

NLP

into to nlp

Natural Language Processing (NLP) is a field of artificial intelligence that focuses on the interaction between computers and human languages.

Real world application

- contextual advertisement
- email spam
- social media filtration
- chatbot

Common task

- Text document classification
- Sentiment analysis
- info retrieval
- parts of speech tagging
- Language detection and machine translation
- Conversational agents
- knowledge graph and QA system
- Text summarization
- Topic modelling
- Text generation
- Spell checking and grammer correction
- speech to text

Approaches to nlp

Heuristic approach's

- Regular Expressions
- Wordnet
- open mind common sense

ML approach's

- Naive bayes
- Logistic Regression
- SVM
- LDA
- Hidden markov model

DL approach's

- RNN
- LSTM
- GRU
- Transformer
- Auto encoders

Challenges in nlp

Ambiguity

i saw the on the beach with my binoculars

Contextual words

i ran to the store because we ran out of milk -> word ran

Slang

pulling your leg

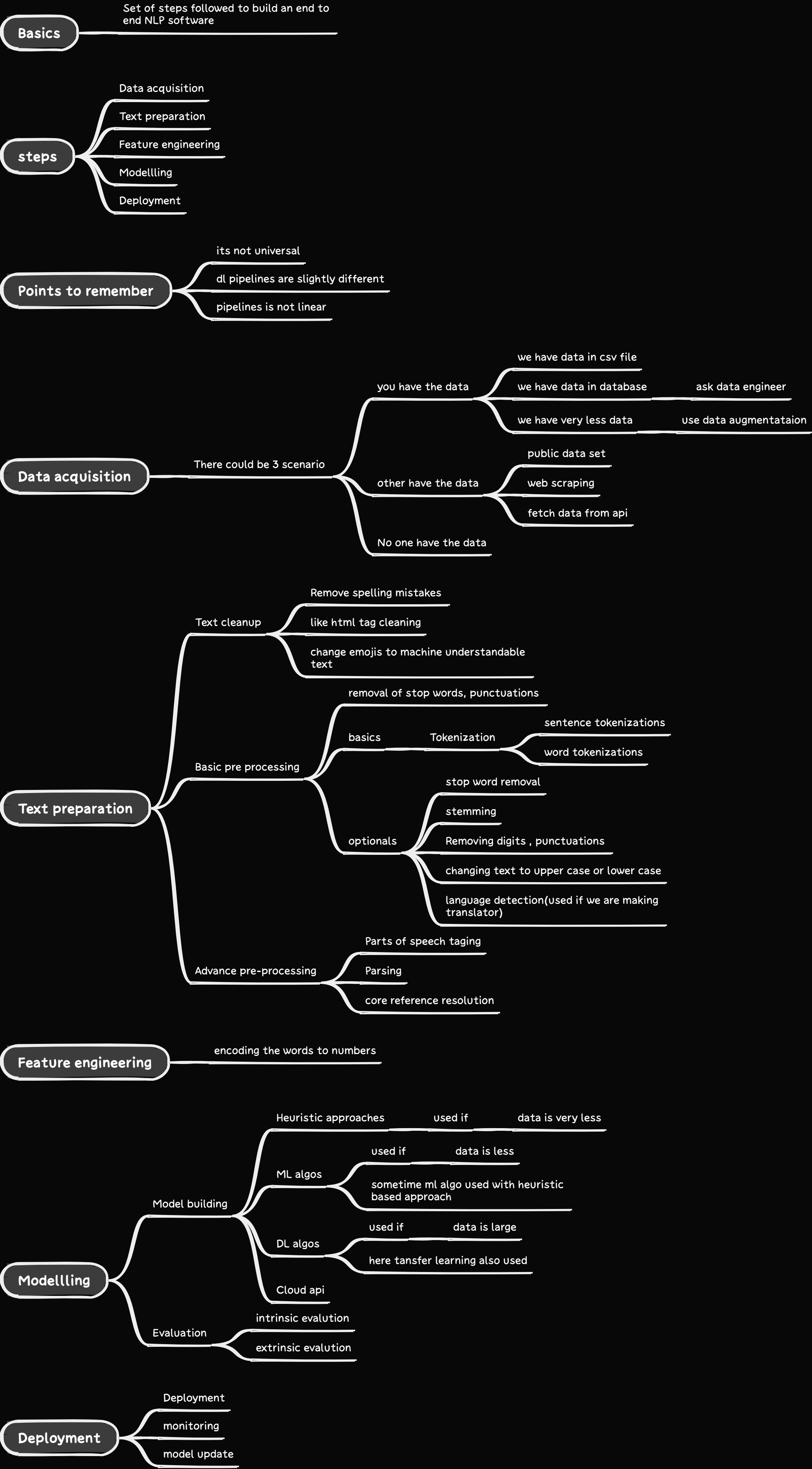
irony, sarcasm

spelling error

creativity

Diversity

NLP Pipeline



Text Pre-processing

Lowercasing

column df['cl'] = df['cl'].str.lower()

