

NLP \Rightarrow Natural Language processing.

Application

- * Auto complete.
- * Spam / ham email filters.
- * Voice Assistant (google, Alexa, Siri)
- * Language translation.
- * Chat bots.
- * Document summarization.

{ NLP is a field of AI, that enables the machine to }
understand human language and assist accordingly. }

Libraries:

→ spacy, Gensim, NLTK, sklearn, tensorflow,
pytorch, hugging face.

Regular Expression → RegEx (python)

regex101.com ⇒ use this for practice (dynamic page) (x)

Steps of NLP:

* Data Preprocessing - 1

- Tokenization
- stop words
- stemming.
- Lemmatization.

* Data preprocessing - 2:

- Bag of words (BOW)
- TFIDF
- Unigrams, Bigrams

* Data preprocessing - 3

- Word2Vec
- Average Word2Vec

* Traditional NLP models.

- Naive Bayes

{ Data preprocessing = Text preprocessing }

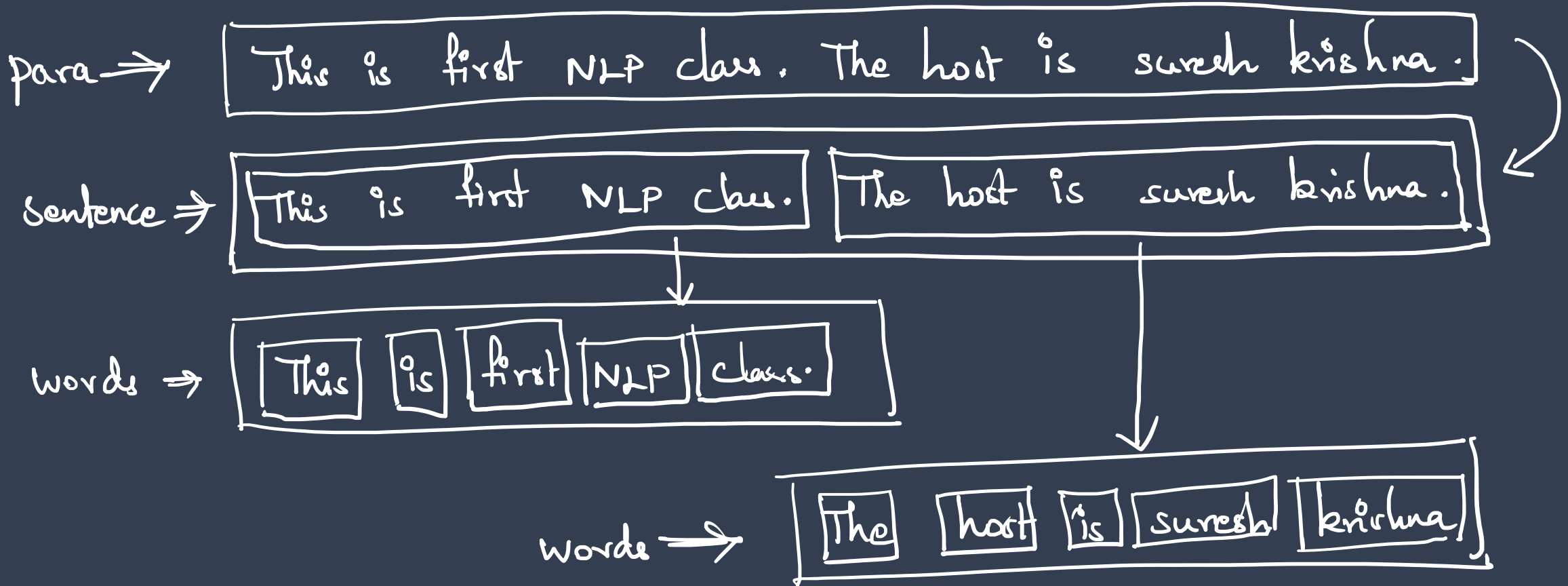
Spam / ham mail filter:

Mail subject	Mail body	class
NLP Doubts	Hi surah, I have a doubt.....	ham
<u>W</u> hurray!! Lottery	Woow!!! You won <u>a</u> lottery. claim	spam

10,00,000
↑
{ size of }
the data }

Data preprocessing - 1 :

① Tokenization — split & convert paragraph \Rightarrow sentence \Rightarrow words.



⑥ stop words:

⇒ Technique by which we eliminate the words that does not contribute to the prediction.

⇒ NLTK has predefined set of stopwords, but we can create our own stopwords.

→ Ex: How is the climate today ?
↳ stopwords

④ stemming \Rightarrow process of finding the root word.

programming / program , programs , programmed , programmer
↳ Rootword = Program

achieve, achieving, achieved

↳ Root word = Achiev

↳ This doesn't make sense

disadvantage
of
stemming
(X)

(X) Advantage ⇒ Faster in computation.

(*) Lemmatization ⇒ overcomes the disadvantage of stemming.

{ Uses the dictionary of predefined words from the
library to compare and find the meaningful
word from its library. }

Advantage → It gives meaningful root word from dictionary
disadvantage → slow in terms of computation.

Application of stemming & lemmatization:

stemming

- spam / ham.
- sentiment /
review classification.

lemmatization

- language translation.
- text summarization.
- chat bots.