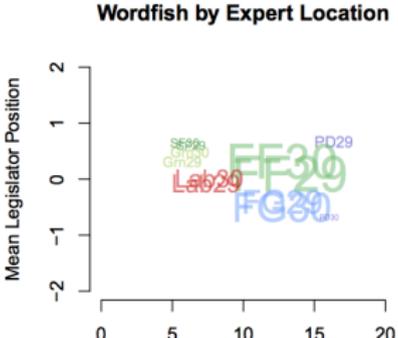
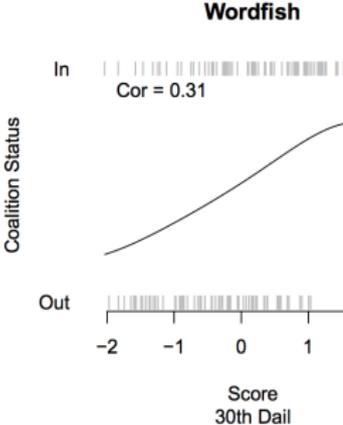
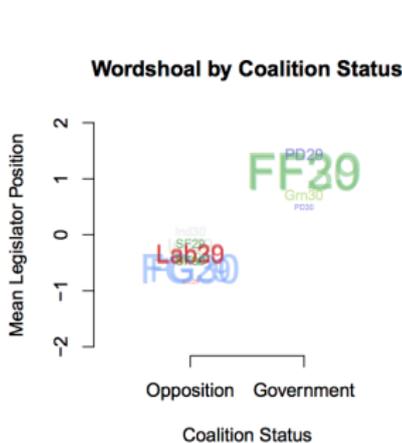
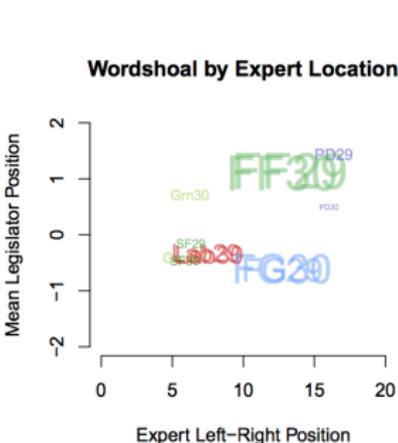
- ▶ **Intuition**: debate-specific estimates are aggregated into a single position using dimensionality reduction

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



Outline

- ▶ Unsupervised scaling of documents
 - ▶ Basics of supervised 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
- ▶ Unsupervised scaling of features
 - ▶ Word embeddings
 - ▶ Examples with word2vec

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 factor analysis for categorical data
- ▶ Following normalization of the marginals, it uses Singular Value Decomposition to reduce the dimensionality of the document-feature matrix
- ▶ This allows projection of the positioning of the words as well as the texts into multi-dimensional space
- ▶ The number of dimensions – as in factor analysis – can be decided based on the eigenvalues from the SVD

Singular Value Decomposition

- ▶ A matrix \mathbf{X} can be represented in a dimensionality equal to its rank d as:
 $n \times k$

$$\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}' \quad (7)$$

$n \times k$ $n \times d$ $d \times d$ $d \times k$

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

Correspondence analysis

1. Compute matrix of standardized residuals, \mathbf{S} :

$$\mathbf{S} = \mathbf{D}_r^{1/2}(\mathbf{P} - \mathbf{r}\mathbf{c}^T)\mathbf{D}_c^{1/2}$$

where $\mathbf{P} = \mathbf{Y} / \sum_{ij} y_{ij}$

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

$\mathbf{D}_r = \text{diag}(\mathbf{r})$, $\mathbf{D}_c = \text{diag}(\mathbf{c})$

