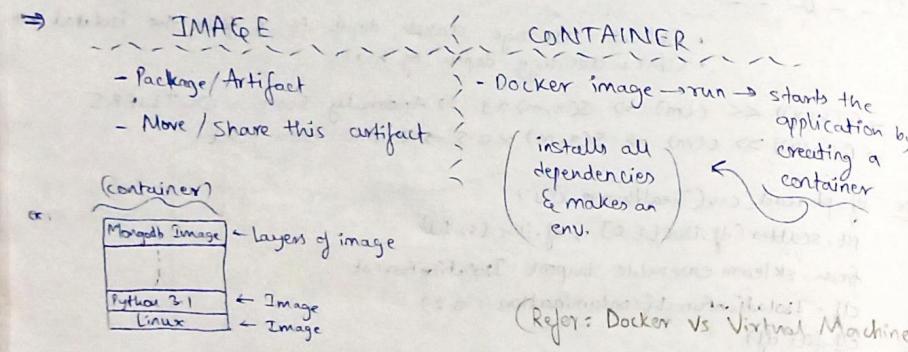


Locality is given by KNN, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors. These are considered outliers.

## DOCKER

### Containers:

- A way to package application with all the necessary dependencies & configuration.
- Portable artifact; easily share & move this package to any env.
- Makes development & deployment more easy & efficient.



⇒ **Docker - Compose**: a tool that helps us define & share multi-container applications.

(docker-compose.yml)

⇒ **Volumes**: preferred mechanism for persisting data generated by & used by docker containers.

## PROJECTS

### Project 1:

1) Created a repository on Github.

- In VS code terminal: → conda create -p venv python==3.8 -y → Creates venv folder
- (conda activate venv)

2) Create .gitignore file.

3) Create setup.py, requirements.txt → (pip install -r requirements.txt)

- setup.py = responsible for creating my ML application as a package. Now we can also install this package in different projects. After making it as a package, we can upload it on Python PyPi and from there anybody can download it.

• Create SRC folder & add this file " \_\_init\_\_.py ", so that when find\_packages() from setup.py searches all the packages needed;

it will consider the SRC folder as a package.

• Write setup.py

• In requirements.txt , adding '-e', so that it automatically triggers setup.py .

• Create components folder & create \_\_init\_\_.py file.

↳ (modules that we are going to specifically use in the project)

• Creating these files inside components folder:  
→ data-ingestion.py  
→ data-transformation.py  
→ model-trainer.py

• Create pipeline folder. Create these files in this folder:  
→ \_\_init\_\_.py  
→ train-pipeline.py  
→ predict-pipeline.py

• Create logger.py, exception.py, utility in src folder. Write own custom exception.py . Write logger.py to log the exceptions which has occurred.

• Create notebook folder & refer the dataset & both the .ipynb files.

• Write → data-ingestion.py → data-transformation.py → model-trainer.py

• Write app.py & predict-pipeline.py .

• Deploying on AWS, Azure.

\* Refer: MLflow, BentoML, DugHub !

## NLP FOR MACHINE LEARNING

### Overview:



### Tokenization:

• CORPUS = Paragraph      • DOCUMENTS = Sentences

• VOCABULARY = Unique Words      • WORDS = words

Tokenization: Process wherein we convert paragraphs into sentences for sentences into words.

↓  
(aka Tokens)

"My name is Ankit. I am a Krishna devotee as well as billionaire by money by age 25." ↗ (CORPUS)

(Tokenization) ↓ Converting to tokens (here: sentences)

1) My name is Ankit.

2) I am a Krishna devotee as well as billionaire by money by age 25.

(Tokenization) ↓ Converting to words (tokens)  
.....

Unique Words = Vocabulary = 16

Ex: ↗

!pip install nltk

corpus = "Currently learning NLP to progress. Watch to become experts in NLP..."

print(corpus)

import nltk

nltk.download('punkt')

#Tokenization {Paragraph → words | Sentence → words}

from nltk.tokenize import sent\_tokenize → (Algo. 1)

documents = sent\_tokenize(corpus)

type(documents) → (list)

for sentence in documents:  
print(sentence) ↗ Currently learning NLP to progress.  
Watch to become experts in NLP.

#Tokenization

from nltk.tokenize import word\_tokenize → (Algo. 2)

word\_tokenize(corpus) → ['Currently', 'learning', 'NLP', 'to', 'progress',  
'Watch', 'to', 'become', 'expert', 'is', 'in',  
'NLP', 'to']

for sentence in documents:

print(word\_tokenize(sentence))

→ ['Currently', 'learning', 'NLP', 'to', 'progress', '.']  
['Watch', 'to', 'become', 'expert', 'is', 'in', 'NLP', '.']

from nltk.tokenize import wordpunct\_tokenize → (Algo. 3)

wordpunct\_tokenize(corpus) → [',', 'become', 'expert', 'is', 'in', 'NLP', '.']

(ID Analyse)

from nltk.tokenize import TreebankWordTokenizer → (Algo. 4)

tokenizer = TreebankWordTokenizer()

tokenizer.tokenize(corpus)

→ [',', 'to', 'progress', '.', 'Watch', '...', 'expert', 'is', 'in',  
'NLP', '.']

- NLTK: • Natural Language Toolkit  
• It's a string processing library & it provides access to many algorithms.  
• Researchers prefer to use NLTK in their work.

