

NLP : technology used by machine to understand, analyze, manipulate, and interpret Human's language

↳ **NLG** : helps machine to understand & Analyze human language

Text Mining : Analyze and extract patterns from the text data.

NLU : it acts as a translator that converts computerized data into human language.

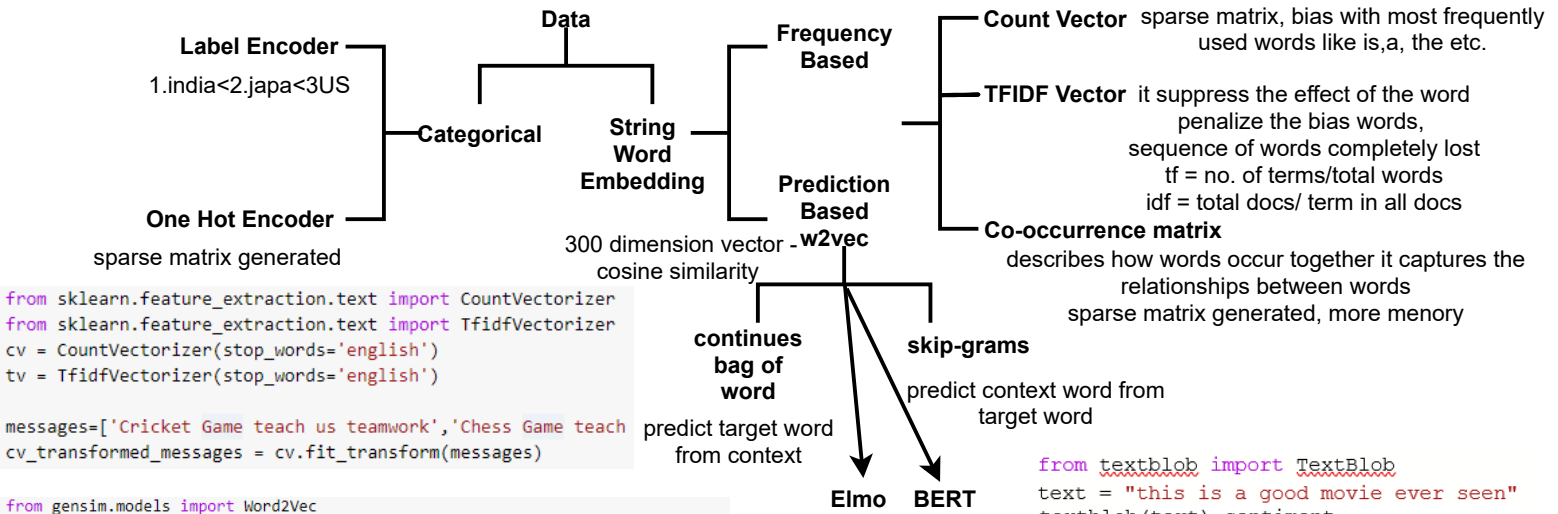
NLP Concepts:

```
1 import spacy
2 nlp = spacy.load('en_core_web_sm')
3 doc = nlp(u'It cost $5 million and Companies Email is ansarijuned4all@gmail.com and i live in india, right now I am eating mango")
4 print("Sentence Token:", [sentence for sentence in doc.sents])
5 # Sentence Token: [It cost $5 million and Companies Email is ansarijuned4all@gmail.com, and i live in india, , right now I am eating mango]
6 print("Token:", [token.text for token in doc])
7 # Token: ['It', 'cost', '$', '5', 'million', 'and', 'Companies', 'Email', 'is', 'ansarijuned4all@gmail.com', 'and', 'i', 'live', 'in', 'india', ',', 'right', 'now', 'I', 'am', 'eating', 'mango']
8 print("Token Index:", [token.i for token in doc])
9 # Token Index: [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21]
10 print("\n")
11 print("Named Entity Recongnisation(NER):")
12 for token in doc.ents:
13     print(token.text, ":", token.label_, ":", token.label_, ":", Explain: ", spacy.explain(token.label_)")
14 """
15 Named Entity Recongnisation(NER):
16 $5 million ::Label(N E R):: MONEY :: Explain:: Monetary values, including unit
17 Companies Email ::Label(N E R):: ORG :: Explain:: Companies, agencies, institutions, etc.
18 ansarijuned4all@gmail.com ::Label(N E R):: ORG :: Explain:: Companies, agencies, institutions, etc.
19 india ::Label(N E R):: GPE :: Explain:: Countries, cities, states
20 """
21 print("\n")
22 print("POS:", [token.pos_ for token in doc])
23 # POS: ['PRON', 'VERB', 'SYM', 'NUM', 'NUM', 'CONJ', 'PROPN', 'PROPN', 'AUX', 'X', 'CONJ', 'PRON', 'VERB', 'ADP', 'PROPN', 'PUNCT', 'ADV', 'ADV', 'PRON', 'AUX', 'VERB', 'NOUN']
24 print("tag:", [token.tag_ for token in doc])
25 # tag: ['PRP', 'VBD', '$', 'CD', 'CD', 'CC', 'NNPS', 'NNP', 'VBZ', 'ADD', 'CC', 'PRP', 'VBP', 'IN', 'NNP', ',', 'RB', 'RB', 'PRP', 'VBP', 'VBG', 'NN']
26 print("tag_:", [token.tag_ for token in doc])
27 # tag: ['pronoun, personal', 'verb, past tense', 'symbol, currency', 'cardinal number', 'cardinal number', 'conjunction, coordinating', 'noun, proper plural', 'noun, proper singular', 'verb, 3rd per
28 print("lemma:", [token.lemma_ for token in doc])
29 # lemma: ['-PRON-', 'cost', '$', '5', 'million', 'and', 'Companies', 'Email', 'be', 'ansarijuned4all@gmail.com', 'and', 'i', 'live', 'in', 'india', ',', 'right', 'now', '-PRON-', 'be', 'eat', 'mang
30 print("is_stop:", [token.is_stop for token in doc])
31 # is_stop: [True, False, False, False, False, True, False, False, True, False, True, True, False, True, False, False, False, True, True, True, False, False]
32 print("is_alpha:", [token.is_alpha for token in doc])
33 # is_alpha: [True, True, True, False, True, True, True, True, True, False, True, True, True, True, True, False, True, True, True, True, True, True]
34 print("is_punct:", [token.is_punct for token in doc])
35 # is_punct: [False, False, False, False, False, False, False, False, False, False, False, False, False, False, True, False, False, False, False, False, False]
36 print("like_num:", [token.like_num for token in doc])
37 # like_num: [False, False, False, True, True, False, False, False, False, False, False, False, False, False, False, False, False, False, False, False, False]
38 print("like_email:", [token.like_email for token in doc])
39 # like_email: [False, False, False, False, False, False, False, True, False, False, False, False, False, False, False, False, False, False, False, False, False]
40 print("doc.ents Length:", len(doc.ents))
41 #doc.ents Length: 4
42 """
```