2. Calculate SVD of \mathbf{S}
3. Project rows and columns onto low-dimensional space:

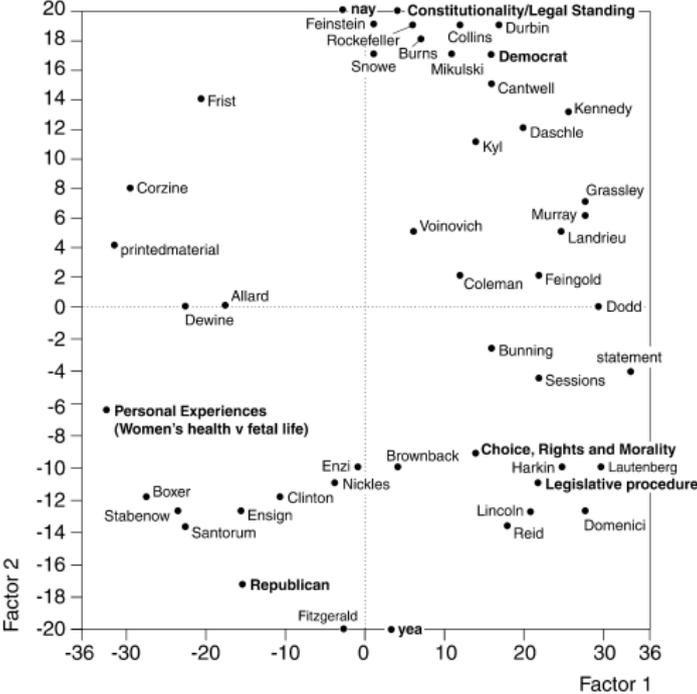
$$\theta = \mathbf{D}_r^{1/2}\mathbf{U} \text{ for rows (documents)}$$

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

Mathematically close to [log-linear poisson regression model](#)

(Lowe, 2008)

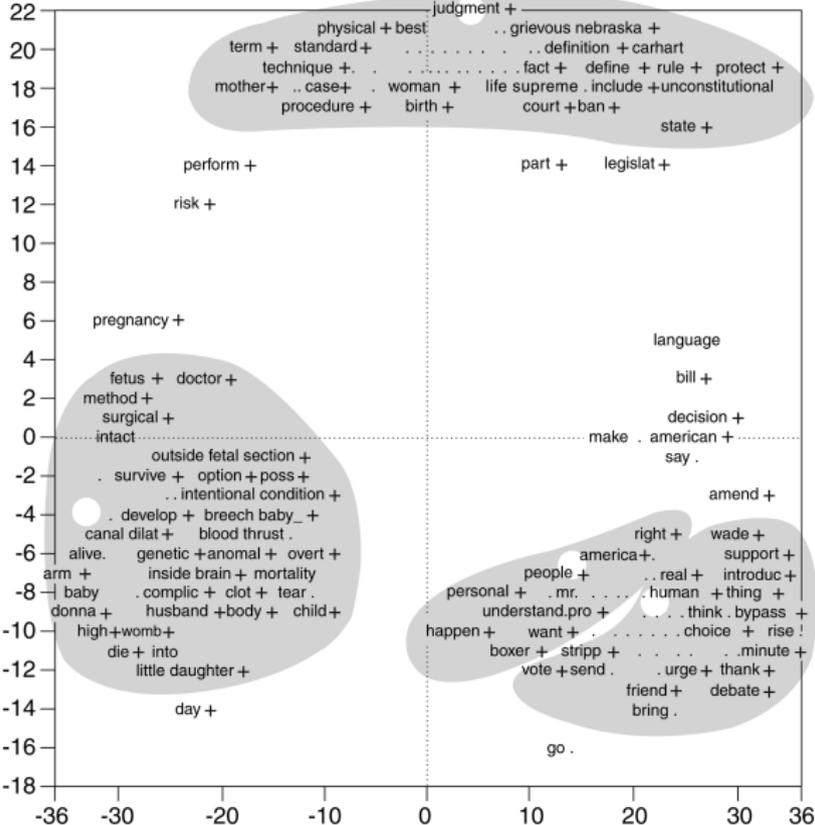
Example: Schonhardt-Bailey (2008) - speakers



	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

- ▶ Unsupervised scaling of documents
 - ▶ Basics of supervised scaling methods
 - ▶ Parametric scaling models: Wordfish and Wordshoal
 - ▶ Non-parametric scaling methods: correspondence analysis
 - ▶ Practical aspects: interpretation, computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Unsupervised scaling of features
 - ▶ Word embeddings
 - ▶ Examples with word2vec

