

NATURAL LANGUAGE PROCESSING

Niharika Krishnan

Before we get started....

NIHARIKA KRISHNAN

- Machine Learning Engineer, TCS
 - Build Chatbots for a living!
- Lead at PyLadies Chennai
 - Community of 100+ women tech enthusiasts
- Tech speaker
 - PyCon India'19, Canada'19
 - Google Women Techmakers,
- AI and NLP enthusiast
- Always willing to learn more!



WHAT TO EXPECT ?

TECH-TALK

Cognitive
Computing

Natural Language
Processing

NLP Around Us

HANDS-ON

NLP Techniques
+
Sentiment Analysis

TECH-TALK

CHATBOTS
Why, What, How?

HANDS-ON

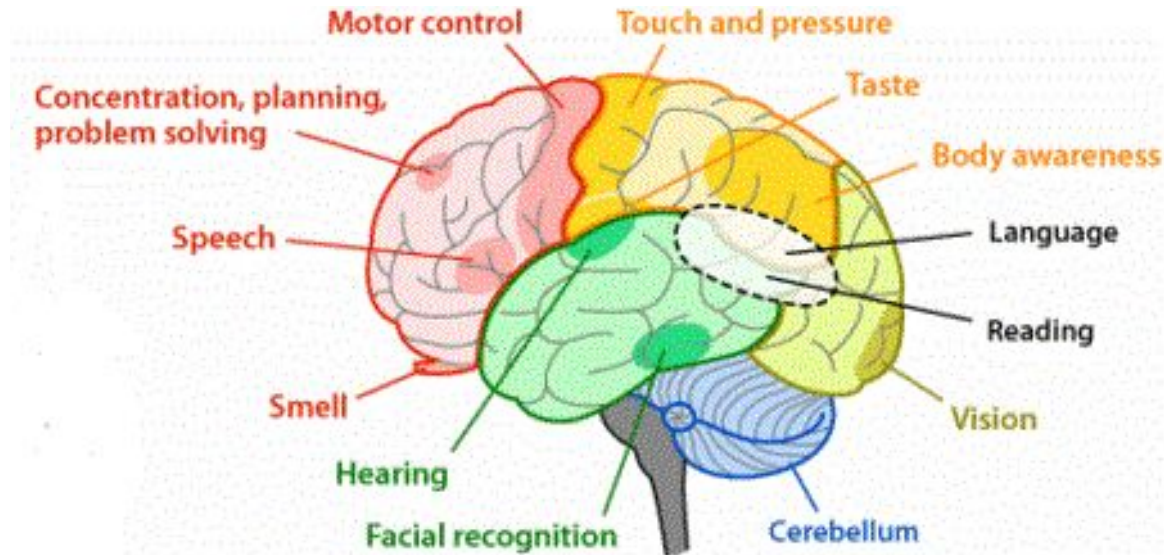
Building a chatbot
using Python

INDUSTRY CONNECT

Upskilling and
staying relevant

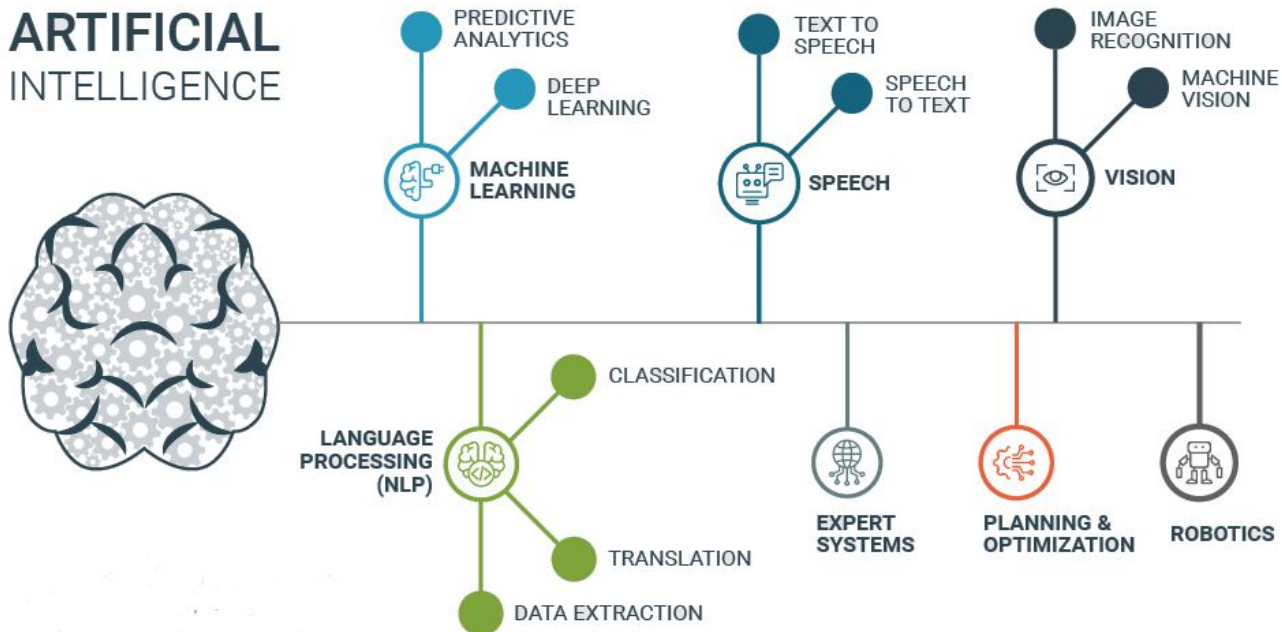
Cognitive Computing

200,000 years and counting....

