

Text as Data

Justin Grimmer

Associate Professor
Department of Political Science
Stanford University

October 7th, 2014

Estimating Word Discrimination

1) Task

- a) **Classification** \rightsquigarrow learn word weights for dictionaries
- b) **Fictitious prediction problem** \rightsquigarrow Identify features that discriminate between groups to learn features that are indicative of some group

2) Objective function

$$f(\theta, \mathbf{X}) = f(\theta, \mathbf{X}, \mathbf{Y})$$

where:

\mathbf{Y} = Document Labels

\mathbf{X} = Document Features

θ = Parameters that measure words discrimination between categories

3) Optimization \rightsquigarrow method specific

4) Validation \rightsquigarrow depends on task

- i) Classification \rightsquigarrow Accuracy, Precision, Recall
- ii) Fictitious prediction \rightsquigarrow Face, convergent, discriminatory, and **confound**

Stylometry ~ Who Wrote Disputed Federalist Papers?

Federalist papers ~ Mosteller and Wallace (1963)

- Persuade citizens of New York State to adopt constitution
- Canonical texts in study of American politics
- 77 essays
 - Published from 1787-1788 in Newspapers
 - And under the name **Publius**, anonymously

Who Wrote the Federalist papers?

- Jay wrote essays 2, 3, 4,5, and 64
- Hamilton: wrote 43 papers
- Madison: wrote 12 papers

Disputed: Hamilton or Madison?

- Essays: 49-58, 62, and 63
- Joint Essays: 18-20

Task: identify authors of the disputed papers.

Task: Classify papers as Hamilton or Madison using dictionary methods

Setting up the Analysis

Training \rightsquigarrow papers Hamilton, Madison are known to have authored

Test \rightsquigarrow unlabeled papers

Preprocessing:

- Hamilton/Madison both discuss similar issues
- Differ in extent they use **stop words**
- Focus analysis on the stop words

Setting up the Analysis

- $\mathbf{Y} = (Y_1, Y_2, \dots, Y_N) = (\text{Hamilton, Hamilton, Madison, } \dots, \text{Hamilton})$
 $N \times 1$ matrix with author labels

- Define the number of words in federalist paper i as num_i

$$\mathbf{X} = \begin{pmatrix} \frac{1}{\text{num}_1} & \frac{2}{\text{num}_1} & \frac{0}{\text{num}_1} & \cdots & \frac{3}{\text{num}_1} \\ \frac{0}{\text{num}_2} & \frac{1}{\text{num}_2} & \frac{0}{\text{num}_2} & \cdots & \frac{0}{\text{num}_2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{0}{\text{num}_N} & \frac{0}{\text{num}_N} & \frac{1}{\text{num}_N} & \cdots & \frac{0}{\text{num}_N} \end{pmatrix}$$

$N \times J$ counting stop word usage rate

- $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_J)$

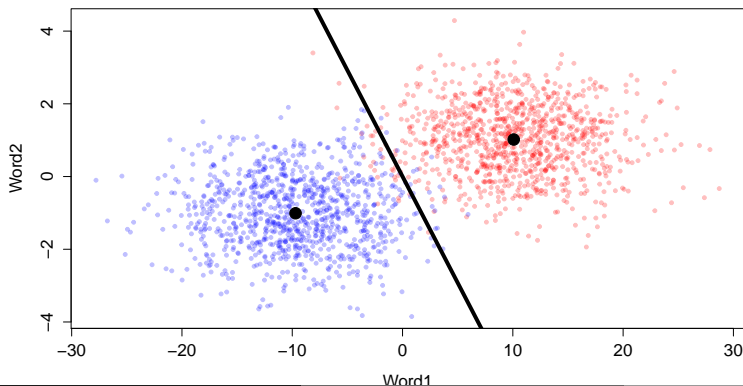
Word weights.

