

PI

- <http://www.youtube.com/watch?v=jG7vhMMXagQ>

WORD SIMILARITY

David Kauchak
CS159 Spring 2011

Class presentations

- IR (3/30)
 - Article 1: Scott and Maksym
 - Article 2: Devin and Dandre
- MT (4/11)
 - Article 1: Jonny, Chysanthia and Daniel M.
 - Article 2: Eric and Benson
- IE (4/18)
 - Article 1: Kathryn and Audrey
 - Article 2: Josh and Michael
- QA (4/25)
 - Article 1: Dustin and Brennen
 - Article 2: Sam and Martin
- Summ (4/27???)
 - Article 1: Andres and Camille
 - Article 2: Jeremy and Dan F.

Admin

- Assignment 5 posted, due next Friday (4/1) at 6pm
 - can turn in by Sunday at 6pm
- Class schedule

Final project

- Read the entire handout
- Groups of 2-3 people
 - e-mail me asap if you're looking for a group
- research-oriented project
 - must involve some evaluation!
 - must be related to NLP
- Schedule
 - Monday, 4/4 project proposal
 - 4/15 status report 1
 - 4/27 status report 2
 - 5/2, 5/4 presentations
 - 5/4 writeup
- There are lots of resources out there that you can leverage

Final project ideas

- pick a text classification task
 - evaluate different machine learning methods
 - implement a machine learning method
 - analyze different feature categories
- n-gram language modeling
 - implement and compare other smoothing techniques
 - implement alternative models
- parsing
 - PCFG-based language modeling
 - lexicalized PCFG (with smoothing)
 - true n-best list generation
 - parse output reranking
 - implement another parsing approach and compare
 - parsing non-traditional domains (e.g. twitter)
- EM
 - word-alignment for text-to-text translation
 - grammar induction

Final project ideas

- spelling correction
- part of speech tagger
- text chunker
- dialogue generation
- pronoun resolution
- compare word similarity measures (more than the ones we're looking at for assign. 5)
- word sense disambiguation
- machine translation
 - compare sentence alignment techniques
- information retrieval
- information extraction
- question answering
- summarization
- speech recognition

Text Similarity

- A common question in NLP is how similar are texts

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

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

Text similarity recapped

- Set based – easy and efficient to calculate
 - word overlap
 - Jaccard
 - Dice
- Vector based
 - create a feature vector based on word occurrences (or other features)
 - Can use any distance measure
 - L1 (Manhattan)
 - L2 (Euclidean)
 - Cosine
 - Normalize the length
 - Feature/dimension weighting
 - inverse document frequency (IDF)

Stoptlists: 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 document
 - 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

Stoplist

I	all-over	around	beneath	due	go
a	almost	as	beside	durin	goddamn
aboard	along	aside	besides	during	goody
about	alongside	astride	between	each	goth
above	alto	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	an	because	ca.	everything	hey
agreed-upon	and	before	de	except	him
ah	another	beforehand	des	far	himself
alas	any	behind	despite	fer	his
albeit	anyone	behynde	do	for	ho
all	anything	below	down	from	how

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

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 stoplist?
 - 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

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)

Word similarity: Exercise

- How could you calculate word similarity if your only resource was:
 1. the words themselves
 2. WordNet
 3. a dictionary
 4. a corpus

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

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
- EDIT(turned, truned) = ?
 EDIT(computer, commuter) = ?
 EDIT(banana, apple) = ?
 EDIT(wombat, worcester) = ?

Edit distance

- EDIT(turned, truned) = 2
 - delete u
 - insert u
- EDIT(computer, commuter) = 1
 - replace p with m
- EDIT(banana, apple) = 5
 - delete b
 - replace n with p
 - replace a with p
 - replace n with l
 - replace a with e
- EDIT(wombat, 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