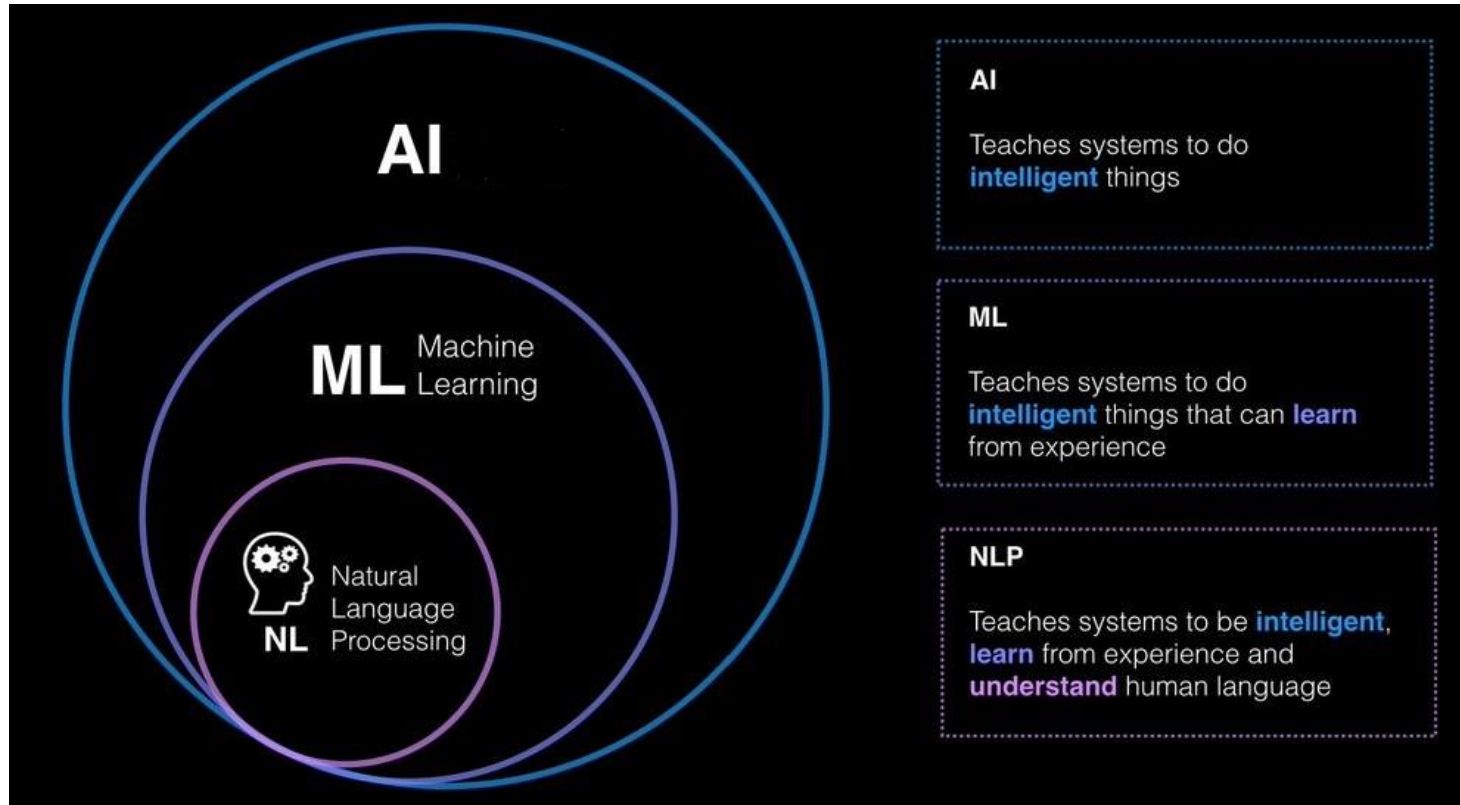


70 years and counting....

ARTIFICIAL INTELLIGENCE



AI vs ML vs NLP ?

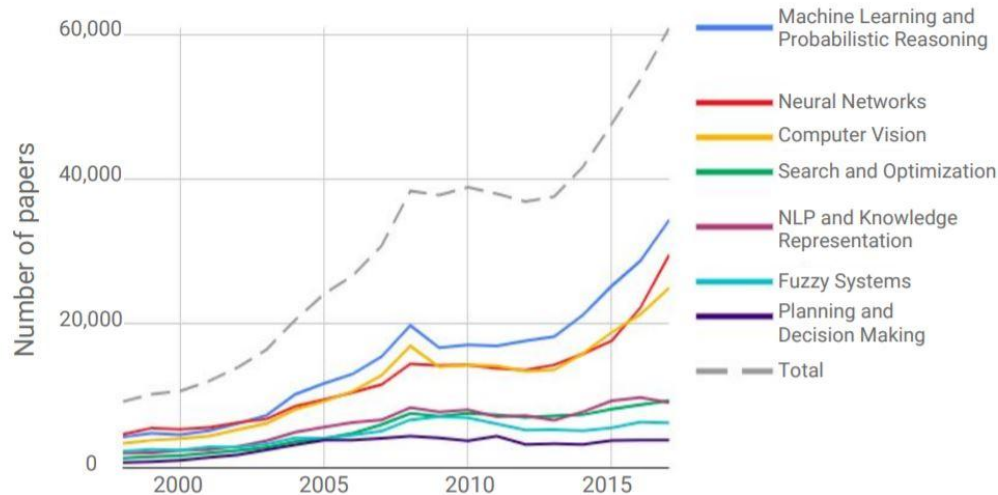


What is Natural Language Processing?

