Computational Text Analysis for Legal Practice (Class 2)

Charles Crabtree & Kevin L. Cope

January 19, 2022

github.com/cdcrabtree/uva-2022

1 Data management with texts.

2 Understanding your text data

- 1 Data management with texts.
- 2 Understanding your text data.
 - Word clouds.

- 1 Data management with texts.
- 2 Understanding your text data.
 - Word clouds.
 - Discriminating terms

- 1 Data management with texts.
- 2 Understanding your text data.
 - Word clouds.
 - Discriminating terms.
 - Word frequencies and associations

- 1 Data management with texts.
- 2 Understanding your text data.
 - Word clouds.
 - Discriminating terms.
 - Word frequencies and associations.
 - DTMs

- 1 Data management with texts.
- 2 Understanding your text data.
 - Word clouds.
 - Discriminating terms.
 - Word frequencies and associations.
 - DTMs.

- 1 Data management with texts.
- 2 Understanding your text data.
 - 3 Word clouds.
 - 4 Discriminating terms.
 - 5 Word frequencies and associations.
 - 6 DTMs.

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches --> Partisan language.

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts --> Different Style/Substantive Focus.

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.
- Reasons:

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.

• Reasons:

1 Interesting in their own right

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.

- 1 Interesting in their own right.
- 2 Create custom dictionaries for a classification task

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.

- 1 Interesting in their own right.
- 2 Create custom dictionaries for a classification task.
- 3 Feature selection: inclusion of features in some subsequent analysis.

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination).
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.

- 1 Interesting in their own right.
- 2 Create custom dictionaries for a classification task.
- 3 Feature selection: inclusion of features in some subsequent analysis.
- Method: Distinctive / Discriminating / Separating word

- Goal: Find words that distinguish one group of texts from another group of texts.
 - Comparing Documents That Are Responsive/Non-Responsive to a Discovery Request → Contain Phrases About a Given Event (e.g., a Termination)
 - Comparing Rep., Dem. speeches → Partisan language.
 - Comparing Judicial Opinions Between Judges/Courts → Different Style/Substantive Focus.

- 1 Interesting in their own right.
- 2 Create custom dictionaries for a classification task.
- 3 Feature selection: inclusion of features in some subsequent analysis.
- Method: Distinctive / Discriminating / Separating word

What does 'distinctive' means?

- Goal: find words (or features) distinctive to each corpus.
- Requires a decision about what 'distinctive' means.

What does 'distinctive' means?

- Goal: find words (or features) distinctive to each corpus.
- Requires a decision about what 'distinctive' means.
- There are a variety of definitions that we might use.

What does 'distinctive' means?

- Goal: find words (or features) distinctive to each corpus.
- Requires a decision about what 'distinctive' means.
- There are a variety of definitions that we might use.

- Distinctive = exclusive.
- If Trump uses the word 'bigly' and Obama never does, we should count 'bigly' as distinctive.

- Distinctive = exclusive.
- If Trump uses the word 'bigly' and Obama never does, we should count 'bigly' as distinctive.
- These words tend not to be terribly interesting or informative.

- Distinctive = exclusive.
- If Trump uses the word 'bigly' and Obama never does, we should count 'bigly' as distinctive.
- These words tend not to be terribly interesting or informative.
- More likely to capture differences in linguistic style than content.

- Distinctive = exclusive.
- If Trump uses the word 'bigly' and Obama never does, we should count 'bigly' as distinctive.
- These words tend not to be terribly interesting or informative.
- More likely to capture differences in linguistic style than content.

- Distinctive = difference in frequency.
- Compare the number of times each author, Trump or Obama, uses a word.

- Distinctive = difference in frequency.
- Compare the number of times each author, Trump or Obama, uses a word.
- Find the largest absolute difference

- Distinctive = difference in frequency.
- Compare the number of times each author, Trump or Obama, uses a word.
- Find the largest absolute difference.
- Doesn't take into account difference in total words.

- Distinctive = difference in frequency.
- Compare the number of times each author, Trump or Obama, uses a word.
- Find the largest absolute difference.
- Doesn't take into account difference in total words.