Objective Function

Heuristically: find $\theta^* = (\theta_1^*, \theta_2^*, \dots, \theta_J^*)$ used to create score

$$p_i = \sum_{j=1}^J \theta_j^* X_{ij}$$

that maximally discriminates between categories



Objective Function

Define:

$$\mu_{\text{Madison}} = \frac{1}{N_{\text{Madison}}} \sum_{i=1}^N I(Y_i = \text{Madison}) \mathbf{X}_i$$
$$\mu_{\text{Hamilton}} = \frac{1}{N_{\text{Hamilton}}} \sum_{i=1}^N I(Y_i = \text{Hamilton}) \mathbf{X}_i$$

Objective Function

We can then define functions that describe the “projected” mean and variance for each author

$$g(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Madison}) = \frac{1}{N_{\text{Madison}}} \sum_{i=1}^N I(Y_i = \text{Madison}) \boldsymbol{\theta}' \mathbf{X}_i = \boldsymbol{\theta}' \boldsymbol{\mu}_{\text{Madison}}$$

$$g(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Hamilton}) = \frac{1}{N_{\text{Hamilton}}} \sum_{i=1}^N I(Y_i = \text{Hamilton}) \boldsymbol{\theta}' \mathbf{X}_i = \boldsymbol{\theta}' \boldsymbol{\mu}_{\text{Hamilton}}$$

$$s(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Madison}) = \sum_{i=1}^N I(Y_i = \text{Madison}) (\boldsymbol{\theta}' \mathbf{X}_i - \boldsymbol{\theta}' \boldsymbol{\mu}_{\text{Madison}})^2$$

$$s(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Hamilton}) = \sum_{i=1}^N I(Y_i = \text{Hamilton}) (\boldsymbol{\theta}' \mathbf{X}_i - \boldsymbol{\theta}' \boldsymbol{\mu}_{\text{Hamilton}})^2$$

Objective Function \rightsquigarrow Optimization

$$\begin{aligned} f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}) &= \frac{(g(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Hamilton}) - g(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Madison}))^2}{s(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Hamilton}) + s(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y}, \text{Madison})} \\ &= \frac{(\boldsymbol{\theta}'(\boldsymbol{\mu}_{\text{Hamilton}} - \boldsymbol{\mu}_{\text{Madison}}))^2}{\text{Scatter}_{\text{Hamilton}} + \text{Scatter}_{\text{Madison}}} \end{aligned}$$

Optimization \rightsquigarrow find $\boldsymbol{\theta}^*$ to maximize $f(\boldsymbol{\theta}, \mathbf{X}, \mathbf{Y})$, assuming independence across dimensions.

(Fisher's) Linear Discriminant Analysis

Optimization \rightsquigarrow Word Weights

For each word j , construct weight θ_j^* ,

$$\mu_{j,\text{Hamilton}} = \frac{\sum_{i=1}^N I(Y_i = \text{Hamilton})X_{ij}}{\sum_{j=1}^J \sum_{i=1}^N I(Y_i = \text{Hamilton})X_{ij}}$$

$$\mu_{j,\text{Madison}} = \frac{\sum_{i=1}^N I(Y_i = \text{Madison})X_{ij}}{\sum_{j=1}^J \sum_{i=1}^N I(Y_i = \text{Madison})X_{ij}}$$

$$\sigma_{j,\text{Hamilton}}^2 = \text{Var}(X_{i,j}|\text{Hamilton})$$

$$\sigma_{j,\text{Madison}}^2 = \text{Var}(X_{i,j}|\text{Madison})$$

We can then generate weight θ_j^* as

$$\theta_j^* = \frac{\mu_{j,\text{Hamilton}} - \mu_{j,\text{Madison}}}{\sigma_{j,\text{Hamilton}}^2 + \sigma_{j,\text{Madison}}^2}$$

Optimization \rightsquigarrow Trimming the Dictionary

- Trimming weights: Focus on discriminating words (very simple **regularization**)
- Cut off: For all $|\theta_j^*| < 0.025$ set $\theta_j^* = 0$.

Classification \rightsquigarrow Determining Authorship

For each disputed document i , compute discrimination statistic

$$p_i = \sum_{j=1}^J \theta_j^* X_{ij}$$

$p_i \rightsquigarrow$ classification (**linear discriminator**)

- Above midpoint in training set \rightarrow Hamilton text
- Below midpoint in training set \rightarrow Madison text

Findings: Madison is the author of the disputed federalist papers.