- Natural Language Processing (NLP) is a branch of artificial intelligence that deals with the interaction between computers and humans using the natural language.
- **Objective:** Read, decipher, understand, and make sense
- $\text{NLP} = \text{Computer Science} + \text{AI} + \text{Computational Linguistics}$

Number of AI papers on Scopus by subcategory (1998–2017)

Source: Elsevier



What makes NLP so
hard?

Standard \ Nonstandard

Examples

Standard

- am not, is not, has not
- very good
- very
- to play a trick
- you all

Nonstandard

- ain't
- cool
- damn
- to pull one's leg
- ya'll



“jaguar” can refer to a car or to an animal (Disambiguation)

I Scream vs Ice Cream (Phonetics)

AMBIGUITY



How we deal with text data.

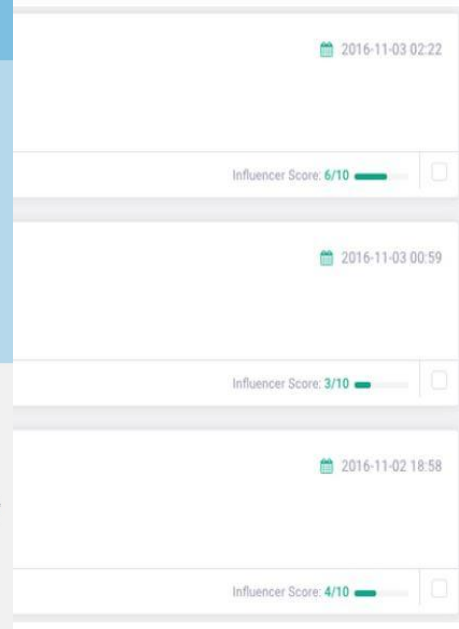
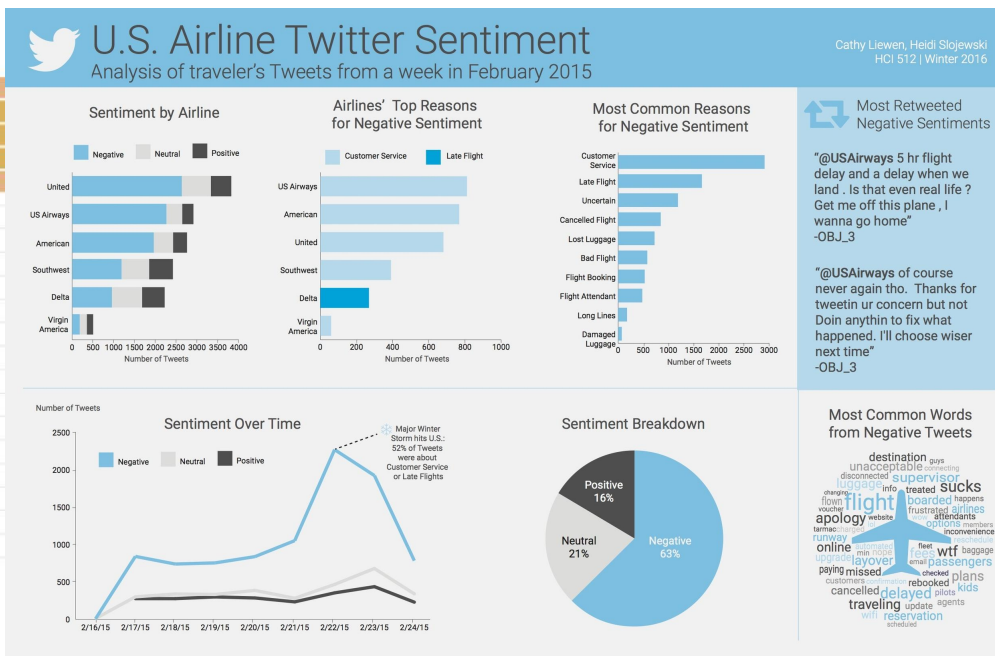
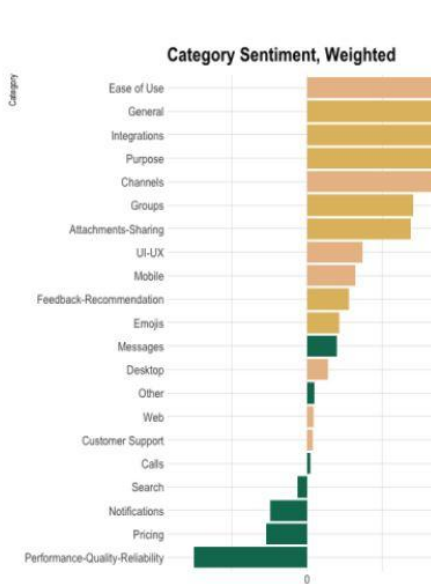
“The key to NLP is data”

The more data you collect, the more you can correct your algorithm's mistakes, and reinforce its correct answers. With unlimited data and unlimited compute, we would have perfect NLP today.

