

# Natural Language Processing 1

## Lecture 4: Lexical and distributional semantics

Katia Shutova

ILLC  
University of Amsterdam

# Outline.

Dependency parsing (finishing off)

Introduction to lexical semantics

Distributional semantics

## Dependency structure

A dependency structure consists of **dependency relations**, which are **binary** and **asymmetric**.

John hit the ball

A relation consists of

- ▶ a head (H) — **hit**
- ▶ a dependent (D) — **John**
- ▶ a label identifying the relation between H and D — **Subject**

## Dependency structure

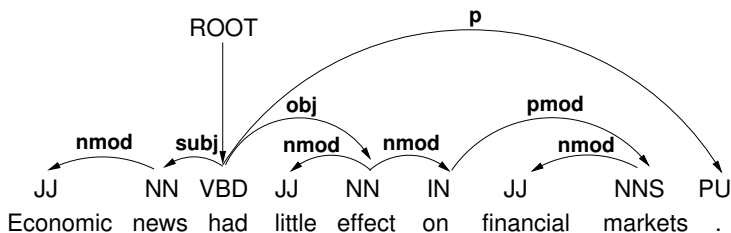
A dependency structure consists of **dependency relations**, which are **binary** and **asymmetric**.

John hit the ball

A relation consists of

- ▶ a head (H) — **hit**
- ▶ a dependent (D) — **ball**
- ▶ a label identifying the relation between H and D — **Object**

## Example dependency structure



[From Joakim Nivre, Dependency Grammar and Dependency Parsing.]

## Dependency parsing

Output a list of dependencies between words in the sentence.

John hit the ball.

(**SUBJ** head=hit dep=**John**)

(**OBJ** head=hit dep=**ball**)

(**DET** head=ball dep=the)



Why is it useful?

- ▶ dependencies provide an interface to semantics

*“Who did what to whom”*

## The cost of parsing errors...

### Incorrect dependencies

(**SUBJ** head=hit dep=**ball**)

(**OBJ** head=hit dep=**John**)

(DET head=ball dep=the)



# Outline.

Dependency parsing (finishing off)

**Introduction to lexical semantics**

Distributional semantics



# Semantics

## Compositional semantics:

- ▶ studies how meanings of phrases are constructed out of the meaning of individual words
- ▶ principle of compositionality: meaning of each whole phrase derivable from meaning of its parts
- ▶ sentence structure conveys some meaning: obtained by syntactic representation

## Lexical semantics:

- ▶ studies how the meanings of individual words can be represented and induced

## What is lexical meaning?

- ▶ recent results in psychology and cognitive neuroscience give us some clues
- ▶ but we don't have the whole picture yet
- ▶ different representations proposed, e.g.
  - ▶ formal semantic representations based on logic,
  - ▶ *or* taxonomies relating words to each other,
  - ▶ *or* distributional representations in statistical NLP
- ▶ but none of the representations gives us a complete account of lexical meaning

## How to approach lexical meaning?

- ▶ **Formal semantics**: set-theoretic approach  
e.g.,  $\text{cat}'$ : the set of all cats;  $\text{bird}'$ : the set of all birds.
- ▶ meaning postulates, e.g.

$$\forall x[\text{bachelor}'(x) \rightarrow \text{man}'(x) \wedge \text{unmarried}'(x)]$$

- ▶ Limitations, e.g. *is the Pope a bachelor?*
  - ▶ Defining concepts through enumeration of all of their features in practice is highly problematic
  - ▶ How would you define e.g. *chair, tomato, thought, democracy?* – impossible for most concepts
  - ▶ **Prototype theory** offers an alternative to set-theoretic approaches

## How to approach lexical meaning?

- ▶ **Formal semantics**: set-theoretic approach  
e.g.,  $\text{cat}'$ : the set of all cats;  $\text{bird}'$ : the set of all birds.
- ▶ meaning postulates, e.g.

$$\forall x[\text{bachelor}'(x) \rightarrow \text{man}'(x) \wedge \text{unmarried}'(x)]$$

