# Text Mining, pt. I

Philip D. Waggoner

MACS 40800: Unsupervised Machine Learning

November 19, 2019

#### Lecture Outline

- Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

#### Lecture Outline

- Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- 4 Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

 Today (only) we will depart from unsupervised learning and focus on supervised learning for text

- Today (only) we will depart from unsupervised learning and focus on supervised learning for text
- The goal is to offer a well-rounded toolbox for text mining, which necessarily includes both supervised and unsupervised techniques

- Today (only) we will depart from unsupervised learning and focus on supervised learning for text
- The goal is to offer a well-rounded toolbox for text mining, which necessarily includes both supervised and unsupervised techniques
- Today: supervised (basics, dictionaries, and sentiment scoring)

- Today (only) we will depart from unsupervised learning and focus on supervised learning for text
- The goal is to offer a well-rounded toolbox for text mining, which necessarily includes both supervised and unsupervised techniques
- Today: supervised (basics, dictionaries, and sentiment scoring)
- Thursday: unsupervised (topic models)

• For most of its history text analysis was qualitative

- For most of its history text analysis was qualitative
- And today?

- For most of its history text analysis was qualitative
- And today? We are still interested in appying subjective judgement to explain and judge that which text is revealing

- For most of its history text analysis was qualitative
- And today? We are still interested in appying subjective judgement to explain and judge that which text is revealing
- Thus, our goal is distant reading, rather than close reading

- For most of its history text analysis was qualitative
- And today? We are still interested in appying subjective judgement to explain and judge that which text is revealing
- Thus, our goal is distant reading, rather than close reading
- The "quantification" of text analysis is a monumental advance in this ancient field in that quantitative work is reliable, replicable, and can easily handle large volumes of material

Obtain some group of texts

- Obtain some group of texts
- 2 Create a document-term matrix

- Obtain some group of texts
- 2 Create a document-term matrix
- Operate/mine ((dis)similarities, sentiments, topics, complexity, classification, etc.)

- Obtain some group of texts
- 2 Create a document-term matrix
- Operate/mine ((dis)similarities, sentiments, topics, complexity, classification, etc.)
- Oraw inferences

• Some set of texts or documents we wish to analyze, which can be

- Some set of texts or documents we wish to analyze, which can be
  - structured: useful document information is known (e.g., beginning and end, authorship, etc.)

- Some set of texts or documents we wish to analyze, which can be
  - structured: useful document information is known (e.g., beginning and end, authorship, etc.)
  - unstructured: desired quantity is unknown (e.g., sentiment)

- Some set of texts or documents we wish to analyze, which can be
  - structured: useful document information is known (e.g., beginning and end, authorship, etc.)
  - unstructured: desired quantity is unknown (e.g., sentiment)
- Further, a corpus may be annotated where metadata data not part
  of the document itself is available, e.g., date, linguistic tagging, etc.

• The corpus is made up of the documents within it, but these may be a sample of the total population of documents available

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity,

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- $\bullet$  Often sample due to time, resources or (<code>legal</code>) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%
- Sometimes, authors claim to have the universe of cases in their corpus...

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%
- Sometimes, authors claim to have the universe of cases in their corpus...
- But, you still need to think about sampling error

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%
- Sometimes, authors claim to have the universe of cases in their corpus...

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%
- Sometimes, authors claim to have the universe of cases in their corpus...
- The point here?

- The corpus is made up of the documents within it, but these may be a sample of the total population of documents available
- Often sample due to time, resources or (legal) necessity, e.g. Twitter gives you  $\sim 1\%$  of all their tweets, but it would presumably be prohibitively expensive to store 100%
- Sometimes, authors claim to have the universe of cases in their corpus...
- The point here? Be cautious as you approach text mining and think carefully about error, where texts came from, and so on, as the corpus should be representative in some sense for inferences to be meaningful

• Language is extraordinarily complex, and involves great subtlety and nuanced interpretation

- Language is extraordinarily complex, and involves great subtlety and nuanced interpretation
- Yet remarkably, we can do very well by simplifying, and representing documents as mathematical objects, making the modeling process much more tractable (more at the end)

- Language is extraordinarily complex, and involves great subtlety and nuanced interpretation
- Yet remarkably, we can do very well by simplifying, and representing documents as mathematical objects, making the modeling process much more tractable (more at the end)
- 'Do very well' → complicated representations add nothing to the quality of inferences, ability to predict outcomes, and model fit

- Language is extraordinarily complex, and involves great subtlety and nuanced interpretation
- Yet remarkably, we can do very well by simplifying, and representing documents as mathematical objects, making the modeling process much more tractable (more at the end)
- 'Do very well' → complicated representations add nothing to the quality of inferences, ability to predict outcomes, and model fit
- Inevitably, the degree to which we simplify is dependent on the task at hand;

- Language is extraordinarily complex, and involves great subtlety and nuanced interpretation
- Yet remarkably, we can do very well by simplifying, and representing documents as mathematical objects, making the modeling process much more tractable (more at the end)
- 'Do very well' → complicated representations add nothing to the quality of inferences, ability to predict outcomes, and model fit
- Inevitably, the degree to which we simplify is dependent on the task at hand; there is no "single best way" to go from texts to numeric data

- Language is extraordinarily complex, and involves great subtlety and nuanced interpretation
- Yet remarkably, we can do very well by simplifying, and representing documents as mathematical objects, making the modeling process much more tractable (more at the end)
- 'Do very well' → complicated representations add nothing to the quality of inferences, ability to predict outcomes, and model fit
- Inevitably, the degree to which we simplify is dependent on the task at hand; there is no "single best way" to go from texts to numeric data
- Though different in means, the end is the same in text mining:
   reduce complexity for inferential clarity

# Quick Note on Terminology

- a **type** is a *unique* sequence of characters that are grouped together in some meaningful way (usually a word)
  - e.g. 'Australia', 'French Revolution', '1984'

# Quick Note on Terminology

- a type is a unique sequence of characters that are grouped together in some meaningful way (usually a word)
  - e.g. 'Australia', 'French Revolution', '1984'
- a token is an instance of type
  - ▶ e.g. "Dog eat dog world", contains three **types**, but four **tokens**

# Quick Note on Terminology

- a type is a unique sequence of characters that are grouped together in some meaningful way (usually a word)
  - e.g. 'Australia', 'French Revolution', '1984'
- a token is an instance of type
  - ▶ e.g. "Dog eat dog world", contains three **types**, but four **tokens**
- a term is a type that the technique recognizes as a type to be recorded
  - e.g. stemmed words like 'motivat' or 'applica'

