

Natural Language Processing

Complex AI Systems

Asmaa Elbadrawy
PhD, Lecturer
IFT Program, ASU

Humans use **language** to
represent knowledge (ideas)
and **communicate** them to
others.

Human language is called **Natural Language** (English, Spanish, etc.) as opposed to computer languages, which are called **Formal Languages** (first order logic, programming, etc.)


Why Computers need to process Natural Languages

- Communicate with humans
- Build KB's from human writings (learn)
- Advance scientific understanding of language and how humans use then.

**Natural Language Processing
(NLP) has many applications.**

Classifying news articles and web pages to facilitate searching.

Features in document classification		
Category: psychology	Category: nutrition	Category: sports
Depression	Calories	Reps
Anxiety	Protein	Tension
Syndrome	Diet	Overload
Mental	Fat	Exercise
Thought	Carbs	Abs



altexsoft
software r&d engineering

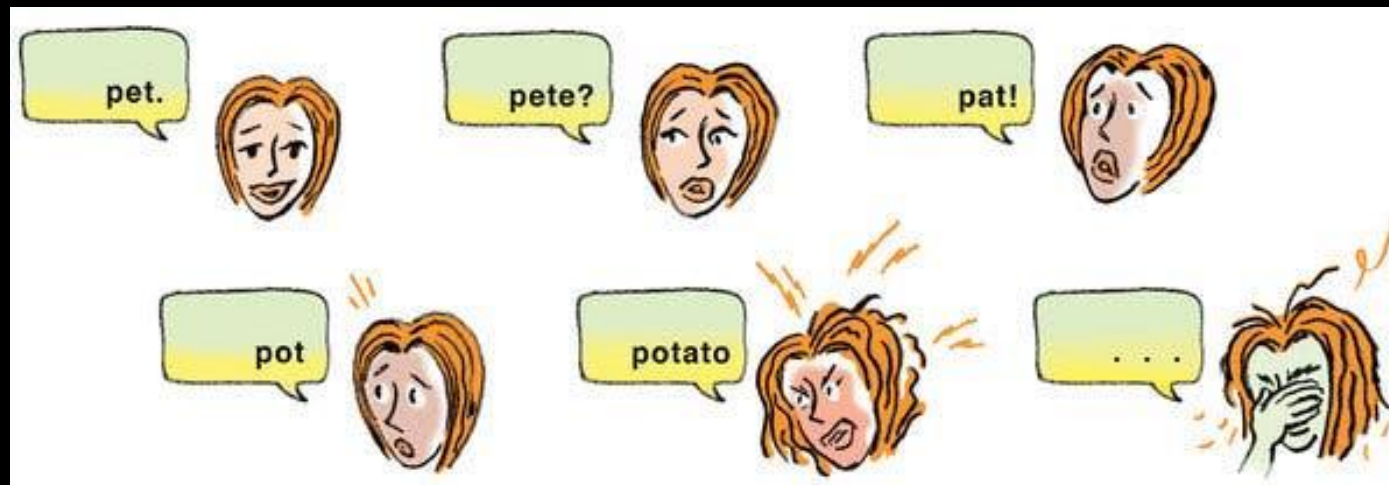
Source: <https://www.altexsoft.com/blog/document-classification/>

Auto-Translation.



Source: <https://www.nytimes.com/2013/05/02/technology/personaltech/the-utility-and-drawbacks-of-translation-apps.html>

Suggesting spelling & grammar corrections.



Source: <https://www.nytimes.com/2015/01/09/style/when-autocorrect-goes-wrong-and-so-so-right.html>

**Suggesting spelling & grammar
corrections.**

**We need a language model that
can tell which word is likely to
come next in a sentence.**

Language Models: **Bag-of-Words**

Represent a set of **documents**
using the **collection (bag)** of
words contained in all of them.

- **Document represented as a set of words**

Class Label

1. Stocks rallied on Monday, with major indexes gaining 1% as optimism persisted over the first quarter earnings season.

Business

2. Heavy rain continued to pound much of the east coast on Monday, with flood warnings issued in New York City and other locations.