- ▶ Limitations, e.g. *is the Pope a bachelor?*
- ▶ Defining concepts through enumeration of all of their features in practice is highly problematic
- ▶ How would you define e.g. *chair, tomato, thought, democracy?* – impossible for most concepts
- ▶ **Prototype theory** offers an alternative to set-theoretic approaches

## Prototype theory

- ▶ introduced the notion of **graded semantic categories**
- ▶ no clear boundaries
- ▶ no requirement that a property or set of properties be shared by all members
- ▶ certain members of a category are more central or **prototypical** (i.e. instantiate the prototype)

*furniture: chair is more prototypical than stool*

Eleanor Rosch 1975. *Cognitive Representation of Semantic Categories* (J Experimental Psychology)

## Prototype theory (continued)

- ▶ Categories form around prototypes; new members added on basis of resemblance to prototype
- ▶ Features/attributes generally graded
- ▶ Category membership a matter of degree
- ▶ Categories do not have clear boundaries

# Semantic relations

## Hyponymy: IS-A

*dog* is a **hyponym** of *animal*  
*animal* is a **hypernym** of *dog*

- ▶ hyponymy relationships form a **taxonomy**
- ▶ works best for concrete nouns

## Other semantic relations

**Meronymy: PART-OF** e.g., *arm* is a **meronym** of *body*, *steering wheel* is a meronym of *car*

**Synonymy** e.g., *aubergine/eggplant*.

**Antonymy** e.g., *big/little*

Also:

**Near-synonymy/similarity** e.g., *exciting/thrilling*  
e.g., *slim/slender/thin/skinny*



# WordNet

- ▶ large scale, open source resource for English
- ▶ hand-constructed
- ▶ wordnets being built for other languages
- ▶ organized into **synsets**: synonym sets (near-synonyms)
- ▶ synsets connected by semantic relations

S: (v) interpret, construe, see (make sense of; assign a meaning to) - "How do you interpret his behavior?"

S: (v) understand, read, interpret, translate (make sense of a language) "She understands French";  
"Can you read Greek?"

## Polysemy and word senses

The children **ran** to the store

If you see this man, **run!**

Service **runs** all the way to Cranbury

She is **running** a relief operation in Sudan

the story or argument **runs** as follows

Does this old car still **run** well?

Interest rates **run** from 5 to 10 percent

Who's **running** for treasurer this year?

They **ran** the tapes over and over again

These dresses **run** small

## Polysemy

- ▶ **homonymy**: unrelated word senses. *bank* (raised land) vs *bank* (financial institution)
- ▶ *bank* (financial institution) vs *bank* (in a casino): related but distinct senses.
- ▶ **regular polysemy** and sense extension
  - ▶ metaphorical senses, e.g. *swallow* [food], *swallow* [information], *swallow* [anger]
  - ▶ metonymy, e.g. he played *Bach*; he drank his *glass*.
  - ▶ zero-derivation, e.g. *tango* (N) vs *tango* (V)
- ▶ vagueness: *nurse*, *lecturer*, *driver*
- ▶ cultural stereotypes: *nurse*, *lecturer*, *driver*

No clearcut distinctions.

## Word sense disambiguation

- ▶ Needed for many applications
- ▶ relies on context, e.g. *striped bass* (the fish) vs *bass guitar*.

### Methods:

- ▶ **supervised** learning:
  - ▶ Assume a predefined set of word senses, e.g. WordNet
  - ▶ Need a large sense-tagged training corpus (difficult to construct)
- ▶ **semi-supervised** learning
  - ▶ bootstrap from a few examples
- ▶ **unsupervised** sense induction
  - ▶ e.g. cluster contexts in which a word occurs

# Outline.

Dependency parsing (finishing off)

Introduction to lexical semantics

Distributional semantics

## Distributional hypothesis

*You shall know a word by the company it keeps* (Firth)

*The meaning of a word is defined by the way it is used*  
(Wittgenstein).

it was authentic **scrumpy**, rather sharp and very strong

we could taste a famous local product — **scrumpy**

spending hours in the pub drinking **scrumpy**

Cornish **Scrumpy** Medium Dry. £19.28 - Case

## Distributional hypothesis

*You shall know a word by the company it keeps* (Firth)