 Collect raw text in machine-readable form; decide what constitutes a document

- Collect raw text in machine-readable form; decide what constitutes a document
- 2 Strip away extraneous material, e.g., ASCII characters, capitalization, punctuation, stop words, and so on

- Collect raw text in machine-readable form; decide what constitutes a document
- Strip away extraneous material, e.g., ASCII characters, capitalization, punctuation, stop words, and so on
- Out document up into useful elementary pieces: tokenization

- Collect raw text in machine-readable form; decide what constitutes a document
- Strip away extraneous material, e.g., ASCII characters, capitalization, punctuation, stop words, and so on
- Out document up into useful elementary pieces: tokenization
- Map tokens to common form: lemmatization, stemming

- Collect raw text in machine-readable form; decide what constitutes a document
- Strip away extraneous material, e.g., ASCII characters, capitalization, punctuation, stop words, and so on
- Out document up into useful elementary pieces: tokenization
- Map tokens to common form: lemmatization, stemming
- (Sometimes) Add descriptive annotations that preserve text context: tagging

- Collect raw text in machine-readable form; decide what constitutes a document
- 2 Strip away extraneous material, e.g., ASCII characters, capitalization, punctuation, stop words, and so on
- Out document up into useful elementary pieces: tokenization
- Map tokens to common form: lemmatization, stemming
- Sometimes Add descriptive annotations that preserve text context: tagging
- Operate/model/mine

• We can compress these steps a bit...

- We can compress these steps a bit...
  - Collect raw text in machine readable form; decide what constitutes a document

- We can compress these steps a bit...
  - Collect raw text in machine readable form; decide what constitutes a document
  - Preprocess Stage

- We can compress these steps a bit...
  - Collect raw text in machine readable form; decide what constitutes a document
  - Preprocess Stage
  - Operate/model/mine

• So you've collected some text and got it into electronic format

- So you've collected some text and got it into electronic format
- Now, we have the first preprocessing step

- So you've collected some text and got it into electronic format
- Now, we have the first preprocessing step
- Generally think control characters here → if they do not contribute to substantive importance then remove them

- So you've collected some text and got it into electronic format
- Now, we have the first preprocessing step
- Generally think control characters here → if they do not contribute to substantive importance then remove them
- Punctuation may also be unhelpful, e.g.,

- So you've collected some text and got it into electronic format
- Now, we have the first preprocessing step
- Generally think control characters here → if they do not contribute to substantive importance then remove them
- Punctuation may also be unhelpful, e.g.,
- are, wash, wash., wash, wash) really different words?

 stopwords serve as linguistic connectors and can be removed at this first preprocessing stage

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines
- So what are stopwords?

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines
- So what are stopwords? compiled lists in any NLP package (e.g., in NLTK, https://gist.github.com/sebleier/554280 or the stopwords library in R)

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines
- So what are stopwords? compiled lists in any NLP package (e.g., in NLTK, https://gist.github.com/sebleier/554280 or the stopwords library in R)
- But stopwords can also vary by application/domain-specific needs

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines
- So what are stopwords? compiled lists in any NLP package (e.g., in NLTK, https://gist.github.com/sebleier/554280 or the stopwords library in R)
- But stopwords can also vary by application/domain-specific needs
- e.g. with Congressional speech data, representative might be a stop word; in British Parliamentary data, honourable might be

- stopwords serve as linguistic connectors and can be removed at this first preprocessing stage
- Removing stopwords simplifies the document with little information loss, e.g., they are ignored by search engines
- So what are stopwords? compiled lists in any NLP package (e.g., in NLTK, https://gist.github.com/sebleier/554280 or the stopwords library in R)
- But stopwords can also vary by application/domain-specific needs
- e.g. with Congressional speech data, representative might be a stop word; in British Parliamentary data, honourable might be
- And inversely, there might be unique cases where stopwords should actually be retained, such as studying linguistic complexity and/or

 After basic cleaning, we want to pull out the meaningful subunits, which are the tokens

- After basic cleaning, we want to pull out the meaningful subunits, which are the tokens
- In other words, we want to turn human-readable text into machine-readable text

- After basic cleaning, we want to pull out the meaningful subunits, which are the tokens
- In other words, we want to turn human-readable text into machine-readable text
- While tokens are usually words, they might include numbers or punctuation too

- After basic cleaning, we want to pull out the meaningful subunits, which are the tokens
- In other words, we want to turn human-readable text into machine-readable text
- While tokens are usually words, they might include numbers or punctuation too
- A common rule for a tokenizer is to use whitespace as the marker, but some applicatins might require more subtlety

- After basic cleaning, we want to pull out the meaningful subunits, which are the tokens
- In other words, we want to turn human-readable text into machine-readable text
- While tokens are usually words, they might include numbers or punctuation too
- A common rule for a tokenizer is to use whitespace as the marker, but some applicatins might require more subtlety
- e.g., "Brown vs Board of Education" 
   → 'Brown', 'vs', 'Board', 'of', 'Education'

- After basic cleaning, we want to pull out the meaningful subunits, which are the tokens
- In other words, we want to turn human-readable text into machine-readable text
- While tokens are usually words, they might include numbers or punctuation too
- A common rule for a tokenizer is to use whitespace as the marker, but some applicatins might require more subtlety
- e.g., "Brown vs Board of Education" 
   → 'Brown', 'vs', 'Board', 'of', 'Education' ???

• Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are *many* tokenizers

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- **Word tokenizers**: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)
- N-gram tokenizers: some contiguous sequence of words (e.g., he makes no sense)

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)
- N-gram tokenizers: some contiguous sequence of words (e.g., he makes no sense)
- Sentence and paragraph tokenizers: single or multiple sentences, or even paragraphs

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, more in a bit)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)
- N-gram tokenizers: some contiguous sequence of words (e.g., he makes no sense)
- Sentence and paragraph tokenizers: single or multiple sentences, or even paragraphs
- "Text chunking": similar to paragraph tokenizers, but can be even longer; stipulation is chunks of equal size (usually for large texts)

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)
- N-gram tokenizers: some contiguous sequence of words (e.g., he makes no sense)
- Sentence and paragraph tokenizers: single or multiple sentences, or even paragraphs
- "Text chunking": similar to paragraph tokenizers, but can be even longer; stipulation is chunks of equal size (usually for large texts)
- And so on...