### Stemming:

It's the process of reducing a word to its word stem that affixes to suffixes & prefixes or to the roots of words known as a 'lemma'. Stemming is imp. in NLU & NLP.

Ex: ↗

words = ["eating", "eaten", "history", "congratulations", "finally", "fairly"]  
#PorterStemmer → (Algo. 1)

from nltk.stem import PorterStemmer

stemming = PorterStemmer()

for word in words:

print(stemming.stem(word))

→ op: eat  
eat  
congratul  
final  
fairly

### #RegexpStemmer Class → (Algo. 2)

#NLTK has (↓) class with the help of which we can easily implement Regular Expression Stemmer Algorithms. It basically takes a single regular expression & removes any prefix or suffix that matches the expression. Let us see an example:

from nltk.stem import RegexpStemmer

reg\_stemmer = RegexpStemmer('ing\\$|s\\$|able\\$', min=4)

reg\_stemmer.stem('eating') → eat

reg\_stemmer.stem('ingeating') → ingeat

min. len. of string to do stemming

→ If last letters = ing, s, able

Removing '\$' will remove whole occurrence of that word/character

### #Snowball Stemmer → (Algo. 3)

#It's a stemming algo. which is also known as the Porter2 stemming algorithm as it is a better version of the Porter Stemmer since some issues of it were fixed in this stemmer.

from nltk.stem import SnowballStemmer  
snowball\_stemmer = SnowballStemmer('english') → supports many languages

for word in words:

print(snowball\_stemmer.stem(word))

→ op: eat  
eat  
histori  
congratul  
final  
fair

Disadvantage: ⇒ Words are not coming in correct form. The meaning of word is basically changing.

⇒ In order to prevent this, we use something called as 'lemmatizer'.

- Spacy: • Need to download packages related to English text.  
• Provides most efficient NLP algo. for a given task.  
• New library & popular amongst developers.

## Lemmatizer:

Lemmatization technique is like stemming. The op we will get after lemmatization is called 'lemma' which is a root word rather than a valid word that means the same thing.

### Ex: # WordNet Lemmatizer

#NLTK provides WordNetLemmatizer class which is a thin wrapper around the wordnet corpus. This class uses morph() function to the WordNet CorpusReader class to find a lemma.  
Use/case: QnA, chatbots, text summarization.

```
import nltk
nltk.download('wordnet')
from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
" pos = POS Tag
    Noun - n (means, 'noun' denoted as 'n')
    Verb - v
    adjective - a
    adverb - r
    " " "
    lemmatizer.lemmatize("going", pos='v') → go
for word in words:
    print(lemmatizer.lemmatize(word, pos='v'))
```

eat  
eat  
history  
congratulations  
finally  
fairly

## Stopwords:

Stopwords are crucial in NLP to remove common words that add little semantic value, improving the efficiency & effectiveness of text analysis. (Refer Notebook for output)

Ex: Paragraph = "I have three visions for India. (-) ..."

from nltk.corpus import stopwords

import nltk

nltk.download('stopwords')

stopwords.words('english') → (available in diff. languages)

stopwords.words('arabic')

from nltk.stem import SnowballStemmer

snowballstemmer = SnowballStemmer('english')

# Apply Stopwords → filter → apply snowball stemming

for i in range(len(sentences)):

words = nltk.word\_tokenize(sentences[i])

words = [snowballstemmer.stem(word) for word in words if word not in set(stopwords.words('english'))]

sentences[i] = ' '.join(words) # converting all the list of words into sentences

sentences # snowball also ensures that all the words are in small case & converts any, if found in uppercase.

from nltk.stem import WordNetLemmatizer

lemmatizer = WordNetLemmatizer()

for i in range(len(sentences)):

words = nltk.word\_tokenize(sentences[i])

words = [lemmatizer.lemmatize(word.lower(), pos='v') for word in words if word not in set(stopwords.words('english'))]

sentences[i] = ' '.join(words)

sentences,

## Parts of Speech Tags :

Ex: CC → coordinating conjunction

MD → modal (ex: could, will)

NN → noun, singular (ex: desk)

NNPS → proper noun, plural (ex: Indians)

PRP → personal pronoun (ex: I, he, she)

RBR → adverb, comparative (ex: better)

RBS → adverb, superlative (ex: best)

(Refer Google / Notebook for all POS Tags)

Ex: import nltk

from nltk.corpus import stopwords

print(nltk.pos\_tag("Taj Mahal is a beautiful Monument".split()))

[('Taj', 'NNP'), ('Mahal', 'NNP'), ('is', 'VBZ'), ('a', 'DT'), ('beautiful', 'JJ'), ('Monument', 'NN')]

NNP = proper noun

VBZ = verb, 3rd person singular present tense

JJ = adjective

# If we have paragraph & we want to find POS Tag of each word

for i in range(len(sentences)):

words = nltk.word\_tokenize(sentences[i])

words = [word for word in words if word not in set(stopwords.words('english'))]

pos\_tag = nltk.pos\_tag(words)

print(pos\_tag)

## Named Entity Recognition:

Types of tag → (ex: person, place/location, date, time, money, percent, ...)