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

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

Vector character-based word similarity

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

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

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

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

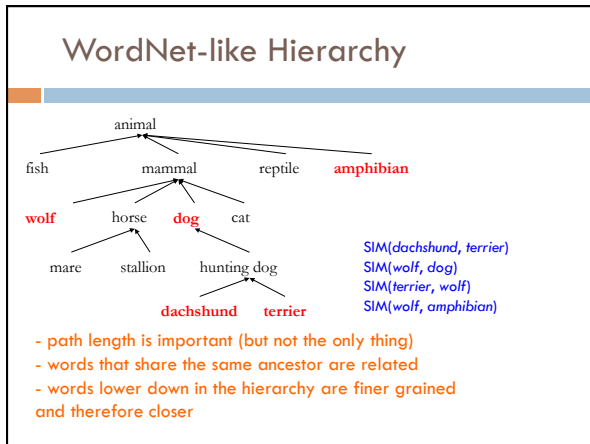
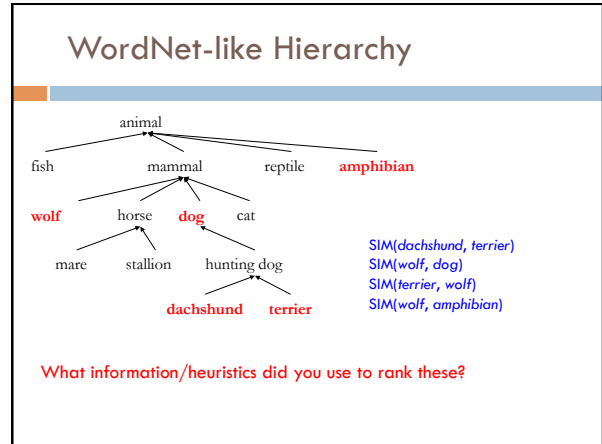
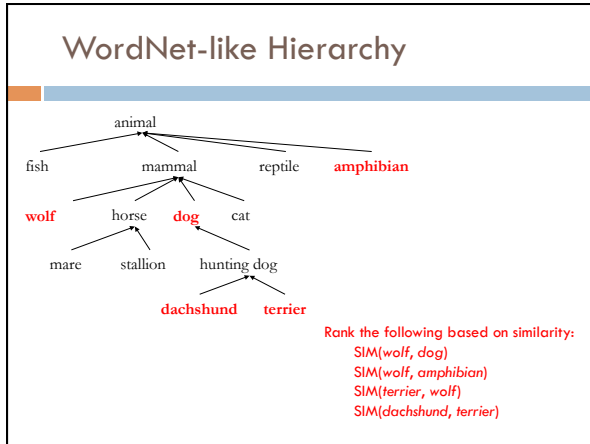
- Four general categories
 - Character-based
 - turned vs. truned
 - cognates (night, nacht, nicht, natt, nat, noc, noch)
 - Semantic web-based (e.g. WordNet)
 - Dictionary-based
 - Distributional similarity-based
 - similar words occur in similar contexts

WordNet-like Hierarchy

```

graph TD
    animal --> fish
    animal --> mammal
    animal --> reptile
    animal --> amphibian
    mammal --> wolf
    mammal --> horse
    mammal --> dog
    mammal --> cat
    horse --> mare
    horse --> stallion
    dog --> hunting_dog
    dog --> dachshund
    dog --> terrier
    
```

To utilize WordNet, we often want to think about some graph-based measure.



- ### 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:
 - ▣ utilize the probability of a word based on the corpus frequency counts of the word and all children of that word (-log of this is the information content)
 - 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)

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

Dictionary-based similarity

Word

armadillo

beagle

dog

Dictionary blurb

a large, nocturnal, burrowing mammal, *Oryzomys latipes*, of central and southern Africa, feeding on ants and termites and having a long, extensible tongue, strong claws, and long ears.

One of a breed of small hounds having long ears, short legs, and a usually black, tan, and white coat.

Any carnivore of the family Canidae, having prominent canine teeth and, in the wild state, a long and slender muzzle, a deep-chested muscular body, a bushy tail, and large, erect ears. Compare canid.

Dictionary-based similarity

Utilize our text similarity measures

$\text{sim}(\text{dog}, \text{beagle}) =$

$\text{sim}(\text{One of a breed of small hounds having long ears, short legs, and a usually black, tan, and white coat.},$

$\text{Any carnivore of the family Canidae, having prominent canine teeth and, in the wild state, a long and slender muzzle, a deep-chested muscular body, a bushy tail, and large, erect ears. Compare canid.})$

Dictionary-based similarity

- noun**
1. a domesticated canid, *Canis familiaris*, bred in many varieties.
 2. any carnivore of the dogfamily *Canidae*, having prominent canine teeth and, in the wild state, a long and slender muzzle, a deep-chested muscular body, a bushy tail, and large, erect ears. Compare *catfel*.
 3. the male of such an animal.
 4. any of various animals resembling a dog.
 5. a despicable man or youth.
 6. *Informal* - a fellow in general: a lucky dog.
 7. dogs; slang - food.
 8. Slang -
 - a. something worthless or of extremely poor quality: That used car *dog* bought is a dog.
 - b. an utter failure: *Top*: *Cricket* say his new play is a dog.
 9. Slang - an ugly, boring, or crude person.
 10. Slang - *dog-dog*.
 11. [*initial capital letter*] Astronomy - either of two constellations, *Canis Major* or *Canis Minor*.
 12. *Adjective* -
 - any of various mechanical devices, as for gripping or holding something.
 - a projection on a moving part for moving steadily or for tripping another part with which it engages.
 13. Also called **grasper**, **ripper**, *Metallworking* - a device on a drawbench for drawing the work through the die.
 14. a crane binding together two timbers.
 15. an iron bar driven into a stone or timber to provide a means of lifting it.
 16. an andiron; firedog.
 17. *Meteorology* - a sundog or fogdog.
 18. a word formerly used in communications to represent the letter D.