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

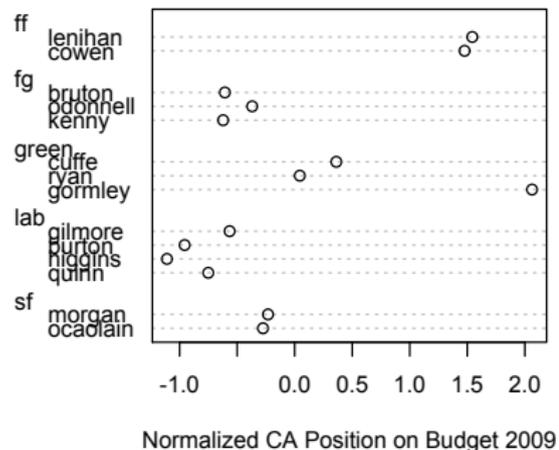
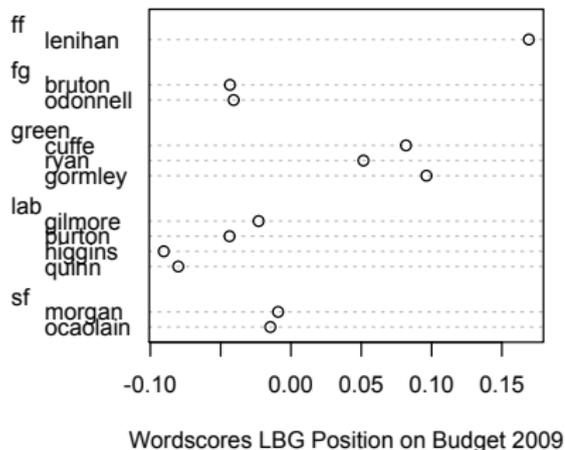
- ▶ There are problems with bootstrapping results from correspondence analysis (Milan and Whittaker 2004):
 - ▶ rotation of the principal components
 - ▶ inversion of singular values
 - ▶ reflection in an axis
- ▶ Ignore the problem and hope it will go away?
 - ▶ SVD-based methods (e.g. correspondence 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?



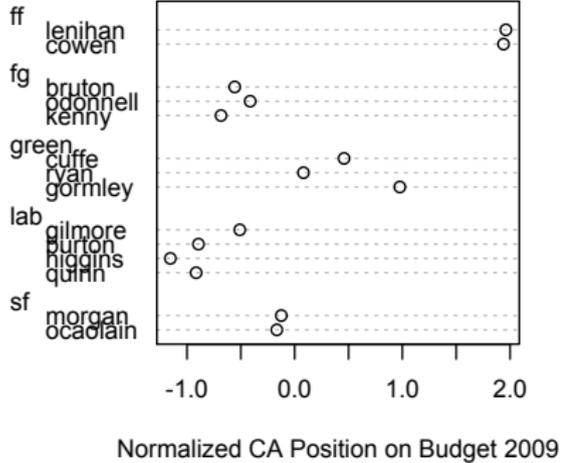
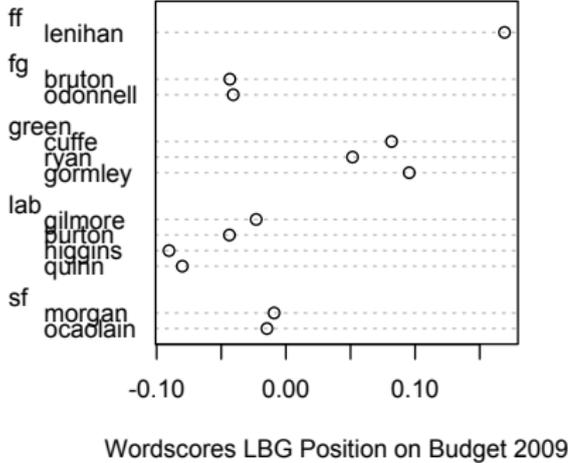
John Gormley: leader of the Green Party and Minister for the Environment, Heritage and Local Government

“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

- ▶ Unsupervised scaling of documents
 - ▶ Basics of supervised scaling methods
 - ▶ Parametric scaling models: Wordfish and Wordshoal
 - ▶ Non-parametric scaling methods: correspondence analysis
 - ▶ Practical aspects: interpretation, computing uncertainty, multiple dimensions, sensitivity to inclusion of irrelevant text
- ▶ Unsupervised scaling of features
 - ▶ Word embeddings
 - ▶ Examples with word2vec

Beyond bag-of-words

Most applications of text analysis rely on a **bag-of-words** representation of documents

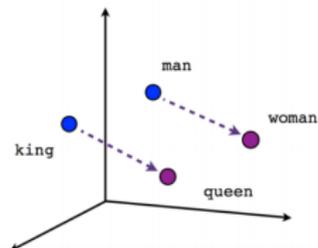
- ▶ Only relevant feature: frequency of features
- ▶ Ignores context, grammar, word order...
- ▶ Wrong but often irrelevant

One alternative: **word embeddings**

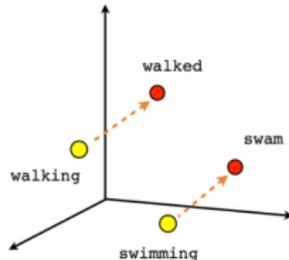
- ▶ Represent words as **real-valued vector** in a multidimensional space (often 100–500 dimensions), common to all words
- ▶ Distance in space captures syntactic and semantic regularities, i.e. words that are close in space have similar meaning
 - ▶ How? Vectors are learned based on context similarity
 - ▶ Distributional hypothesis: words that appear in the same context share semantic meaning
- ▶ Operations with vectors are also meaningful