NLP is all around us

Sentiment Analysis

- **3.5 billion people** are active social media users - 45% of the world's population
- Every single minute of the day, **500,000 Tweets** and **510,000 Facebook comments**
- Checked Zomato before dining out? Read amazon reviews before buying?



Spell Check / Auto Complete

- Spellchecker points to spelling errors and possibly suggests alternatives.
- Autocorrector automatically picks the most likely word
- Auto complete: Understands context + grammar + provides suggestions
- Levenshtein Distance Algorithm - Insertion, Deletion, Substitution (Edit Distance)
- Python Packages: Symspell vs Janspell vs KenLM
- Grammarly, Google Doc, Plagiarism Detection

natural lanugageee processing

All News Images Books Videos More

About 53,40,00,000 results (0.64 seconds)

Showing results for **natural language processing**
Search instead for natural lanugageee processing

New Message

Recipients

Subject

Hi,

Looking forward to hearing from you

Text Mining

- Process of exploring sizeable textual data and find patterns generate valuable insights, enabling companies to make data-driven decisions
- Unstructured Data to Structured Data
- Packages: PyPDF, pydocx
- Information Retrieval
- Summarization - Semantic Map

Manufacturers

- Identify root causes of product issues quicker
- Identify trends in market segments
- Understand competitors' products

Government

- Identify fraud
- Understand public sentiments about unmet needs
- Find emerging concerns that can shape policy

Financial Institutions

- Use contact center transcriptions understand customers
- Identify money laundering or other fraudulent situations

Retail

- Identify profitable customers and understand the reasons for their loyalty
- Manage the brand on social media

Legal

- Identify topics and keywords in discovery documents
- Find patterns in defendant's communications

Healthcare

- Find similar patterns in doctor's reports
- Use social media to detect disease outbreaks earlier
- Identify patterns in patient claims data

Telecommunications

- Prevent customer churn
- Suggest up-sell/cross-sell opportunities by understanding customer comments

Life Sciences

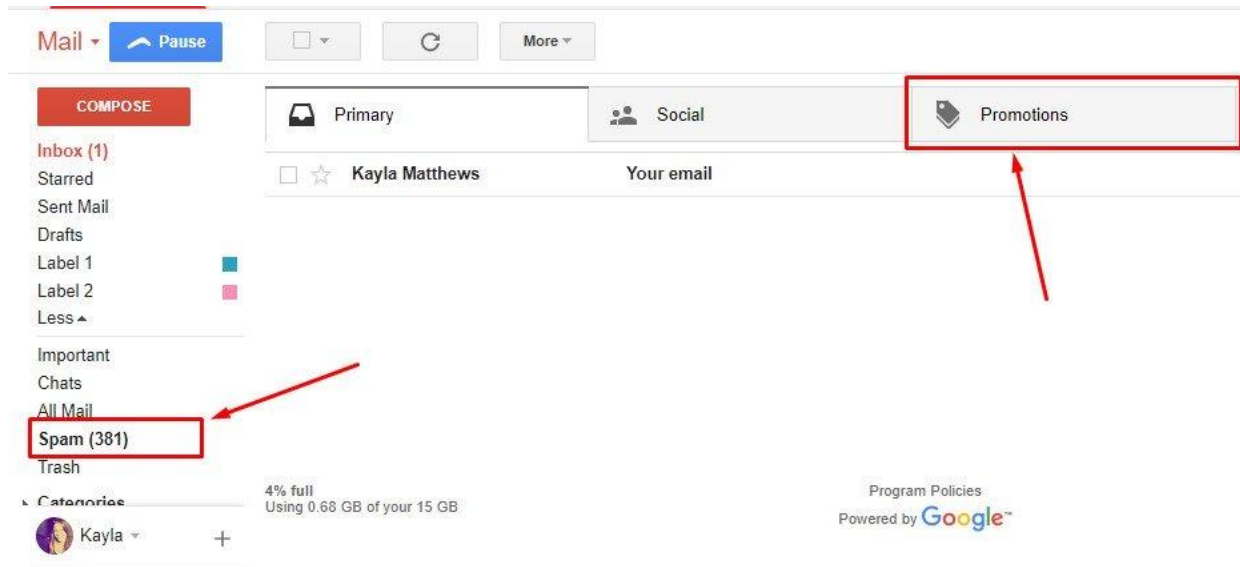
- Identify adverse events in medicines or vaccines
- Recommend appropriate research materials

Insurance

- Identify fraudulent claims
- Track competitive intelligence
- Manage the brand on social media