Removing

To remove something we use regular expression

use this tool to make regular expression

```
import re
def remove_html_tags(text):
    pattern = re.compile('<.*?>')
    return pattern.sub(r'', text)
```

df['cl'] = df['cl'].apply(remove_html_tags)

```
def remove_url(text):
    pattern = re.compile(r'https?://\S+|www\.\S+')
    return pattern.sub(r'', text)
```

df['review'] = df['review'].apply(remove_url)

Removing Punctuations

Punctuations

exclude = string.punctuation

```
def remove_punct(text):
    return text.translate(str.maketrans('', '', exclude))
```

df['cl'] = df['cl'].apply(remove_punct)

Chat word treatment

spelling corrections

from textblob import TextBlob

```
def spelling_correction(text):
    textBlb = TextBlob(incorrect_text)
    return textBlb.correct().string
```

df['cl'] = df['cl'].apply(spelling_correction)

Removing stop words

If you want to do parts of speech tagging then do not remove stop words

from nltk.corpus import stopwords

def remove_stopwords(text):

new_text = []

```
for word in text.split():
    if word in stopwords.words('english'):
        new_text.append("")
    else:
        new_text.append(word)
x = new_text[:]
new_text.clear()
return " ".join(x)
```

df['cl'] = df['cl'].apply(remove_stopwords)

Handling Emojis

Ways

Remove

```
import re
def remove_emoji(text):
    emoji_pattern = re.compile("["
        u"\U0001F600-\U0001F64F"
        # emoticons
        u"\U0001F300-\U0001F5FF"
        # symbols & pictographs
        u"\U0001F680-\U0001F6FF"
        # transport & map symbols
        u"\U0001F1E0-\U0001F1FF"
        # flags (iOS)
        u"\U00002702-\U000027B0"
        u"\U000024C2-\U0001F251"
    ]+", flags=re.UNICODE)
    return emoji_pattern.sub(r'', text)
```

df['cl'] = df['cl'].apply(remove_emoji)

Replace it with its meaning

import emoji

import re

```
def replace_emojis_with_meanings(text):
    # First, demojize the text to get :emoji_name: format
    demojized_text = emoji.demojize(text)
```

Use regex to find all :emoji_name: patterns and replace them with readable text

```
def convert_match(match):
    # Get the matched emoji name without colons
    emoji_name = match.group(1)
    # Replace underscores with spaces
    readable_name = emoji_name.replace('_', ' ')
    return readable_name
```

```
# Replace all :emoji_name: with readable text
readable_text = re.sub(r':([a-zA-Z0-9_]+):', convert_match, demojized_text)
```

return readable_text

df['cl'] = df['cl'].apply(replace_emojis_with_meanings)

Stemming and lemmatization

changing walking , walks , walked --> walk

Stemming is the process of reducing inflection in words to their root form such as mapping a group of words to the same stem even if the stem itself is not a valid word in the language

that why when you need to show text to user we do not use stemming -> in that case we use lemmatization(lemmatization also get the root of word but it make sure the root word is valid in language)

lemmatization is slow as comapre to stemming

from nltk.stem.porter import PorterStemmer

```
ps = PorterStemmer()
def stem_words(text):
    return " ".join([ps.stem(word) for word in text.split()])
```

df['cl'] = df['cl'].apply(PorterStemmer)

lemmatization

code

```
import nltk
from nltk.stem import WordNetLemmatizer
wordnet_lemmatizer = WordNetLemmatizer()

sentence = "He was running and eating at same time. He has bad habit of swimming after playing long hours in the Sun."
punctuations="?!:.,'"
sentence_words = nltk.word_tokenize(sentence)
for word in sentence_words:
    if word in punctuations:
        sentence_words.remove(word)

sentence_words
print("{}(0:20){}(1:20)".format("Word", "Lemma"))
for word in sentence_words:
    print("{}(0:20){}(1:20)".format(word, wordnet_lemmatizer.lemmatize(word, pos='v')))
```

POS Tagging

Theory

machine learning techniques used for POS tagging

Hidden Markov Models or Conditional Random Fields

used where

name entity recognition

chatbots

word sense diaambiguation

Example

i ran to store because i ran out of milk-> here ran has 2 meaning

code

!pip install spacy

import spacy

nlp = spacy.load('en_core_web_sm')

doc = nlp(u"I will google about facebook")

doc.text 'I will google about facebook'

doc[-1] facebook

doc[2].pos_ 'VERB'

doc[2].tag_ 'VB'

spacy.explain('VB') 'verb, base form'

Tokenization

breaking text document into smaller parts it can be word , sentence , etc (mostly words)

Code

```
import spacy
nlp = spacy.load('en_core_web_sm')
```

doc4 = nlp("I am going to visit delhi")

```
for token in doc4:
    print(token)
```


Text Representation