(eg 1 million dollar)

Ex: sentence = "This is test sentence by Ankit from 2005"

import nltk

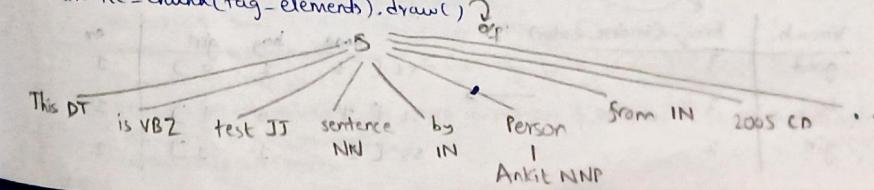
words = nltk.word\_tokenize(sentence)

tag\_elements = nltk.pos\_tag(words)

nltk.download('maxent\_ne\_chunker')

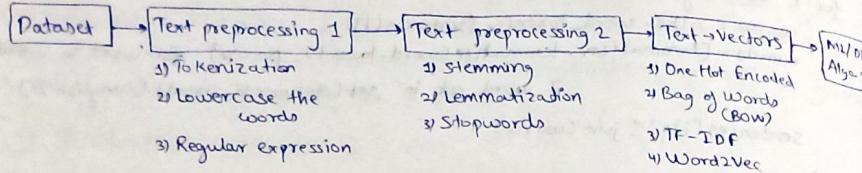
nltk.download('words')

nltk.ne\_chunk(tag\_elements).draw()



## Flow:

We are doing "sentiment analysis":



1) One Hot Encoding: → (ex: sentiment analysis) → aim: convert word to vector

ex:	Text	Op	Vocabulary size = 7
D1: The food is good		1	
D2: The food is bad		0	The food is good bad Pizza Amazing
D3: Pizza is Amazing		1	

$$\begin{aligned} D1 &= [[1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0], \quad D2 = [[1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0], \quad D3 = [[0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0], \\ &\quad [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0], \quad [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0], \quad [0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0], \\ &\quad [0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0], \quad [0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0], \quad [0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0]] \\ &\quad 4 \times 7 \qquad \qquad \qquad 4 \times 7 \qquad \qquad \qquad 3 \times 7 \end{aligned}$$

- Advantages:

• easy to implement with python [sklearn → OneHotEncoder , pd.get\_dummies()]

- Disadvantages:

• Sparse matrix → leads to overfitting

• No semantic meaning is getting captured.

- for ML Algo's, we need fixed size data (here we have  $3 \times 7$ ,  $4 \times 7$ )
- Out of Vocabulary (OOV)
- (i.e. if in test data, a word comes which is not in vocab)

2) Bag of Words: → (ex: sentiment analysis)

ex:	Text	Op	Lowercase all the words → apply stopwords
He is a good boy		1	S1 → good boy
She is a good girl		1	S2 → good girl
Boy and girl are good		1	S3 → boy girl good

Vocab.	Freq. (in desc. order)	good	boy	girl	Op
good	3	S1 [ 1 ]	1	0 ]	1
boy	2	S2 [ 1 ]	0 ]	1 ]	1
girl	2	S3 [ 1 ]	1 ]	1 ]	1

• Binary BOW and BOW (Count will get updated based on freq.)  
 (1 and 0)  
 (prev. ex.)

Advantages:

• fixed size 1<sup>n</sup> (it will help ML Algo's)

Disadvantages:

• Sparse matrix/array → overfitting  
 • Order of the word is getting changed  
 (Based on freq. of word → refer ex.)

• Out of Vocabulary (OOV)

• Semantic meaning is still not captured.

(Ex: The food is good → [1 1 1 0 1] → v<sub>1</sub> {cosine similarity} → {v<sub>1</sub>, v<sub>2</sub>} BUT in actual they're complete opposite)

3) N-grams:

ex: S1 → The food is good } after removing stopword → food not good → (unigram)  
 S2 → The food is not good } food not good → (bigram)

food	not	good	food	good	food	not	not	good
1	0	1	1	0	0	1	1	1

Now our model can clearly distinguish btw both sentences.  
 (As there are many mis-matches)

N-grams = (1,1) → unigrams  
 = (1,2) → unigram, bigram → (comb. of 2 words)  
 = (1,3) → unigram, bigram, trigram → (comb. of 3 words)  
 = (2,3) → bigram, trigram

Ex: (Bag of Words + N-grams) → Refer notebook for more...  
 (separater)

messages = pd.read\_csv('Location', sep='|', names=['label', 'message'])

# Data Cleaning & Preprocessing

```

import re
import nltk
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()
corpus = []
  
```

for i in range(0, len(messages)):

```

review = review.lower()
review = review.replace('[^a-zA-Z]', ' ', messages['message'][i])
review = review.split() # Can try with 'lemmatization' as well
review = [ps.stem(word) for word in review if word not in stopwords.words('english')]
corpus.append(review)
  
```

## #Create BOW

```
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=100, binary=True)
x = cv.fit_transform(corpus).toarray() for binary BOW
```

x.shape  
x  
#N-grams  
(selecting top 100 vocab. wrt frequency)

cv.vocabulary\_ #Gives top 100 words → (unigrams only)

```
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=100, binary=True, ngram_range=(2,2))
x = cv.fit_transform(corpus).toarray()
```

cv.vocabulary\_ → gives top 100 words (Bigram + Trigram)

x

## 4) TF-IDF : (Term Frequency - Inverse Document Frequency)

Ex: S1 → good boy

S2 → good girl

S3 → boy girl good

(from prev. ex. after preprocessing)

$$TF = \frac{\text{No. of repetition of words in sentence}}{\text{No. of words in sentence}}$$

$$IDF = \log_e \left( \frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right)$$

	TF		
	S1	S2	S3
good	½	½	⅓
boy	½	0	⅓
girl	0	½	⅓

$$\begin{aligned} * \quad & \text{IDF} \\ \text{Words} & \quad \text{IDF} \\ \text{good} & \quad \log_e(3/2) = 0 \\ \text{boy} & \quad \log_e(3/2) \\ \text{girl} & \quad \log_e(3/2) \end{aligned}$$

Note: If a word is present in every sentence → less priority.