#### **Tokenizers**

- Thus, the naive whitespace approach implies many forms of "tokenizing" text: there are many tokenizers
- Word tokenizers: whitespace (or with "word-stem," stemmed words distinguished by whitespace, *more in a bit*)
- Character and shingle tokenizers: individual characters or small character-based subsets or "shingles" (e.g., qu or th)
- N-gram tokenizers: some contiguous sequence of words (e.g., he makes no sense)
- Sentence and paragraph tokenizers: single or multiple sentences, or even paragraphs
- "Text chunking": similar to paragraph tokenizers, but can be even longer; stipulation is chunks of equal size (usually for large texts)
- And so on...
- Ultimately, the choice of tokenizer depends on the needs of the specific project

• Stemming is closely related to tokenizing

- Stemming is closely related to tokenizing
- Documents may use different forms of words ('lacked', 'lacking', 'lack'), or words which are similar in concept ('political', 'politician', 'politicizing')

- Stemming is closely related to tokenizing
- Documents may use different forms of words ('lacked', 'lacking', 'lack'), or words which are similar in concept ('political', 'politician', 'politicizing')
- **Stemming** maps these variants to the root of the word using a crude heuristic that chops off the affixes and returns a **stem**

- Stemming is closely related to tokenizing
- Documents may use different forms of words ('lacked', 'lacking', 'lack'), or words which are similar in concept ('political', 'politician', 'politicizing')
- **Stemming** maps these variants to the root of the word using a crude heuristic that chops off the affixes and returns a **stem**
- Lemmatization is similar, but stems based on vocabulary, parts of speech, and mapping rules, and returns a word in the dictionary

- Stemming is closely related to tokenizing
- Documents may use different forms of words ('lacked', 'lacking', 'lack'), or words which are similar in concept ('political', 'politician', 'politicizing')
- **Stemming** maps these variants to the root of the word using a crude heuristic that chops off the affixes and returns a **stem**
- Lemmatization is similar, but stems based on vocabulary, parts of speech, and mapping rules, and returns a word in the dictionary
- e.g. lemmatization would return 'see' or 'saw' if it came across 'saw' (clearly this is subjective, and requires constant quality checking)

Original Word		Stemmed Word
abolish	$\mapsto$	abolish
abolished	$\mapsto$	abolish
abolishing	$\mapsto$	abolish
abolition	$\mapsto$	abolit
abortion	$\mapsto$	abort
abortions	$\mapsto$	abort
abortive	$\mapsto$	abort
treasure	$\mapsto$	treasure
treasured	$\mapsto$	treasure
treasures	$\mapsto$	treasure
treasuring	$\mapsto$	treasure
treasury	$\mapsto$	treasuri

#### NYT

Emergency measures adopted for Beijing's first "red alert" over air pollution left millions of schoolchildren cooped up at home, forced motorists off the roads and shut down factories across the region on Tuesday, but they failed to dispel the toxic air that shrouded the Chinese capital in a soupy, metallic haze.

#### marked up

Emergency measures adopted for Beijing s first red alert over air pollution left millions of schoolchildren cooped up at home, forced motorists off the roads and shut down factories across the region on Tuesday, but they failed to dispel the toxic air that shrouded the Chinese capital in a soupy, metallic haze.

#### NYT

Emergency measures adopted for Beijing's first \red alert" over air pollution left millions of schoolchildren cooped up at home, forced motorists off the roads and shut down factories across the region on Tuesday, but they failed to dispel the toxic air that shrouded the Chinese capital in a soupy, metallic haze.

#### Stemmed

Emergenc measur adopt for Beij s first red alert over air pollut left million of schoolchildren coop up at home forc motorist off the road and shut down factori across the region on Tuesdai but thei fail to dispel the toxic air that shroud the Chines capit in a soupi metal haze.

• To this point, there have been no distinctions drawn between nouns, verbs, and so on regarding our tokens, which makes sense for some applications (e.g., classification)

- To this point, there have been no distinctions drawn between nouns, verbs, and so on regarding our tokens, which makes sense for some applications (e.g., classification)
- But, sometimes we may want to know information about the part-of-speech a word represents

- To this point, there have been no distinctions drawn between nouns, verbs, and so on regarding our tokens, which makes sense for some applications (e.g., classification)
- But, sometimes we may want to know information about the part-of-speech a word represents
- E.g., we are often interested in recording who did what to whom, e.g.,
   'the UK bombing will force ISIS to surrender'

- To this point, there have been no distinctions drawn between nouns, verbs, and so on regarding our tokens, which makes sense for some applications (e.g., classification)
- But, sometimes we may want to know information about the part-of-speech a word represents
- E.g., we are often interested in recording who did what to whom, e.g.,
   'the UK bombing will force ISIS to surrender'
- Here, force is a verb, not a noun

- To this point, there have been no distinctions drawn between nouns, verbs, and so on regarding our tokens, which makes sense for some applications (e.g., classification)
- But, sometimes we may want to know information about the part-of-speech a word represents
- E.g., we are often interested in recording who did what to whom, e.g.,
   'the UK bombing will force ISIS to surrender'
- Here, **force** is a verb, not a noun
- Annotating in this way is called parts-of-speech tagging (e.g., the RDRPOSTagger library in R)

• We have now pre-processed our texts

• We have now pre-processed our texts!!!



We have now preprocessed our texts → \*yay\*

- We have now preprocessed our texts → \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things

- We have now preprocessed our texts → \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things
- And we do (almost) as well without ordering as when we retain it

- We have now preprocessed our texts → \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things
- And we do (almost) as well without ordering as when we retain it
- In such a world, we are treating a document as a "bag of words"

- We have now preprocessed our texts → \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things
- And we do (almost) as well without ordering as when we retain it
- In such a world, we are treating a document as a "bag of words"
- ullet e.g. "The leading Republican presidential candidate has said Muslims should be banned from entering the US."  $\leadsto$

- ullet We have now preprocessed our texts  $\leadsto$  \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things
- And we do (almost) as well without ordering as when we retain it
- In such a world, we are treating a document as a "bag of words"
- ullet e.g. "The leading Republican presidential candidate has said Muslims should be banned from entering the US."  $\leadsto$
- ullet "lead republican presidenti candid said muslim ban enter us"  $\equiv$

- We have now preprocessed our texts → \*yay\*
- Generally, we are willing to ignore the order of the words in a document, which drastically simplifies things
- And we do (almost) as well without ordering as when we retain it
- In such a world, we are treating a document as a "bag of words"
- ullet e.g. "The leading Republican presidential candidate has said Muslims should be banned from entering the US."  $\leadsto$
- "lead republican presidenti candid said muslim ban enter us" ≡
   "us lead said candid presidenti ban muslim republican enter"

• Another way to think about collections of text is a vector space model

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document
- e.g. the document "Philip is short" can be thought of a vector in 3 dimensions:

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document
- e.g. the document "Philip is short" can be thought of a vector in 3 dimensions:  $d_1$  corresponds to how 'Philip'-ish it is,  $d_2$  corresponds to how 'is'-ish it is, and  $d_3$  corresponds to how 'short'-ish it is