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words
- Difference in Liberal, Conservative language \rightsquigarrow Ideological Language

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words
- Difference in Liberal, Conservative language \rightsquigarrow Ideological Language
- Difference in Secret/Not Secret Language \rightsquigarrow Secretive Language (Gill and Spirling 2014)

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words
- Difference in Liberal, Conservative language \rightsquigarrow Ideological Language
- Difference in Secret/Not Secret Language \rightsquigarrow Secretive Language (Gill and Spirling 2014)
- Difference in Toy advertising

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words
- Difference in Liberal, Conservative language \rightsquigarrow Ideological Language
- Difference in Secret/Not Secret Language \rightsquigarrow Secretive Language (Gill and Spirling 2014)
- Difference in Toy advertising
- Difference in Language across groups \rightsquigarrow Labeling output from Clustering/Topic Models

Inferring Separating Words

Classification \rightsquigarrow Custom Dictionaries

- Stylometry \rightsquigarrow Classify Authors
- Dictionary based classification \rightsquigarrow Gentzkow and Shapiro (2010) and measures of media slant
- Dictionary based classification \rightsquigarrow Customized to particular setting

Fictitious Prediction Problem \rightsquigarrow Infer words that are indicative of some class/group

- Difference in Republican, Democratic language \rightsquigarrow **Partisan** words
- Difference in Liberal, Conservative language \rightsquigarrow Ideological Language
- Difference in Secret/Not Secret Language \rightsquigarrow Secretive Language (Gill and Spirling 2014)
- Difference in Toy advertising
- Difference in Language across groups \rightsquigarrow Labeling output from Clustering/Topic Models

Vague and **Difficult** to derive before hand

Congressional Language Across Sources

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?
 - Yes: press releases have different purposes, targets, and need not relate to official business

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?
 - Yes: press releases have different purposes, targets, and need not relate to official business
 - No: press releases are just reactive to floor activity, will follow floor statements

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?
 - Yes: press releases have different purposes, targets, and need not relate to official business
 - No: press releases are just reactive to floor activity, will follow floor statements
- Deeper question: what does it mean for two text collections to be **different?**

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?
 - Yes: press releases have different purposes, targets, and need not relate to official business
 - No: press releases are just reactive to floor activity, will follow floor statements
- Deeper question: what does it mean for two text collections to be **different?**
- One Answer: **texts used for different purposes**

Congressional Language Across Sources

Congressional Press Releases and Floor Speeches

- Collected 64,033 press releases
- Problem: are they **distinct** from floor statements (approx. 52,000 during same time)?
 - Yes: press releases have different purposes, targets, and need not relate to official business
 - No: press releases are just reactive to floor activity, will follow floor statements
- Deeper question: what does it mean for two text collections to be **different?**
- One Answer: **texts used for different purposes**
- Partial answer: identify words that distinguish press releases and floor speeches

A Method for Identifying Distinguishing Words

A Method for Identifying Distinguishing Words

Mutual Information

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty
 - Unrelated: Conditional uncertainty = uncertainty

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty
 - Unrelated: Conditional uncertainty = uncertainty
 - Perfect predictor: Conditional uncertainty = 0

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty
 - Unrelated: Conditional uncertainty = uncertainty
 - Perfect predictor: Conditional uncertainty = 0
- Mutual information(X_j): uncertainty - conditional uncertainty (X_j)

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty
 - Unrelated: Conditional uncertainty = uncertainty
 - Perfect predictor: Conditional uncertainty = 0
- Mutual information(X_j): uncertainty - conditional uncertainty (X_j)
 - Maximum: Uncertainty $\rightarrow X_j$ is perfect predictor

