CS 505: Introduction to Natural Language Processing Wayne Snyder Boston University

Lecture 6: Conclusion to N-Grams; Introduction to Vector Models: BOW, Distance Metrics, TF-IDF, PCA



From last time: Lack of data for all N-Grams

Training set:

... denied the allegations... denied the reports... denied the claims... denied the request

Test set
 ... denied the offer
 ... denied the loan

P("offer" | denied the) = 0

Zero probability bigrams

- Bigrams with zero probability
 - mean that we will assign 0 probability to the test set!
- And hence we cannot compute perplexity (can't divide by 0)!

The intuition of smoothing (from Dan Klein)

When we have sparse statistics:

- P(w | denied the) 3 allegations 2 reports
- 1 claims
- 1 request
- 7 total



Steal probability mass to generalize better

P(w | denied the) 2.5 allegations 1.5 reports 0.5 claims 0.5 request 2 other 7 total



Add-one estimation

- Also called Laplace smoothing
- Pretend we saw each word one more time than we did
- Just add one to all the counts!
- Normal (Most Likely Estimate):
- Add-1 estimate:

$$P_{MLE}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P_{Add-1}(w_i \mid w_{i-1}) = \frac{c(w_{i-1}, w_i) + 1}{c(w_{i-1}) + V}$$

Another solution to 0 N-gram counts: Backoff

Another solution to 0 counts for N-Grams is to use less left context:

This is called Backoff; suppose you have a 4-Gram model:

use quadrigram	if you have "good evidence,"
use trigram	if you have "good evidence,"
otherwise bigram,	if you have "good evidence,"
else unigram (single word),	if you have "good evidence."

If you have a word in test set that does not occur in the training set, use # occurrences of word in corpus / size of corpus

In the testing set:

<s> John likes to play Beethoven symphonies on a harmonica

In the training set:

0 occurrences of "symphonies on a harmonica" 1 occurrence of "on a harmonica" 5 occurrences of "a harmonica" 10 occurrences of "harmonica"

What is "good evidence here"?

Backoff and Interpolation

In the testing set:

<s> John likes to play Beethoven symphonies on a harmonica

In the training set: 0 occurrences of "symphonies on a harmonica" 1 occurrence of "on a harmonica" 5 occurrences of "a harmonica" 10 occurrences of "harmonica"

What is "good evidence" for the probability of "harmonica" following "symphonies on a"? How to calculate this probability using the evidence?

Interpolation: Weighted sum of probabilities for each of trigram, bigram, and unigram; weights can be estimated, or learned from training set.

Stupid Backoff: Multiply (recursively) by a constant "discount factor" 0.4:

P("symphonies on a harmonica") = 0.4 * P("on a harmonica") (on HW 03!)

N-gram Smoothing Summary

Used to deal with missing data

- Add-1 smoothing:
 - OK for text categorization, not for language modeling
- o Backoff and Interpolation
 - Use thresholds for counts
 - Learn weights for interpolation
- o Combination approaches
 - Extended Interpolated Kneser-Ney (state of the art, covered in the text)
- BUT: when you have enough data, doesn't matter which you use...

Vector Models

A BOW model represents a text by a vector of frequency counts over all the vocabulary:

	BC	W vector for				
	tex	t Julius Caesa	ar	Texts		
Vocabulary						
П						
4	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

BOW = Term frequency vector

- The term frequency $TF_{t,d}$ of term *t* in document *d* is defined as the number of times that *t* occurs in *d*.
- We want to use TF when using BOW vectors in NLP tasks such as
 - Queries: Which play talks about religion the most?
 - Classification: Which plays are tragedies and which are comedies?
 - Similarity: Which play is most similar to Julius Caesar?
- Raw term frequency is not what we want:
 - A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.
 - But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency.

Document frequency

- Rare terms are more informative than frequent terms
 - Recall stop words
- Consider a term in the query that is rare in the collection (e.g., arachnocentric)
- A document containing this term is very likely to be relevant to the query "What plays are *arachnocentric*?"
- \rightarrow We want a high weight for rare terms like *arachnocentric*.

Document frequency, continued

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., *high, increase, line*)
- A document containing such a term is more likely to be relevant than a document that doesn't
- But it's not a sure indicator of relevance.
- $\circ \rightarrow$ For frequent terms, we want high positive weights for words like *high*, *increase, and line*
- But lower weights than for rare terms.
- We will use document frequency DF to capture this.

