

# DD1418

## Språkteknologi med introduktion till maskininlärning

### 1a: Intro

Johan Boye, KTH

# What this course is about

This course is about processing human languages.

- English, Swedish, Chinese, Russian, Arabic, ...

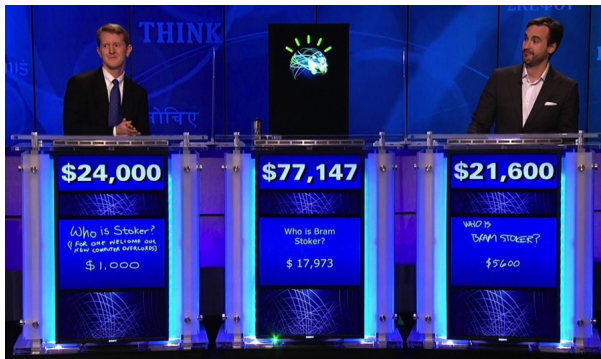
What is “processing”?

- One-liner: “Making computers understand language”
- More precisely: Extracting useful information from natural language, and generating (correct, useful) language.

# What's the point?

- A lot of human knowledge is represented as text.
- Lots of useful applications.
- Understanding language is a key aspect in the AI programme.
- We can gain insights about human language from the study of computational models.
- Language processing is one of the main drivers of machine learning (together with image recognition).
- Getting an interesting and well-paying job!

# Question answering: Watson



On February 16, 2011, IBM's computer system Watson defeated the world's best human Jeopardy champions

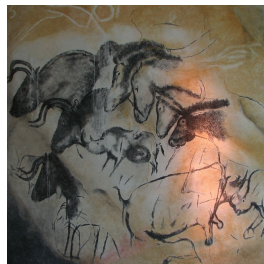
<https://vimeo.com/222234104>

# Language engineering applications

- Question-answering
- Speech recognition
- Machine translation
- Search engines
- Natural-language interfaces
- Text summarization
- Text classification
- Word prediction
- Spell checking
- Grammar checking
- ...

# Some things we don't know about language

- How old is it? (dunno, but perhaps around 100,000 years).
- How did it arise?
- Why do children learn it so effortlessly?
- How does the brain generate and analyse it?



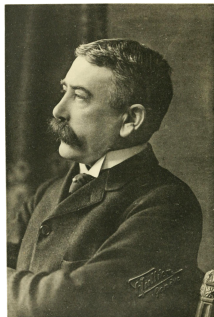
# Some things we do know about language

- It's uniquely human.
- There are 5,000-7,000 known languages.
- The relation between words and meaning is arbitrary.
- Language is rule-bound.
- Language is productive.



# 'Arbitrariness of the sign'

- The relation between words and meaning is arbitrary.
  - The word *dog* doesn't look like a dog or bark like a dog...
  - ... but means 'dog' all the same.
  - Seems obvious but is an important prerequisite for the effectiveness of language.
- This observation is usually attributed to *Ferdinand de Saussure* (1857-1913).





# Language is rule-bound and productive

Language can be broken down to smallest units (words), which are combined using the rules of the language.

These rules are a *naturally occurring phenomenon*, not something we learn in school.

Using these sub-conscious rules, we can produce and understand an infinite number of sentences, e.g.:

- He went skiing with a kangaroo and a watermelon in his left shoe.

Even if words are arbitrary, language structure definitely is not!

# Why is language understanding hard?

Language is *underspecified*

Language interpretation *requires knowledge about the world.*

Language is *ambiguous*

- “John made the pie in the fridge”

+ many more (smaller) challenges...

# Processing language by computers (1)

## Rule-based approaches

- The earliest work strived to uncover the subconscious rules that guide our language understanding.
- Grammars, automata, logic formulas, etc.
- Great when you don't have much data
- Predominant approach beginnings-1990's
- Rationalism

# Processing language by computers (2)

## Statistical / “traditional” machine-learning approaches

- Learns patterns directly from text data (from *corpora*)
- Great if you have lots of data
- Predominant approach late 1990’s-2014 (about)
- Empiricism

## Neural approaches

- Great if you have huge amounts of data
- End-to-end systems

# Machine learning

Machine learning is a way of **finding patterns** and **making predictions** from empirical data.

Useful when there is no known algorithm to solve a problem, e.g. ...