*The meaning of a word is defined by the way it is used*  
(Wittgenstein).

it was authentic **scrumpy**, rather sharp and very strong

we could taste a famous local product — **scrumpy**

spending hours in the pub drinking **scrumpy**

Cornish **Scrumpy** Medium Dry. £19.28 - Case

## Distributional hypothesis

*You shall know a word by the company it keeps* (Firth)

*The meaning of a word is defined by the way it is used*  
(Wittgenstein).

it was authentic **scrumpy**, rather sharp and very strong

we could taste a famous local product — **scrumpy**

spending hours in the pub drinking **scrumpy**

Cornish **Scrumpy** Medium Dry. £19.28 - Case



## Distributional hypothesis

*You shall know a word by the company it keeps* (Firth)

*The meaning of a word is defined by the way it is used*  
(Wittgenstein).

it was authentic **scrumpy**, rather sharp and very strong

we could taste a famous local product — **scrumpy**

spending hours in the pub drinking **scrumpy**

Cornish **Scrumpy** Medium Dry. £19.28 - Case

## Distributional hypothesis

*You shall know a word by the company it keeps* (Firth)

*The meaning of a word is defined by the way it is used*  
(Wittgenstein).

it was authentic **scrumpy**, rather sharp and very strong

we could taste a famous local product — **scrumpy**

spending hours in the pub drinking **scrumpy**

Cornish **Scrumpy** Medium Dry. £19.28 - Case

# Scrumpy



## Distributional hypothesis

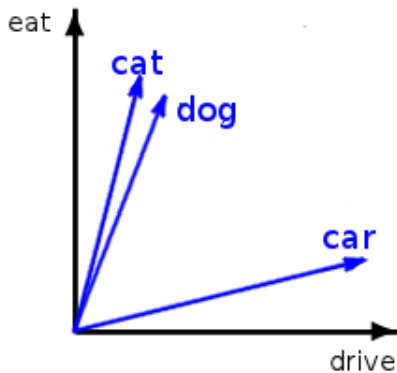
This leads to the **distributional hypothesis** about word meaning:

- ▶ the context surrounding a given word provides information about its meaning;
- ▶ words are similar if they share similar linguistic contexts;
- ▶ semantic similarity  $\approx$  distributional similarity.

## The general intuition

- ▶ **Distributions** are vectors in a multidimensional semantic space.
- ▶ The **semantic space** has dimensions which correspond to possible contexts – **features**.
- ▶ For our purposes, a distribution can be seen as a point in that space (the vector being defined with respect to the origin of that space).
- ▶ *scrumpy* [...pub 0.8, drink 0.7, strong 0.4, joke 0.2, mansion 0.02, zebra 0.1...]

# Vectors



# Feature matrix

	feature <sub>1</sub>	feature <sub>2</sub>	...	feature <sub>n</sub>
word <sub>1</sub>	$f_{1,1}$	$f_{2,1}$		$f_{n,1}$
word <sub>2</sub>	$f_{1,2}$	$f_{2,2}$		$f_{n,2}$
...				
word <sub>m</sub>	$f_{1,m}$	$f_{2,m}$		$f_{n,m}$

## The notion of context

- 1 Word windows (unfiltered):  $n$  words on either side of the lexical item.

**Example:**  $n=2$  (5 words window):

*| The prime **minister** acknowledged the |  
question.*

*minister* [ the 2, prime 1, acknowledged 1, question 0 ]



## Context

- 2 Word windows (filtered):  $n$  words on either side removing some words (e.g. function words, some very frequent content words). Stop-list or by POS-tag.

**Example:**  $n=2$  (5 words window), stop-list:

*| The prime **minister** acknowledged the |  
question.*

*minister* [ prime 1, acknowledged 1, question 0 ]

## Context

- 3 Lexeme window (filtered or unfiltered); as above but using stems.