IDF: Inverse Document Frequency

- DF_t is the document frequency of a token t = the number of documents
 that contain t
 - DF_t is an inverse measure of the informativeness of t
 - $DF_t \leq N$ (length of corpus in tokens)
- We define the IDF (inverse document frequency) of t by

$$IDF_t = log_{10}\left(\frac{N}{DF_t}\right)$$

• We use log (N/DF_t) instead of N/DF_t to "dampen" the effect of IDF.

IDF Example: suppose N = 1 million

calpurnia	10	1,000,000
animal	100	500,000
sunday	1,000	333,333.33
fly	10,000	250,000
under	100,000	200,000
the	1,000,000	166666.66

There is one idf value for each term *t* in a collection.

TF-IDF weighting

 The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$\mathbf{w}_{t,d} = \log(1 + \mathrm{tf}_{t,d}) \times \log_{10}(N/\mathrm{df}_t)$$

- Best known weighting scheme in information retrieval
 - Note: the "-" in TF-IDF is a hyphen, not a minus sign!
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection

TF-IDF Vector Models

Each document is now represented by a real-valued vector of tf-idf weights $\in \mathbb{R}^{|V|}$

	Julius Caesar
	73
	157
DF vector:	227
	10
	0
	0
	0

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	5.25	3.18	0	0	0	0.35
Brutus	1.21	6.1	0	1	0	0
Caesar	8.59	2.54	0	1.51	0.25	0
Calpurnia	0	1.54	0	0	0	0
Cleopatra	2.85	0	0	0	0	0
mercy	1.51	0	1.9	0.12	5.25	0.88
worser	1.37	0	0.11	4.15	0.25	1.95

V = Vocabulary

|V| = Size of V

Documents as vectors (TF BOW, TF-IDF, etc....)

- So we have a |V|-dimensional vector space
- o Terms are axes of the space
- o Documents are points or vectors in this space
- Very high-dimensional: Number of dimensions = size of vocabulary,
 often 10,000 or more.
- o These are very sparse vectors
 - most entries are zero.



Formalizing vector space proximity

- Basic idea: Two documents are similar if their vectors are "close" to each other = but how do we define "closeness"?
- First cut: distance between two points
 - (= distance between the end points of the two vectors)
- Euclidean distance?
- Euclidean distance is a bad idea . . .
- ... because Euclidean distance is large for vectors of different lengths.

Why distance is a bad idea

The Euclidean distance between qand d_2 is large even though the distribution of terms in the query q and the distribution of terms in the document d_2 are very similar.



Use angle instead of distance

- Thought experiment: take a document *d* and append it to itself. Call this document *d*'.
- "Semantically" d and d' have the same content
- The Euclidean distance between the two documents can be quite large
- The angle between the two documents is 0, corresponding to maximal similarity.



From angles to cosines

Instead of angle in degrees, we use cosine of the angle:

 $\cos 0 = 0$ angle



Cosine Similarity

$$\operatorname{cosine \ similarity} = S_{C}(A, B) := \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_{i}B_{i}}{\sqrt{\sum_{i=1}^{n} A_{i}^{2} \cdot \sum_{i=1}^{n} B_{i}^{2}}},$$

$$\operatorname{Cos} \quad \operatorname{Angle} \quad \operatorname{Meaning}$$

$$1 \quad 0 \quad \operatorname{Same} \quad 1 \quad \sqrt{\sum_{i=1}^{n} A_{i}^{2} \cdot \sum_{i=1}^{n} B_{i}^{2}},$$

$$0 \quad 90/270 \quad \operatorname{Orthogonal}$$

$$-1 \quad 180 \quad \operatorname{Opposite} \quad \sqrt{\theta} \quad \sqrt{\psi(q)} \quad \sqrt{\psi(d_{2})},$$

$$\psi(d_{3}) \quad \psi(d_{3}) \quad \psi(d_{3})$$

0

RICH

1

Vector Space Models for Shakespeare



Vector Space Models for Words



Vector Space Models for Queries in Corpus



Vector Space Models for Words



Vector Space Models for Quantitative Models



Figure 11: 3D plot - angle 2