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document
- e.g. the document "Philip is short" can be thought of a vector in 3 dimensions:  $d_1$  corresponds to how 'Philip'-ish it is,  $d_2$  corresponds to how 'is'-ish it is, and  $d_3$  corresponds to how 'short'-ish it is
- Features are typically unigram frequencies of the tokens in the document

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document
- e.g. the document "Philip is short" can be thought of a vector in 3 dimensions:  $d_1$  corresponds to how 'Philip'-ish it is,  $d_2$  corresponds to how 'is'-ish it is, and  $d_3$  corresponds to how 'short'-ish it is
- Features are typically unigram frequencies of the tokens in the document
- ullet e.g. "the dog chased the cat" becomes (2,1,1,1) for dimensions defined as simple counts across the, dog, chased, cat

- Another way to think about collections of text is a vector space model
- We have some document, which is considered a collection of W terms/features (words, tokens, etc.), where each term is a dimension
- Documents (comprising these terms) are the linear combinations of vectors along the axes, implying document  $W \rightsquigarrow \mathbb{R}^W$
- The angle between vectors (cosine of the angle) captures the "\_ness" of the document
- e.g. the document "Philip is short" can be thought of a vector in 3 dimensions:  $d_1$  corresponds to how 'Philip'-ish it is,  $d_2$  corresponds to how 'is'-ish it is, and  $d_3$  corresponds to how 'short'-ish it is
- Features are typically unigram frequencies of the tokens in the document
- e.g. "the dog chased the cat" becomes (2,1,1,1) for dimensions defined as simple counts across the, dog, chased, cat
- Thus, the vector space model represents a document vector in d-dimensional feature space

ullet  $d=1,\ldots,D \leadsto$  indexes documents in the corpus

- $d = 1, \dots, D \leadsto$  indexes documents in the corpus
- $w = 1, ..., W \rightsquigarrow$  indexes features in the documents

- $d=1,\ldots,D \leadsto$  indexes documents in the corpus
- $w = 1, ..., W \rightsquigarrow$  indexes features in the documents
- $oldsymbol{ iny y}_d \in \mathbb{R}^{\mathcal{W}} \leadsto$  representation of document d in a some feature space

- $d = 1, \dots, D \leadsto$  indexes documents in the corpus
- $w = 1, ..., W \rightsquigarrow$  indexes features in the documents
- $oldsymbol{ iny y}_d \in \mathbb{R}^{W} \leadsto$  representation of document d in a some feature space
- So each document is now a vector, with each entry representing the frequency of a particular token or feature

## Vector Space Models: Notation

- $d = 1, \dots, D \leadsto \text{indexes documents in the corpus}$
- $w = 1, ..., W \rightsquigarrow$  indexes features in the documents
- $oldsymbol{ iny y}_d \in \mathbb{R}^{W} \leadsto$  representation of document d in a some feature space
- So each document is now a vector, with each entry representing the frequency of a particular token or feature
- Stacking these vectors on top of each other gives the document term matrix (DTM) (sometimes called the document feature matrix (DFM))

 The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)

- The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)
- The frequency of a term is adjusted for how often and/or rarely it is used

- The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)
- The frequency of a term is adjusted for how often and/or rarely it is used
  - Calculate a term's frequency (tf)

- The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)
- The frequency of a term is adjusted for how often and/or rarely it is used
  - Calculate a term's frequency (tf)
  - ② Calculate a term's inverse document frequency (idf,  $In(\frac{N_{docs}}{N_{docs,containing,term}}))$

- The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)
- The frequency of a term is adjusted for how often and/or rarely it is used
  - Calculate a term's frequency (tf)
  - 2 Calculate a term's inverse document frequency (idf,  $ln(\frac{N_{docs}}{N_{docs}}, containing, term))$
  - tf-idf → product of tf and idf

- The most common approach to weighting word importance across documents is the tf-idf weighting scheme, based on Zipf's law (frequency a word appears is inversely proportional to its rank)
- The frequency of a term is adjusted for how often and/or rarely it is used
  - Calculate a term's frequency (tf)
  - ② Calculate a term's inverse document frequency (idf,  $In(\frac{N_{docs}}{N_{docs,containing,term}}))$
  - tf-idf → product of tf and idf
- Functionally, the tf-idf captures a decrease in the weight for commonly used words and, for more rarely used words, an increase in the weight for words use more rarely in some set of documents, D, that are not used very much in a collection of documents

 We may also be interested in mining the diversity of some document or set of documents 
 → lexical diversity

- Recall the elementary components of some text are called tokens (usually words, but could be numbers, etc.)

- Recall the elementary components of some text are called tokens (usually words, but could be numbers, etc.)
- The **types** in a document are the set of *unique* tokens

- Recall the elementary components of some text are called tokens (usually words, but could be numbers, etc.)
- The **types** in a document are the set of *unique* tokens
- Thus we typically have many more tokens than types, because authors repeat tokens

- Recall the elementary components of some text are called tokens (usually words, but could be numbers, etc.)
- The types in a document are the set of unique tokens
- Thus we typically have many more tokens than types, because authors repeat tokens
- To measure lexical diversity, we can use the type-to-token ratio TTR

$$TTR = \frac{N_{types}}{N_{tokens}}$$

- Recall the elementary components of some text are called tokens (usually words, but could be numbers, etc.)
- The types in a document are the set of unique tokens
- Thus we typically have many more tokens than types, because authors repeat tokens
- To measure lexical diversity, we can use the type-to-token ratio TTR

$$TTR = \frac{N_{types}}{N_{tokens}}$$

• This may provide increased clarity of the text, e.g., authors with limited vocabularies might have a low lexical diversity

#### Lecture Outline

- 1 Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- 4 Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

 We've covered the basics of text mining, simplification, document representation, and so on

- We've covered the basics of text mining, simplification, document representation, and so on
- We can build on this to begin to think about documents as members of categories or classes

- We've covered the basics of text mining, simplification, document representation, and so on
- We can build on this to begin to think about documents as members of categories or classes
- Integral to classification is the dictionary

- We've covered the basics of text mining, simplification, document representation, and so on
- We can build on this to begin to think about documents as members of categories or classes
- Integral to classification is the dictionary
- This allows us to move on to supervised learning problems

• What is supervised learning?

• What is supervised learning? labels are key

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y = 1), some that were negative (y = 0);

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y=1), some that were negative (y=0); whether some statements that were liberal, some that were conservative; etc.