**decoding handwritten text,**

An attempt to get more information about the Admiralty House meeting will be made in the House of Commons this afternoon. Labour M.P.s already have many questions to the Prime Minister asking for a statement. President Kennedy flew from London Airport last night to arrive in Washington this morning. He is to make a 30-minute nation-wide broadcast and television report on his talks with Mr. Khrushchev this evening.

**or sentiment analysis**



For some problems, it's just terribly difficult to write an algorithm that solves it.

E.g. sentiment analysis: Are these reviews positive or negative?

- *Uncanny, haunting, I must have read this novel at the right time for me as it found a sure spot under my skin and disturbed my normally peaceful sleep.*
- *If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut.*

# Machine learning

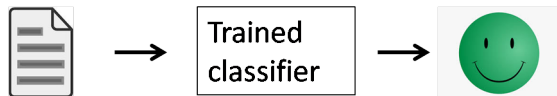
It's hard to write rules to decide whether or not a review is positive (the rationalist approach).

The empiricist approach would instead be to **train** a learning algorithm on **examples** of positive and negative reviews, and let it **generalize** from these examples.

- *Uncanny, haunting, I must have read this novel at the right time for me as it found a sure spot under my skin and disturbed my normally peaceful sleep.* 😊
- *If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut.* 😡
- *A wonderful book that I will remember for a long time.* 😊
- *I hated this movie. Hated, hated, hated, hated, hated this movie. Hated it. Hated every simpering stupid vacant audience-insulting moment of it.* 😡
- ...

# Machine learning

The trained algorithm can then be used for classifying reviews:



We'll see in the course how this is achieved.

Recent advances in language engineering has made it possible to solve even more difficult problems (like translation). We'll see some of that in the course as well.



# DD1418

## Lecture 1b: Words

Johan Boye, KTH

# Levels of linguistic analysis

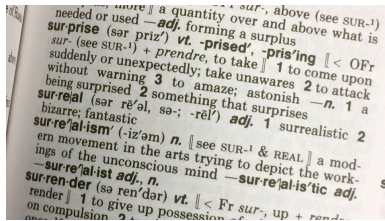
Words	Morphology, Phonology
Sentences	Syntax
Meaning	Semantics
Language use	Pragmatics

# Some terminology

A *corpus* (pl. *corpora*) is a collection of texts, possibly annotated.

A *lexicon* or *vocabulary* is a list of all the unique words in a corpus.

The *lemma* (pl. *lemmata*) of a word is what you look up in a dictionary.



# Words and tokens

*A hat and a hat make two hats*

contain 8 *tokens* (sw. *löpord*), but 5 words (not counting morphological variants), or 6 words (counting morphological variants).

In language engineering, we see two variants:

**full-form lexicon** Save all forms of every word in the lexicon (*go*, *goes*, *gone*, *went*, *going*, etc.)

**morphological analysis** Store only lemmas + rules to generate all the other forms.

# Word classes

Words can be divided into classes depending on their use in the language.

- Noun, verb, adjective, adverb, preposition, pronoun, conjunction, interjection, determiner, etc.
- These classes are often called *parts-of-speech* (or *POS tags*).

Lots of debate in linguistics about their nature and generality.

The idea of word classes can be traced back to Aristotle (384-322 BC) and Dionysius Thrax (about 170-90 BC).

# Quiz

Determine the word classes in the following sentence:

Jag borde åka hem nu

# Quiz

Determine the word classes in the following sentence:

Jag	borde	åka	hem	nu
-----	-------	-----	-----	----

Pronoun	Auxiliary verb	Verb	Adverb	Adverb
---------	----------------	------	--------	--------

# Open and closed classes

Closed classes, e.g. :

- determiners (a, an, the, some, ...)
- pronouns (she, her, I, you, me, ...)
- prepositions (on, to, under, from, ...)