Option 3: Difference in averages

- Distinctive = difference in rates.
- Compare the average rate each author uses a word.

Option 3: Difference in averages

- Distinctive = difference in rates.
- Compare the average rate each author uses a word.

Option 3: Difference in averages

- Normalize DTM from counts to proportions.
- For each word *p* in an arbitrary corpus *c*:

$$\mu_{\mathcal{D}} = rac{\sum_{i=1}^{N} \mathcal{D}_i}{T}$$

where p_i is the number of times a p appears in document i, N is the total number of documents in c and T is the total number of words in c.

 Take the difference between one author's proportion of a word and another's proportion of the same word.

$$heta_{
ho} = \mu_{
ho, Trum
ho} - \mu_{
ho, Obama}$$

• Find words with highest absolute difference.

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.
- Word 2: 5/1000 (Trump); 1/1000 (Obama) → Score 4.9/1000.

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.
- Word 2: 5/1000 (Trump); 1/1000 (Obama) → Score: 4.9/1000.
- Ignores cases when one author uses a word frequently and a another author barely uses it.

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.
- Word 2: 5/1000 (Trump); 1/1000 (Obama) → Score: 4.9/1000.
- Ignores cases when one author uses a word frequently and a another author barely uses it.
- More generally: Differences in rates of frequent words > Differences in rates of rare words.

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.
- Word 2: 5/1000 (Trump); 1/1000 (Obama) → Score: 4.9/1000.
- Ignores cases when one author uses a word frequently and a another author barely uses it.
- More generally: Differences in rates of frequent words > Differences in rates of rare words.
- Adjustment: Divide the difference in authors' average rates by the average rate across all authors.

Differences in averages: Problems

- Favors more frequent words.
- Word 1: 30/1000 (Trump); 25/1000 (Obama) → Score: 5/1000.
- Word 2: 5/1000 (Trump); 1/1000 (Obama) → Score: 4.9/1000.
- Ignores cases when one author uses a word frequently and a another author barely uses it.
- More generally: Differences in rates of frequent words > Differences in rates of rare words.
- Adjustment: Divide the difference in authors' average rates by the average rate across all authors.

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?
 - Depends on context, goal

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?
 - Depends on context, goal.
 - Classification → accuracy, precision, recall.

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?
 - Depends on context, goal.
 - Classification → accuracy, precision, recall.
 - Qualitative inference → face validity, convergence, etc

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?
 - Depends on context, goal.
 - Classification → accuracy, precision, recall.
 - Qualitative inference → face validity, convergence, etc.
 - More on this later (at the end of slides).

- Other metrics for 'distinctiveness':
 - Standardized mean difference (account for variability).
 - Standard Log Odds (used in Monroe, Colaresi, and Quinn, 2009).
 - Many more!
- How do we choose?
 - Depends on context, goal.
 - Classification
 Accuracy, precision, recall.

 - More on this later (at the end of slides).

Why do we care?

- Qualitative inference comparing 2 groups.
- Create custom dictionaries for classification task.

Why do we care?

- Qualitative inference comparing 2 groups.
- Create custom dictionaries for classification task.

• Federalist Papers:

1 Canonical texts in study of American politics

- Federalist Papers:
 - 1 Canonical texts in study of American politics.
 - 2 Designed to persuade citizens of New York to adopt constitution

- Federalist Papers:
 - 1 Canonical texts in study of American politics.
 - 2 Designed to persuade citizens of New York to adopt constitution.
 - 3 77 essays, published from 1787-1799 in newspapers, published anonymously under the name Publius.
- Who wrote the Federalist papers? (Hostler and Wallace, 1963)

- Federalist Papers:
 - 1 Canonical texts in study of American politics.
 - 2 Designed to persuade citizens of New York to adopt constitution.
 - 3 77 essays, published from 1787-1799 in newspapers, published anonymously under the name Publius.
- Who wrote the Federalist papers? (Hostler and Wallace, 1963)
 - Jay: wrote 5 essays
 - Hamilton: wrote 43 papers.
 - Madison: wrote 12 papers.
 - Disputed (Hamilton or Madison?): Essays 49-58, 62, and 63.