[good boy girl]
S1    0 $\frac{1}{2} \log_e(3/2)$ 0
S2    0    0 $\frac{1}{2} \log_e(3/2)$
S3    0 $\frac{1}{3} \log_e(3/2)$ $\frac{1}{3} \log_e(3/2)$

Op → (given)

### -Advantages:

- Intuitive
- Fixed Size → Vocab Size

- Word Importance is getting captured

### -Disadvantages:

- Sparcity still exists

. OOV

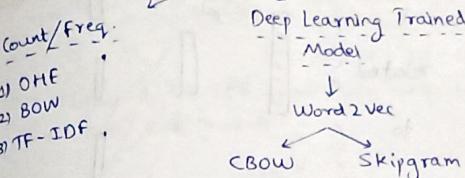
Ex. 14 → Lemmatize instead of stemming

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer(max_features=100)
x = tfidf.fit_transform(corpus).toarray() stop 100 vocabs wrt frequency
```

### #N-Grams

```
tfidf = TfidfVectorizer(max_features=100, ngram_range=(2,2))
x = tfidf.fit_transform(corpus).toarray()
tfidf.vocabulary_
```

## Word Embeddings



term used for representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning.

### 5) Word2Vec:

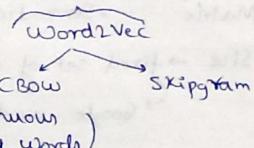
This algo. uses a neural network model to learn word associations from a large corpus of text. As the name implies, word2vec represents each distinct word with a particular list of numbers called a vector. Each word in vocabulary is converted to 'Feature representation', means each word is converted into vector, based on some features.

Ex: Boy apple mango

Gender	→	0.01	0.05
Food	→	-0.02	0.02
Age	→	0.03	0.09
:	⋮	⋮	⋮

(300 dimensions) → used by 'Google'

(is a unsupervised ML Algo.)



### 5.1) C BOW (Continuous Bag of Words)

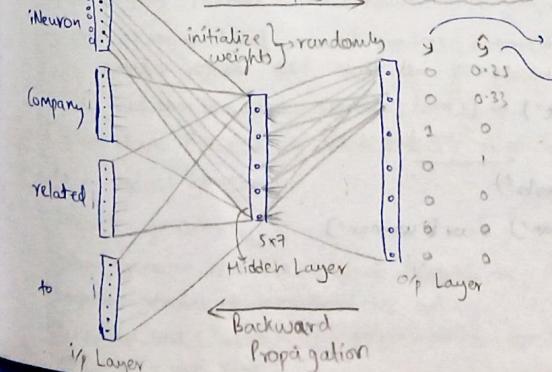
corpus = dataset → [iNeuron Company is related to data science]  
Let window size = 5 & now taking its 'middle' word.

iP → [iNeuron, company, related, to] is  
[company, is, to, data] related  
[is, related, data, science] to

(Vocab.)

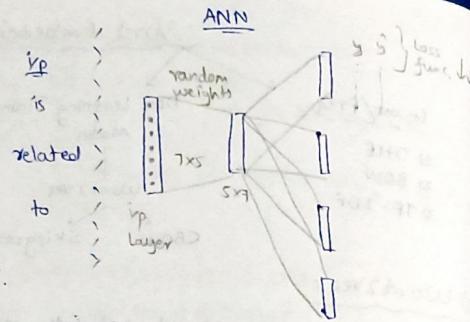
(While sending it to process in ML Algo., we convert it into ONE)  
Ex: iNeuron [1 0 0 0 0 0 0]  
Company [0 1 0 0 0 0 0]  
related [0 0 0 1 0 0 0]  
to [0 0 0 0 1 0 0]

Fully Connected NN: → Each ie every node will be connected to other node.



## 5.2 Skipgram

[iNeuron, company, related, to]  
 [company, is, to, data]  
 [is, related, data, science]



\* We apply 'softmax' func. to get 'g'.

- Improving: → Increase the training data.  
→ Increase the window size = vector dimension is also increasing.
- When to apply: → Small dataset = CBOW  
→ Large dataset = Skipgram

### - Advantages of Word2Vec:

- Dense Matrix → Semantic info. is getting captured. • OOV is also solved.
- Vocab. Size → fixed set of dimension vectors  
 $\hookrightarrow$  Google word2vec (300 dimensions)

### 6) Average Word2Vec:

Ex: [The food is good] = Sentence 1  
 $\begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$  → avg. of prev. vector values  
 $\underbrace{\quad\quad\quad\quad}_{300 \text{ dim.}}$   
 $\hookrightarrow$  (vector representation of whole sentence)

We will use "gensim" & "glove" libraries to implement.

### Ex: (Word2Vec)