Stemming	Lemmatization
<ul style="list-style-type: none">Produced by “stemmers”Produces a word's “stem”	<ul style="list-style-type: none">Produced by “lemmatizers”Produces a word's “lemma”
<ul style="list-style-type: none">am- amthe goin - the gohaving - hav	<ul style="list-style-type: none">am- bethe going - the goinghaving - have

Feature Extraction:

Corpus: it is a collection of documents **Vocabulary**: collection of unique words present in corpus **BOW**: total unique words in all the documents



```
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
cv = CountVectorizer(stop_words='english')
tv = TfidfVectorizer(stop_words='english')

messages=['Cricket Game teach us teamwork','Chess Game teach
cv_transformed_messages = cv.fit_transform(messages)
```

```
from gensim.models import Word2Vec
mymodel_cbow = Word2Vec(tokenized_sentences, min_count=1, size = 100)
mymodel_skipgram = Word2Vec(tokenized_sentences, min_count=1, size = 100,sg=1)
mymodel_cbow.similarity("python","Java")
mymodel_cbow.most_similar("Data")
```

NLP Libraries

Scikit-learn: It provides a wide range of algorithms for building machine learning models in Python.

Natural language Toolkit (NLTK): NLTK is a complete toolkit for all NLP techniques.

Pattern: It is a web mining module for NLP and machine learning.

TextBlob: It provides an easy interface to learn basic NLP tasks like sentiment analysis, noun phrase extraction, or pos-tagging, spelling correction.

Query: Query is used to transform natural language questions into queries in a database query language.
tga = pipeline(task="table-question-answering", model="google/tapas-base-finetuned-wtq")

SpaCy: SpaCy is an open-source NLP library which is used for Data Extraction, Data Analysis, Sentiment Analysis, and Text Summarization.

Gensim: Gensim works with large datasets and processes data streams.

GPT1.2

GPT-Neo: clone of gpt-3 but let dimension.

TS(Text to Text Transfer Transformer): Context + Answer → Generate Question

- AutomaticSpeechRecognitionPipeline
- ConversationalPipeline
- FeatureExtractionPipeline
- FillMaskPipeline
- ImageClassificationPipeline
- QuestionAnsweringPipeline
- SummarizationPipeline
- TextClassificationPipeline
- TextGenerationPipeline
- TokenClassificationPipeline
- TranslationPipeline
- ZeroShotClassificationPipeline
- Text2TextGenerationPipeline
- TableQuestionAnsweringPipeline

```
from textblob import TextBlob
text = "this is a good movie ever seen"
textblob(text).sentiment
```

```
from transformers import pipeline
```

```
!pip install autoviml
from autoviml.Auto_NLP import Auto_NLP
nlp_column = 'reviews'
target = 'target'
train_nlp, test_nlp, nlp_transformer, preds = Auto_NLP(
    nlp_column, train, test, target, score_type='balanced_accuracy',
    modeltype='Classification', top_num_features=50, verbose=2,
    build_model=True)
nlp_transformer.predict(test[nlp_column])
```