

# Natural Language Processing 1

## Lecture 6: Distributional semantics: generalisation and word embeddings

Katia Shutova

ILLC  
University of Amsterdam

15 November 2018

## Experimental corpus

- ▶ Dump of entire **English Wikipedia**, parsed with the English Resource Grammar producing dependencies.
- ▶ **Dependencies** include:
  - ▶ **For nouns**: head verbs (+ any other argument of the verb), modifying adjectives, head 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; head prepositions (+ any other argument of the preposition)  
*black: 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

## Corpus choice

- ▶ As much data as possible?
  - ▶ British National Corpus (BNC): 100 m words
  - ▶ Wikipedia: 897 m words
  - ▶ UKWac: 2 bn words
  - ▶ ...
- ▶ In general preferable, *but*:
  - ▶ More data is not necessarily the data you want.
  - ▶ More data is not necessarily realistic from a psycholinguistic point of view. We perhaps encounter 50,000 words a day. BNC = 5 years' text exposure.

## 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

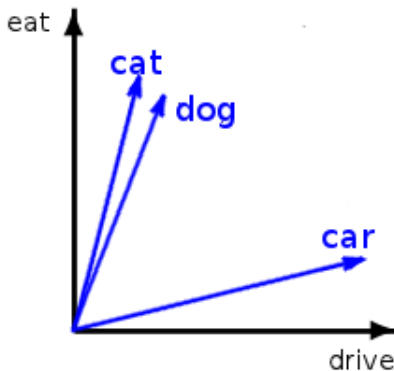
0.37::space\_n+and\_c

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 Jaccard, 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	0.28::incompetent_a
0.48::and_c+civilian_n	0.28::pron_rel_+shoot_v
0.42::soldier_n+and_c	0.28::hat_n+poss_rel
0.41::and_c+soldier_n	0.28::terrorist_n+and_c
0.38::secret_a	0.27::and_c+crowd_n
0.37::people_n+include_v	0.27::military_a
0.37::corrupt_a	0.27::helmet_n+poss_rel
0.36::uniformed_a	0.27::father_n+be_v
0.35::uniform_n+poss_rel	0.26::on_p+duty_n
0.35::civilian_n+and_c	0.25::salary_n+poss_rel
0.31::iraqi_a	0.25::on_p+horseback_n
0.31::lot_n+poss_rel	0.25::armed_a
0.31::chechen_a	0.24::and_c+nurse_n
0.30::laugh_v	0.24::job_n+as_p
0.29::and_c+criminal_n	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.



# Clustering

- ▶ clustering techniques group objects into clusters
- ▶ similar objects in the same cluster, dissimilar objects in different clusters
- ▶ allows us to obtain generalisations over the data
- ▶ widely used in various NLP tasks:
  - ▶ semantics (e.g. word clustering);
  - ▶ summarization (e.g. sentence clustering);
  - ▶ text mining (e.g. document clustering).

## Distributional word clustering

We will:

- ▶ cluster words based on the contexts in which they occur
- ▶ assumption: words with similar meanings occur in similar contexts, i.e. are distributionally similar
- ▶ we will consider noun clustering as an example
- ▶ cluster 2000 nouns – most frequent in the British National Corpus
- ▶ into 200 clusters

## Clustering nouns

truck lorry path  
bike highway way  
car taxi street  
bicycle driver road avenue  
mechanic lab building house  
engineer scientist office flat shack  
plumber writer office flat shack  
journalist proceedings dwelling  
book journal  
newspaper magazine

# Clustering nouns



## Feature vectors

- ▶ can use different kinds of context as features for clustering
  - ▶ window based context
  - ▶ parsed or unparsed
  - ▶ syntactic dependencies
- ▶ different types of context yield different results
- ▶ **Example experiment:** use verbs that take the noun as a direct object or a subject as features for clustering
- ▶ **Feature vectors:** verb lemmas, indexed by dependency type, e.g. subject or direct object
- ▶ **Feature values:** corpus frequencies

## Extracting feature vectors: Examples

### tree (Dobj)

85 plant\_v  
 82 climb\_v  
 48 see\_v  
 46 cut\_v  
 27 fall\_v  
 26 like\_v  
 23 make\_v  
 23 grow\_v  
 22 use\_v  
 22 round\_v  
 20 get\_v  
 18 hit\_v  
 18 fell\_v  
 18 bark\_v  
 17 want\_v  
 16 leave\_v  
 ...

### crop (Dobj)

76 grow\_v  
 44 produce\_v  
 16 harvest\_v  
 12 plant\_v  
 10 ensure\_v  
 10 cut\_v  
 9 yield\_v  
 9 protect\_v  
 9 destroy\_v  
 7 spray\_v  
 7 lose\_v  
 6 sell\_v  
 6 get\_v  
 5 support\_v  
 5 see\_v  
 5 raise\_v  
 ...