```
!pip install gensim
import gensim
from gensim.models import Word2Vec, KeyedVectors
import gensim.downloader as api
wv = api.load('word2vec-google-news-300')
vec_king = wv['King']
vec_king.shape → (300, )
wv['cricket'] → gives values of all those 300 vector values.
wv.most_similar('cricket') → [(most_similar_words, % matching)]
wv.similarity('hockey', 'sports')
vec = wv['King'] - wv['man'] + wv['woman']
vec
wv.most_similar([vec])
```

### Ex: (Bow, TFIDF ML Algo's)

- Flow → 1) Preprocessing & cleaning  
 2) Train Test Split → We want our model to have no info. on test-data, so before changing sentences into vector, we split it hence avoiding data leakage.  
 3) BOW & TF-IDF {sentences → vectors}  
 4) Train Our models.

(14)

### # Create BOW

```
y = pd.get_dummies(messages['label'])
y = y.loc[:, 0].values # Just taking 'ham' values
```

① → train\_test\_split (corpus, y, test\_size=0.25, random\_state=42)

from sklearn.feature\_extraction.text import CountVectorizer

cv = CountVectorizer(max\_features=2500, ngram\_range=(1,2))

x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.25, random\_state=42)

len(x\_train), len(y\_train)

# Independent Features

x\_train = cv.fit\_transform(x\_train).toarray()

x\_test = cv.transform(x\_test).toarray()

x\_train

cv.vocabulary\_

from sklearn.naive\_bayes import MultinomialNB

spam\_detect\_model = MultinomialNB().fit(x\_train, y\_train)

y\_pred = spam\_detect\_model.predict(x\_test)

⑥

### # Creating TF-IDF Model

(→)

from sklearn.feature\_extraction.text import TfidfVectorizer

tv = TfidfVectorizer(max\_features=2500, ngram\_range=(1,2))

x\_train = tv.fit\_transform(x\_train).toarray()

x\_test = tv.transform(x\_test).toarray()

x\_train

tv.vocabulary\_

from sklearn.naive\_bayes import MultinomialNB

spam\_tfidf\_model = MultinomialNB().fit(x\_train, y\_train)

y\_pred = spam\_tfidf\_model.predict(x\_test)

⑥

### Ex: (Spam Ham Project = Word2Vec, AvgWord2Vec)

!pip install gensim

import gensim

from gensim.models import Word2Vec, KeyedVectors

import gensim.downloader as api

wv = api.load('word2vec-google-news-300')

vec\_king = wv['King']

vec\_king

⑯ → Lemmatize instead of stemming

[ $i, j, k$ ] for  $i, j, k$  in zip(list(map(lambda x: len(x) > 0, corpus)), corpus, messages['message'])  
if  $k < 1$

```
from nltk import sent_tokenize
from gensim.utils import simple_preprocess
words = []
for sent in corpus:
    sent_token = sent_tokenize(sent)
    for word in sent_token:
        words.append(simple_preprocess(word))
words
```

# Training Word2Vec from Scratch

```
model = gensim.models.Word2Vec(words)
model.wv.index_to_key # To get all the vocabulary
model.corpus_count # Vocabulary Size
model.epochs # Higher epochs = better performance of model
model.wv.similar_by_word('good')
model.wv['good'].shape
```

Word for word in doc if word in  
model.wv.index\_to\_key

```
words[0]
def avg_word2vec(doc):
    return np.mean([model.wv[word] for word in doc if word in
                    model.wv.index_to_key], axis=0)
```

!pip install tqdm

```
from tqdm import tqdm
import numpy as np
x = []
for i in tqdm(range(len(words))):
    x.append(avg_word2vec(words[i]))
```

len(x) → Refer why 3 values are missing!

# Independent Features

```
x_new = np.array(x)
messages.shape
X[1] # Vectors of 2nd sentence in paragraph
x_new.shape
```

x\_new[0].shape

```
y = messages[[lambda x: len(x) > 0, corpus]]
```

$X[0]$  shape = (100, )  
Now to add it as a row, we reshape  
it to (1, -1) → so new shape = (1, 100)

```
y = pd.get_dummies(y['label'])
y = y.iloc[:, 0].values
y.shape
X[0].reshape(1, -1).shape
```

Final independent features

```
# Final independent features
df = pd.DataFrame()
for i in range(0, len(x)):
    df = df.append(pd.DataFrame(x[i].reshape(1, -1), ignore_index=True))
df.head()
df['Output'] = y
df.head()
df.dropna(inplace=True)
df.isnull().sum()
# Independent Features
X = df
X.isnull().sum()
y = df['Output']
```

⑮