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y=1), some that were negative (y=0); whether some statements that were liberal, some that were conservative; etc.
- Next, we train a machine (model) on these examples, using the features (DTM) as the independent variable,

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y=1), some that were negative (y=0); whether some statements that were liberal, some that were conservative; etc.
- Next, we train a machine (model) on these examples, using the features (DTM) as the independent variable, e.g. does the commentator use the word "fetus" or "baby" in discussing abortion law?

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y=1), some that were negative (y=0); whether some statements that were liberal, some that were conservative; etc.
- Next, we train a machine (model) on these examples, using the features (DTM) as the independent variable, e.g. does the commentator use the word "fetus" or "baby" in discussing abortion law? Gives a clue as to the "ideology" of the speaker/author

- What is supervised learning? labels are key
- We start by labelling some examples of each category, e.g. some reviews that were positive (y=1), some that were negative (y=0); whether some statements that were liberal, some that were conservative; etc.
- Next, we train a machine (model) on these examples, using the features (DTM) as the independent variable, e.g. does the commentator use the word "fetus" or "baby" in discussing abortion law? Gives a clue as to the "ideology" of the speaker/author
- The **goal**, then, is to classify by using the learned relationship between labels and features to predict the outcomes of *future* documents (e.g.,  $y \in \{0,1\}$ , sentiment) *not* in the training set

#### Lecture Outline

- 1 Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

• The most fundamental concept in supervised text classification:

• The most fundamental concept in supervised text classification: dictionary

- The most fundamental concept in supervised text classification: dictionary
- Set of pre-defined words with specific connotations that allow us to classify documents automatically, quickly and accurately

- The most fundamental concept in supervised text classification: dictionary
- Set of pre-defined words with specific connotations that allow us to classify documents automatically, quickly and accurately
- Very common in opinion mining/sentiment analysis, and also in coding events or party platforms

- The most fundamental concept in supervised text classification: dictionary
- Set of pre-defined words with specific connotations that allow us to classify documents automatically, quickly and accurately
- Very common in opinion mining/sentiment analysis, and also in coding events or party platforms
- Often used in supervised learning problems, as a starting point

- The most fundamental concept in supervised text classification: dictionary
- Set of pre-defined words with specific connotations that allow us to classify documents automatically, quickly and accurately
- Very common in opinion mining/sentiment analysis, and also in coding events or party platforms
- Often used in supervised learning problems, as a starting point
- As such, we will cover them today in this context

- The most fundamental concept in supervised text classification: dictionary
- Set of pre-defined words with specific connotations that allow us to classify documents automatically, quickly and accurately
- Very common in opinion mining/sentiment analysis, and also in coding events or party platforms
- Often used in supervised learning problems, as a starting point
- As such, we will cover them today in this context → sentiment analysis

#### Classification with Dictionaries

• Goal: usually, we are trying to do one of two closely related things:

#### Classification with Dictionaries

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements
- A dictionary guides us, the (weighted) presence of which helps us with either of these goals

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements
- A dictionary guides us, the (weighted) presence of which helps us with either of these goals
- Some weights are binary (either +1 or -1), while others are continuous (e.g., ranging from -5 to +5)

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements
- A dictionary guides us, the (weighted) presence of which helps us with either of these goals
- Some weights are binary (either +1 or -1), while others are continuous (e.g., ranging from -5 to +5)
- Some dictionaries capture broad sentiment (e.g., positive/negative),

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements
- A dictionary guides us, the (weighted) presence of which helps us with either of these goals
- Some weights are binary (either +1 or -1), while others are continuous (e.g., ranging from -5 to +5)
- Some dictionaries capture broad sentiment (e.g., positive/negative), others capture emotions (e.g., anger, joy, disgust),

- Goal: usually, we are trying to do one of two closely related things:
  - Categorize documents in unique classes
    - ★ This review is 'positive'; this speech is 'liberal'
  - Measure extent to which document is associated with a category
    - ★ This review is generally 'positive', but has some negative elements
- A dictionary guides us, the (weighted) presence of which helps us with either of these goals
- Some weights are binary (either +1 or -1), while others are continuous (e.g., ranging from -5 to +5)
- Some dictionaries capture broad sentiment (e.g., positive/negative), others capture emotions (e.g., anger, joy, disgust), still others are non-sentiment (e.g., politics, food, places)

• In solving the classification problem, have a set of key words with scores, e.g., "terrible" =-1; "fantastic" =+1

- In solving the classification problem, have a set of key words with scores, e.g., "terrible" =-1; "fantastic" =+1
- The relative rate of occurrence of these terms tells us about the overall tone/class the document should be placed in

- In solving the classification problem, have a set of key words with scores, e.g., "terrible" =-1; "fantastic" =+1
- The relative rate of occurrence of these terms tells us about the overall tone/class the document should be placed in
- The tone, Y, based on document i and words m = 1, ..., M in the **dictionary**,

$$Y_i = \sum_{m=1}^{M} \frac{s_m w_{im}}{N_i}$$

#### where

- $ightharpoonup s_m$  is the score of the word m
- w<sub>im</sub> is the number of occurrences of the mth dictionary word in document i
- $\triangleright$   $N_i$  is the total number of all **dictionary** words in the document

 To classify, simply add up the number of times the dictionary words appear and multiply by the score (normalizing by document dictionary presence)

- To classify, simply add up the number of times the dictionary words appear and multiply by the score (normalizing by document dictionary presence)
- Dependent on the dictionary, the score, Y ("tone") impacts document classification,

- To classify, simply add up the number of times the dictionary words appear and multiply by the score (normalizing by document dictionary presence)
- Dependent on the dictionary, the score, Y ("tone") impacts document classification,
  - ▶  $Y > 0 \rightsquigarrow positive$

- To classify, simply add up the number of times the dictionary words appear and multiply by the score (normalizing by document dictionary presence)
- Dependent on the dictionary, the score, Y ("tone") impacts document classification,
  - ▶  $Y > 0 \rightsquigarrow positive$
  - ▶  $Y < 0 \rightsquigarrow \text{negative}$

- To classify, simply add up the number of times the dictionary words appear and multiply by the score (normalizing by document dictionary presence)
- Dependent on the dictionary, the score, Y ("tone") impacts document classification,
  - ▶  $Y > 0 \rightsquigarrow positive$
  - ▶  $Y < 0 \rightsquigarrow \text{negative}$
  - ►  $Y \approx 0 \rightsquigarrow \text{ambiguous}$

## For example... The Big Short (newsreview.com)

Director and co-screenwriter Adam McKay (Step Brothers) bungles a great opportunity to savage the architects of the 2008 financial crisis in The Big Short, wasting an A-list ensemble cast in the process. Steve Carell, Brad Pitt, Christian Bale and Ryan Gosling play various tenuously related members of the finance industry. men who made made a killing by betting against the housing market, which at that point had superficially swelled to record highs. All of the elements are in place for a lacerating satire, but almost every aesthetic choice in the film is bad, from the U-Turn-era Oliver Stone visuals to Carell's sketch-comedy performance to the cheeky cutaways where Selena Gomez and Anthony Bourdain explain complex financial concepts. After a brutal opening half, it finally settles into a groove, and there's a queasy charge in watching a credit-drunk America walking towards that cliff's edge, but not enough to save the film.