A Method for Identifying Distinguishing Words

Mutual Information

- Unconditional uncertainty (entropy):
 - Randomly sample a press release
 - Guess press release/floor statement
 - Uncertainty about guess
 - Maximum: No. press releases = No. floor statements
 - Minimum : All documents in one category
- Conditional uncertainty (X_j) (conditional entropy)
 - Condition on presence of word X_j
 - Randomly sample a press release
 - Guess press release/floor statement
 - Word presence reduces uncertainty
 - Unrelated: Conditional uncertainty = uncertainty
 - Perfect predictor: Conditional uncertainty = 0
- Mutual information(X_j): uncertainty - conditional uncertainty (X_j)
 - Maximum: Uncertainty $\rightarrow X_j$ is perfect predictor
 - Minimum: 0 $\rightarrow X_j$ fails to separate speeches and floor statements

A Method for Identifying Distinguishing Words

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech
- Define **entropy** $H(\text{Doc})$

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech
- Define **entropy** $H(\text{Doc})$

$$H(\text{Doc}) = - \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t) \log_2 \Pr(t)$$

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech
- Define **entropy** $H(\text{Doc})$

$$H(\text{Doc}) = - \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t) \log_2 \Pr(t)$$

- $\log_2?$ Encodes bits

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech
- Define **entropy** $H(\text{Doc})$

$$H(\text{Doc}) = - \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t) \log_2 \Pr(t)$$

- \log_2 ? Encodes bits
- Maximum: $\Pr(\text{Press}) = \Pr(\text{Speech}) = 0.5$

A Method for Identifying Distinguishing Words

- $\Pr(\text{Press}) \equiv$ Probability selected document press release
- $\Pr(\text{Speech}) \equiv$ Probability selected document speech
- Define **entropy** $H(\text{Doc})$

$$H(\text{Doc}) = - \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t) \log_2 \Pr(t)$$

- \log_2 ? Encodes bits
- Maximum: $\Pr(\text{Press}) = \Pr(\text{Speech}) = 0.5$
- Minimum: $\Pr(\text{Press}) \rightarrow 0$ (or $\Pr(\text{Press}) \rightarrow 1$)

A Method for Identifying Distinguishing Words

A Method for Identifying Distinguishing Words

- Consider presence/absence of word X_j

A Method for Identifying Distinguishing Words

- Consider presence/absence of word X_j
- Define **conditional entropy** $H(\text{Doc}|X_j)$

A Method for Identifying Distinguishing Words

- Consider presence/absence of word X_j
- Define **conditional entropy** $H(\text{Doc}|X_j)$

$$H(\text{Doc}|X_j) = - \sum_{s=0}^1 \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t, X_j = s) \log_2 \Pr(t|X_j = s)$$

A Method for Identifying Distinguishing Words

- Consider presence/absence of word X_j
- Define **conditional entropy** $H(\text{Doc}|X_j)$

$$H(\text{Doc}|X_j) = - \sum_{s=0}^1 \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t, X_j = s) \log_2 \Pr(t|X_j = s)$$

- Maximum: X_j unrelated to Press Releases/Floor Speeches

A Method for Identifying Distinguishing Words

- Consider presence/absence of word X_j
- Define **conditional entropy** $H(\text{Doc}|X_j)$

$$H(\text{Doc}|X_j) = - \sum_{s=0}^1 \sum_{t \in \{\text{Pre}, \text{Spe}\}} \Pr(t, X_j = s) \log_2 \Pr(t|X_j = s)$$

- Maximum: X_j unrelated to Press Releases/Floor Speeches
- Minimum: X_j is a perfect predictor of press release/floor speech

A Method for Identifying Distinguishing Words

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

$$\text{Mutual Information}(X_j) = H(\text{Doc}) - H(\text{Doc}|X_j)$$

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

$$\text{Mutual Information}(X_j) = H(\text{Doc}) - H(\text{Doc}|X_j)$$

- Maximum: entropy $\Rightarrow H(\text{Doc}|X_j) = 0$

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

$$\text{Mutual Information}(X_j) = H(\text{Doc}) - H(\text{Doc}|X_j)$$

- Maximum: entropy $\Rightarrow H(\text{Doc}|X_j) = 0$
- Minimum: $0 \Rightarrow H(\text{Doc}|X_j) = H(\text{Doc})$.

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

$$\text{Mutual Information}(X_j) = H(\text{Doc}) - H(\text{Doc}|X_j)$$

- Maximum: entropy $\Rightarrow H(\text{Doc}|X_j) = 0$
- Minimum: 0 $\Rightarrow H(\text{Doc}|X_j) = H(\text{Doc})$.