```
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier()
classifier.fit(X_train, y_train)
y_pred = classifier.predict(X_test)
```

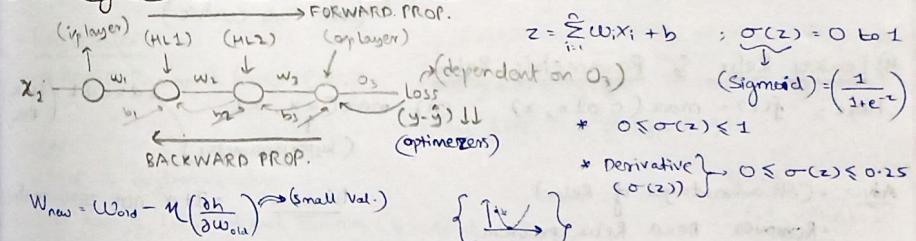
⑯

Ex: (Refer Notebook → Kindle Review Sentiment Analysis)

## DEEP LEARNING

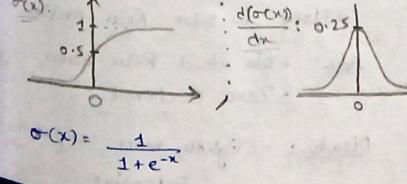
Refer "Application of Derivatives in Deep Learning Neural Network" → PREV. CHAPTER

Vanishing Gradient Problem & Activation functions.



The "Vanishing Gradient" occurs when gradients used for updating weights during training becomes very small. This makes it hard for the network to learn & update its weights effectively.

1) Sigmoid Activation Function: (transforms value btw 0 to 1)

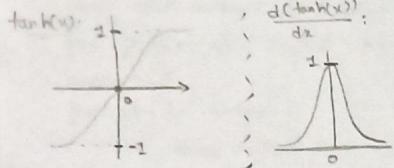


Adv.:

- Binary classification suitable
- Clear prediction (1 or 0)

- Disadv.:
- Prone to vanishing gradient problem.
  - Func. output is not zero centered  $\rightarrow$  (efficient weight update)
  - Mathematical operations are relatively time consuming

## 2) Tanh Activation Function:



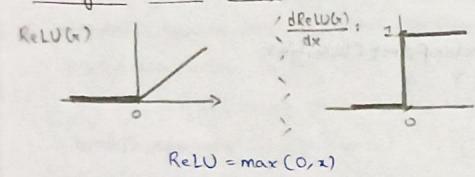
$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Adv.: • zero centric  $\rightarrow$  weight update is efficient

Disadv.: • prone to vanishing gradient prob.  
• Time complexity

## 3) ReLU Activation Function:

### (Rectified Linear Unit)



$$\text{ReLU} = \max(0, x)$$

- If derivative of ReLU

either  $\frac{dy}{dx} = 1$  or  $\frac{dy}{dx} = 0$

(wt. update) (dead neuron)

will happen  $\downarrow$  (When  $x=0$ )

- Adv.:
- Solving vanishing gradient problem
  - $\max(0, x) \rightarrow$  calculation is superfast. The ReLU func. has a linear relationship.
  - It is much faster than sigmoid or Tanh

Disadv.:

- Dead Neuron
- ReLU func. op  $\rightarrow$  either '0' or ' $x$ ';  $x$  is a positive number.  $\rightarrow$  (Not zero centric)

## 4) Leaky ReLU & Parametric ReLU:

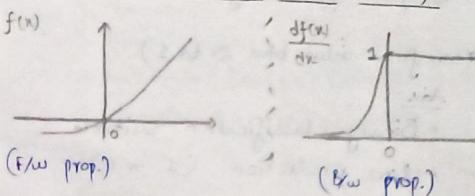
Here:  $f(x) = \max(0.01x, x)$  (or)  $\max(\alpha x, x)$

$\alpha = 0.01$  (Leaky ReLU)  
(hyperparameter)

- Adv.:
- All advantages of ReLU
  - Removes Dead ReLU problem.  $\rightarrow$  (Dead Neuron)

Disadv.: • It is not zero-centric.

## 5) ELU (Exponential Linear Units)



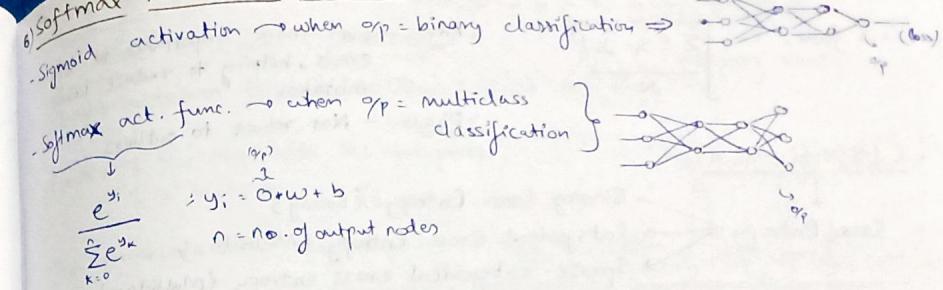
$$f(x) = \begin{cases} x & : \text{if } x > 0 \\ \alpha(e^x - 1) & : \text{otherwise} \end{cases}$$

$\rightarrow$  Used to solve ReLU problem.

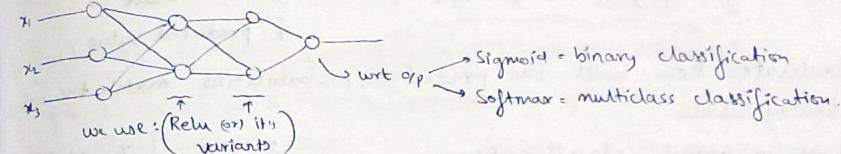
- Adv.:
- No dead ReLU issues.
  - Zero centered.

Disadv.: • Slightly more computationally intensive.

## Softmax Activation Function:



## Which Activation Func. to use when?



## Loss func. & Cost func.:

$\underline{\text{op}}: x_1, x_2, x_3, \text{op}$

$\underline{\text{Loss func.}}: \text{MSE} = (y - \hat{y})^2 \rightarrow$  We take every single datapoint & do fw-bw propagation.

$\underline{\text{Cost func.}}: \text{We take all the datapoints at once & do fw-bw propagation.}$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

ANN  $\curvearrowright$  Classification  
 $\curvearrowright$  Regression

## Regression:

### 1) Mean Squared Error (MSE):

Adv.:

- MSE is differentiable

- It has 1 local/global minima (gradient descent)

- It converges faster.

### 2) Mean Absolute Error (MAE):

Adv.: - Robust to outliers

- Loss func. =  $|y - \hat{y}|$

- Cost func. =  $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

Disadv.: - convergence usually takes time in MAE