## For example... The Big Short (newsreview.com)

```
great
                       savage
crisis
                        wasting
         tenuously
                   killing
                         superficially swelled
                                   bad
 complex
                                       brutal
        drunk
enough
```

• Using the Bing dictionary, we can see that there are 11 negative words, and 2 positive words

- Using the Bing dictionary, we can see that there are 11 negative words, and 2 positive words
- Thus, the tone could be calculated as,

$$=\frac{2-11}{13}$$

- Using the Bing dictionary, we can see that there are 11 negative words, and 2 positive words
- Thus, the tone could be calculated as,

$$=\frac{2-11}{13}$$

$$=\frac{-9}{13}$$

- Using the Bing dictionary, we can see that there are 11 negative words, and 2 positive words
- Thus, the tone could be calculated as,

$$=\frac{2-11}{13}$$

$$=\frac{-9}{13}$$

-0.6923

- Using the Bing dictionary, we can see that there are 11 negative words, and 2 positive words
- Thus, the tone could be calculated as,

$$=\frac{2-11}{13}$$

$$=\frac{-9}{13}$$

-0.6923

bad

• Dictionaries are context dependent

- Dictionaries are context dependent
- Most dictionaries do not take into account qualifiers (e.g. "no good")

- Dictionaries are context dependent
- Most dictionaries do not take into account qualifiers (e.g. "no good")
- All ignore sarcasm, irony, nuance

- Dictionaries are context dependent
- Most dictionaries do not take into account qualifiers (e.g. "no good")
- All ignore sarcasm, irony, nuance
- Ultimately context matters, and any dictionary used should be clearly justified

• Common sentiment dictionaries:

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words
  - ▶ **Bing**: Binary classification for positive/negative sentiments

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words
  - ▶ **Bing**: Binary classification for positive/negative sentiments
  - ▶ NRC: Binary classification for a variety of categories like positive, negative, disgust, fear, joy, trust, surprise

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words
  - ▶ **Bing**: Binary classification for positive/negative sentiments
  - ▶ NRC: Binary classification for a variety of categories like positive, negative, disgust, fear, joy, trust, surprise
  - ► LIWC: 2300 words grouped into 70 classes on positive/negative emotions

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words
  - ▶ **Bing**: Binary classification for positive/negative sentiments
  - ► NRC: Binary classification for a variety of categories like positive, negative, disgust, fear, joy, trust, surprise
  - LIWC: 2300 words grouped into 70 classes on positive/negative emotions
  - ► **AFINN**: Continuous scores between −5 and 5 for negative and positive sentiment, respectively

- Common sentiment dictionaries:
  - ► **ANEW** (Affective Norms for English Words): "Happiness" dictionary with 1034 words, measuring affective reactions to words
  - ▶ **Bing**: Binary classification for positive/negative sentiments
  - ► NRC: Binary classification for a variety of categories like positive, negative, disgust, fear, joy, trust, surprise
  - LIWC: 2300 words grouped into 70 classes on positive/negative emotions
  - ► **AFINN**: Continuous scores between −5 and 5 for negative and positive sentiment, respectively
  - ► **General Inquirer Database**: 3627 negative and positive word strings (widely used across domains)

### Lecture Outline

- 1 Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

### Creating Dictionaries

 These dictionaries are great, and very widely used, especially these days

### Creating Dictionaries

- These dictionaries are great, and very widely used, especially these days
- But what if the domain requirements limits their utility?

### Creating Dictionaries

- These dictionaries are great, and very widely used, especially these days
- But what if the domain requirements limits their utility?
- In other words, what if we are interested in weighting word frequencies with non-sentiment-based words, like, political language, gender language, racist language, and so on...

- These dictionaries are great, and very widely used, especially these days
- But what if the domain requirements limits their utility?
- In other words, what if we are interested in weighting word frequencies with non-sentiment-based words, like, political language, gender language, racist language, and so on...
- In this case, we might consider creating our own dictionaries

• Though there are a several ways, here are the three most widely used/recognized:

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods
  - ② Manual generation → careful thought about useful words

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods
  - ② Manual generation → careful thought about useful words
  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods
  - ② Manual generation → careful thought about useful words
  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)
    - ★ Undergraduates

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods

  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)
    - ★ Undergraduates: Pizza ~> Research Output

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods
  - Manual generation → careful thought about useful words
  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)
    - ★ Undergraduates: Pizza ~> Research Output
    - ★ Mechanical turkers

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods

  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)
    - ★ Undergraduates: Pizza ~> Research Output
    - ★ Mechanical turkers, e.g., ask turkers: how happy is, elevator, car, pretty, young

- Though there are a several ways, here are the three most widely used/recognized:
  - Separating methods

  - Populations of people who are (surprisingly) willing to perform ill-defined tasks (i.e., crowd-sourcing)
    - ★ Undergraduates: Pizza ~> Research Output
    - ★ Mechanical turkers, e.g., ask turkers: how happy is, elevator, car, pretty, young ~ Output as dictionary

• What is the goal here?

 What is the goal here? Find words that discriminate between groups of texts, i.e., words distinctive to each document

- What is the goal here? Find words that discriminate between groups of texts, i.e., words distinctive to each document
- There are three main ways to separate, and thus manually locate distinctive words in text

- What is the goal here? Find words that discriminate between groups of texts, i.e., words distinctive to each document
- There are three main ways to separate, and thus manually locate distinctive words in text
  - Exclusive/unique use: words used in a simgle document

- What is the goal here? Find words that discriminate between groups of texts, i.e., words distinctive to each document
- There are three main ways to separate, and thus manually locate distinctive words in text
  - Exclusive/unique use: words used in a simgle document
  - 2 Difference in frequencies: absolute differences frequency of use

- What is the goal here? Find words that discriminate between groups of texts, i.e., words distinctive to each document
- There are three main ways to separate, and thus manually locate distinctive words in text
  - Exclusive/unique use: words used in a simgle document
  - ② Difference in frequencies: *absolute* differences frequency of use
  - Oifference in rates/average usage: difference between proportions of same word use across documents (where high difference = good/distinct)

# Separating Methods in R

• Quick demo of each in R

• Dictionary methods are context invariant

- Dictionary methods are context invariant
  - ► Optimization  $\leadsto$  allows you to incorporate information specific to context

- Dictionary methods are context invariant
  - ▶ Optimization ~ allows you to incorporate information specific to context
  - ▶ No optimization → same word weights regardless of texts or contexts

- Dictionary methods are context invariant
  - ▶ Optimization → allows you to incorporate information specific to context
  - ▶ No optimization → same word weights regardless of texts or contexts
  - Without optimization, its unclear about dictionaries' performances (e.g., did they correctly classify?)

