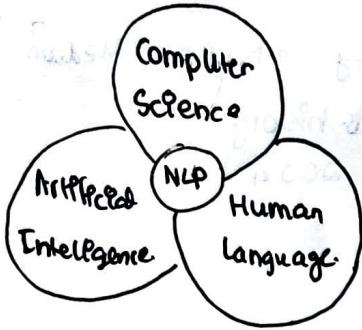


1

NATURAL LANGUAGE PROCESSING

NLP is a subfield of linguistics, computer science & AI concerned with interactions b/w computers & human language. In particular how to program computers to process & analyse large amounts of natural language data.



Need for NLP

- * Natural language evolved naturally in humans through use & repetition
- * Ordinary language without conscious planning or premeditation.

REAL WORLD APPLICATIONS

- * Contextual Advertisements
- * Email clients - spam filtering, smart reply.
- * Social Media - removing adult content, opinion mining
- * Search engines.
- * Chatbots

Common NLP Tasks

- * Text / document classification (News)
- * Sentiment Analysis (Reviews) (Customer review)
- * Information retrieval.
- * Parts of Speech Tagging (noun, verb, adjective).
- * Language Detection & Machine translation.
- * Conversational Agents.
- * Knowledge graphs & QA Systems (Who is prime min of India?)
- * Text summarization (Inshorts)
- * Topic modelling
- * Text generation (Keyboard)
- * Spell Checking & Grammar correction
- * Text parsing
- * Speech-to-Text

APPROACHES TO NLP

- * Heuristic Methods
- * Machine learning Based Methods
- * Deep learning Based Methods

Heuristic Method: Hard rules

Example Not efficient

- * Regular Expressions
 - Removing HTML tags
 - searching for salutation
- * Wordnet
 - Lexical dictionary
- * Open Mind common sense.

Advantages

- * Quick approaches.
- * Current scenario

Machine learning approach:

The Big advantage.

ML Workflow

Algorithms used

- Naive Bayes
- Logistic Regression
- SVM
- LDA
- Hidden Markov Models

Deep Learning Approach

* Architectures used

RNN

LSTM

GRU / CNN

Transformers (BERT)

Auto encoders

Challenges in NLP

* Ambiguity.

I saw the boy on the beach with my binoculars.

I have never taste a cake quite like that one before.

* Contextual Words

I ran to the store because we ran out of milk.

* Colloquialisms & Slang

Piece of cake, pulling your legs.

* Synonyms

* Irony, Sarcasm and tonal difference

That's just what I needed today!

* Spelling Errors

Poems, dialogue, scripts.

* Diversity

Assignment

1) make a blog on NLP at Medium

→ NLP & its history.