### 3) Huber Loss: ( $\in$ MAE + MSE)

Cost func. =  $\begin{cases} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 & : \text{if } |y_i - \hat{y}_i| \leq \delta \rightarrow \text{(hyperparameter)} \\ \delta |y_i - \hat{y}_i| - \frac{1}{2} \delta^2 & : \text{otherwise} \end{cases}$

## Disadv.:

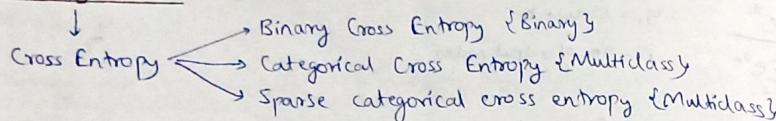
- Not robust to outliers

#### 4) RMSE (Root Mean Squared Error):

$$\text{Cost func.} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{N}}$$

- Adv.: Gives more importance to errors, helping to reduce big mistakes
- Disadv.: Not robust to outliers

#### - Classification:



#### 1) Binary Cross Entropy:

- $\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$ ;  $y$  = actual value  
 $(\log \text{loss}) \rightarrow \text{same used in Logistic regression}$   $\hat{y}$  = predicted value
- Evaluates how well the predicted probabilities match the actual binary labels.
- Used in binary classification

$$\hat{y} = \frac{1}{1+e^{-z}} \quad \text{Sigmoid Act. func.}$$

#### 2) Categorical Cross Entropy:

- Loss func. =  $L(x_i, y_i) = -\sum_{j=1}^C y_{ij} \ln(\hat{y}_{ij})$ ;  $C$  = no. of categories
- Converts  $\sigma_p$  into OHE  $\rightarrow$   $\sigma_p \Rightarrow$  Good Bad Neutral
- $\sigma_p \sim \hat{y}_{ij}$  = probabilities  
 $\in [0.2, 0.3, 0.5]$

$\therefore$  This also gives the prob. of other categories.

#### 3) Sparse Categorical Cross Entropy:

- Disadv.: - Loosing info. about the prob. of other category.

Final  $\sigma_p$  = index of category

$$\begin{matrix} & 1 & 0 & 0 \\ \sigma_p = & 0 & 1 & 0 \\ & 0 & 0 & 1 \end{matrix}$$

$\downarrow$   
 $\% = 2^{\text{nd}}$  index

#### \* Right Combination:

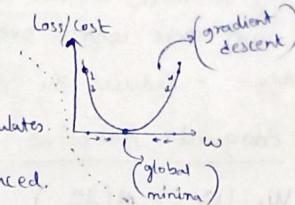
Hidden Layers	Op Layer	Problem Statement	Loss Function
1) ReLU or its variants	Sigmoid	Binary Classification	Binary Cross Entropy
2) ReLU or its variants	Softmax	Multi class	Categorical/Sparse CE
3) ReLU or its variants	Linear	Regression	MSE, MAE, Huber Loss, RMSE

Optimizers: (responsible for reducing the loss func. in the back prop.)  
 (changes 'wt' basically)

#### 1) Gradient Descent Optimizer:

$$\text{Weight updation formula} \Rightarrow W_{\text{new}} = W_{\text{old}} - \eta \left( \frac{\partial h}{\partial w_{\text{old}}} \right)$$

Learning rate



Epoch = 1 f/w prop. with ALL datapoints

& 1 bw prop. {1 epoch = 1 iteration}

This optimizer takes all datapoints at once & calculates.

Epochs done until loss func. is not reduced.

Adv.: - Convergence will happen

Disadv.: - Huge resource (RAM, GPU) needed.

#### 2) Stochastic Gradient Descent (SGD):

Takes 1 record at a time (Similar to gradient descent opt.)

Adv.: - solves resource issue

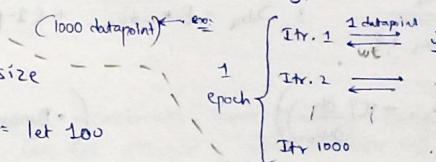


// Since we are solving 1 datapoint at a time, weight updation doesn't happen efficiently & we get noise (zig-zag type).

#### 3) Mini Batch SGD:

3 terms: epoch, iteration, batch-size

$$\text{No. of itr's} = \frac{\text{total datapoints}}{\text{batch\_size}} = \text{let } 100$$



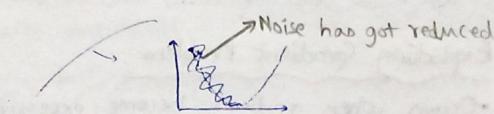
1 epoch = 100 itr.

Adv.: - Convergence speed will increase

- Noise will be less when compared to SGD.

- Efficient Resource Usage

Disadv.: - Noise still exists.



#### 4) SGD with Momentum:

(smoothing)

$$W_t = W_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}$$

$$W_{\text{new}} = W_{\text{old}} - \eta \left( \frac{\partial h}{\partial w_{\text{old}}} \right)$$

(bias update)