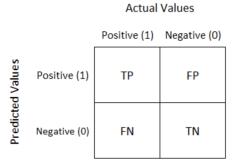
- Dictionary methods are context invariant
  - ► Optimization → allows you to incorporate information specific to context
  - ▶ No optimization → same word weights regardless of texts or contexts
  - Without optimization, its unclear about dictionaries' performances (e.g., did they correctly classify?)
- Importantly, just because we get a positive or negative score, does
   NOT mean that these are accurate measures in our text

- Dictionary methods are context invariant
  - ▶ Optimization ~→ allows you to incorporate information specific to context
  - ▶ No optimization → same word weights regardless of texts or contexts
  - Without optimization, its unclear about dictionaries' performances (e.g., did they correctly classify?)
- Importantly, just because we get a positive or negative score, does
   NOT mean that these are accurate measures in our text
- This points to the need for validation

- Classification validity (requires hand coded documents):
  - ► Training: build dictionary on subset of documents with known labels
  - Testing: apply dictionary method to other documents with known labels
    - ★ Is the classification scheme well defined for your texts?
    - \* Can humans accomplish the coding task with consistency (e.g., Cronbach's α)?
    - ★ Is the dictionary appropriate?
- Replicate classification exercise
  - How well does our method perform on held out documents?
  - Why "held out"? Over-fitting
  - (Cross)validation
  - ► Can also use off-the-shelf dictionaries to compare

 The most widely use method for assessing classification is a confusion matrix

 The most widely use method for assessing classification is a confusion matrix



• Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:

- Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:
  - ► Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN}$

- Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:
  - ▶ Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN}$
  - ▶ Specificity:  $\frac{TN}{TN+FP}$

- Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:
  - ▶ Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN}$
  - ▶ Specificity:  $\frac{TN}{TN+FP}$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP}$

- Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:
  - ▶ Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN}$
  - ▶ Specificity:  $\frac{TN}{TN+FP}$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP}$
  - Accuracy:  $\frac{TP+TN}{TP+TN+FP+FN}$

- Confusion matrices help us visualize our classification results, while also gauging the precision of our estimate and determining error:
  - ▶ Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN}$
  - ▶ Specificity:  $\frac{TN}{TN+FP}$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP}$
  - Accuracy:  $\frac{TP+TN}{TP+TN+FP+FN}$
  - ▶ F1 score ("F-measure"):  $\frac{2TP}{2TP+FP+FN}$

 Suppose our classification of fraudulent documents yielded the following results

	fraud (1)	genuine (0)
$\widehat{\textit{fraud}}$ (1)	4	1
genuine (0)	0	2

fraud (1) genuine (0)
$$\widehat{fraud} (1) \qquad \qquad 4 \qquad \qquad 1$$

$$\widehat{genuine} (0) \qquad \qquad 0 \qquad \qquad 2$$

- Now we can check our accuracy by simply plugging in the values
  - ► Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN} \rightsquigarrow \frac{4}{4+0} = 1.0$

fraud (1) genuine (0)
$$\widehat{fraud} (1) \qquad \qquad 4 \qquad \qquad 1$$

$$\widehat{genuine} (0) \qquad \qquad 0 \qquad \qquad 2$$

- Now we can check our accuracy by simply plugging in the values
  - ► Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN} \rightsquigarrow \frac{4}{4+0} = 1.0$
  - Specificity:  $\frac{TN}{TN+FP} \rightsquigarrow \frac{2}{2+1} = 0.67$

fraud (1) genuine (0)
$$\widehat{fraud} (1) \qquad \qquad 4 \qquad \qquad 1$$

$$\widehat{genuine} (0) \qquad \qquad 0 \qquad \qquad 2$$

- Now we can check our accuracy by simply plugging in the values
  - ► Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN} \rightsquigarrow \frac{4}{4+0} = 1.0$
  - Specificity:  $\frac{TN}{TN+FP} \rightsquigarrow \frac{2}{2+1} = 0.67$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP} \rightsquigarrow \frac{4}{4+1} = 0.80$

#### The Confusion Matrix

fraud (1) genuine (0) 
$$\widehat{fraud}$$
 (1) 4 1  $\widehat{genuine}$  (0) 0 2

- Now we can check our accuracy by simply plugging in the values
  - ▶ Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN} \rightsquigarrow \frac{4}{4+0} = 1.0$
  - ► Specificity:  $\frac{TN}{TN+FP} \rightsquigarrow \frac{2}{2+1} = 0.67$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP} \rightsquigarrow \frac{4}{4+1} = 0.80$
  - Accuracy:  $\frac{TP+TN}{TP+TN+FP+FN} \leftrightarrow \frac{4+2}{4+2+1+0} = 0.86$

#### The Confusion Matrix

fraud (1) genuine (0) 
$$\widehat{fraud}$$
 (1) 4 1  $\widehat{genuine}$  (0) 0 2

- Now we can check our accuracy by simply plugging in the values
  - ► Sensitivity (recall or "hit rate"):  $\frac{TP}{TP+FN} \rightsquigarrow \frac{4}{4+0} = 1.0$
  - ► Specificity:  $\frac{TN}{TN+FP} \rightsquigarrow \frac{2}{2+1} = 0.67$
  - ▶ Precision ("positive predicted value"):  $\frac{TP}{TP+FP} \rightsquigarrow \frac{4}{4+1} = 0.80$
  - Accuracy:  $\frac{TP+TN}{TP+TN+FP+FN} \rightsquigarrow \frac{4+2}{4+2+1+0} = 0.86$
  - ► F1 score ("F-measure"):  $\frac{2TP}{2TP+FP+FN} \rightsquigarrow \frac{2(4)}{2(4)+1+0} = 0.89$
- In R: table(), e.g., table(predicted\_classes, actual\_classes)

#### Lecture Outline

- 1 Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- 4 Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone
- Set of (human) hand-coded documents
  - ▶ Training Set: documents we'll use to learn how to code
  - ► Test/Validation Set: documents we'll use to learn how well we code

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone
- Set of (human) hand-coded documents
  - ► Training Set: documents we'll use to learn how to code
  - ► Test/Validation Set: documents we'll use to learn how well we code
- Obtain a new set of unlabeled documents

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone
- Set of (human) hand-coded documents
  - ► Training Set: documents we'll use to learn how to code
  - ► Test/Validation Set: documents we'll use to learn how well we code
- Obtain a new set of unlabeled documents
- Rinse and repeat

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone
- 2 Set of (human) hand-coded documents
  - ► Training Set: documents we'll use to learn how to code
  - ► Test/Validation Set: documents we'll use to learn how well we code
- 3 Obtain a new set of unlabeled documents
- Rinse and repeat
- But how do we classify?