Open classes

- nouns, verbs, adjectives, adverbs



# SUC (Stockholm-Umea Corpus) tag set

Code	Swedish category	Example	English translation
AB	Adverb	<i>inte</i>	Adverb
DT	Determinerare	<i>denna</i>	Determiner
HA	Frågande/relativt adverb	<i>när</i>	Interrogative/Relative Adverb
HD	Frågande/relativ determinerare	<i>vilken</i>	Interrogative/Relative Determiner
HP	Frågande/relativt pronomen	<i>som</i>	Interrogative/Relative Pronoun
HS	Frågande/relativt possessivt pronomen	<i>vars</i>	Interrogative/Relative Possessive
IE	Infinitivmärke	<i>att</i>	Infinitive Marker
IN	Interjektion	<i>ja</i>	Interjection
JJ	Adjektiv	<i>glad</i>	Adjective
KN	Konjunktion	<i>och</i>	Conjunction
NN	Substantiv	<i>pudding</i>	Noun
PC	Particip	<i>utsänd</i>	Participle
PL	Partikel	<i>ut</i>	Particle
PM	Egennamn	<i>Mats</i>	Proper Noun
PN	Pronomen	<i>hon</i>	Pronoun
PP	Preposition	<i>av</i>	Preposition
PS	Possessivt pronomen	<i>hennes</i>	Possessive
RG	Grundtal	<i>tre</i>	Cardinal number
RO	Ordningstal	<i>tredje</i>	Ordinal number
SN	Subjunktion	<i>att</i>	Subjunction
UO	Utländskt ord	<i>the</i>	Foreign Word
VB	Verb	<i>kasta</i>	Verb

# Penn Treebank tag set

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			

# Universal Dependencies tag set

Universal Dependencies is an initiative to create a multilingual tagset and a set of multilingual analysis tools.

Open class words	Closed class words	Other
<u>ADJ</u>	<u>ADP</u>	<u>PUNCT</u>
<u>ADV</u>	<u>AUX</u>	<u>SYM</u>
<u>INTJ</u>	<u>CCONJ</u>	<u>X</u>
<u>NOUN</u>	<u>DET</u>	
<u>PROPN</u>	<u>NUM</u>	
<u>VERB</u>	<u>PART</u>	
	<u>PRON</u>	
	<u>SCONJ</u>	

# Part-of-Speech tagging

The *part-of-speech tagging* problem is to assign a POS tag to each word in a text.

Why is this useful?

- Speech synthesis (**insult** (noun) vs **insult** (verb))
- Syntactic analysis (“*Time flies like an arrow*”)
- Finding content words in text  
ADJ NOUN “*linear function*”  
NOUN NOUN “*regression coefficient*”  
ADJ ADJ NOUN “*Gaussian random variable*”  
NOUN PREP NOUN “*degrees of freedom*”
- Back-off model for machine learning (when we have sparse data).

# Ambiguities

Some words can belong to more than one class, e.g. *like*:

VERB “*I like her.*”

NOUN “*He got a like on Facebook.*”

ADJ “*The portrait is very like.*”

ADV “*This is, like, crazy!*”

ADP “*It looks like an accident*”

SCONJ “*He acted like he was all alone*”

# Ambiguities

Swedish also has this kind of words, e.g. så:

ADV “Så gick det till.”

PRON “På så sätt!”

SCONJ “Han åt så han blev mätt.”

INTJ “Så, det var intressant att höra!

VERB “Man måste så innan man kan skörda.”

NOUN “Grisarna drack ur en så.”

# New words

Even with a large dictionary, you will always encounter new words.

Names

Neologisms like *metrosexual*, to *google*, etc.

New concepts (*bogvisir*)



# Part-of-speech tagging

Retrieve all possible tags from a dictionary, then decide which ones are the most likely, e.g. :

*I like plays about bolliwogs*

	VERB/NOUN			
	ADJ/ADV			
PRON	/ADP/CONJ	VERB/NOUN	ADV/ADP/ADJ	?



PRON	VB	NOUN	PREP	NOUN
------	----	------	------	------



# Difficulty of POS tagging

POS tagging is considered to be an “easy” problem.

State-of-the-art POS taggers have about 97% accuracy, which is as good as it gets (what human experts can agree on, the *gold standard*).

However, a straightforward method (a *baseline*) already gives about 90% accuracy.

- Tag every word with its most frequent tag.
- Tag unknown words as nouns.

The high accuracy is partly due to:

- Many words are unambiguous (including *a*, *an*, *the*).
- Punctuation (“.”, “!”, “?”) is unambiguous.

# Difficulty of POS tagging

Some POS tagging decisions are difficult for computers:

*ADP “He turned off the road”*

*PART “He turned off the radio”*

# DD1418

## Lecture 1c: The structure of words

Johan Boye, KTH

# The structure of words

*Morphology* is the study of how *words* are built from smaller units called *morphemes* (smallest *meaningful* unit)

Two kinds of morphemes:

**stem** the core unit

**affixes** small units signalling various grammatical functions

**un- fortun -ate -ly**  
prefix stem suffixes

Note that the stem doesn't have to be a word! Stem  $\neq$  lemma!