What about words that have multiple senses/parts of speech?

Dictionary-based similarity

- noun**
1. a domesticated canid, *Canis familiaris*, bred in many varieties.
 2. any carnivore of the dogfamily *Canidae*, having prominent canine teeth and, in the wild state, a long and slender muzzle, a deep-chested muscular body, a bushy tail, and large, erect ears. Compare *catfel*.
 3. the male of such an animal.
 4. any of various animals resembling a dog.
 5. a despicable man or youth.
 6. *Informal* - a fellow in general: a lucky dog.
 7. dogs; slang - food.
 8. Slang -
 - something worthless or of extremely poor quality: That used car *dog* bought is a dog.
 - b. an utter failure: *Top*: *Cricket* say his new play is a dog.
 9. Slang - an ugly, boring, or crude person.
 10. Slang - *dog-dog*.
 11. [*initial capital letter*] Astronomy - either of two constellations, *Canis Major* or *Canis Minor*.
 12. *Adjective* -
 - any of various mechanical devices, as for gripping or holding something.
 - a projection on a moving part for moving steadily or for tripping another part with which it engages.
 13. Also called **grasper**, **ripper**, *Metallworking* - a device on a drawbench for drawing the work through the die.
 14. a crane binding together two timbers.
 15. an iron bar driven into a stone or timber to provide a means of lifting it.
 16. an andiron; firedog.
 17. *Meteorology* - a sundog or fogdog.
 18. a word formerly used in communications to represent the letter D.

1. part of speech tagging
2. word sense disambiguation
3. most frequent sense
4. average similarity between all senses
5. max similarity between all senses
6. sum of similarity between all senses

Dictionary + WordNet

- WordNet also includes a “gloss” similar to a dictionary definition
- Other variants include the overlap of the word senses as well as those word senses that are related (e.g. hypernym, hyponym, etc.)
 - incorporates some of the path information as well
 - Banerjee and Pedersen, 2003

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

Corpus-based approaches

Word

cardvark

beagle

dog

ANY blurb



Ideas?

Corpus-based

The **Beagle** is a breed of small to medium-sized dog. A member of the Hound Group, it is similar in appearance to the Foxhound but smaller, with shorter leg

Beagles are intelligent, and are popular as pets because of their size, even temper, and lack of inherited health problems.

Dogs of similar size and purpose to the modern **Beagle** can be traced in Ancient Greece[2] back to around the 5th century B.C.

From medieval times, **beagle** was used as a generic description for the smaller hounds, though these dogs differed considerably from the modern breed.

In the 1840s, a standard **Beagle** type was beginning to develop: the distinction between the North Country Beagle and Southern

Corpus-based: feature extraction

The **Beagle** is a breed of small to medium-sized dog. A member of the Hound Group, it is similar in appearance to the Foxhound but smaller, with shorter leg

- We'd like to utilize or vector-based approach
- How could we create a vector from these occurrences?
 - collect word counts from all documents with the word in it
 - collect word counts from all sentences with the word in it
 - collect all word counts from all words within X words of the word
 - collect all words counts from words in specific relationship: subject-object, etc.

Word-context co-occurrence vectors

The **Beagle** is a **breed** of small to medium-sized dog. A member of the Hound Group, it is similar in appearance to the Foxhound but smaller, with shorter leg

Beagles are intelligent, and are popular as pets because of their size, even temper, and lack of inherited health problems.

Dogs of similar size and purpose **to the modern Beagle can be traced** in Ancient Greece[2] back to around the 5th century B.C.

From medieval times, **beagle was used** as a generic description for the smaller hounds, though these dogs differed considerably from the modern breed.

In the 1840s, a standard **Beagle type was beginning** to develop: the distinction between the North Country Beagle and Southern

Word-context co-occurrence vectors

The **Beagle** is a breed
Beagles are intelligent, and
to the modern **Beagle** can be traced
From medieval times, **beagle** was used as
1840s, a standard **Beagle** type was beginning

the: 2
is: 1
a: 2
breed: 1
are: 1
intelligent: 1
and: 1
to: 1
modern: 1
...

