

WORD SIMILARITY

David Kauchak
CS159 Fall 2014

Admin

Assignment 4

Quiz #2 Thursday

- ▣ Same rules as quiz #1
 - First 30 minutes of class
 - Open book and notes

Assignment 5 out on Thursday

Quiz #2

Topics

- ▣ Linguistics 101
- ▣ Parsing
 - Grammars, CFGs, PCFGs
 - Top-down vs. bottom-up
 - CKY algorithm
 - Grammar learning
 - Evaluation
 - Improved models
- ▣ Text similarity
 - Will also be covered on Quiz #3, though

Text Similarity

A common question in NLP is how similar are texts

score: $\text{sim}(\text{document}_1, \text{document}_2) = ?$

rank: $\text{rank}(\text{document}_1, \text{document}_2, \text{document}_3) = ?$

Bag of words representation

For now, let's ignore word order:

Obama said banana repeatedly last week on tv, "banana, banana, banana"

(4, 1, 1, 0, 0, 1, 0, 0, ...)

banana	4
obama	1
said	1
california	0
across	0
tv	1
wrong	0
capital	0

Frequency of word occurrence

"Bag of words representation": multi-dimensional vector, one dimension per word in our vocabulary

Vector based word

A

a ₁ : When	1
a ₂ : the	2
a ₃ : defendant	1
a ₄ : and	1
a ₅ : courthouse	0
...	

B

b ₁ : When	1
b ₂ : the	2
b ₃ : defendant	1
b ₄ : and	0
b ₅ : courthouse	1
...	

Multi-dimensional vectors, one dimension per word in our vocabulary

How do we calculate the similarity based on these vectors?

Normalized distance measures

Cosine

$$sim_{cos}(A,B) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

L2

$$dist_{L2}(A,B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

L1

$$dist_{L1}(A,B) = \sum_{i=1}^n |a_i - b_i|$$

a' and b' are length normalized versions of the vectors

Our problems

So far...

- word order
- length
- synonym
- spelling mistakes
- word importance
- word frequency

Word importance

Include a weight for each word/feature

A

a_1 : When	1	w_1
a_2 : the	2	w_2
a_3 : defendant	1	w_3
a_4 : and	1	w_4
a_5 : courthouse	0	w_5
...

B

b_1 : When	1	w_1
b_2 : the	2	w_2
b_3 : defendant	1	w_3
b_4 : and	0	w_4
b_5 : courthouse	1	w_5
...

Distance + weights

We can incorporate the weights into the distances

Think of it as either (both work out the same):

- preprocessing the vectors by multiplying each dimension by the weight
- incorporating it directly into the similarity measure

$$sim_{cos}(A,B) = A \cdot B = \sum_{i=1}^n a_i b_i = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

$$= \frac{\sum_{i=1}^n w_i a_i w_i b_i}{\sqrt{\sum_{i=1}^n (w_i a_i)^2} \sqrt{\sum_{i=1}^n (w_i b_i)^2}} \quad \text{with weights}$$

Document vs. overall frequency

The overall frequency of a word is the number of occurrences in a dataset, counting multiple occurrences

Example:

Word	Overall frequency	Document frequency
insurance	10440	3997
try	10422	8760

Which word is a more informative (and should get a higher weight)?

Document frequency

Word	Collection frequency	Document frequency
insurance	10440	3997
try	10422	8760

Document frequency is often related to word importance, but we want an actual weight. Problems?

$$sim_{cos}(A,B) = A \cdot B = \frac{\sum_{i=1}^n w_i a_i w_i b_i}{\sqrt{\sum_{i=1}^n (w_i a_i)^2} \sqrt{\sum_{i=1}^n (w_i b_i)^2}}$$

From document frequency to weight

Word	Collection frequency	Document frequency
insurance	10440	3997
try	10422	8760

weight and document frequency are **inversely** related
 □ higher document frequency should have lower weight and vice versa

document frequency is unbounded

document frequency will change depending on the size of the data set (i.e. the number of documents)