### tree (Subj)

131 grow\_v  
 49 plant\_v  
 40 stand\_v  
 26 fell\_v  
 25 look\_v  
 23 make\_v  
 22 surround\_v  
 21 show\_v  
 20 seem\_v  
 20 overhang\_v  
 20 fall\_v  
 19 cut\_v  
 18 take\_v  
 18 go\_v  
 18 become\_v  
 17 line\_v  
 ...

### crop (Subj)

78 grow\_v  
 23 yield\_v  
 10 sow\_v  
 9 fail\_v  
 8 plant\_v  
 7 spray\_v  
 7 come\_v  
 6 produce\_v  
 6 feed\_v  
 6 cut\_v  
 5 sell\_v  
 5 make\_v  
 5 include\_v  
 5 harvest\_v  
 4 follow\_v  
 3 ripen\_v  
 ...

## Feature vectors: Examples

### tree

131 grow\_v\_Subj  
85 plant\_v\_Dobj  
82 climb\_v\_Dobj  
49 plant\_v\_Subj  
48 see\_v\_Dobj  
46 cut\_v\_Dobj  
40 stand\_v\_Subj  
27 fall\_v\_Dobj  
26 like\_v\_Dobj  
26 fell\_v\_Subj  
25 look\_v\_Subj  
23 make\_v\_Subj  
23 make\_v\_Dobj  
23 grow\_v\_Dobj  
22 use\_v\_Dobj  
22 surround\_v\_Subj  
22 round\_v\_Dobj  
20 overhang\_v\_Subj

...

### crop

78 grow\_v\_Subj  
76 grow\_v\_Dobj  
44 produce\_v\_Dobj  
23 yield\_v\_Subj  
16 harvest\_v\_Dobj  
12 plant\_v\_Dobj  
10 sow\_v\_Subj  
10 ensure\_v\_Dobj  
10 cut\_v\_Dobj  
9 yield\_v\_Dobj  
9 protect\_v\_Dobj  
9 fail\_v\_Subj  
9 destroy\_v\_Dobj  
8 plant\_v\_Subj  
7 spray\_v\_Subj  
7 spray\_v\_Dobj  
7 lose\_v\_Dobj  
6 feed\_v\_Subj

...

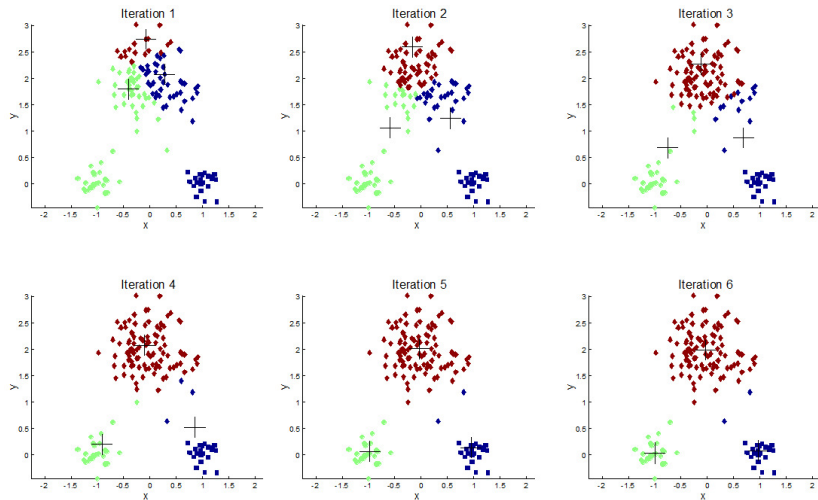
## Clustering algorithms, K-means

- ▶ many clustering algorithms are available
- ▶ example algorithm: K-means clustering
  - ▶ given a set of  $N$  data points  $\{x_1, x_2, \dots, x_N\}$
  - ▶ partition the data points into  $K$  clusters  $C = \{C_1, C_2, \dots, C_K\}$
  - ▶ minimize the sum of the squares of the distances of each data point to the cluster mean vector  $\mu_i$ :

$$\arg \min_C \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \mu_i\|^2 \quad (2)$$



# K-means clustering



## Noun clusters

tree crop flower plant root leaf seed rose wood grain stem forest garden

consent permission concession injunction licence approval

lifetime quarter period century succession stage generation decade phase  
interval future

subsidy compensation damages allowance payment pension grant

carriage bike vehicle train truck lorry coach taxi

official officer inspector journalist detective constable police policeman re-  
porter

girl other woman child person people

length past mile metre distance inch yard

tide breeze flood wind rain storm weather wave current heat

sister daughter parent relative lover cousin friend wife mother husband  
brother father