Often do some preprocessing like lowercasing
and removing stop words

Corpus-based similarity

$\text{sim}(\text{dog}, \text{beagle}) =$

$\text{sim}(\text{context_vector}(\text{dog}), \text{context_vector}(\text{beagle}))$

the:	5	the:	2
is:	1	is:	1
a:	4	a:	2
breeds:	2	breed:	1
are:	1	are:	1
intelligent:	5	intelligent:	1
...		and:	1
		to:	1
		modern:	1
		...	

Another feature weighting

- TFIDF weighting takes into account the general importance of a feature
- For distributional similarity, we have the feature (f_i), but we also have the word itself (w) that we can use for information
- This is different from traditional text similarity where we only have f_i
- Another feature weighting idea
 - ▣ don't use raw co-occurrence
 - ▣ count *how likely* feature f_i and word w are to occur together
 - incorporates co-occurrence
 - but also incorporates how often w and f_i occur in other instances

Mutual information

- A bit more probability ☺

$$I(X,Y) = \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

When will this be high and when will this be low?

Mutual information

- A bit more probability ☺

$$I(X,Y) = \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

- if x and y are independent (i.e. one occurring doesn't impact the other occurring) $p(x,y) = p(x)p(y)$ and the sum is 0

- if they're dependent then $p(x,y) = p(x)p(y|x) = p(y)p(x|y)$ then we get $p(y|x)/p(y)$ (i.e. how much more likely are we to see y given x has a particular value) or vice versa $p(x|y)/p(x)$

Pointwise mutual information

Mutual information

$$I(X,Y) = \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

How related are two variables (i.e. over all possible values/events)

Pointwise mutual information

$$PMI(x,y) = \log \frac{p(x,y)}{p(x)p(y)}$$

How related are two events/values

PMI weighting

- Mutual information is often used for features selection in many problem areas
- PMI weighting weights co-occurrences based on their correlation (i.e. high PMI)

context_vector(beagle)

the: 2
is: 1
a: 2
breed: 1
are: 1
intelligent: 1
and: 1
to: 1
modern: 1
...

$$\log \frac{p(\text{beagle,the})}{p(\text{beagle})p(\text{the})}$$

this would likely be lower

$$\log \frac{p(\text{beagle,breed})}{p(\text{beagle})p(\text{breed})}$$

this would likely be higher

Web-based similarity

beagle



beagle

Beagle - Wikipedia, the free encyclopedia
The Beagle is a breed of small to medium-sized dog. A member of the Hound Group, it is similar in appearance to the Foxhound but smaller, with shorter legs ...
History · Description · Variations · Temperament
en.wikipedia.org/wiki/Beagle - Cached · Similar

Beagle (software) - Wikipedia, the free encyclopedia
Beagle is a search system for Linux and other such modern UNIX-like systems ...
en.wikipedia.org/wiki/Beagle_(software) - Cached · Similar

Beagle Information and Pictures, Beagles
All about the Beagle, info, pictures, breeders, rescues, care, temperament, health, puppy pictures and much more.
www.dogbreedinfo.com/beagle.htm - Cached · Similar

Beagles & Buddies: PET ADOPTION, BEAGLE SHELTER, DOG RESCUE
Rescue shelter for Beagles, as well as other small dogs, from pounds, humane societies & off the street. We have a no-kill policy, our rescue facility keeps ...
www.beaglesandbuddies.com/ - Cached · Similar

How can we make a document/blurb from this?

Web-based similarity

The image shows a screenshot of search results for the term "Beagle". The results are as follows:

- Beagle - Wikipedia, the free encyclopedia**
The Beagle is a breed of small to medium-sized dog. A member of the Hound Group, it is similar in appearance to the Foxhound but smaller, with shorter legs...
en.wikipedia.org/wiki/Beagle - Cached - Similar
- Beagle (software) - Wikipedia, the free encyclopedia**
Beagle is a search system for Linux and other such modern Unix-like systems ...
en.wikipedia.org/wiki/Beagle_(software) - Cached - Similar
- Beagle Information and Pictures, Beagles**
All about the Beagle, info, pictures, breeders, rescues, care, temperament, health, puppy pictures and much more.
www.dogbreedinfo.com/beagle.htm - Cached - Similar
- Beagles & Buddies: PET ADOPTION, BEAGLE SHELTER, DOG RESCUE**
Rescue shelter for Beagles, as well as other small dogs, from pounds, humane societies & off the street. We have a no-kill policy, our rescue facility keeps ...
www.beagleandbuddies.com/ - Cached - Similar

Annotations on the right side of the screenshot:

- An arrow points from the first two results to the text: "Concatenate the snippets for the top N results".
- An arrow points from the last two results to the text: "Concatenate the web page text for the top N results".