Bigger mutual information \Rightarrow better discrimination

A Method for Identifying Distinguishing Words

- Define **Mutual Information**(X_j) as

$$\text{Mutual Information}(X_j) = H(\text{Doc}) - H(\text{Doc}|X_j)$$

- Maximum: entropy $\Rightarrow H(\text{Doc}|X_j) = 0$
- Minimum: $0 \Rightarrow H(\text{Doc}|X_j) = H(\text{Doc})$.

Bigger mutual information \Rightarrow better discrimination

Objective function and optimization \rightsquigarrow estimate probabilities that we then place in mutual information

A Method for Identifying Distinguishing Words

Formula for mutual information

(based on ML estimates of probabilities)

n_p = Number Press Releases

n_s = Number of Speeches

D = $n_p + n_s$

n_j = $\sum_{i=1}^D X_{i,j}$ (No. docs X_j appears)

n_{-j} = No. docs X_j does not appear

$n_{j,p}$ = No. press and X_j

$n_{j,s}$ = No. speech and X_j

$n_{-j,p}$ = No. press and not X_j

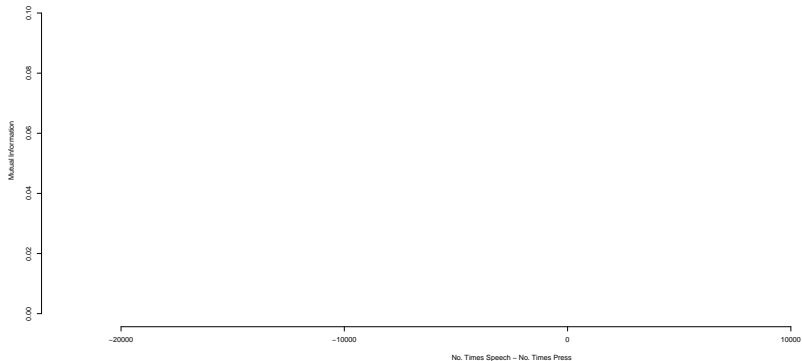
$n_{-j,s}$ = No. speech and not X_j

A Method for Identifying Distinguishing Words

Formula for Mutual Information

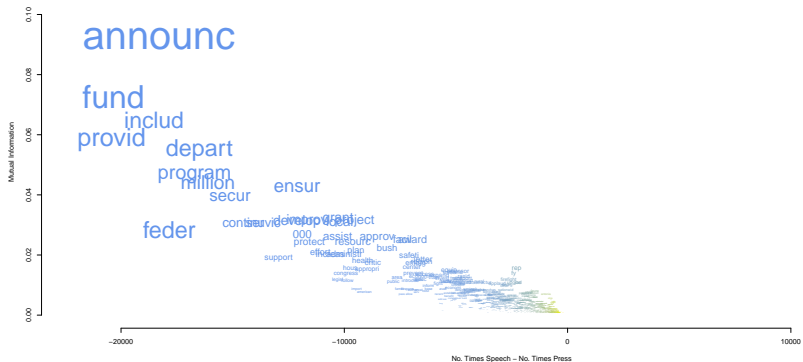
$$\begin{aligned} \text{MI}(X_j) = & \frac{n_{j,p}}{D} \log_2 \frac{n_{j,p}D}{n_j n_p} + \frac{n_{j,s}}{D} \log_2 \frac{n_{j,s}D}{n_j n_s} \\ & + \frac{n_{-j,p}}{D} \log_2 \frac{n_{-j,p}D}{n_{-j} n_p} + \frac{n_{-j,s}}{D} \log_2 \frac{n_{-j,s}D}{n_{-j} n_s}. \end{aligned}$$

What's Different About Press Releases



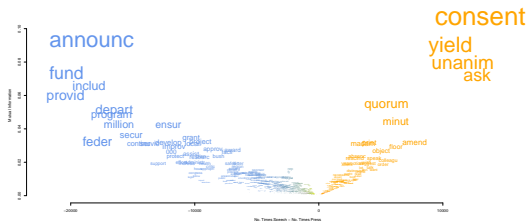
What's Different?