## Different senses of *run*

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

## Subject arguments of *run*

0.2125 drop tear sweat paint blood water juice

0.1665 technology architecture program system product version interface

software tool computer network processor chip package

0.1657 tunnel road path trail lane route track street bridge

0.1166 carriage bike vehicle train truck lorry coach taxi

0.0919 tide breeze flood wind rain storm weather wave current heat

0.0865 tube lock tank circuit joint filter battery engine device disk furniture

machine mine seal equipment machinery wheel motor slide disc instrument

0.0792 ocean canal stream bath river waters pond pool lake

0.0497 rope hook cable wire thread ring knot belt chain string

0.0469 arrangement policy measure reform proposal project programme

scheme plan course

0.0352 week month year

0.0351 couple minute night morning hour time evening afternoon

## Subject arguments of *run* (continued)

0.0341 criticism appeal charge application allegation claim objection  
suggestion case complaint

0.0253 championship open tournament league final round race match  
competition game contest

0.0218 desire hostility anxiety passion doubt fear curiosity enthusiasm  
impulse instinct emotion feeling suspicion

0.0183 expenditure cost risk expense emission budget spending

0.0136 competitor rival team club champion star winner squad county player  
liverpool partner leads

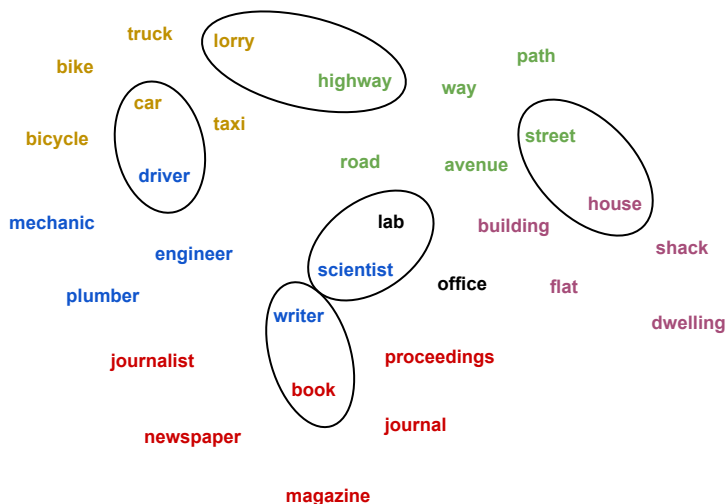
0.0102 being species sheep animal creature horse baby human fish male  
lamb bird rabbit female insect cattle mouse monster

...

# Clustering nouns



## Clustering nouns



## We can also cluster verbs...

sparkle glow widen flash flare gleam darken narrow flicker shine blaze  
bulge

gulp drain stir empty pour sip spill swallow drink pollute seep flow drip  
purify ooze pump bubble splash ripple simmer boil tread

polish clean scrape scrub soak

kick hurl push fling throw pull drag haul

rise fall shrink drop double fluctuate dwindle decline plunge decrease  
soar tumble surge spiral boom

initiate inhibit aid halt trace track speed obstruct impede accelerate  
slow stimulate hinder block

work escape fight head ride fly arrive travel come run go slip move



## Uses of word clustering in NLP

Widely used in NLP as a source of lexical information:

- ▶ Word sense induction and disambiguation
- ▶ Modelling predicate-argument structure (e.g. semantic roles)
- ▶ Identifying figurative language and idioms
- ▶ Paraphrasing and paraphrase detection
- ▶ Used in applications directly, e.g. machine translation, information retrieval etc.

# Distributional semantic models

## 1. Count-based models:

- ▶ Explicit vectors: dimensions are elements in the context
- ▶ **long sparse** vectors with **interpretable** dimensions

## 2. Prediction-based models:

- ▶ Train a model to predict plausible contexts for a word
- ▶ learn word representations in the process
- ▶ **short dense** vectors with **latent** dimensions

## Sparse vs. dense vectors

### Why dense vectors?

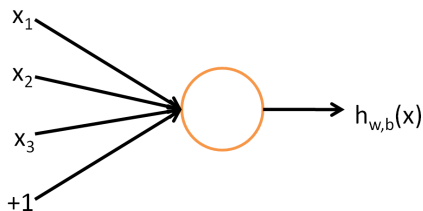
- ▶ easier to use as features in machine learning (less weights to tune)
- ▶ may generalize better than storing explicit counts
- ▶ may do better at capturing synonymy:
  - ▶ e.g. *car* and *automobile* are distinct dimensions in count-based models
  - ▶ will not capture similarity between a word with *car* as a neighbour and a word with *automobile* as a neighbour

## Brief introduction to neural networks

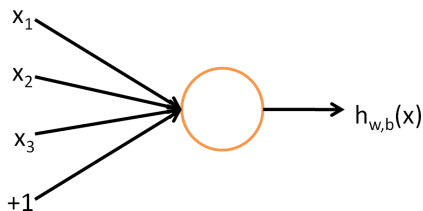
Supervised learning framework.