**Example:**  $n=2$  (5 words window), stop-list:

*| The prime **minister** acknowledged the |  
question.*

*minister* [ prime 1, acknowledge 1, question 0 ]

## Context

- 4 Dependencies (directed links between heads and dependents). Context for a lexical item is the dependency structure it belongs to (various definitions).

**Example:**

*The prime **minister** acknowledged the question.*

*minister* [ prime\_a 1, acknowledge\_v 1 ]

*minister* [ prime\_a\_mod 1, acknowledge\_v\_subj 1 ]

*minister* [ prime\_a 1, acknowledge\_v+question\_n 1 ]

## Parsed vs unparsed data: examples

### **word (unparsed)**

meaning\_n  
derive\_v  
dictionary\_n  
pronounce\_v  
phrase\_n  
latin\_j  
ipa\_n  
verb\_n  
mean\_v  
hebrew\_n  
usage\_n  
literally\_r

### **word (parsed)**

or\_c+phrase\_n  
and\_c+phrase\_n  
syllable\_n+of\_p  
play\_n+on\_p  
etymology\_n+of\_p  
portmanteau\_n+of\_p  
and\_c+deed\_n  
meaning\_n+of\_p  
from\_p+language\_n  
pron\_rel\_+utter\_v  
for\_p+word\_n  
in\_p+sentence\_n

## Context weighting

- ▶ Binary model: if context  $c$  co-occurs with word  $w$ , value of vector  $\vec{w}$  for dimension  $c$  is 1, 0 otherwise.

... [a long long long **example** for a distributional semantics] model... ( $n=4$ )

... {a 1} {dog 0} {long 1} {sell 0} {semantics 1}...

- ▶ Basic frequency model: the value of vector  $\vec{w}$  for dimension  $c$  is the number of times that  $c$  co-occurs with  $w$ .

... [a long long long **example** for a distributional semantics] model... ( $n=4$ )

... {a 2} {dog 0} {long 3} {sell 0} {semantics 1}...

## Characteristic model

- ▶ Weights given to the vector components express how *characteristic* a given context is for word  $w$ .
- ▶ Pointwise Mutual Information (PMI)

$$PMI(w, c) = \log \frac{P(w, c)}{P(w)P(c)} = \log \frac{P(w)P(c|w)}{P(w)P(c)} = \log \frac{P(c|w)}{P(c)}$$

$$P(c) = \frac{f(c)}{\sum_k f(c_k)}, \quad P(c|w) = \frac{f(w, c)}{f(w)},$$

$$PMI(w, c) = \log \frac{f(w, c) \sum_k f(c_k)}{f(w) f(c)}$$

$f(w, c)$ : frequency of word  $w$  in context  $c$

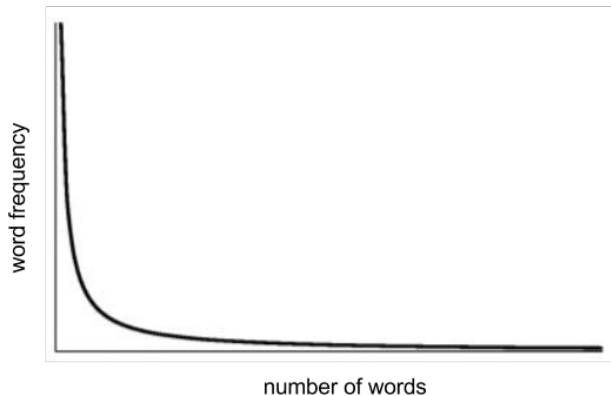
$f(w)$ : frequency of word  $w$  in all contexts

$f(c)$ : frequency of context  $c$

## What semantic space?

- ▶ Entire vocabulary.
  - ▶ + All information included – even rare contexts
  - ▶ - Inefficient (100,000s dimensions). Noisy (e.g. *002.png/thumb/right/200px/graph\_n*). **Sparse**
- ▶ Top  $n$  words with highest frequencies.
  - ▶ + More efficient (2000-10000 dimensions). Only ‘real’ words included.
  - ▶ - May miss out on infrequent but relevant contexts.