# Affixes

Affixes come in four varieties:

prefix *un-familiar*

suffix *quick-ly*

infix (sv.) *korru-m-pera*

circumfix (ge.) *ge-sag-t*

# Word formation

Words can be formed from other words by:

inflection (sv. böjning) *cat* - *cats*

derivation (sv. avledning) *friend* - *friendly* - *friendliness*

compounding (sv. sammansättning) *smartphone*, *anti-missile*

# Inflections: Verbs

A typical English verb has 4 or 5 forms.

- *ask - asks - asked - asking*
- *go - goes - gone - went - going*

Swedish: about 10 forms

- *äta - äter - åt - ätit - ätande - ät - ätas - äts - åts - ätits - äten*

French: >40 forms

Classic Greek: 350 forms\*

Turkish: 3 million forms\*

\*S. Pinker (1997) *The language instinct*, Penguin.

# Inflections: Verbs

## Parler

The verb *parler* "to speak", in French orthography and IPA transcription

	Indicative				Subjunctive		Conditional	Imperative
	Present	Simple past	Imperfect	Simple future	Present	Imperfect	Present	Present
<b>Je</b>	parl-e /paʁl/	parl-ai /paʁle/	parl-ais /paʁtɛ/	parl-erai /paʁlɛʁe/	parl-e /paʁl/	parl-asse /paʁlas/	parl-erais /paʁlɛʁɛ/	
<b>tu</b>	parl-es /paʁl/	parl-as /paʁla/	parl-ais /paʁtɛ/	parl-eras /paʁlɛʁa/	parl-es /paʁl/	parl-asses /paʁlas/	parl-erais /paʁlɛʁɛ/	parl-e /paʁl/
<b>il</b>	parl-e /paʁl/	parl-a /paʁla/	parl-ait /paʁtɛ/	parl-era /paʁlɛʁa/	parl-e /paʁl/	parl-ât /paʁla/	parl-erait /paʁlɛʁɛ/	
<b>nous</b>	parl-ons /paʁlɔ̃/	parl-âmes /paʁlam/	parl-ions /paʁljɔ̃/	parl-erons /paʁlɛʁɔ̃/	parl-ions /paʁljɔ̃/	parl-assions /paʁlasjɔ̃/	parl-erions /paʁlɛʁjɔ̃/	parl-ons /paʁlɔ̃/
<b>vous</b>	parl-ez /paʁle/	parl-âtes /paʁlat/	parl-iez /paʁlje/	parl-erez /paʁlɛʁe/	parl-iez /paʁlje/	parl-assiez /paʁlasje/	parl-eriez /paʁlɛʁje/	parl-ez /paʁle/
<b>ils</b>	parl-ent /paʁl/	parl-èrent /paʁlɛːʁ/	parl-aient /paʁtɛ/	parl-eront /paʁlɛʁɔ̃/	parl-ent /paʁl/	parl-assent /paʁlas/	parl-eraient /paʁlɛʁɛ/	



# Inflections: Nouns

**English** A typical noun has 2 forms: *cat*, *cats*

- 1 feature: *Number* with 2 values: *Singular*, *Plural*.

**Swedish** typically 8 forms: *stol*, *stolen*, *stolar*, *stolarna*, *stols*, *stolens*, *stolars*, *stolarnas*  
*äpple*, *äpplet*, *äpplen*, *äpplena*, *äpples*, *äpplets*,  
*äpplens*, *äpplenas*

- *Number* with 2 values: *Singular*, *Plural*.
- *Species* with 2 values: *Indefinite*, *Definite*
- *Case* with 2 values: *Nominative*, *Genitive*
- *Gender* with 3 values: *Utrum*, *Neutrum*, *Masculine*