- Federalist Papers:
 - 1 Canonical texts in study of American politics.
 - 2 Designed to persuade citizens of New York to adopt constitution.
 - 3 77 essays, published from 1787-1799 in newspapers, published anonymously under the name Publius.
- Who wrote the Federalist papers? (Hostler and Wallace, 1963)
 - Jay: wrote 5 essays.
 - Hamilton: wrote 43 papers.
 - Madison: wrote 12 papers.
 - Disputed (Hamilton or Madison?): Essays 49-58, 62, and 63.

- Tasks: Identify authors of disputed papers.
- Method: Classify papers as Hamilton or Madison using dictionary methods.
- Training data → Hamilton, Madison are known to have authored.
- Test data → disputed (i.e. unlabeled) papers.
- Preprocessing:
 - Hamilton/Madison discuss similar themes.
 - Differ on the extent they use stop words.
 - Focus analysis on the stop words.

Word Weights: Standardized Mean Difference

- For each word p, construct weight θ_p , $\mu_{p,Hamilton} = \text{Rate}(p) \text{ in subcorpus of Hamilton docs}$ $\mu_{p,Madison} = \text{Rate}(p) \text{ in subcorpus of Madison docs}$ $\sigma_{p,Hamilton}^2 = \text{Var}(p) \text{ in subcorpus of Hamilton docs}$ $\sigma_{p,Madison}^2 = \text{Var}(p) \text{ in subcorpus of Madison docs}$
- We can then generate weight θ_D as

$$heta_{
m p}=rac{\mu_{
m p,Hamilton}-\mu_{
m p,Madison}}{\sigma_{
m p,Hamilton}^2+\sigma_{
m p,Madison}^2}$$

Trimming the dictionary

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

Trimming the dictionary

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

Classification for determining authorship

 For each disputed document i, compute discrimination statistic.

$$Y_i = \sum_{p=1}^{p} \theta_p X_{ip}$$

- $Y_i \rightsquigarrow$ classification (linear discriminator)
 - Above midpoint in training set \rightsquigarrow Hamilton text.
 - Below midpoint in training set → Madison text.
- Findings: Madison is the author of the disputed federalist papers.

- How do we choose between different 'distinctive word' metrics?
- How do we chose between dictionaries?

- How do we choose between different 'distinctive word' metrics?
- How do we chose between dictionaries?
- How do we evaluate our findings?

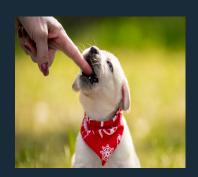
- How do we choose between different 'distinctive word' metrics?
- How do we chose between dictionaries?
- How do we evaluate our findings?
- Three evaluation strategies:

- How do we choose between different 'distinctive word' metrics?
- How do we chose between dictionaries?
- How do we evaluate our findings?
- Three evaluation strategies:
 - Face validity (do these results make sense?)
 - Convergence (do different metrics lead to the same result?)
 - 'Gold Standard' (do our results align with human coding?)

- How do we choose between different 'distinctive word' metrics?
- How do we chose between dictionaries?
- How do we evaluate our findings?
- Three evaluation strategies:
 - Face validity (do these results make sense?)
 - Convergence (do different metrics lead to the same result?)
 - 'Gold Standard' (do our results align with human coding?)

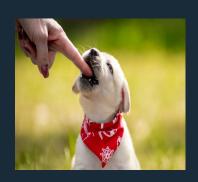
Application

- Suppose you received a request for document production in a dogbite case:
 - "Any and all documents in your possession regarding the care and keeping of the subject dog."



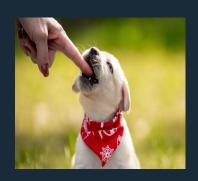
Application

- Suppose you received a request for document production in a dogbite case:
 - "Any and all documents in your possession regarding the care and keeping of the subject dog."
- "What terms would you imagine."



Application

- Suppose you received a request for document production in a dogbite case:
 - "Any and all documents in your possession regarding the care and keeping of the subject dog."
- "What terms would you imagine."



Let's Program

Launch RStudio!

kcope@law.virginia.edu|crabtree@dartmouth.edu