## Word frequency: Zipfian distribution





## What semantic space?

- ▶ Entire vocabulary.
  - ▶ + All information included – even rare contexts
  - ▶ - Inefficient (100,000s dimensions). Noisy (e.g. *002.png/thumb/right/200px/graph\_n*). **Sparse**.
- ▶ Top  $n$  words with highest frequencies.
  - ▶ + More efficient (2000-10000 dimensions). Only ‘real’ words included.
  - ▶ - May miss out on infrequent but relevant contexts.

## What semantic space?

- ▶ Singular Value Decomposition (SVD): the number of dimensions is reduced by exploiting redundancies in the data.
  - ▶ + Very efficient (200-500 dimensions). Captures generalisations in the data.
  - ▶ - SVD matrices are not interpretable.
- ▶ Non-negative matrix factorization (NMF)
  - ▶ Similar to SVD in spirit, but performs factorization differently

## Our reference text

### Douglas Adams, *Mostly harmless*

The major difference between a thing that might go wrong and a thing that cannot possibly go wrong is that when a thing that cannot possibly go wrong goes wrong it usually turns out to be impossible to get at or repair.

- ▶ **Example:** Produce distributions using a word window, PMI-based model

## The semantic space

### Douglas Adams, *Mostly harmless*

The major difference between a thing that might go wrong and a thing that cannot possibly go wrong is that when a thing that cannot possibly go wrong goes wrong it usually turns out to be impossible to get at or repair.

- ▶ Assume only keep open-class words.
- ▶ **Dimensions:**

difference  
get  
go  
goes

impossible  
major  
possibly  
repair

thing  
turns  
usually  
wrong

## Frequency counts...

### Douglas Adams, *Mostly harmless*

The major difference between a thing that might go wrong and a thing that cannot possibly go wrong is that when a thing that cannot possibly go wrong goes wrong it usually turns out to be impossible to get at or repair.

#### ► **Counts:**

difference 1

get 1

go 3

goes 1

impossible 1

major 1

possibly 2

repair 1

thing 3

turns 1

usually 1

wrong 4

## Conversion into 5-word windows...

### Douglas Adams, *Mostly harmless*

The major difference between a thing that might go wrong and a thing that cannot possibly go wrong is that when a thing that cannot possibly go wrong goes wrong it usually turns out to be impossible to get at or repair.

- ▶  $\emptyset$   $\emptyset$  **the** major difference
- ▶  $\emptyset$  the **major** difference between
- ▶ the major **difference** between a
- ▶ major difference **between** a thing
- ▶ ...

## Distribution for *wrong*

Douglas Adams, *Mostly harmless*

The major difference between a thing that [might go wrong and a] thing that cannot [possibly go wrong is that] when a thing that cannot [possibly go [wrong goes wrong] it usually] turns out to be impossible to get at or repair.

► **Distribution (frequencies):**

difference 0  
get 0  
go 3  
goes 2

impossible 0  
major 0  
possibly 2  
repair 0

thing 0  
turns 0  
usually 1  
wrong 2

## Distribution for *wrong*

Douglas Adams, *Mostly harmless*

The major difference between a thing that [might go wrong and a] thing that cannot [possibly go wrong is that] when a thing that cannot [possibly go [wrong goes wrong] it usually] turns out to be impossible to get at or repair.

► **Distribution (PPMIs):**

difference 0  
get 0  
go 0.70  
goes 1

impossible 0  
major 0  
possibly 0.70  
repair 0

thing 0  
turns 0  
usually 0.70  
wrong 0.40



## Experimental corpus

- ▶ Dump of entire **English Wikipedia**, parsed with the English Resource Grammar producing dependencies.
- ▶ **Dependency contexts** include:
  - ▶ **For nouns:** verbs (+ any other argument of the verb), modifying adjectives, prepositions (+ any other argument of the preposition).  
*cat: chase\_v+mouse\_n, black\_a, of\_p+neighbour\_n*
  - ▶ **For verbs:** arguments (NPs and PPs), adverbial modifiers.  
*eat: cat\_n+mouse\_n, in\_p+kitchen\_n, fast\_a*
  - ▶ **For adjectives:** modified nouns; prepositions (+ any other argument of the preposition)  
*angry: cat\_n, at\_p+dog\_n*

## System description

- ▶ Semantic space: top 100,000 contexts.
- ▶ Weighting: pointwise mutual information (PMI).

