# Text Mining

# Tweets

- Storm pushes Presidential Race from spotlight
- Romney holds slim 1-point lead among likely voters in CO race
- Did Pennsylvania run misleading voter ID ad?
- Mitt Romney is just not that into Federal Disaster Relief

# Modified Tweets

- Storm pushes Presidential Race from spotlight
- Romney leads  likely voters in race
- Run misleading voter ad?
- Obama, Romney run race close president race

# Reduction

- A: storm push president race spotlight
- B: romney lead like vote race
- C: run mislead vote ad
- D: obama romey run race close president race

# Bag of words

- Unique words across all documents
- storm, push, president, race, spotlight, romney, lead, like, vote, run, mislead, ad, obama, close

| | A | B | C | D | NumDocs |
|---|---|---|---|---|---|
| storm | 1 | 0 | 0 | 0 | 1 |
| push | 1 | 0 | 0 | 0 | 1 |
| president | 1 | 0 | 0 | 1 | 2 |
| race | 1 | 1 | 0 | 2 | 3 |
| spotlight | 1 | 0 | 0 | 0 | 1 |
| romney | 0 | 1 | 0 | 1 | 2 |
| lead | 0 | 1 | 0 | 0 | 1 |
| like | 0 | 1 | 0 | 0 | 1 |
| vote | 0 | 1 | 1 | 0 | 2 |
| run | 0 | 0 | 1 | 1 | 2 |
| mislead | 0 | 0 | 1 | 0 | 1 |
| ad | 0 | 0 | 1 | 0 | 1 |
| obama | 0 | 0 | 0 | 1 | 1 |
| close | 0 | 0 | 0 | 1 | 1 |
| TOTAL | 5 | 5 | 4 | 7 | |

# Similarity between documents

- Do documents use the same terms?
- Don't care about common terms
- Want to control for the length of the document

# Similarity between documents

- **Term frequency**: fraction of the words in a document are this term

- **Document frequency**: fraction of the documents contain this term

- Normalized vector:

 V = term freq * inverse document freq

= TF/DF

| | A | B | C | D | NumDocs |
|---|---|---|---|---|---|
| storm | 0.20 | 0 | 0 | 0 | 4 |
| push | 0.20 | 0 | 0 | 0 | 4 |
| president | 0.20 | 0 | 0 | 0.143 | 2 |
| race | 0.20 | 0.20 | 0 | 0.287 | 1.33 |
| spotlight | 0.20 | 0 | 0 | 0 | 4 |
| romney | 0 | 0.20 | 0 | 0.143 | 2 |
| lead | 0 | 0.20 | 0 | 0 | 4 |
| like | 0 | 0.20 | 0 | 0 | 4 |
| vote | 0 | 0.20 | 0.25 | 0 | 2 |
| run | 0 | 0 | 0.25 | 0.143 | 2 |
| mislead | 0 | 0 | 0.25 | 0 | 4 |
| ad | 0 | 0 | 0.25 | 0 | 4 |
| obama | 0 | 0 | 0 | 0.143 | 4 |
| close | 0 | 0 | 0 | 0.143 | 4 |
| TOTAL | 1 | 1 | 1 | 1 | |

# Distance between documents

- Dist(V, W) = ½( KL(V, AVG) + KL(W, AVG))

- where: KL stands for Kulback-Leibler measure
    KL(V, AVG) = sum( log(V/AVG) * AVG)

- and V = TF* IDF

# Similarity Matrix

|   | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 1.80 | 2.10 | 1.44 |
| B | 1.80 | 0 | 1.65 | 1.30 |
| C | 2.10 | 1.65 | 0 | 1.61 |
| D | 1.44 | 1.30 | 1.61 | 0 |

# Multi-dimensional Scaling

- Information visualization technique for high-dimensional data.

- Consider the matrix of dis-similarities above for the four documents.

- Assign locations in 2 dimensions so that the distances between documents is roughly preserved.

# Example

|   | A | B | C |
|---|---|---|---|
| A | 0 | 3 | 4 |
| B | 3 | 0 | 5 |
| C | 4 | 5 | 0 |

Could represent
as a triangle in
two dimensions

# Our Documents

# MDS

- Doesn't produce unique representations of the data,

- Does give you the opportunity to compare objects (documents in our case)

- Look for clusters and gaps

# Hierarchical clustering

- Build a binary tree that successively merges similar groups.

- This implies that we need a metric or measure of similarity between groups of points.

- There are various algorithms that can be used to create the binary tree.

# Agglomerative Clustering

1. Start with each point in its own group.

2. Merge the two most similar groups.

3. Repeat step 2 until all groups have been merged into one

- Note that the similarity between two groups being merged at any stage must, by design, be decreasing because we merge less and less similar groups.

# Measure of similarity between groups

- Single linkage: smallest distance between any point in one group and a point in the other group.

- Complete linkage: largest distance between any point in one group and a point in the other group.