Weather

- **The words appearing in each document represent its features!**
- **Build a dataset in which each document has a class label.**
- **Train a model (DT, Deep Learning, etc.) to classify a document based on its words.**

Language Models: **N-Grams**

While some words can be common among different subjects, some phrases (that consist of n consecutive words) are more specific to one subject.

Quarter

This is a common word in economics and in sports!

We can use a sequence of 4 words instead of a single word

- **First quarter earnings report**

More specific to business

- **Fourth quarter touchdown passes**

More specific to sports

This 4-word sequence is called an n-gram (4-gram).

The n-gram model can tell us if a phrase/short sentence is likely a correct English phrase.

- **A black cat.**

**Is likely proper English
because it appears in many
training examples.**

- **Black cat a.**

**Is likely not proper English
because it appears in zero
training examples.**

Language Models: **Part-of-Speech Tagging**

**A native English speaker can directly say that “a black cat” is sound because it follows a familiar pattern
(article-adjective-noun)
but “Black cat a” does not follow any familiar pattern.**

- **A black cat.**

- **Black cat a.**

If we can tag words with their POS, then define sound sentence patterns, this can help us identify proper sentences in a more generic way without having to have an example of each possible sentence in our data corpus.

- A black cat.**

- Black cat a.**

This is similar to building a KB of language rules. More sophisticated language models consider the grammar rules.

POS Tags

Tag	Word	Description	Tag	Word	Description
CC	<i>and</i>	Coordinating conjunction	PRP\$	<i>your</i>	Possessive pronoun
CD	<i>three</i>	Cardinal number	RB	<i>quickly</i>	Adverb
DT	<i>the</i>	Determiner	RBR	<i>quicker</i>	Adverb, comparative
EX	<i>there</i>	Existential there	RBS	<i>quickest</i>	Adverb, superlative
FW	<i>per se</i>	Foreign word	RP	<i>off</i>	Particle
IN	<i>of</i>	Preposition	SYM	<i>+</i>	Symbol
JJ	<i>purple</i>	Adjective	TO	<i>to</i>	to
JJR	<i>better</i>	Adjective, comparative	UH	<i>eureka</i>	Interjection
JJS	<i>best</i>	Adjective, superlative	VB	<i>talk</i>	Verb, base form
LS	<i>I</i>	List item marker	VBD	<i>talked</i>	Verb, past tense
MD	<i>should</i>	Modal	VBG	<i>talking</i>	Verb, gerund
NN	<i>kitten</i>	Noun, singular or mass	VCN	<i>talked</i>	Verb, past participle
NNS	<i>kittens</i>	Noun, plural	VBP	<i>talk</i>	Verb, non-3rd-sing
NNP	<i>Ali</i>	Proper noun, singular	VBZ	<i>talks</i>	Verb, 3rd-sing
NNPS	<i>Fords</i>	Proper noun, plural	WDT	<i>which</i>	Wh-determiner
PDT	<i>all</i>	Predeterminer	WP	<i>who</i>	Wh-pronoun
POS	<i>'s</i>	Possessive ending	WP\$	<i>whose</i>	Possessive wh-pronoun
PRP	<i>you</i>	Personal pronoun	WRB	<i>where</i>	Wh-adverb
\$	<i>\$</i>	Dollar sign	#	<i>#</i>	Pound sign
“	<i>‘</i>	Left quote	”	<i>,</i>	Right quote
(<i>[</i>	Left parenthesis)	<i>]</i>	Right parenthesis
,	<i>,</i>	Comma	.	<i>!</i>	Sentence end
:	<i>;</i>	Mid-sentence punctuation			

Fig 23.1, Russell & Norvig's Textbook

Considering **POS** and familiar **sentence patterns** is similar to building **a KB of language rules**.
More sophisticated language models consider **grammar rules**.

This approach can facilitate tasks such as **auto translation**.



English <**adjective then noun**>

Spanish <**noun then adjective**>