## An example noun

► *language*:

0.54::other+than\_p+English\_n

0.53::English\_n+as\_p

0.52::English\_n+be\_v

0.49::english\_a

0.48::and\_c+literature\_n

0.48::people\_n+speak\_v

0.47::French\_n+be\_v

0.46::Spanish\_n+be\_v

0.46::and\_c+dialects\_n

0.45::grammar\_n+of\_p

0.45::foreign\_a

0.45::germanic\_a

0.44::German\_n+be\_v

0.44::of\_p+instruction\_n

0.44::speaker\_n+of\_p

0.42::pron\_rel\_+speak\_v

0.42::colon\_v+English\_n

0.42::be\_v+English\_n

0.42::language\_n+be\_v

0.42::and\_c+culture\_n

0.41::arabic\_a

0.41::dialects\_n+of\_p

0.40::percent\_n+speak\_v

0.39::spanish\_a

0.39::welsh\_a

0.39::tonal\_a

## An example adjective

► *academic*:

0.52::Decathlon\_n

0.51::excellence\_n

0.45::dishonesty\_n

0.45::rigor\_n

0.43::achievement\_n

0.42::discipline\_n

0.40::vice\_president\_n+for\_p

0.39::institution\_n

0.39::credentials\_n

0.38::journal\_n

0.37::journal\_n+be\_v

0.37::vocational\_a

0.37::student\_n+achieve\_v

0.36::athletic\_a

0.36::reputation\_n+for\_p

0.35::regalia\_n

0.35::program\_n

0.35::freedom\_n

0.35::student\_n+with\_p

0.35::curriculum\_n

0.34::standard\_n

0.34::at\_p+institution\_n

0.34::career\_n

0.34::Career\_n

0.33::dress\_n

0.33::scholarship\_n

0.33::prepare\_v+student\_n

0.33::qualification\_n

## Data sparsity

- ▶ Distribution for *unicycle*, as obtained from Wikipedia.

0.45::motorized_a	0.17::slip_v
0.40::pron_rel_+ride_v	0.16::and_c+1_n
0.24::for_p+entertainment_n	0.16::autonomous_a
0.24::half_n+be_v	0.16::balance_v
0.24::unwieldy_a	0.13::tall_a
0.23::earn_v+point_n	0.12::fast_a
0.22::pron_rel_+crash_v	0.11::red_a
0.19::man_n+on_p	0.07::come_v
0.19::on_p+stage_n	0.06::high_a
0.19::position_n+on_p	

# Polysemy

- ▶ Distribution for *pot*, as obtained from Wikipedia.

0.57::melt_v	0.32::boil_v
0.44::pron_rel_+smoke_v	0.31::bowl_n+and_c
0.43::of_p+gold_n	0.31::ingredient_n+in_p
0.41::porous_a	0.30::plant_n+in_p
0.40::of_p+tea_n	0.30::simmer_v
0.39::player_n+win_v	0.29::pot_n+and_c
0.39::money_n+in_p	0.28::bottom_n+of_p
0.38::of_p+coffee_n	0.28::of_p+flower_n
0.33::amount_n+in_p	0.28::of_p+water_n
0.33::ceramic_a	0.28::food_n+in_p
0.33::hot_a	

# Polysemy

- ▶ Some researchers incorporate word sense disambiguation techniques.
- ▶ But most assume a single space for each word: can perhaps think of subspaces corresponding to senses.
- ▶ Graded rather than absolute notion of polysemy.

## Idiomatic expressions

- ▶ Distribution for *time*, as obtained from Wikipedia.