What's Different About Press Releases



What's Different?

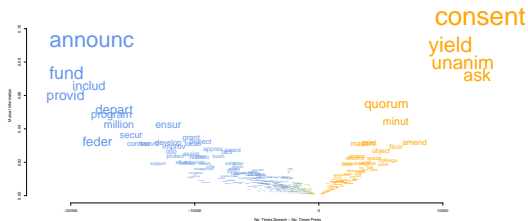
What's Different About Press Releases



What's Different?

- Press Releases: Credit Claiming
- Floor Speeches: Procedural Words
- Validate: Manual Classification

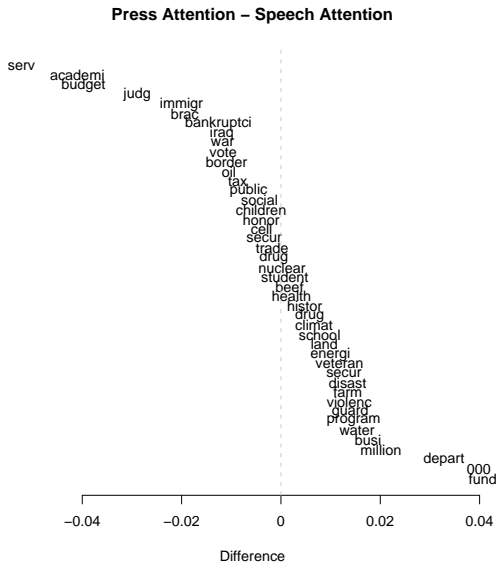
What's Different About Press Releases



What's Different?

- Press Releases: Credit Claiming
- Floor Speeches: Procedural Words
- Validate: Manual Classification
- Sample 500 Press Releases, 500 Floor Speeches

What's Different About Press Releases



Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

$$\text{Odds}(E) = \frac{P(E)}{1 - P(E)}$$

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

$$\text{Odds}(E) = \frac{P(E)}{1 - P(E)}$$

$$\text{Odds Ratio}(E, F) = \frac{\frac{P(E)}{1 - P(E)}}{\frac{P(F)}{1 - P(F)}}$$

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

$$\text{Odds}(E) = \frac{P(E)}{1 - P(E)}$$

$$\text{Odds Ratio}(E, F) = \frac{\frac{P(E)}{(1 - P(E))}}{\frac{P(F)}{1 - P(F)}}$$

$$\text{Log Odds Ratio}(E, F) = \log\left(\frac{P(E)}{1 - P(E)}\right) - \log\left(\frac{P(F)}{1 - P(F)}\right)$$

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

$$\text{Odds}(E) = \frac{P(E)}{1 - P(E)}$$

$$\text{Odds Ratio}(E, F) = \frac{\frac{P(E)}{(1 - P(E))}}{\frac{P(F)}{1 - P(F)}}$$

$$\text{Log Odds Ratio}(E, F) = \log\left(\frac{P(E)}{1 - P(E)}\right) - \log\left(\frac{P(F)}{1 - P(F)}\right)$$

Strategy \rightsquigarrow Construct objective function on **proportions** (and then calculate log-odds)

Fightin' Words \rightsquigarrow An Introduction to Regularization

Monroe, Colaresi, and Quinn (2009) \rightsquigarrow what makes a word partisan?

Argue for using **Log Odds Ratio**, weighted by variance

Recall: For some event E and F

$$P(E) = 1 - P(E^c)$$

$$\text{Odds}(E) = \frac{P(E)}{1 - P(E)}$$

$$\text{Odds Ratio}(E, F) = \frac{\frac{P(E)}{(1 - P(E))}}{\frac{P(F)}{1 - P(F)}}$$

$$\text{Log Odds Ratio}(E, F) = \log\left(\frac{P(E)}{1 - P(E)}\right) - \log\left(\frac{P(F)}{1 - P(F)}\right)$$

Strategy \rightsquigarrow Construct objective function on **proportions** (and then calculate log-odds)

Objective Function

Suppose we're interested in how a word separates partisan speech.