Word embeddings example

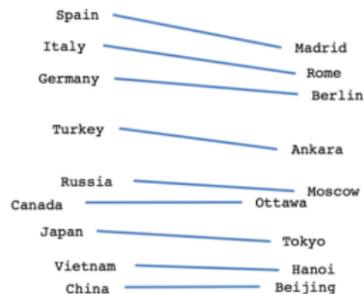
word	D_1	D_2	D_3	...	D_N
man	0.46	0.67	0.05
woman	0.46	-0.89	-0.08
king	0.79	0.96	0.02
queen	0.80	-0.58	-0.14



Male-Female



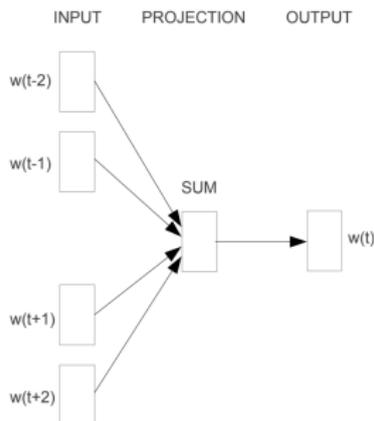
Verb tense



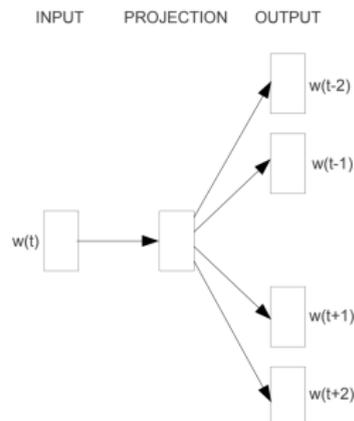
Country-Capital

word2vec (Mikolov 2013)

- ▶ Statistical method to efficiently learn word embeddings from a corpus, developed by Google engineer
- ▶ Most popular, in part because pre-trained vectors are available
- ▶ Two models to learn word embeddings:



CBOW



Skip-gram

Example: Pomeroy et al 2018

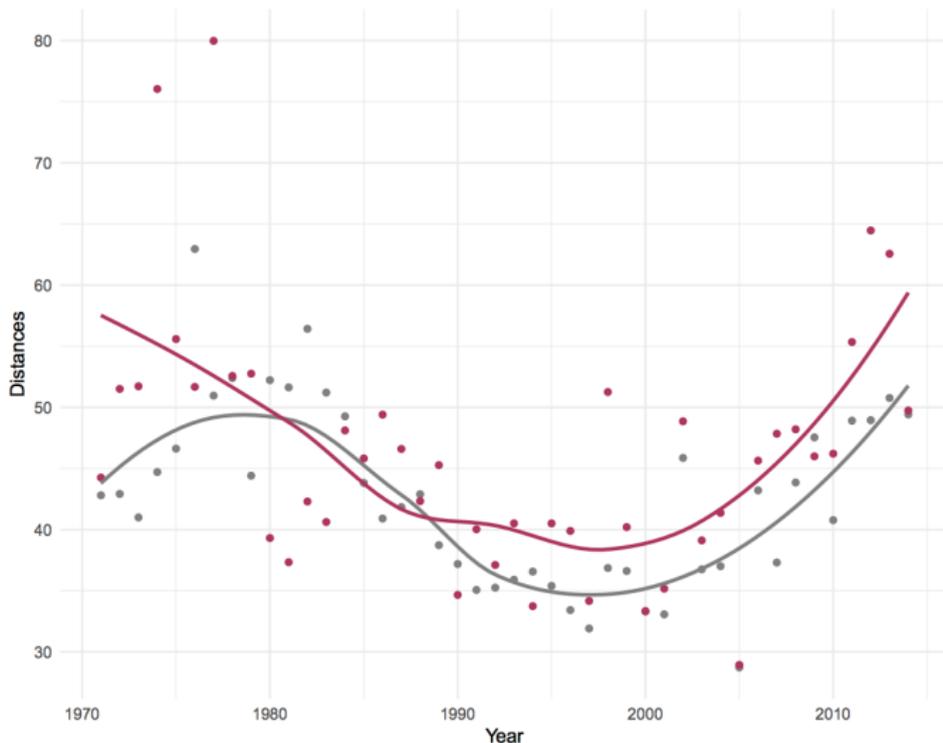


Figure 4: *Distances by core countries*. Plot of Euclidian distances between US and Russia (gray), and US and China (maroon).