0.46::of\_p+death\_n

0.45::same\_a

0.45::1\_n+at\_p(temp)

0.45::Nick\_n+of\_p

0.42::spare\_a

0.42::playoffs\_n+for\_p

0.42::of\_p+retirement\_n

0.41::of\_p+release\_n

0.40::pron\_rel\_+spend\_v

0.39::sand\_n+of\_p

0.39::pron\_rel\_+waste\_v

0.38::place\_n+around\_p

0.38::of\_p+arrival\_n

0.38::of\_p+completion\_n

0.37::after\_p+time\_n

0.37::of\_p+arrest\_n

0.37::country\_n+at\_p

0.37::age\_n+at\_p

0.37::space\_n+and\_c

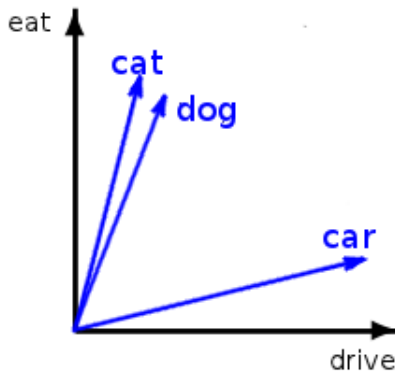
0.37::in\_p+career\_n

0.37::world\_n+at\_p



## Calculating similarity in a distributional space

- ▶ Distributions are vectors, so distance can be calculated.



## Measuring similarity

- ▶ Cosine:

$$\cos(\theta) = \frac{\sum v1_k * v2_k}{\sqrt{\sum v1_k^2} * \sqrt{\sum v2_k^2}} \quad (1)$$

- ▶ The cosine measure calculates the angle between two vectors and is therefore length-independent. This is important, as frequent words have longer vectors than less frequent ones.
- ▶ Other measures include Euclidean distance etc.

## The scale of similarity: some examples

house – building 0.43  
gem – jewel 0.31  
capitalism – communism 0.29  
motorcycle – bike 0.29  
test – exam 0.27  
school – student 0.25  
singer – academic 0.17  
horse – farm 0.13  
man – accident 0.09  
tree – auction 0.02  
cat – county 0.007

## Words most similar to *cat*

as chosen from the 5000 most frequent nouns in Wikipedia.

1 cat	0.29 human	0.25 woman	0.22 monster
0.45 dog	0.29 goat	0.25 fish	0.22 people
0.36 animal	0.28 snake	0.24 squirrel	0.22 tiger
0.34 rat	0.28 bear	0.24 dragon	0.22 mammal
0.33 rabbit	0.28 man	0.24 frog	0.21 bat
0.33 pig	0.28 cow	0.23 baby	0.21 duck
0.31 monkey	0.26 fox	0.23 child	0.21 cattle
0.31 bird	0.26 girl	0.23 lion	0.21 dinosaur
0.30 horse	0.26 sheep	0.23 person	0.21 character
0.29 mouse	0.26 boy	0.23 pet	0.21 kid
0.29 wolf	0.26 elephant	0.23 lizard	0.21 turtle
0.29 creature	0.25 deer	0.23 chicken	0.20 robot

## But what is similarity?

- ▶ In distributional semantics, very broad notion: synonyms, near-synonyms, hyponyms, taxonomical siblings, antonyms, etc.
- ▶ Correlates with a psychological reality.
- ▶ Test via correlation with human judgments on a test set:
  - ▶ Miller & Charles (1991)
  - ▶ WordSim
  - ▶ MEN
  - ▶ SimLex

## Miller & Charles 1991

3.92 automobile-car	3.05 bird-cock	0.84 forest-graveyard
3.84 journey-voyage	2.97 bird-crane	0.55 monk-slave
3.84 gem-jewel	2.95 implement-tool	0.42 lad-wizard
3.76 boy-lad	2.82 brother-monk	0.42 coast-forest
3.7 coast-shore	1.68 crane-implement	0.13 cord-smile
3.61 asylum-madhouse	1.66 brother-lad	0.11 glass-magician
3.5 magician-wizard	1.16 car-journey	0.08 rooster-voyage
3.42 midday-noon	1.1 monk-oracle	0.08 noon-string
3.11 furnace-stove	0.89 food-rooster	
3.08 food-fruit	0.87 coast-hill	

- ▶ Distributional systems, reported correlations 0.8 or more.