Inverse document frequency

$$idf_w = \log \frac{N}{df_w}$$

← # of documents in dataset
 ← document frequency of w

IDF is inversely correlated with DF
 □ higher DF results in lower IDF

N incorporates a dataset dependent normalizer

log dampens the overall weight

IDF example, suppose N=1 million

word	df _w	idf _w
calpurnia	1	
animal	100	
sunday	1,000	
fly	10,000	
under	100,000	
the	1,000,000	

What are the IDFs assuming log base 10?

$$idf_w = \log \frac{N}{df_w}$$

IDF example, suppose N=1 million

word	df _w	idf _w
calpurnia	1	6
animal	100	4
sunday	1,000	3
fly	10,000	2
under	100,000	1
the	1,000,000	0

There is one idf value/weight for each word

$$idf_w = \log \frac{N}{df_w}$$

IDF example, suppose $N=1$ million

word	df_w	idf_w
calpurnia	1	
animal	100	
sunday	1,000	
fly	10,000	
under	100,000	
the	1,000,000	

What if we didn't use the log to dampen the weighting?

$$idf_w = \log \frac{N}{df_w}$$

IDF example, suppose $N=1$ million

word	df_w	idf_w
calpurnia	1	1,000,000
animal	100	10,000
sunday	1,000	1,000
fly	10,000	100
under	100,000	10
the	1,000,000	1

Tends to overweight rare words!

TF-IDF

One of the most common weighting schemes

TF = term frequency

IDF = inverse document frequency

$$a'_i = a_i \times \log N / df_i$$

TF
IDF (word importance weight)

We can then use this with any of our similarity measures!

Stoplists: extreme weighting

Some words like 'a' and 'the' will occur in almost every document

- IDF will be 0 for any word that occurs in all documents
- For words that occur in almost all of the documents, they will be nearly 0

A **stoplist** is a list of words that should **not** be considered (in this case, similarity calculations)

- Sometimes this is the n most frequent words
- Often, it's a list of a few hundred words manually created

Stoptlist

I	all-over	around	beneath	due	go
a	almost	as	beside	durin	goddamn
aboard	along	aside	besides	during	goody
about	alongside	astride	between	each	gosh
above	alho	at	between	eh	half
across	although	atop	beyond	either	have
after	amid	avec	bi	en	he
afterwards	amidst	away	both	every	hell
against	among	back	but	ever	her
agin	amongst	be	by	everyone	herself
ago	on	because	ca.	everything	hey
agreed-upon	and	before	de	except	him
ah	another	beforehand	des	far	himself
alias	any	behind	despite	fer	his
alib	anyone	belhnde	do	for	ho
all	anything	below	down	from	how

If most of these end up with low weights anyway, why use a stoptlist?

Stoptlists

Two main benefits

- More fine grained control: some words may not be frequent, but may not have any content value (alas, teh, gosh)
- Often does contain many frequent words, which can drastically reduce our storage and computation

Any downsides to using a stoptlist?

- For some applications, some stop words may be important

Our problems

Which of these have we addressed?

- word order
- length
- synonym
- spelling mistakes
- word importance
- word frequency

A model of word similarity!

Word overlap problems

A: When the defendant and his lawyer walked into the court, some of the victim supporters turned their backs to him.

B: When the defendant walked into the courthouse with his attorney, the crowd truned their backs on him.

Word similarity

How similar are two words?

score: $\text{sim}(w_1, w_2) = ?$ rank: $w \ ?$ w_1
 w_2
 w_3

list: w_1 and w_2 are synonyms

applications?

Word similarity applications

- General text similarity
- Thesaurus generation
- Automatic evaluation
- Text-to-text
 - ▣ paraphrasing
 - ▣ summarization
 - ▣ machine translation
- information retrieval (search)

Word similarity

How similar are two words?

score: $\text{sim}(w_1, w_2) = ?$ rank: $w \ ?$ w_1
 w_2
 w_3

list: w_1 and w_2 are synonyms

ideas? useful
resources?