→ 1960s to 2024

End to End NLP Pipeline

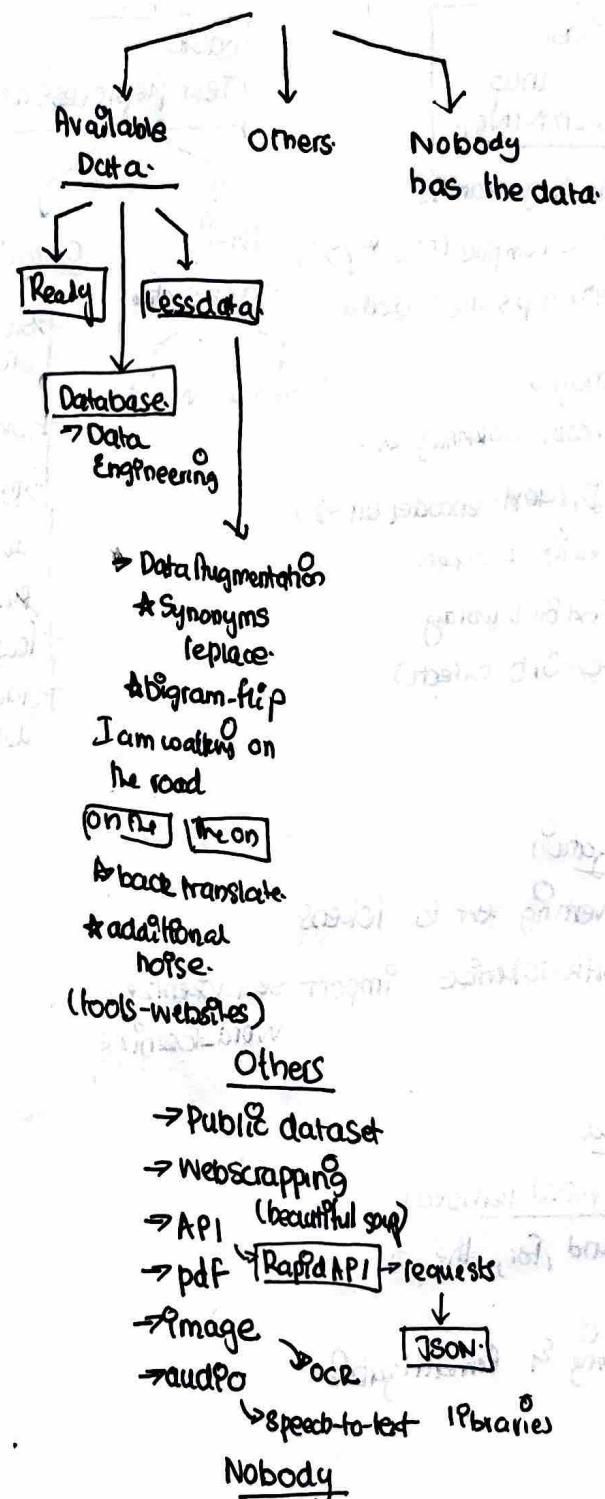
- Set of steps followed to build end-to-end NLP software
- NLP Software consists of following Steps:
 - ★ Data Acquisition
 - ★ Text Preparation
 - Text cleanup
 - Basic preprocessing
 - Advance preprocessing
 - ★ Feature Engineering
 - Bag of words
 - TFIDF
 - ★ Modelling
 - Applying algorithm - Model Building
 - Evaluation
 - ★ Deployment
 - Deployment
 - Monitoring
 - Model Update.

Points to remember:

- ★ It's not universal.
- ★ Deep Learning Pipelines are slightly different
- ★ Pipeline is Non-linear.

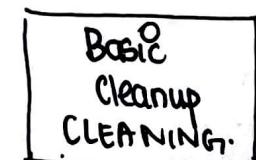
Data Acquisition

(2)



Text preparation

HTML tag
<p>Hello</p>

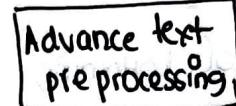
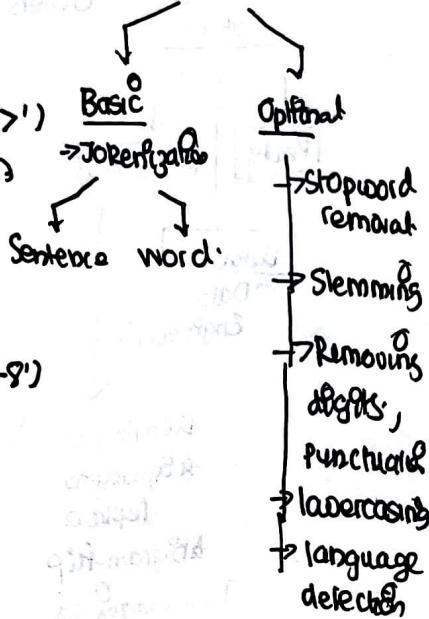
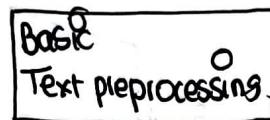


- HTML tag cleaning


```
p=re.compile(r'<.*?>')
return p.sub("",data)
```
- emoji's (UNICODE NORMALIZATION)


```
emojis=re.compile('utf-8')
```
- spelling checker


```
TextBlob library
textBlob.correct()
```



- POS tagging
- Parsing
- coreference resolution

Tokenization

Converting text to tokens
from nltk.tokenize import sent_tokenize,
word_tokenize.

Optional

* Stop word removal

and, for, the, an.