- ▶ **Input:** a set of labelled training examples  $(x^{(i)}, y^{(i)})$
- ▶ **Output:** hypotheses  $h_{W,b}(x)$  with parameters  $W, b$  which we fit to our data

The simplest possible neural network — single **neuron**



## Neuron as a computational unit



$$h_{W,b}(x) = f(W^T x + b) = f\left(\sum_{i=1}^3 W_i x_i + b\right)$$

where  $f : \mathbb{R} \mapsto \mathbb{R}$  is the **activation function**,  
 $W$  is a matrix of trainable weights,  
 $b$  is the bias term.

## Activation functions (common choices)

**Sigmoid** function

$$f(z) = \frac{1}{1 + e^{-z}}$$

*output in range [0,1]*

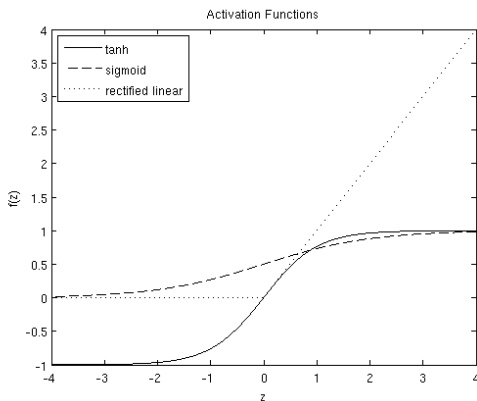
**Hyperbolic tangent (tanh):**

$$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

*output in range [-1,1]*

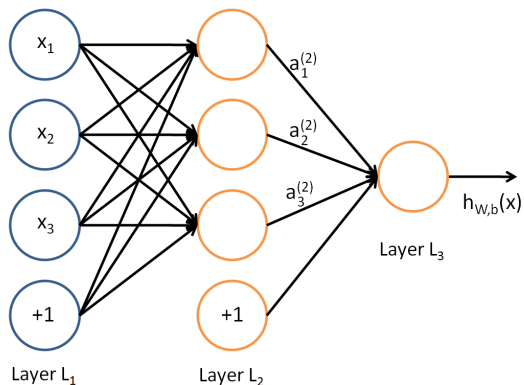
**Rectified linear (ReLU):**

$$f(z) = \max(0, z)$$



# Multi-layer neural network

Feed-forward architecture

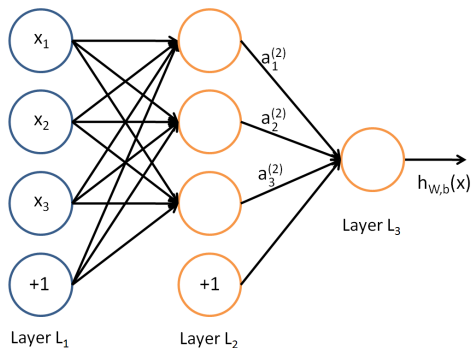


Input layer

Hidden layer

Output layer

## Multi-layer neural network

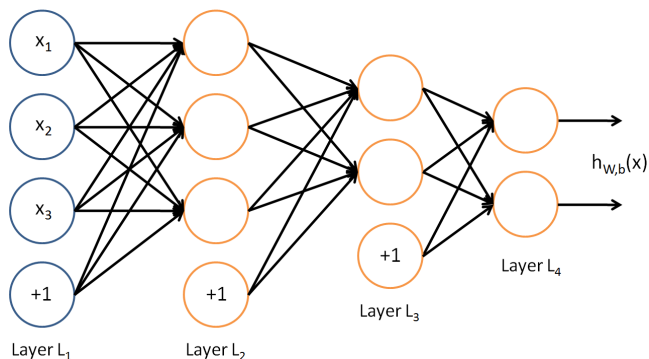


$$z^{(2)} = W^{(1)}x + b^{(1)}$$
$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$a^{(2)} = f(z^{(2)})$$
$$h_{W,b}(x) = a^{(3)} = f(z^{(3)})$$



## Deep neural networks and multi-class classification



## Softmax function

Used in multi-class classification problems.

- ▶ Takes a vector of real values and squashes them into the range  $[0,1]$ , so that they add up to 1
- ▶ use this as a **probability distribution** over output classes

$$\text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^d e^{z_k}}$$

$d$  is the dimensionality of the output layer

## Acknowledgement

*Some slides were adapted from Aurelie Herbelot*

*The introduction to neural networks is based on this helpful tutorial:*

`http://ufldl.stanford.edu/tutorial/supervised/  
MultiLayerNeuralNetworks/`