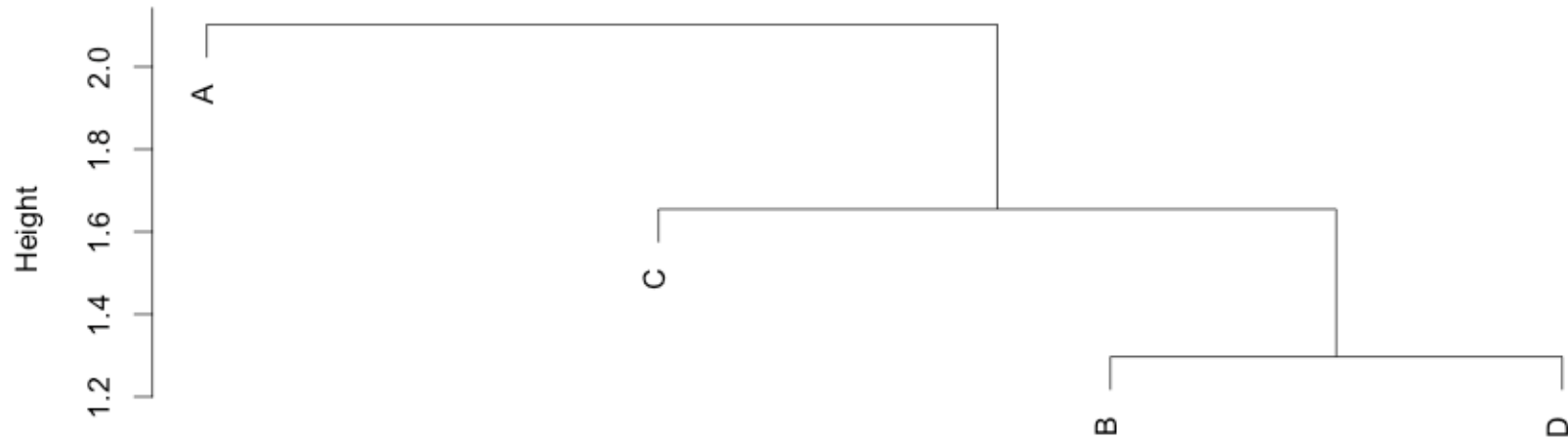- Average linkage: average distance between each point in one group and every point in the other group

- Single linkage tends to result in chaining, where you successively add on one point to a group

- Complete linkage tends not to merge close groups when one point in one group is far from the other group.

# Dendrogram

- Useful visualization of the clustering process.
- Typically the tree is drawn such that the heights of the branches proportional to the dissimilarity between the two groups.
- This visual helps you see where a good place to "cut" the tree might be and create clusters
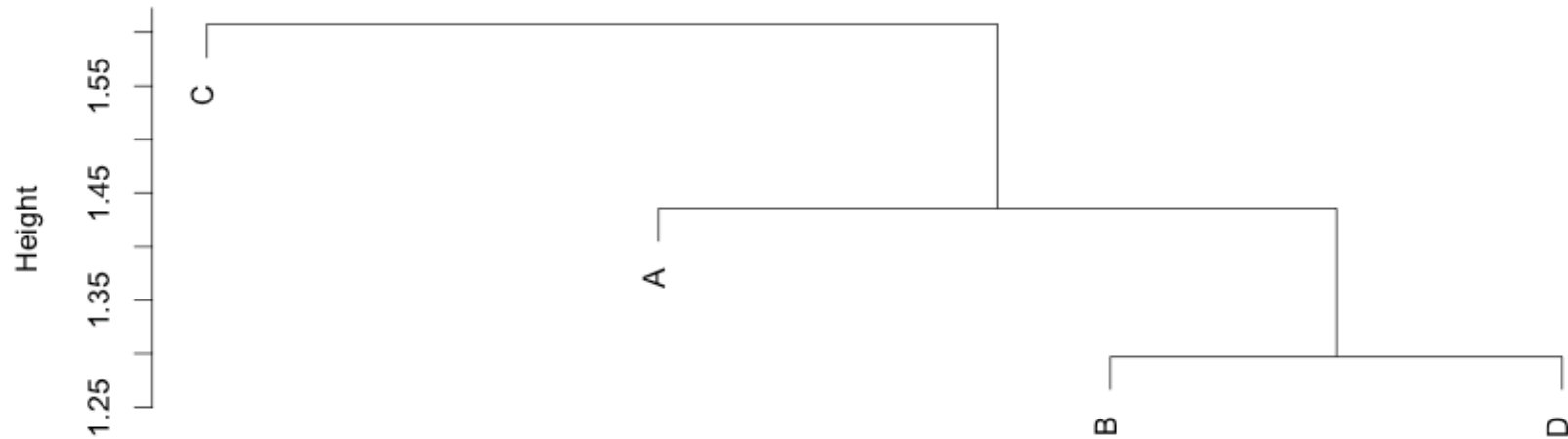
# Complete linkage



**Cluster Dendrogram**

documents
hclust (*, "complete")

# Single linkage



**Cluster Dendrogram**

documents
hclust (*, "single")

# Dendrogram

- Different definitions of similarity can give very different trees.

- The algorithm imposes a hierarchy on a set of data, even if there isn't one.

# Your Turn

State of the Union addresses

# State of the union speeches

- Use readLines() to read in the speeches
- Return value: character vector with one element/character string per line in the file
- Regular expressions to find ***
- Use *** to identify the date of the speech
- Use regular expressions to extract the year
- Use regular expressions to extract the month
- Use *** to extract the name of the president

# State of the union speeches

- Chop the speeches up into a list there is one element for each speech.
- Each element is a character vector.
-  Each element of the vector is a character string corresponding to a sentence in the speech

# Word Vectors

- Eliminate apostrophes, numbers, and the phrase: (Applause.) from the text.
- Make all the characters lower case.
- Split the sentences up where there are blanks and punctuation
- Drop any empty words that resulted from this split
- Load the library Rstem and use the function wordStem() to stem words

- Find the bag of words
- Create a word vector for each speech
- Normalize the word vectors to get term frequencies

# Analysis

- Exploratory analysis of the data:
  - Number of sentences, long words, political party

- Multidimensional scaling

- Hierarchical clustering