* Stemming & lemmatization

POSTAGGING

→ Parts of speech tagging

Feature Engineering

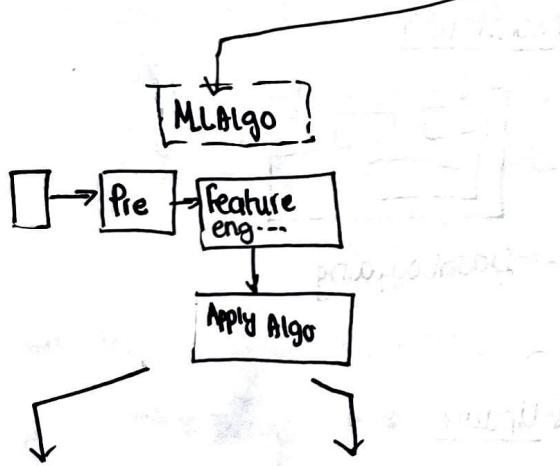
→ Converting text to numbers.

Sentiment analysis of movie reviews

<u>review text</u>	<u>sentiment</u>
1	0 Negative Sentiment
2	1 Positive Sentiment
3	
4	
5	
6	
7	
8	
9	
10	

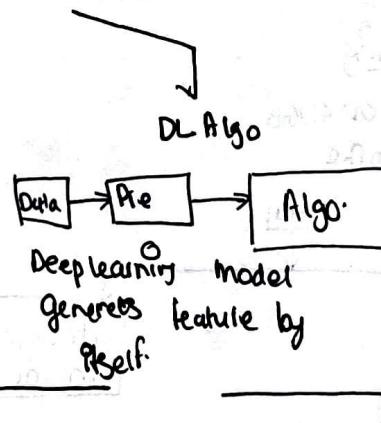
(50,000, 4) - Text vectorization -

# of five words	# of -ve words	# of Neutral words	sentiment
5	2	6	0
3	1	6	1



Advantage.
Can justify the interpretability

Disadvantage
Domain knowledge required to find features



Advantage
★ features are generated automatically

Disadvantage
★ Cannot justify why something is working.

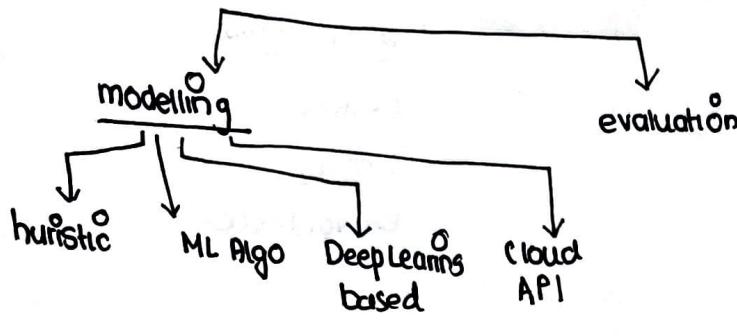
Bag of words

↳ Tfidf

↳ OHE

→ word2vec.

Modelling



which approach to apply?

→ Amount of data

→ Nature of problem.

Email Spam classifier

Less data → heuristic approach.

Increased data → ML algo.

Huge data → DL approach.

# ve	# -ve	from spam nr
3	6	0
4	6	1

Deployment

Deployment Monitoring Update

* Deployment

→ Deploy model on any cloud service as API (microservice)

Chatbot

Cloud API: Comes at a cost

DL → Transfer Learning

Eg: BERT → Trained on 40GB data on the Internet

Evaluation

Intrinsic evaluation

Eg: Text generation

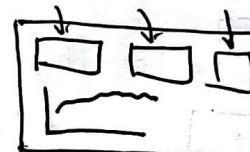
Next word while typing in keypad.

Metric: perplexity.

Tells how confused is NLP model

Extrinsic evaluation

* Monitoring



→ Dashboarding

* Update

NLP Pipeline Questions

1. Data Acquisition

a. From where would you acquire the data?

2. Text preparation

a. What kind of cleaning steps would you perform?

b. What kind of text processing would you apply?

c. Is advanced text processing required?

3 Feature Engineering

a. What kind of features would you create?

4. Modelling

a. What algorithm would you use to solve the problem at hand?

b. What intrinsic evaluation metrics would you use?

c. What extrinsic evaluation metrics would you use?

5 Deployment

a) ~~what~~ How would you deploy your solution into the entire product?

b) How & what things will you monitor?

c) What would be your model update strategy?