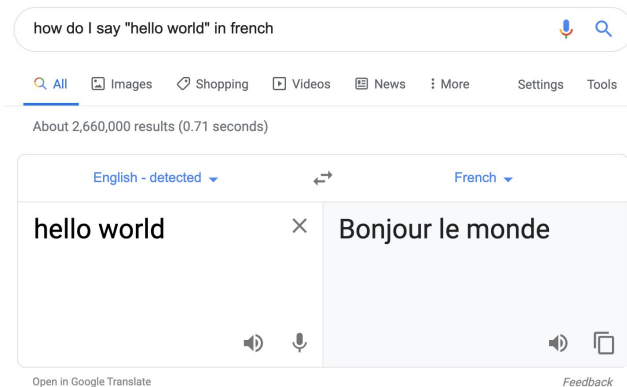
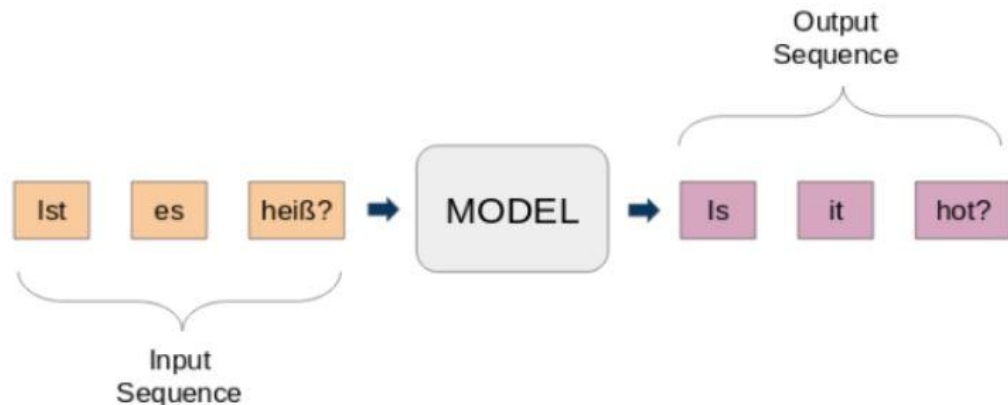
Text Classification

- Process of assigning tags or categories to text according to its content
- Document Classification, Sentence Classification
- Customer Support, Sentiment Analysis, Content Recommendation
- Python Packages: Fasttext, Flair



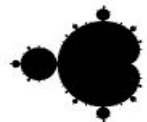
Machine Translation

- Translating one source language or text into another language
- Interpret & analyze all of the elements in the text and how each word may influence another. Consider grammar, syntax (sentence structure), semantics (meanings) in the source and target languages, as well as familiarity with each local region.
- Concepts of Deep Neural Networks - Recurrent Neural Networks - LSTM
- Seq2seq: Encoder ~ Decoder



Python Packages

Natural
Language
ToolKit



TextBlob

For NLP Preprocessing



flair

*fast*Text

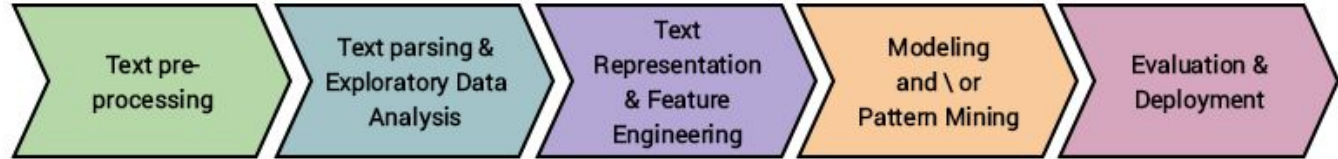
spaCy



GENSIM
topic modelling for humans



Lifecycle



**“A machine learning model
is only as good as the data it
is fed.”**

Text Preprocessing

Text Normalization

NLP is very Interesting to learn!!!

- Converting all letters to lower or upper case - `str.lower()` / `str.toupper()`
- Converting numbers into words or removing numbers - `replace/regex`
- Removing punctuations, accent marks and other diacritics - `regex`
- Removing white spaces - `str.strip()`
- Expanding abbreviations - `str.replace()`

Natural language processing is very interesting to learn

Stopword Removal

- Common language articles, pronouns and prepositions such as “and”, “the” or “to”
- Has no impact towards NLP
- There is no Universal list of stopwords → Domain specific
- Beware of NOTs

Sample text with Stop Words	Without Stop Words
GeeksforGeeks – A Computer Science Portal for Geeks	GeeksforGeeks , Computer Science, Portal ,Geeks
Can listening be exhausting?	Listening, Exhausting
I like reading, so I read	Like, Reading, read

Tokenization

- count all words in a piece of text
- occurrence matrix for the sentence or document, disregarding grammar and word order
- word frequencies or occurrences are then used as features for training a classifier.

```
from nltk.tokenize import sent_tokenize, word_tokenize

data = "All work and no play makes jack a dull boy, all work and no play"
print(word_tokenize(data))
```

['All', 'work', 'and', 'no', 'play', 'makes', 'jack', 'dull', 'boy', ',', 'all', 'work', 'and', 'no', 'play']

```
from nltk.tokenize import sent_tokenize, word_tokenize

data = "All work and no play makes jack dull boy. All work and no play makes jack a dull boy."
print(sent_tokenize(data))
```

['All work and no play makes jack dull boy.', 'All work and no play makes jack a dull boy.']

- End up removing punctuations: Dr. → Dr
- Hyphens? Parenthesis? → PROBLEM!!

