

DISTRIBUTIONAL WORD SIMILARITY

David Kauchak
CS159 Fall 2014

Admin

Assignment 5 out

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

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

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

cardvark

beagle

dog

Dictionary blurb

a large, nocturnal, burrowing mammal, *Orycteropus afer*, of central and southern Africa, feeding on ants and termites and having a long, extensile 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 canid.
 3. the male of such an animal.
 4. any of various *Atalapha* resembling a dog.
 5. a despicable man or youth.
 6. *Informal* - a fellow in general: a lucky dog.
 7. dog; Stamp - *See*.
 8. Stamp -
 - a. something worthless or of extremely poor quality: That new car *dog* bought is a dog.
 - b. an utter failure: *Bob*: *Chris* say his new play is a dog.
 9. Stamp - an ugly, boring, or crude person.
 10. Stamp - *See* *DOG*.
 11. (initial capital letter) Astronomy - either of two constellations, *Canis Major* or *Canis Minor*.
 12. Machinery -
 - a. any of various mechanical devices, as for gripping or holding something.
 - b. a projection on a moving part for moving steadily or for tripping another part with which it engages.
 13. Also called *gripbar*, *ripper*. Metallurgy - a device on a drawbench for drawing the work through the die.
 14. a clamp binding together two timbers.
 15. an iron bar driven into a stone or timber to provide a means of lifting it.
 16. an andree, freedog.
 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-throated muscular body, a bushy tail, and large, erect ears. Compare *catfish*.
 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 to a dog.
 - b. an utter failure; Stop. **Cricket** day his new play is a dog.
 9. Slang - an ugly, boring, or crude person.
 10. Slang - **hot dog**.
 11. (initial capital letter) Astronomy - either of two constellations, *Canis Major* or *Canis Minor*.
 12. **adjutery** -
 - any of various mechanical devices, as for gripping or holding something.
 - a projection on a moving part for moving steadily or for bringing another part with which it engages.
 13. Also called **grasper**, **ripper**, **metalworking** - 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

ANY blurb with the word

aardvark



beagle



Ideas?

dog



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 BC.

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 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 BC.

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	the:	2
	is:	1
Beagles are intelligent, and	a:	2
	breed:	1
to the modern Beagle can be traced	are:	1
	intelligent:	1
From medieval times, beagle was used as	and:	1
	to:	1
1840s, a standard Beagle type was beginning	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
		...	

Web-based similarity

The image shows a Google search interface with the search term 'beagle' entered in the search box. Below the search box are buttons for 'Google Search' and 'I'm Feeling Lucky'. The text 'Ideas?' is visible below the search box.

Web-based similarity

The diagram illustrates the process of web-based similarity. It starts with the word 'beagle' on the left. An arrow points to a Google search box containing 'beagle'. Another arrow points down to a list of search results, including Wikipedia entries and other web pages.

Web-based similarity

The diagram shows search results for 'beagle' with annotations. Blue arrows point from specific search results to text boxes on the right. One arrow points to the text 'Concatenate the snippets for the top N results', and another points to 'Concatenate the web page text for the top N results'.

Another feature weighting

TF- IDF 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

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

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

Feature weighting ideas given this additional information?

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

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

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

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

Does IDF capture this?

Not really. IDF only accounts for f_i regardless of w

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) then:

$$p(x,y) =$$

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) then:

$$p(x,y) = p(x)p(y)$$

What does this do to the sum?

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 they are dependent then:

$$p(x,y) = p(x)p(y|x) = p(y)p(x|y)$$



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

Mutual information

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

What is this asking?

When is this high?

How much more likely are we to see y given x has a particular value!

Point-wise 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)

Point-wise mutual information

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

How related are two particular events/values

PMI weighting

Mutual information is often used for feature 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
...

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

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

How do we calculate these?