

Natural Language Processing (NLP)

Natural Language? : Human language . English, Hindi, French etc.

NLP :- Processing the natural language so that we can learn useful features.

NLU : Natural Language Understanding ; syntax, grammar, construction

NLG : Natural Language Generation ; Given an incomplete text we can complete it.

Use cases : Lots of use cases

1) Sentiment analysis : Given a review / text / article what is the sentiment? (+ve / -ve / neutral) } classification

2) Classification of documents

3) Document cluster } clustering Problem.

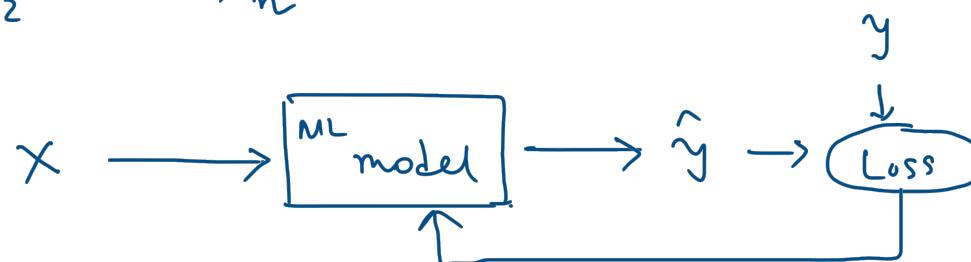
Revisiting the traditional ML:

features

x_1, x_2, \dots, x_n

targets

y



features are numerical / categorical
(oHE / IE)

(Supervised ML model)

Sentiment classification:- This can also be framed as supervised ML task.

→ I liked the product but it needs improvement.

(neutral)

→ I am absolutely delighted by the product

(+ve)

→ The product is not working after using for 2 days.

(-ve)

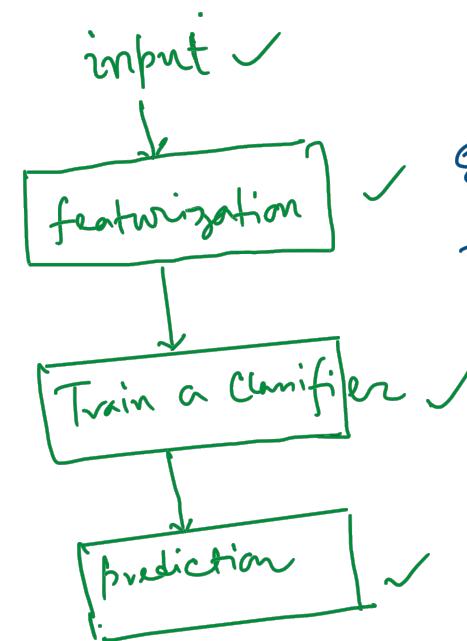
1-2 * (-ve) , 3 * (neutral) , 4+ * (+ve)

Training input

1. Review - 1
2. Review - 2
3. Review - 3
- :
- :
- m. Review - m

Output (Sentiment)

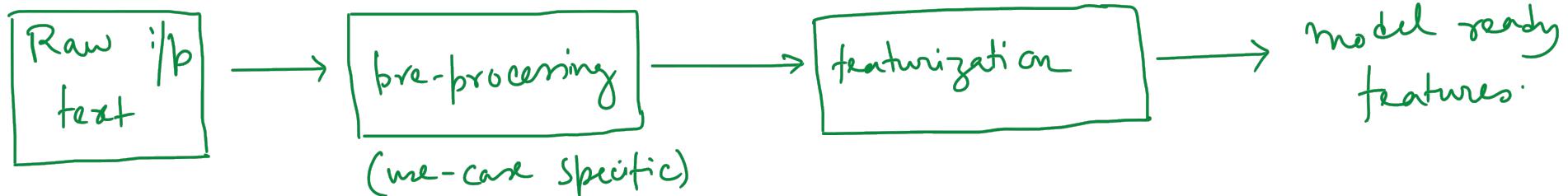
+ve
-ve
-ve
:
:
+ve



Extract features
from raw inputs

Given an unseen review can you predict the sentiment?

The first step of any NLP problem is to do the featurization



Text Preprocessing:

@skarmakar is teaching at @LogicMojo. Please sign up: <https://logicmojo.com>

Preprocessing Step-1: Removing punctuations and stopwords.

punctuations

;
:
,

:/ :) : (: :

Stopwords

the
and
on
one
is

English stopwords.

Domain Specific Stopwords: Suppose you are creating a sentiment classifier for the review of air-fryers.

Total set of stopwords : Language specific \cup Domain Specific.

Example:

I am sad, because the movie was bad. [raw]
→ (sad because movie bad)

The product is great. I am very happy to use the product. [raw]
→ product great very happy use product

Pre-processing step-2: Remove URLs, Handles.

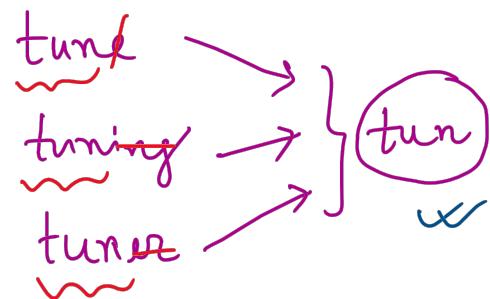
Pre-processing step-3: Lowercasing

Pre-processing step-4: Stemming | Lemmatization

Both Stemming and Lemmatization are text normalization techniques.

Stemming: → It is a crude heuristic process that chips off word endings (prefixes / suffixes) to get the root form.

- It doesn't care about grammar / meaning.
- It often produces non-dictionary words (stems)
- Porter Stemmer algorithm is the most famous stemming algo for english.



<u>Raw</u>	<u>Stemmed</u>
Caring	car
Cars	car
Studies	Studi
better	better

Lemmatization :

- uses vocabulary and morphological analysis of words to return the "lemma" : the bare / root word. (dictionary word)
- Takes into account the parts of speech (POS) and the context.

<u>Raw</u>	<u>Lemmatized word</u>
caring	care
studies	study
better	good

Corpus: Set of texts, intended to use for specific NLP task.

Ex:- Set of wikipedia articles.

Set of reviews of a product in Amazon.

Financial times articles published between 2020 - 2025.

Corpus are used to train NLP models

One of the most common technique to collect corpus is to do
web-scraping (crawling different website and dump textual data)

Most of the sophisticated NLP models are trained on publicly
available corpus (wikipedia article).

Tokenization

Corpus: $\{ \text{document-1}, \text{document-2}, \dots, \text{document-n} \}$

document: $\{ \text{sentence-1}, \text{sentence-2}, \dots \}$

Sentence: I liked the movie @Oppenheimer .

Tokens: Building blocks of a sentence.

In the strict sense tokens are not simple words but
tokens can consist a part of word.

(liked) \rightarrow lik ed

Vocabulary :

Corpus →

Apply preprocessing → processed text → tokenization (breaking into words / tokens)

Create a set of
unique words in the
pre-processed corpus

Set of unique words : Vocabulary (∇)

Ex:-

- 1) I love to read horror stories. } [love, read, horror, story]
- 2) I hate my engineering stream. } [hate, engineer, stream]
- 3) I love my country. } [love, country]

$V = [\underline{\text{love}}, \text{read}, \text{horror}, \text{story}, \text{hate}, \text{engineer}, \text{stream}, \text{country}]$