Stemming

- Reducing inflection in words to their root forms such as mapping a group of words to the same stem even if the stem itself is not a valid word in the Language
- Integral to search queries and information retrieval
- Stem words, sentences, documents
- English Stemmers, Non English Stemmers
- Shorten the lookup & normalize

```
# stemming
from nltk.stem.porter import PorterStemmer
porter = PorterStemmer()

def tokenizer_porter(text):
    return [porter.stem(word) for word in text.split()]

tokenizer_porter('It has become appallingly obvious \
that our technology has exceeded our humanity.')
```



Lemmatization

- Reduces the inflected words properly ensuring that the root word belongs to the language. Root word is called Lemma.
- Consideration the context of the word
- NLTK: WordNetLemmatizer built on WordNet Database
- Part-of-speech parameter to a word (whether it is a noun, a verb, and so on) it's possible to define a role for that word in the sentence

```
# import these modules
from nltk.stem import WordNetLemmatizer

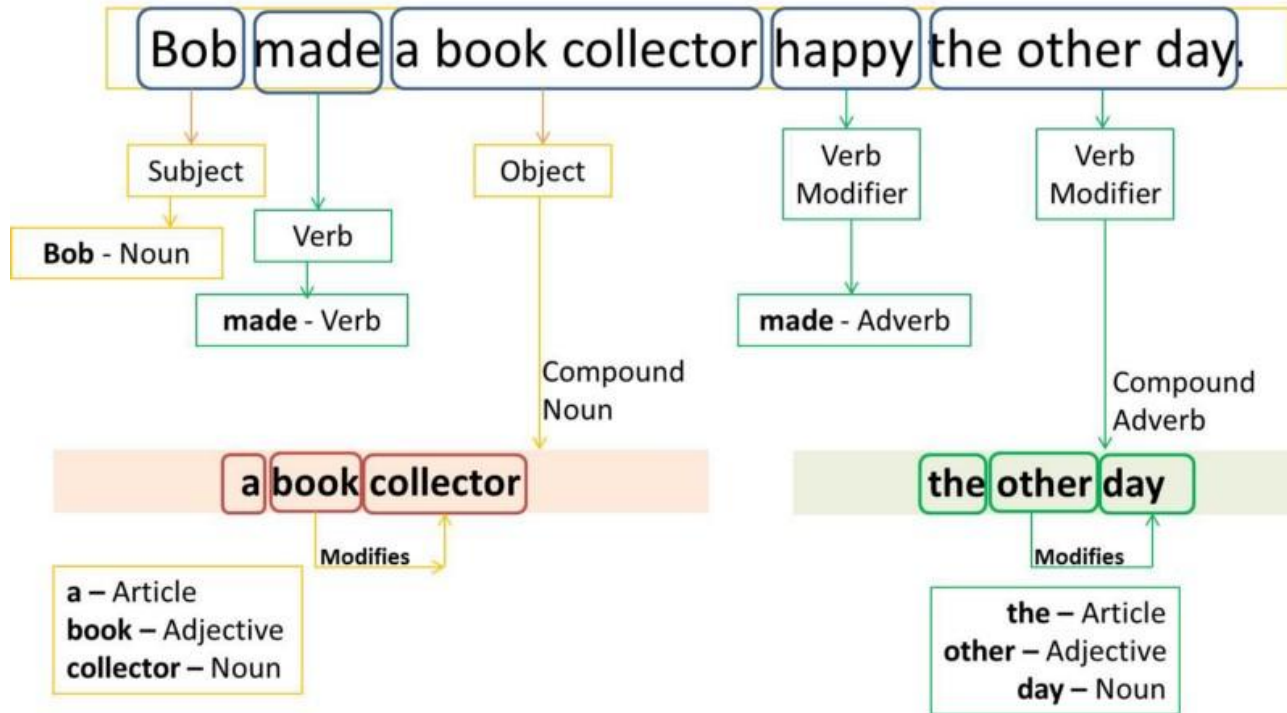
lemmatizer = WordNetLemmatizer()

print("rocks :", lemmatizer.lemmatize("rocks"))
print("corpora :", lemmatizer.lemmatize("corpora"))
```



POS Tagging

- Part-of-speech tagging aims to assign parts of speech to each word of a given text (such as nouns, verbs, adjectives, and others) based on its definition and its context



N-Grams

- n-gram is a contiguous sequence of n items from a given sample of text or speech.
- ngram model models sequence, i.e., predicts next word (n) given previous words (1, 2, 3, ..., n-1). multiple gram (bigram and above) captures context
- **Rule of thumb:** trigram is a common choice with large training corpora (millions of words), whereas a bigram is often used with smaller ones.

This is Big Data AI Book

Uni-Gram

This

Is

Big

Data

AI

Book

Bi-Gram

This Is

Is Big

Big Data

Data AI

AI Book

Tri-Gram

This Is Big

Is Big Data

Big Data AI

Data AI Book

FEATURE ENGINEERING

Bag of Words/One Hot Encoding/Document Term Matrix

- count all words in a piece of text
- occurrence matrix for the sentence or document, disregarding grammar and word order
- word frequencies or occurrences are then used as features for training a classifier

	I	love	dogs	hate	and	knitting	is	my	hobby	passion
Doc 1	1	1	1							
Doc 2	1		1	1	1	1				
Doc 3					1	1	1	2	1	1

- Absence of semantic meaning & context
- What if “the”, “a” are the most repeated?
- **Solution:** Stop word + TF-IDF

Term Frequency - Inverse Document Frequency

- Numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.
- TF = Bag of Words; IDF = Measures the importance of the term across the document.
Closer to 0, the more common the word is
- TF-IDF = TF * IDF (Float value)
- 83% of text-based recommender systems in digital libraries use tf-idf.