Word similarity

Four categories of approaches (maybe more)

- ▣ Character-based
 - turned vs. truned
 - cognates (night, nacht, nicht, natt, nat, noc, noch)
- ▣ Semantic web-based (e.g. WordNet)
- ▣ Dictionary-based
- ▣ Distributional similarity-based
 - similar words occur in similar contexts

Character-based similarity

$$\text{sim}(\textit{turned}, \textit{truned}) = ?$$

How might we do this using only the words (i.e. no outside resources?)

Edit distance (Levenshtein distance)

The edit distance between w_1 and w_2 is the minimum number of operations to transform w_1 into w_2

Operations:

- insertion
- deletion
- substitution

$$\text{EDIT}(\textit{turned}, \textit{truned}) = ?$$

$$\text{EDIT}(\textit{computer}, \textit{commuter}) = ?$$

$$\text{EDIT}(\textit{banana}, \textit{apple}) = ?$$

$$\text{EDIT}(\textit{wombat}, \textit{worcester}) = ?$$

Edit distance

$$\text{EDIT}(\textit{turned}, \textit{truned}) = 2$$

- delete u
- insert u

$$\text{EDIT}(\textit{computer}, \textit{commuter}) = 1$$

- replace p with m

$$\text{EDIT}(\textit{banana}, \textit{apple}) = 5$$

- delete b
- replace n with p
- replace a with l
- replace n with l
- replace a with e

$$\text{EDIT}(\textit{wombat}, \textit{worcester}) = 6$$

Better edit distance

Are all operations equally likely?

- No

Improvement, give different weights to different operations

- replacing a for e is more likely than z for y

Ideas for weightings?

- Learn from actual data (known typos, known similar words)
- Intuitions: phonetics
- Intuitions: keyboard configuration

Vector character-based word similarity

$\text{sim}(\text{turned}, \text{truned}) = ?$

Any way to leverage our vector-based similarity approaches from last time?

Vector character-based word similarity

$\text{sim}(\text{turned}, \text{truned}) = ?$

Generate a feature vector based on the characters (or could also use the set based measures at the character level)

a:	0	a:	0
b:	0	b:	0
c:	0	c:	0
d:	1	d:	1
e:	1	e:	1
f:	0	f:	0
g:	0	g:	0
...		...	

problems?

Vector character-based word similarity

$\text{sim}(\text{restful}, \text{fluster}) = ?$

Character level loses a lot of information

a:	0	a:	0
b:	0	b:	0
c:	0	c:	0
d:	1	d:	1
e:	1	e:	1
f:	0	f:	0
g:	0	g:	0
...		...	

ideas?

Vector character-based word similarity

$\text{sim}(\text{restful}, \text{fluster}) = ?$

Use character bigrams or even trigrams

aa:	0	aa:	0
ab:	0	ab:	0
ac:	0	ac:	0
...		...	
es:	1	er:	1
...		...	
fu:	1	fl:	1
...		...	
re:	1	lu:	1
...		...	

Word similarity

Four general categories

- Character-based
 - turned vs. truned
 - cognates (night, nacht, nicht, natt, nat, noc, noch)
- Semantic web-based (e.g. WordNet)
- Dictionary-based
- Distributional similarity-based
 - similar words occur in similar contexts

WordNet

Lexical database for English

- 155,287 words
- 206,941 word senses
- 117,659 synsets (synonym sets)
- ~400K relations between senses
- Parts of speech: nouns, verbs, adjectives, adverbs

Word graph, with word senses as nodes and edges as relationships

Psycholinguistics

- WN attempts to model human lexical memory
- Design based on psychological testing

Created by researchers at Princeton

- <http://wordnet.princeton.edu/>

Lots of programmatic interfaces

WordNet relations

- synonym
- antonym
- hypernyms
- hyponyms
- holonym
- meronym
- troponym
- entailment
- (and a few others)