- Set of categories, e.g. sentiments
  - ▶ Positive Tone, Negative Tone
- Set of (human) hand-coded documents
  - ► Training Set: documents we'll use to learn how to code
  - ► Test/Validation Set: documents we'll use to learn how well we code
- Obtain a new set of unlabeled documents
- Rinse and repeat
  - But how do we classify? many ways, but we will cover basic regression as an example

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}$ 

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$  We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$ We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

$$f(\beta, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} (y_i - \beta' \boldsymbol{x}_i)^2$$

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$  We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2$$

$$\widehat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2 \right\}$$
(1)

Problem:

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$  We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2$$

$$\widehat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2 \right\}$$
(1)

#### Problem:

- J will likely be large (perhaps J > N)

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$ We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2$$

$$\widehat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2 \right\}$$
(1)

#### Problem:

- J will likely be large (perhaps J > N)
- There many correlated variables

Suppose we have N documents, with each document i having label  $y_i \in \{-1,1\} \leadsto \{\text{negative, positive}\}\$  We represent each document i is  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$ .

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2$$

$$\widehat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^{N} (y_i - \beta' \mathbf{x}_i)^2 \right\}$$
(1)

#### Problem:

- J will likely be large (perhaps J > N)
- There many correlated variables
- Predictions will be variable

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y})$$

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2$$

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^{J} \beta_j^2}_{\text{Penalty}}$$

Penalty for model complexity

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^{J} \beta_j^2}_{\text{Penalty}}$$

where:

Penalty for model complexity

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^{J} \beta_j^2}_{\text{Penalty}}$$

where:

-  $\beta_0 \rightsquigarrow \text{intercept}$ 

Penalty for model complexity

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^{J} \beta_j^2}_{\text{Penalty}}$$

#### where:

- $\beta_0 \rightsquigarrow \text{intercept}$
- $\lambda \leadsto$  penalty parameter

Penalty for model complexity

$$f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \underbrace{\lambda \sum_{j=1}^{J} \beta_j^2}_{\text{Penalty}}$$

#### where:

- $\beta_0 \rightsquigarrow \text{intercept}$
- $\lambda \leadsto$  penalty parameter
- Standardized **X** (coefficients on same scale)

## Ridge Regression → Optimization

$$\boldsymbol{\beta}^{\mathsf{Ridge}} = \arg\min_{\boldsymbol{\beta}} \left\{ f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) \right\}$$

## Ridge Regression <>> Optimization

$$\begin{split} \boldsymbol{\beta}^{\mathsf{Ridge}} &= \operatorname{arg\ min}_{\boldsymbol{\beta}} \left\{ f(\boldsymbol{\beta}, \boldsymbol{X}, \boldsymbol{Y}) \right\} \\ &= \operatorname{arg\ min}_{\boldsymbol{\beta}} \left\{ \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{J} \beta_j^2 \right\} \end{split}$$

## Other Penalized Objective Functions

Different Penalty for Model Complexity: LASSO

$$f(\beta, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{J} \underbrace{|\beta_j|}_{\text{Penalty}}$$

## Other Penalized Objective Functions

Different Penalty for Model Complexity: LASSO

$$f(\beta, \boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{J} \underbrace{|\beta_j|}_{\mathsf{Penalty}}$$

And combining the two criteria --> Elastic-Net

$$f(\beta, \mathbf{X}, \mathbf{Y}) = \frac{1}{2N} \sum_{i=1}^{N} \left( y_i - \beta_0 - \sum_{j=1}^{J} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{J} \left( \frac{1}{2} (1 - \alpha) \beta_j^2 + \alpha |\beta_j| \right)$$

• Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, not isolated term frequencies

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, not isolated term frequencies
- Further, we can do so in a statistically principled manner with penalities, though we can reach a similar conclusion with any number of classification techniques (e.g., naive Bayes)

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, *not* isolated term frequencies
- Further, we can do so in a statistically principled manner with penalities, though we can reach a similar conclusion with any number of classification techniques (e.g., naive Bayes)
- Always think very carefully about the context of texts, assumptions of classifiers, and validation of any initial patterns by testing,

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, not isolated term frequencies
- Further, we can do so in a statistically principled manner with penalities, though we can reach a similar conclusion with any number of classification techniques (e.g., naive Bayes)
- Always think very carefully about the context of texts, assumptions of classifiers, and validation of any initial patterns by testing, re-testing,

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, not isolated term frequencies
- Further, we can do so in a statistically principled manner with penalities, though we can reach a similar conclusion with any number of classification techniques (e.g., naive Bayes)
- Always think very carefully about the context of texts, assumptions of classifiers, and validation of any initial patterns by testing, re-testing, re-testing,

- Ultimately, via more complex supervised techniques, we can bypass the lack of optimization problem inherent in basic dictionary methods
- We can optimize soem cost function with respect to context, by incorporating explanatory features (all word frequencies) in the probability of assignment of some document to a class
- The result is more realistic estimated classes of text based on the entirety of the document, not isolated term frequencies
- Further, we can do so in a statistically principled manner with penalities, though we can reach a similar conclusion with any number of classification techniques (e.g., naive Bayes)
- Always think very carefully about the context of texts, assumptions of classifiers, and validation of any initial patterns by testing, re-testing, re-testing, and... you get the idea

#### Lecture Outline

- 1 Text Mining
- 2 A Crash Course in Supervised Learning
- 3 Dictionaries
- 4 Manually Locating Distinctive Words
- 5 Putting It All Together: Parametric Supervised Classification
- 6 Some useful packages and functions

#### R

- ▶ tm package (e.g., Corpus(), DocumentTermMatrix(), etc.)
- tidytext package (e.g., get\_sentiments(), etc.)
- wordcloud package
- ▶ (more for next class, but still for your problem set) topicmodels
- ▶ stm package for structural topic models (may or may not get there)

#### Python:

- ► NI.TK
- ▶ spaCy
- ▶ Scikit-Learn
- ▶ For viz: Matplotlib and Seaborn