$\mathbf{Y} = (\text{Republican}, \text{Republican}, \text{Democrat}, \dots, \text{Republican})$

$\mathbf{X} =$ Unnormalized matrix of word counts $N \times J$

Define

$$\mathbf{x}_{\text{Republican}} = \left(\sum_{i=1}^N I(Y_i = \text{Republican})X_{i1}, \sum_{i=1}^N I(Y_i = \text{Republican})X_{i2}, \dots, \sum_{i=1}^N I(Y_i = \text{Republican})X_{iJ} \right)$$

with $N_{\text{Republican}} =$ Total number of Republican words

Objective Function

Objective Function

$$\pi_{\text{Republican}} \sim \text{Dirichlet}(\alpha)$$

Objective Function

$$\begin{aligned}\boldsymbol{\pi}_{\text{Republican}} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})\end{aligned}$$

Objective Function

$$\boldsymbol{\pi}_{\text{Republican}} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} \sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})$$

This implies an objective function on $\boldsymbol{\pi}$,

Objective Function

$$\boldsymbol{\pi}_{\text{Republican}} \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} \sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})$$

This implies an objective function on $\boldsymbol{\pi}$,

$$p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y}) \propto p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}, \mathbf{Y})$$

Objective Function

$$\begin{aligned}\boldsymbol{\pi}_{\text{Republican}} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})\end{aligned}$$

This implies an objective function on $\boldsymbol{\pi}$,

$$\begin{aligned}p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y}) &\propto p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}, \boldsymbol{\alpha}, \mathbf{Y}) \\ &\propto \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \pi_j^{x_{\text{Republican},j}}\end{aligned}$$

Objective Function

$$\begin{aligned}\boldsymbol{\pi}_{\text{Republican}} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})\end{aligned}$$

This implies an objective function on $\boldsymbol{\pi}$,

$$\begin{aligned}p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y}) &\propto p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}, \mathbf{X}, \mathbf{Y}) \\ &\propto \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \pi_j^{x_{\text{Republican},j}}\end{aligned}$$

$p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y})$ is a Dirichlet distribution:

Objective Function

$$\begin{aligned}\boldsymbol{\pi}_{\text{Republican}} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})\end{aligned}$$

This implies an objective function on $\boldsymbol{\pi}$,

$$\begin{aligned}p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y}) &\propto p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}, \mathbf{X}, \mathbf{Y}) \\ &\propto \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \pi_j^{x_{\text{Republican},j}}\end{aligned}$$

$p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y})$ is a Dirichlet distribution:

$$\pi_{\text{Republican},j}^* = \frac{x_{\text{Republican},j} + \alpha_j}{N_{\text{Republican}} + \sum_{j=1}^J \alpha_j}$$

Objective Function

$$\begin{aligned}\boldsymbol{\pi}_{\text{Republican}} &\sim \text{Dirichlet}(\boldsymbol{\alpha}) \\ \mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})\end{aligned}$$

This implies an objective function on $\boldsymbol{\pi}$,

$$\begin{aligned}p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y}) &\propto p(\boldsymbol{\pi} | \boldsymbol{\alpha}) p(\mathbf{x}_{\text{Republican}} | \boldsymbol{\pi}, \mathbf{X}, \mathbf{Y}) \\ &\propto \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \pi_j^{x_{\text{Republican},j}}\end{aligned}$$

$p(\boldsymbol{\pi} | \boldsymbol{\alpha}, \mathbf{X}, \mathbf{Y})$ is a Dirichlet distribution:

$$\pi_{\text{Republican},j}^* = \frac{x_{\text{Republican},j} + \alpha_j}{N_{\text{Republican}} + \sum_{j=1}^J \alpha_j}$$

Calculating Log Odds Ratio

Define log Odds Ratio_j as

$$\text{log Odds Ratio}_j = \log \left(\frac{\pi_{\text{Republican},j}}{1 - \pi_{\text{Republican},j}} \right) - \log \left(\frac{\pi_{\text{Democratic},j}}{1 - \pi_{\text{Democratic},j}} \right)$$

$$\text{Var}(\text{log Odds Ratio}_j) \approx \frac{1}{x_{jD} + \alpha_j} + \frac{1}{x_{jR} + \alpha_j}$$

$$\text{Std. Log Odds}_j = \frac{\text{log Odds Ratio}_j}{\sqrt{\text{Var}(\text{log Odds Ratio}_j)}}$$

Applying the Model

How do Republicans and Democrats differ in debate?

