

# WORD SIMILARITY

David Kauchak  
CS159 Spring 2019

## Admin

### Assignment 4

### Quiz #2 Wednesday

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

Assignment 5 out soon

## 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 obama said  
california across tv  
wrong capital

Frequency of word occurrence

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

## Vector based word

A

a1: When 1  
a2: the 2  
a3: defendant 1  
a4: and 1  
a5: courthouse 0  
...

B

b1: When 1  
b2: the 2  
b3: defendant 1  
b4: and 0  
b5: courthouse 1  
...

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

## TF-IDF

One of the most common weighting schemes

TF = term frequency

IDF = inverse document frequency

$$a'_i = \underbrace{a_i}_{TF} \times \underbrace{\log N / df_i}_{IDF \text{ (word importance weight)}}$$

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

## Normalized distance measures

Cosine

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

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

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:  $w \quad ?$

$w_1$

$w_2$

$w_3$

applications?

list:  $w_1$  and  $w_2$  are synonyms

## 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:  $w \ ?$

$w_1$  ideas? useful  
 $w_2$  resources?  
 $w_3$

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

## 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  
...

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)

problems?

## Vector character-based word similarity

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

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

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

Character level loses a lot of information

ideas?

## Vector character-based word similarity

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

aa: 0  
ab: 0  
ac: 0  
...  
es: 1  
...  
fu: 1  
...  
re: 1  
...

aa: 0  
ab: 0  
ac: 0  
...  
er: 1  
...  
fl: 1  
...  
lu: 1  
...

Use character bigrams or even trigrams

## 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: dog

### Noun

- **S: (n) dog, domestic dog, Canis familiaris** (a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) *"the dog barked all night"*
- **S: (n) frump, dog** (a dull unattractive unpleasant girl or woman) *"she got a reputation as a frump"; "she's a real dog"*
- **S: (n) dog** (informal term for a man) *"you lucky dog"*
- **S: (n) cad, boaster, blackguard, dog, hoard, heel** (someone who is morally reprehensible) *"you dirty dog"*
- **S: (n) frank, frankfurter, hotdog, hot dog, dog, wiener, wienerwurst, weenie** (a smooth-textured sausage of minced beef or pork usually smoked; often served on a bread roll)
- **S: (n) pawl, detent, click, dog** (a hinged catch that fits into a notch of a ratchet to move a wheel forward or prevent it from moving backward)
- **S: (n) andiron, firelog, dog, dog-iron** (metal supports for logs in a fireplace) *"the andirons were too hot to touch"*

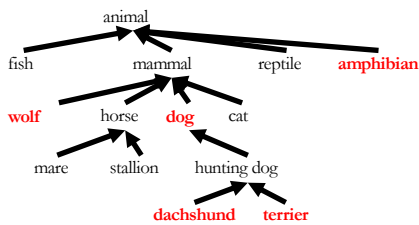
### Verb

- **S: (v) chase, chase after, trail, tail, tag, give chase, dog, go after, track** (go after with the intent to catch) *"The policeman chased the mugger down the alley"; "the dog chased the rabbit"*

## WordNet: dog

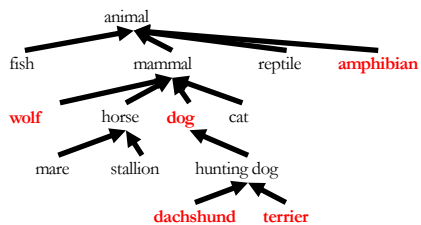
- **S: (n) dog, domestic dog, Canis familiaris** (a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) *"the dog barked all night"*
  - **direct hyponym / full hyponym**
  - **part meronym**
  - **member holonym**
  - **direct hypernym / inherited hypernym / sister term**
- **direct hyponym / full hyponym**
  - **S: (n) puppy** (a young dog)
  - **S: (n) pooch, doggie, doggy, barker, bow-wow** (informal terms for dogs)
  - **S: (n) cur, mongrel, mutt** (an inferior dog or one of mixed breed)
  - **S: (n) lapdog** (a dog small and tame enough to be held in the lap)
  - **S: (n) toy dog, toy** (any of several breeds of very small dogs kept purely as pets)
  - **S: (n) hunting dog** (a dog used in hunting game)
  - **S: (n) working dog** (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
  - **S: (n) dalmatian, coach dog, carriage dog** (a large breed having a smooth white coat with black or brown spots; originated in Dalmatia)
  - **S: (n) basenji** (small smooth-haired breed of African origin having a tightly curled tail and the inability to bark)
  - **S: (n) pug, pug-dog** (small compact smooth-coated breed of Asiatic origin having a tightly curled tail and broad flat wrinkled muzzle)