$$\text{tfidf}_{i,j} = \text{tf}_{i,j} \times \log\left(\frac{N}{\text{df}_i}\right)$$

$\text{tf}_{i,j}$ = total number of occurrences of i in j

df_i = total number of documents (speeches) containing i

N = total number of documents (speeches)

	I	love	dogs	hate	and	knitting	is	my	hobby	passion
Doc 1	0.18	0.48	0.18							
Doc 2	0.18		0.18	0.48	0.18	0.18				
Doc 3					0.18	0.18	0.48	0.95	0.48	0.48

Word Vectors

- Convert input into numerical form that deep neural network can process as inputs

“How to plot dataframe bar graph”

- CBOW - Features: {how,to,plot,bar,graph}, Predict: {dataframe}
 - Predicts current word with the help of neighbouring words
- Skip Gram: Features: {dataframe}, Predict: {how,to,plot,bar,graph}
 - Predicts neighbouring words based on current word



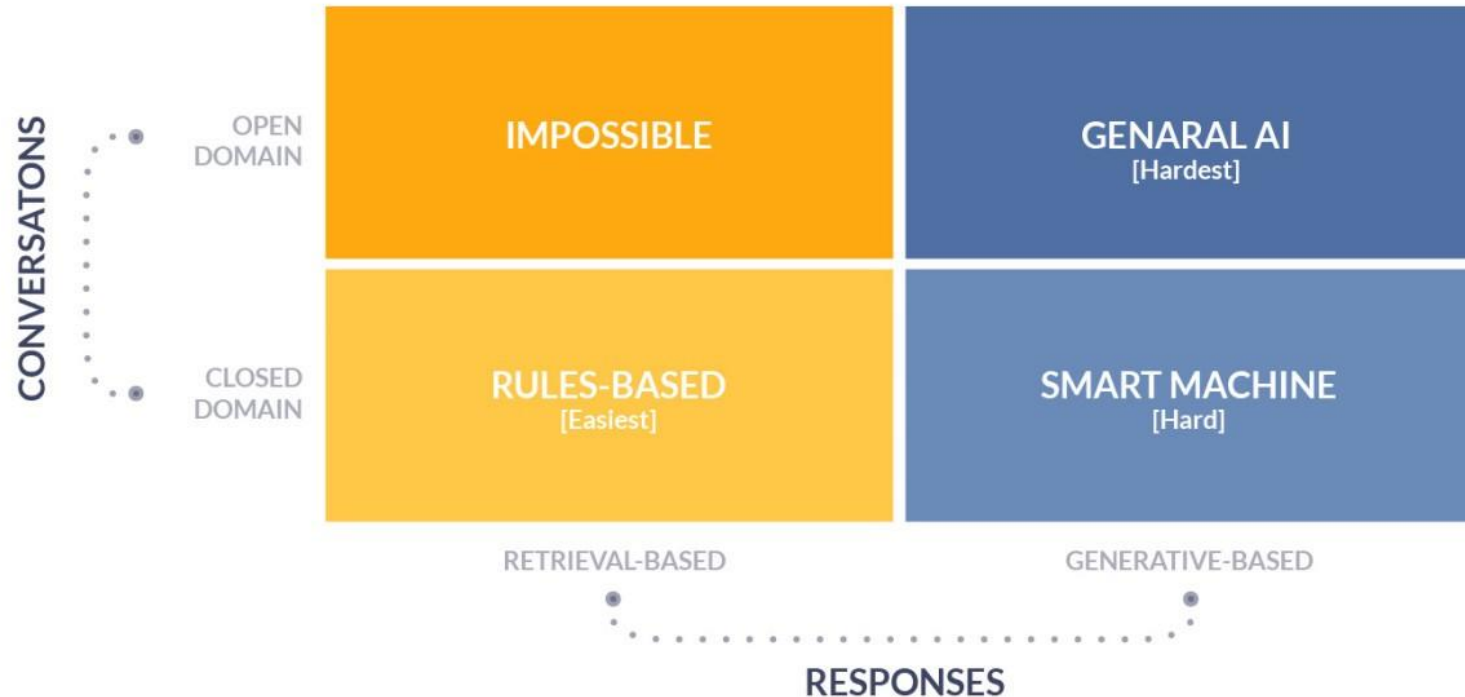
TUTORIAL

CHATBOTS

Chatbots - What, Why, Who, How?

- Artificial intelligence (AI) software that can simulate a conversation with a user
- FAQ/Simple Information Lookup
 - User: Asks Question, Chatbot: Gives Answers, No Machine Learning
- Guided Conversation Form
 - Chatbot: Asks Question, User: Gives Answers, Machine Learning → Extract Entities
- Information Lookup
 - User: Asks Question, Chatbot: Gives Answers, Information from DB/ Web services
 - Machine Learning → Extract Entities
- Smart Home / Simple Information Lookup
 - Performs Predefined Tasks

CHATBOT CONVERSATION FRAMEWORK



Terminologies

Utterance

"Show me yesterday's financial news"