WordNet relations

synonym – X and Y have similar meaning

antonym – X and Y have opposite meanings

hypernyms – subclass

- beagle is a hypernym of dog

hyponyms – superclass

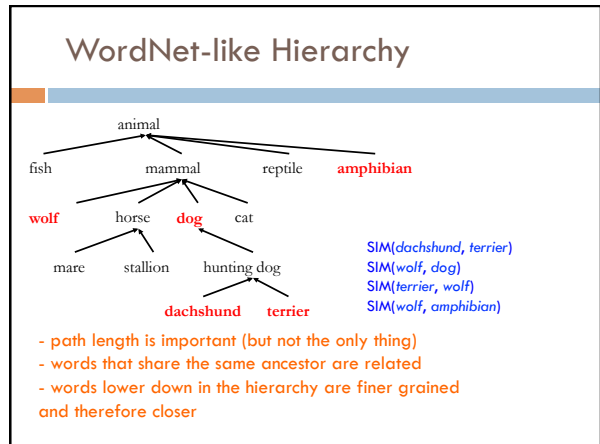
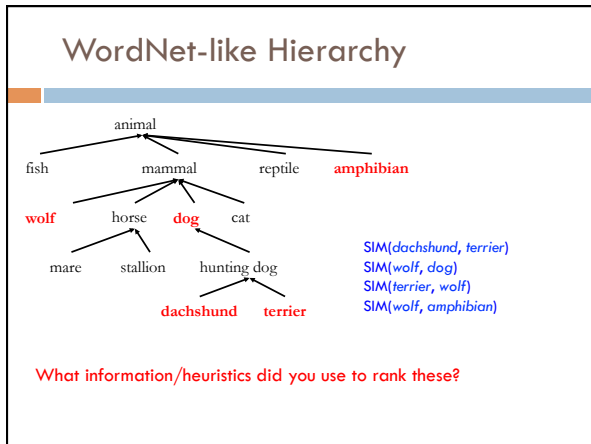
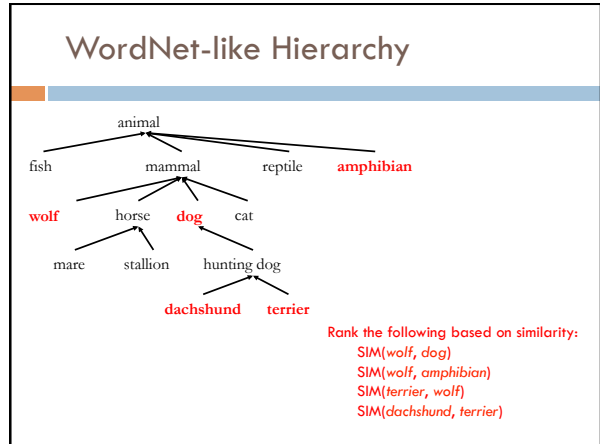
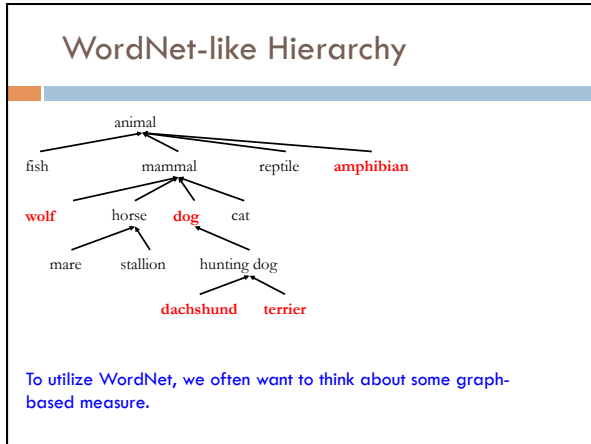
- dog is a hyponym of beagle

holonym – contains part

- car is a holonym of wheel

meronym – part of

- wheel is a meronym of car



WordNet similarity measures

path length doesn't work very well

Some ideas:

- path length scaled by the depth (Leacock and Chodorow, 1998)

With a little cheating:

- Measure the "information content" of a word using a corpus: how specific is a word?
 - words higher up tend to have less information content
 - more frequent words (and ancestors of more frequent words) tend to have less information content

WordNet similarity measures

Utilizing information content:

- information content of the lowest common parent (Resnik, 1995)
- information content of the words minus information content of the lowest common parent (Jiang and Conrath, 1997)
- information content of the lowest common parent divided by the information content of the words (Lin, 1998)