# Natural Language Processing

(NLP)

# Objectives for the morning

- Explain the steps of turning raw text data into useful features
- Explain the difference between stemming and lemmatization
- Compute the TF-IDF of documents

# Motivation

- You run Google News and you want to group news articles by topic

- You run a legal tech firm and you need to sift through 1,000,000 pages of legal documents to find the relevant ones

# Text Featurization Pipeline

- Tokenization
- Lower case conversion
- Punctuation removal
- Stop words removal
- Stemming / Lemmatization
- Bag of words / N-grams

# Terminology

- Corpus:
    - A dataset of text
    - E.g. newspaper articles, tweets
- Document:
    - A single entry from our corpus
    - E.g. article, sentence, tweet
- Vocabulary:
    - All the words that appear in the corpus
- Token:
    - An entity
    - E.g. a word

# Sample Sentence

"Students are learning from other students"

# Tokenization

Take the document and split it into a list of tokens

"Students are learning from other students"

⬇

["Students", "are", "learning", "from", "other", "students"]

# Lower case conversion

["Students", "are", "learning", "from", "other", "students"]

["students", "are", "learning", "from", "other", "students"]

# Stop word removal

- Remove the words we ignore in our analysis that are too common to be useful (since they are too common, they do not provide predictive power)

- Sklearn and nltk have standard list of stop words

["students", "are", "learning", "from", "other", "students"]

["students", "learning", "other", "students"]

# Stemming/Lemmatization

- Stemming:
  - Removes morphical suffixes
    - speaking -> speak
    - actually -> actual
  - Does it without context
  - Sometimes weird
  - Pretty fast

- Lemmatization:
  - Replace words with canonical form
    - better -> good
    - speaking -> speak
  - Uses hash tables
  - Slower

# Stemming/Lemmatization

["students", "learning", "other", "students"]

⬇

["student", "learn", "other", "student"]

# Corpus with 3 documents

Doc 1: "Students are learning from other students"

- ["student", "learn", "other", "student"]

Doc 2: "I am teaching at Galvanize"

- ["teach", "galvanize"]

Doc 3: "There are students learning at Galvanize"

- ["student", "learn", "galvanize"]

# Bag of words

A document represented as a vector of word counts is called "bag of words"

Vector for our corpus: (galvanize, learn, other, student, teach)

| document | galvanize | learn | other | student | teach |
|----------|-----------|-------|-------|---------|-------|
| Doc 1 | | | | | |
| Doc 2 | | | | | |
| Doc 3 | | | | | |

# Bag of words

A document represented as a vector of word counts is called "bag of words"

Vector for our corpus: (galvanize, learn, other, student, teach)

| document | galvanize | learn | other | student | teach |
|----------|-----------|-------|-------|---------|-------|
| Doc 1 | 0 | 1 | 1 | 2 | 0 |
| Doc 2 | 1 | 0 | 0 | 0 | 1 |
| Doc 3 | 1 | 1 | 0 | 1 | 0 |

# Bag of words

- Issues with bag of words:
    - Word counts
        - Counts emphasize results from longer documents
    - Every word has equal weighting
        - "other" and "student" have different predictive power.

# Term Frequency - Inverse Document Frequency

- Measures the relevance of a word
- Adjusts word count by document length and how often it shows up in the corpus

# Term Frequency

Normalize counts within a document to frequency

$$tf(t, d) = \frac{total\ count\ of\ term\ t\ in\ document\ d}{total\ count\ of\ all\ terms\ in\ document\ d}$$

| document | galvanize | learn | other | student | teach |
|---|---|---|---|---|---|
| Doc 1 | | | | | |
| Doc 2 | | | | | |
| Doc 3 | | | | | |

# Term Frequency

Normalize counts within a document to frequency

$$tf(t, d) = \frac{total\ count\ of\ term\ t\ in\ document\ d}{total\ count\ of\ all\ terms\ in\ document\ d}$$

| document | galvanize | learn | other | student | teach |
|---|---|---|---|---|---|
| Doc 1 | 0 | ¼ = 0.25 | ¼ = 0.25 | 2/4 = 0.5 | 0 |
| Doc 2 | ½ = 0.5 | 0 | 0 | 0 | ½ = 0.5 |
| Doc 3 | ⅓ = 0.33 | ⅓ = 0.33 | 0 | ⅓ = 0.33 | 0 |

# Term Frequency

Issues with term frequency:

- Words found only in one document should have highest weighting
- Words found in every document should have lowest weighting

# Inverse Document Frequency

$$idf(t, D) = log\frac{total\ number\ of\ document\ incorpus\ D}{count\ of\ document\ containing\ term\ t}$$

| document | galvanize | learn | other | student | teach |
|---|---|---|---|---|---|
| Doc 1 | | X | X | X | |
| Doc 2 | X | | | | X |
| Doc 3 | X | X | | X | |
| *idf(t,D)* | log(3/2) | log(3/2) | log(3/1) | log(3/2) | log(3/1) |

# TF-IDF

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

| document | galvanize | learn | other | student | teach |
|---|---|---|---|---|---|
| Doc 1 | 0 | 0.25xlog(3/2) = 0.101 | 0.25xlog(3/1) = 0.275 | 0.5xlog(3/2) = 0.203 | 0 |
| Doc 2 | 0.5xlog(3/2) = 0.203 | 0 | 0 | 0 | 0.5xlog(3/1) = .549 |
| Doc 3 | 0.33xlog(3/2) = 0.135 | 0.33xlog(3/2) = 0.135 | 0 | 0.33xlog(3/2) = 0.135 | 0 |

# Comparing TF-IDF vectors of documents

Cosine similarity:  $similarity = cos\theta = \dfrac{A \cdot B}{\|A\| \, \|B\|}$

- Doc 1 vs. Doc 2:
    - ["student", "learn", "other", "student"] vs ["teach", "galvanize"]
    
    $(0, 0.101, 0.275, 0.203, 0) \; vs. \; (0.203, 0, 0, 0, 0.275)$
    
    $similarity = \dfrac{0}{0.36 \times 0.34} = 0$

- Doc 1 vs. Doc 3:
    - ["student", "learn", "other", "student"] vs ["student", "learn", "galvanize"]
    
    $(0, 0.101, 0.275, 0.203, 0) \; vs. \; (0.135, 0.135, 0, 0.135, 0)$
    
    $similarity = \dfrac{0.041}{0.36 \times 0.23} = 0.34$

# N-Grams & Skip Grams

- Sliding window on text
- Handle phrases or words that go together (to capture the relationship between the consecutive words)

"The rain in Spain falls mainly on the plain"

Bi-gram: ["the rain", "rain in", "in spain", "spain falls", … ]

1-skip-bi-gram: ["the in", "rain spain", "in falls", ... ]

# Morning Assignment