Condition on **topic** and examine word usage

- Press Releases (64,033)
- Topic Coded (Structural Topic Model)
- Given press release is about topic, what are the features that distinguish Republican and Democratic language?

Multinomial Inverse Regression

- In classification we're generally interested in:

$$E[Y|\mathbf{X}] = g(X_1, X_2, \dots, X_J)$$

- Problem: J might be very, very big.
- Potential solution \rightsquigarrow invert regression

$$E[\mathbf{X}|Y] = g(Y)$$

- Inversion is particularly useful for **feature selection**

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\mathbf{x}_{\text{Republican}} \sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}})$$

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]}\end{aligned}$$

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j)\end{aligned}$$

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j) \\ \lambda_j &\sim \text{Gamma}(s, r)\end{aligned}$$

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j) \\ \lambda_j &\sim \text{Gamma}(s, r)\end{aligned}$$

Laplace priors \rightsquigarrow **regularize** or **shrink** estimates toward zero

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j) \\ \lambda_j &\sim \text{Gamma}(s, r)\end{aligned}$$

Laplace priors \rightsquigarrow regularize or shrink estimates toward zero

Laplace priors \rightsquigarrow Equivalent to $L1$ or lasso penalization

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j) \\ \lambda_j &\sim \text{Gamma}(s, r)\end{aligned}$$

Laplace priors \rightsquigarrow regularize or shrink estimates toward zero

Laplace priors \rightsquigarrow Equivalent to L_1 or lasso penalization

Gamma-Lasso prior

Multinomial Inverse Regression: Objective Function (Taddy 2014)

As before, $\mathbf{x}_{\text{Republican}}$ to be the Republican count vector.

$$\begin{aligned}\mathbf{x}_{\text{Republican}} &\sim \text{Multinomial}(N_{\text{Republican}}, \boldsymbol{\pi}_{\text{Republican}}) \\ \pi_{\text{Republican},j} &= \frac{\exp[\alpha_j + I(\text{Republican})\phi_j]}{\sum_{l=1}^J \exp[\alpha_l + I(\text{Republican})\phi_l]} \\ \phi_j &\sim \text{Laplace}(\lambda_j) \\ \lambda_j &\sim \text{Gamma}(s, r)\end{aligned}$$

Laplace priors \rightsquigarrow regularize or shrink estimates toward zero

Laplace priors \rightsquigarrow Equivalent to $L1$ or lasso penalization

Gamma-Lasso prior

Optimization \rightsquigarrow Coordinate descent (paper is great!) \rightsquigarrow textir package

Applying Multinomial Inverse Regression: Objective Function

Taddy (2014) considers speeches made on Congressional floor in 2005
“Most” Republican words

Applying Multinomial Inverse Regression: Objective Function

Taddy (2014) considers speeches made on Congressional floor in 2005
“Most” Republican words
un.official,term.care.insurance, weapons.grade.plutonium
million.illegal.immigrant, **grand ole opry**, ...,
personal.injury.lawyer

Applying Multinomial Inverse Regression: Objective Function

Taddy (2014) considers speeches made on Congressional floor in 2005

“Most” Republican words
un.official,term.care.insurance, weapons.grade.plutonium
million.illegal.immigrant, **grand ole opry**, ...,
personal.injury.lawyer

“Most” Democratic Words

Applying Multinomial Inverse Regression: Objective Function

Taddy (2014) considers speeches made on Congressional floor in 2005

“Most” Republican words

un.official,term.care.insurance, weapons.grade.plutonium
million.illegal.immigrant, grand ole opry, ...,
personal.injury.lawyer

“Most” Democratic Words

wild.bird, death.penalty.system, record.budget.deficit
security.private.account, able.buy.gun