## WordNet-like Hierarchy



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

## WordNet-like Hierarchy

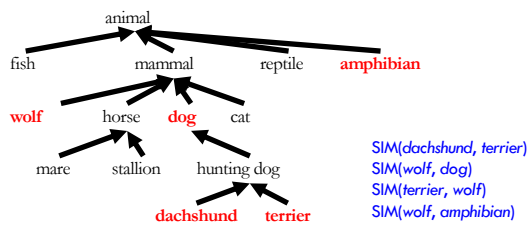


Rank the following based on similarity:

- SIM(wolf, dog)
- SIM(wolf, amphibian)
- SIM(terrier, wolf)
- SIM(dachshund, terrier)

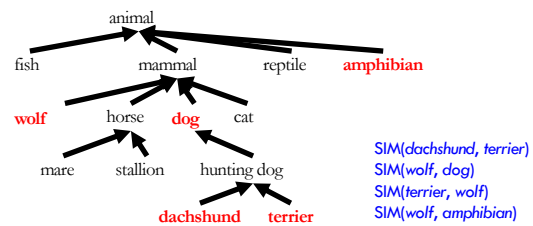


## WordNet-like Hierarchy



What information/heuristics did you use to rank these?

## WordNet-like Hierarchy



- path length is important (but not the only thing)
- words that share the same ancestor are related
- words lower down in the hierarchy are finer grained and therefore closer

## 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)

## 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

aardvark

beagle

dog

### Dictionary blurb

a large, nocturnal, burrowing mammal, *Oryzomys afer*, 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 canid.
  3. the male of such an animal.
  4. any of various domestic animals resembling a dog.
  5. a despicable man or youth.
  6. [Informal] a fellow in general; a lucky dog.
  7. dog; slang - fast.
  8. Slang
    - a. something worthless or of extremely poor quality: That used car you bought is a dog.
    - b. an utter failure; flop: Critics say his new play is a dog.
  9. Slang - an oily, boring, or crude person.
  10. Slang - hot 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 **grasper**, **ripper**. Metalworking - a device on a lathebench for drawing the work through the die.
  14. a cramp binding together two timbers.
  15. an iron bar driven into a stone or timber to provide a means of lifting it.
  16. an andiron; fire dog.
  17. Meteorology - a sundog or fog dog.
  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 *canid*.
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* - *testis*.
8. *Slang* .
  - a. something worthless or of extremely poor quality: That used garage bought is a dog.
  - b. an utter failure; flop: Critics say 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. *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 *grasper*, *ripper*, *metalworking* - a device on a drawbench for drawing the work through the die.
14. a cramp 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 our 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 <b>Beagle</b> is a breed	the:	2
	is:	1
<b>Beagles</b> are intelligent, and	a:	2
	breed:	1
to the modern <b>Beagle</b> can be traced	are:	1
	intelligent:	1
From medieval times, <b>beagle</b> was used as	and:	1
	to:	1
1840s, a standard <b>Beagle</b> 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



beagle  Advanced search Language tools

Ideas?

## Web-based similarity

beagle



beagle  Advanced search Language tools



**Beagle** - Wikipedia, the free encyclopedia <sup>↗</sup>  
 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 <sup>↗</sup>  
 Beagle is a search system for Linux and other such modern Unix-like systems ...  
en.wikipedia.org/wiki/Beagle\_(software) · Cached · Similar

⚡ Show more results from wikipedia.org

**Beagle Information and Pictures, Beagles** <sup>↗</sup>  
 All about the Beagle, info, pictures, breeders, rescue, care, temperament, health, puppy pictures and much more.  
[www.dogsofindia.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

## Web-based similarity

**Beagle** - Wikipedia, the free encyclopedia <sup>↗</sup>  
 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 <sup>↗</sup>  
 Beagle is a search system for Linux and other such modern Unix-like systems ...  
en.wikipedia.org/wiki/Beagle\_(software) · Cached · Similar

⚡ Show more results from wikipedia.org

**Beagle Information and Pictures, Beagles** <sup>↗</sup>  
 All about the Beagle, info, pictures, breeders, rescue, care, temperament, health, puppy pictures and much more.  
[www.dogsofindia.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



Concatenate the snippets for the top N results



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	$\rightarrow \log \frac{p(\text{beagle, the})}{p(\text{beagle})p(\text{the})}$	How do we calculate these?
is:	1		
a:	2		
breed:	1	$\rightarrow \log \frac{p(\text{beagle, breed})}{p(\text{beagle})p(\text{breed})}$	
are:	1		
intelligent:	1		
and:	1		
to:	1		
modern:	1		
...			