## TOEFL synonym test

Test of English as a Foreign Language: task is to find the best match to a word:

Prompt: levied

Choices: (a) imposed  
(b) believed  
(c) requested  
(d) correlated

Solution: (a) imposed

- ▶ Non-native English speakers applying to college in US reported to average 65%
- ▶ Best corpus-based results are 100%

## Distributional methods are a usage representation

- ▶ Distributions are a good conceptual representation if you believe that ‘the meaning of a word is given by its usage’.
- ▶ Corpus-dependent, culture-dependent, register-dependent.  
Example: similarity between *policeman* and *cop*: 0.23



## Distribution for *policeman*

### **policeman**

0.59::ball\_n+poss\_rel

0.48::and\_c+civilian\_n

0.42::soldier\_n+and\_c

0.41::and\_c+soldier\_n

0.38::secret\_a

0.37::people\_n+include\_v

0.37::corrupt\_a

0.36::uniformed\_a

0.35::uniform\_n+poss\_rel

0.35::civilian\_n+and\_c

0.31::iraqi\_a

0.31::lot\_n+poss\_rel

0.31::chechen\_a

0.30::laugh\_v

0.29::and\_c+criminal\_n

0.28::incompetent\_a

0.28::pron\_rel+shoot\_v

0.28::hat\_n+poss\_rel

0.28::terrorist\_n+and\_c

0.27::and\_c+crowd\_n

0.27::military\_a

0.27::helmet\_n+poss\_rel

0.27::father\_n+be\_v

0.26::on\_p+duty\_n

0.25::salary\_n+poss\_rel

0.25::on\_p+horseback\_n

0.25::armed\_a

0.24::and\_c+nurse\_n

0.24::job\_n+as\_p

0.24::open\_v+fire\_n

## Distribution for *cop*

### **cop**

0.45::crooked\_a

0.45::corrupt\_a

0.44::maniac\_a

0.38::dirty\_a

0.37::honest\_a

0.36::uniformed\_a

0.35::tough\_a

0.33::pron\_rel\_+call\_v

0.32::funky\_a

0.32::bad\_a

0.29::veteran\_a

0.29::and\_c+robot\_n

0.28::and\_c+criminal\_n

0.28::bogus\_a

0.28::talk\_v+to\_p+pron\_rel\_

0.27::investigate\_v+murder\_n

0.26::on\_p+force\_n

0.25::parody\_n+of\_p

0.25::Mason\_n+and\_c

0.25::pron\_rel\_+kill\_v

0.25::racist\_a

0.24::addicted\_a

0.23::gritty\_a

0.23::and\_c+interference\_n

0.23::arrive\_v

0.23::and\_c+detective\_n

0.22::look\_v+way\_n

0.22::dead\_a

0.22::pron\_rel\_+stab\_v

0.21::pron\_rel\_+evade\_v

## The similarity of synonyms

- ▶ Similarity between *eggplant/aubergine*: 0.11  
Relatively low cosine. Partly due to frequency (222 for *eggplant*, 56 for *aubergine*).
- ▶ Similarity between *policeman/cop*: 0.23
- ▶ Similarity between *city/town*: 0.73

In general, true synonymy does not correspond to higher similarity scores than near-synonymy.

## Similarity of antonyms

- ▶ Similarities between:
  - ▶ cold/hot 0.29
  - ▶ dead/alive 0.24
  - ▶ large/small 0.68
  - ▶ colonel/general 0.33

## Identifying antonyms

- ▶ Antonyms have high distributional similarity: hard to distinguish from near-synonyms purely by distributions.
- ▶ Identification by heuristics applied to pairs of highly similar distributions.
- ▶ For instance, antonyms are frequently coordinated while synonyms are not:
  - ▶ a selection of cold and hot drinks
  - ▶ wanted dead or alive

## Distributions and knowledge

What kind of information do distributions encode?

- ▶ lexical knowledge
- ▶ world knowledge
- ▶ boundary between the two is blurry
- ▶ no perceptual knowledge

Distributions are partial lexical semantic representations, but useful and theoretically interesting.

# Acknowledgement

*Some slides were adapted from Ann Copestake*