**Finnish** 2253 forms of *kauppa* (shop) listed at

<http://www.ling.helsinki.fi/~fkarlssso/genkau2.html>

# Inflections: Nouns

## Some forms of the Hungarian noun *ablak* (window):

<i>ablaka</i>	its window
<i>ablakában</i>	in its window
<i>ablakából</i>	from its window
<i>ablakai</i>	its windows
<i>ablakaik</i>	their windows
<i>ablakaikkal</i>	with their windows
<i>ablakainak</i>	for their windows
<i>ablakán</i>	on its window
<i>ablakát</i>	its window (accusativus)
<i>ablakba</i>	into the window
<i>ablakhoz</i>	towards the window
<i>ablakkal</i>	with the window
<i>ablakok</i>	windows
<i>ablakokat</i>	windows (accusativus)
<i>ablakokba</i>	into the windows
<i>ablakokkal</i>	with the windows
<i>ablakokon</i>	on the windows
<i>ablakon</i>	on the window
<i>ablakot</i>	window (accusativus)

# Derivations

By *derivation*, new words can be formed (often from a word of another class)

- *black*  $\Rightarrow$  *blackness*
- *affect*  $\Rightarrow$  *affection*  $\Rightarrow$  *affectionate*
- *compute*  $\Rightarrow$  *computation*  $\Rightarrow$  *computational*
- (but not *eat*  $\Rightarrow$  *eatation*\*)

Some Swedish examples:

- *motiv*  $\Rightarrow$  *motivera*  $\Rightarrow$  *motivering*
- *god*  $\Rightarrow$  *godhet*
- *äta*  $\Rightarrow$  *ätbar*

# Compounds

Swedish can form long compound words, e.g.  
*hårdvarukompatibilitetsproblem*

Some variations:

- Vowel changes: *hårdvarukompatibilitetsproblem*
- Vowel drop: *läkarmottagning*
- Extra *s*: *hårdvarukompatibilitetsproblem*
- Nothing at all: *tidrapport*

Some compound words have been *lexicalized* (e.g. *football*).

The longest lexicalized Swedish compound word (according to *SAOL*) is *realisationsvinstbeskattning*. (28 letters)

# Morphological analysis and generation

Morphological analysis Word form  $\Rightarrow$  lemma + features

- *cats* = NOUN:cat + NUMBER:plural
- *stolarnas* = NOUN:stol +  
GENDER:utrum +  
NUMBER:plural +  
SPECIES:definite +  
CASE:genitive

Morphological generation Lemma + features  $\Rightarrow$  word form

# SUC (Stockholm-Umea Corpus) tag set

Feature	Value	Legend	Parts-of-speech where feature applies
Gender	UTR	Uter (common)	DT, HD, HP, JJ, NN, PC, PN, PS, (RG, RO)
	NEU	Neuter	
	MAS	Masculine	
Number	SIN	Singular	DT, HD, HP, JJ, NN, PC, PN, PS, (RG, RO)
	PLU	Plural	
Definiteness	IND	Indefinite	DT, (HD, HP, HS), JJ, NN, PC, PN, (PS, RG, RO)
Case	DEF	Definite	JJ, NN, PC, PM, (RG, RO)
	NO	Nominative	
	M		
Tense	GEN	Genitive	VB
	PRS	Present	
	PRT	Preterite	
	SUP	Supinum	
	INF	Infinite	
Voice	AKT	Active	
	SFO	S-form (passive or deponential)	
Mood	KON	Subjunctive (Sw. konjunktiv)	
Participle form	PRS	Present	PC
Degree	PRF	Perfect	(AB), JJ
	POS	Positive	
	KO	Comparative	
	M		
Pronoun form	SUV	Superlative	PN
	SUB	Subject form	
	OBJ	Object form	
	SMS	Compound (Sw. sammansättningsform)	All parts-of-speech

# What's the point?

Spell checking (what is a correct word?)

Grammar checking (of *agreement*)

- The cat **s** **were** hungry. (number agreement)
- **Den** svart **a** svan **en** (species agreement)
- De **n** svarta svan **en** (gender agreement)
- De **n** svarta svan **en** (number agreement)

Information retrieval By *splitting compound words* and *removing suffixes*, more relevant documents can be found.

- *hårdvarukompatibilitetsproblemen* ⇒  
*hårdvaru-kompatibilitets-problem-en*

Translation, text generation

# Lemmatization

*Lemmatization* is the process of rewriting words into their lemmata.

*The boys are taller than the girls.* →  
*The boy be tall than the girl.*

Often useful in machine learning contexts where we want to reduce the number of dimensions.



# Lemmatization

For English, a lemmatizer can simply be a look-up table.

...

ask	ask
-----	-----

asked	ask
-------	-----

asking	ask
--------	-----

asks	ask
------	-----

...

However, in many languages this solution is not sufficient.