Entity

Entity

Intent: **showNews**

Verb

Noun

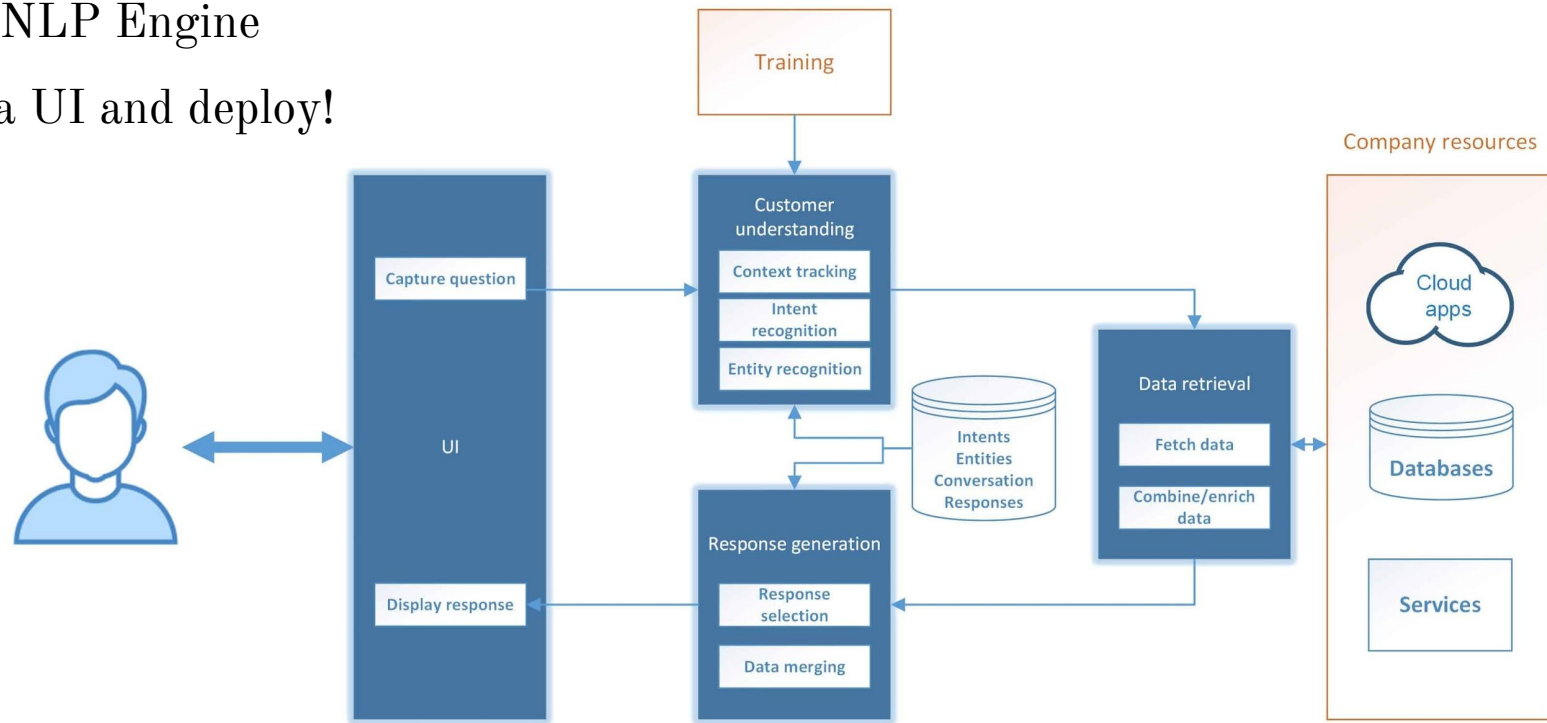
Utterances, Intents, Entities

- Utterances: Query Entered by the user
- Intents: An intent is the user's intention
 - Fasttext: Facebook AI Research Lab
 - Intent Classification: Greeting, Schedule, Resource, Thanking
- Entities: An entity modifies an intent
 - SpaCy: Named Entity Recognition

PERSON	People, including fictional.	LANGUAGE	Any named language.
NORP	Nationalities or religious or political groups.	DATE	Absolute or relative dates or periods.
FAC	Buildings, airports, highways, bridges, etc.	TIME	Times smaller than a day.
ORG	Companies, agencies, institutions, etc.	PERCENT	Percentage, including "%".
GPE	Countries, cities, states.	MONEY	Monetary values, including unit.
LOC	Non-GPE locations, mountain ranges, bodies of water.	QUANTITY	Measurements, as of weight or distance.
PRODUCT	Objects, vehicles, foods, etc. (Not services.)	ORDINAL	"first", "second", etc.

Steps to build a chatbot

- Build corpus / knowledge base
- Intents, entities, actions
- Train NLP Engine
- Build a UI and deploy!



amazon

Google

Microsoft



SAMSUNG




amazon alexa

Google

ASSISTANT

 Cortana


Siri

 Bixby

Let's build a FDP ML
chatbot!



“ NLP made it possible for machines to become our friend.

BILL GATES
MICROSOFT CO-FOUNDER
AT HUNTER COLLEGE IN NYC

INDUSTRY CONNECT

Staying relevant

Meetups



Conferences



Open Source Contribution

Hackathons



SMART INDIA
HACKATHON
2020



Google
Summer of Code

Women in Tech



QUESTIONS

Want to explore further?
Let's connect!



[linkedin.com/in/niharikakrishnan](https://www.linkedin.com/in/niharikakrishnan)



[@Nihaaarika](https://twitter.com/Nihaaarika)



[niharikakrishnan](https://github.com/niharikakrishnan)

Slide Deck: <https://github.com/niharikakrishnan/Talks>