$$b_{\text{new}} = b_{\text{old}} - \eta \left( \frac{\partial h}{\partial b_{\text{old}}} \right)$$

Smoothing happens with the help of "Exponential Weight Average (EWA)"

Time Series  $\Rightarrow$  Time =  $t_1, t_2, t_3, \dots, t_n$ ; Value at time  $t_1 \Rightarrow V_{t_1} = a_1, a_2, a_3, \dots, a_n$

$$V_{t+1} = \beta * V_t + (1-\beta) * q_t$$

$\downarrow$   
(controls smoothening func.)

If ' $\beta$ ' is more, we are saying ' $q_t$ ' should have more control over the next value.

- Adv.: - Reduces the noise - Quick convergence

### 5) Adagrad: Adaptive Gradient Descent:

$$W_t = W_{t-1} - \eta \left( \frac{\partial h}{\partial W_t} \right)$$

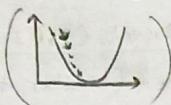
(fixed)  $\rightarrow$  Prev. Learning rate = fixed

$$\eta' = \frac{\eta}{\sqrt{S_{W_t} + \epsilon}}$$

(small value) ( $\epsilon$ )

$$\alpha_t = \sum_{i=1}^t \left( \frac{\partial h}{\partial W_t} \right)^2$$

(derivative of loss wrt  $W_t$ )  
(traversing b/w hidden layers)



but here, we have dynamic learning rate

i.e. As the convergence happens, the learning rate should change.

Disadv.:  $\eta' \rightarrow$  possibility to become a very small value  $\approx 0$

(wt updation) { $W_t \approx W_{t-1}$ }  
(not occurring)

### 6) Adadelta & RMSProp:

$$\eta' = \frac{\eta}{\sqrt{S_{W_t} + \epsilon}} ; S_{W_t} = \beta * S_{W_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial W_t} \right)^2$$

Basically (smoothing 'x\_t')

$$(W_t = W_{t-1} - \eta' \left( \frac{\partial h}{\partial W_t} \right))$$

• Dynamic LR + Smoothing EWA

### 7) Adam Optimizer: {Best Optimizer}

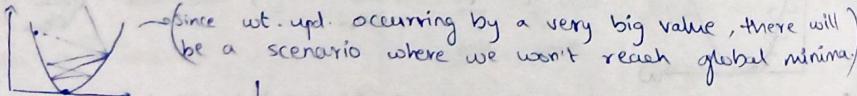
SGD with momentum + RMSProp (Dynamic LR + Smoothing)

All these updates happen in 'weight updation formula' & 'bias updation formula'

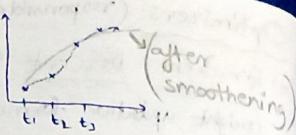
### Exploding Gradient Problem:

Occurs when gradients become excessively large during training the DNN.  
(slope of a func.)

Happens due to repeated multiplication of large weights during backpropagation, leading to exponentially increasing gradients.

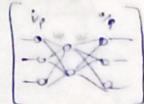


Must ensure that weights are initialised properly, else, if initialised with very high value  $\rightarrow$  above prob. can come



### Initialising Techniques:

- Weights should be small
- Weights should not be same
- Weights should have good variance.



### Uniform Distribution:

We initialize the weights based on some uniform distribution.

$$w_{ij} \approx \text{uniform distribution } \left[ \frac{-1}{\sqrt{\text{input}}} , \frac{1}{\sqrt{\text{input}}} \right]$$

(always give range)

### Xavier / Glorot Initialization:

$$\text{Xavier Uniform Initialization: } w_{ij} \approx \text{uniform distribution } \left[ \frac{-\sqrt{6}}{\sqrt{\text{input} + \text{output}}} , \frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}} \right]$$

$$\text{Xavier Normal Initialization: } w_{ij} \approx N(0, \sigma) ; \sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

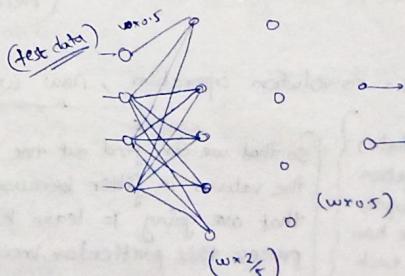
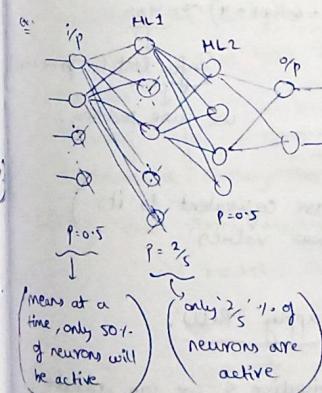
(weights initialised from a normal distribution where mean = 0 & std. dev =  $\sigma$ )

### Training the Initialization:

$$\text{He Uniform: } w_{ij} \approx \text{uniform distribution } \left[ -\sqrt{\frac{6}{\text{input}}} , \sqrt{\frac{6}{\text{input}}} \right]$$

$$\text{He Normal: } w_{ij} \approx N(0, \sigma) ; \sigma = \sqrt{\frac{2}{\text{input}}}$$

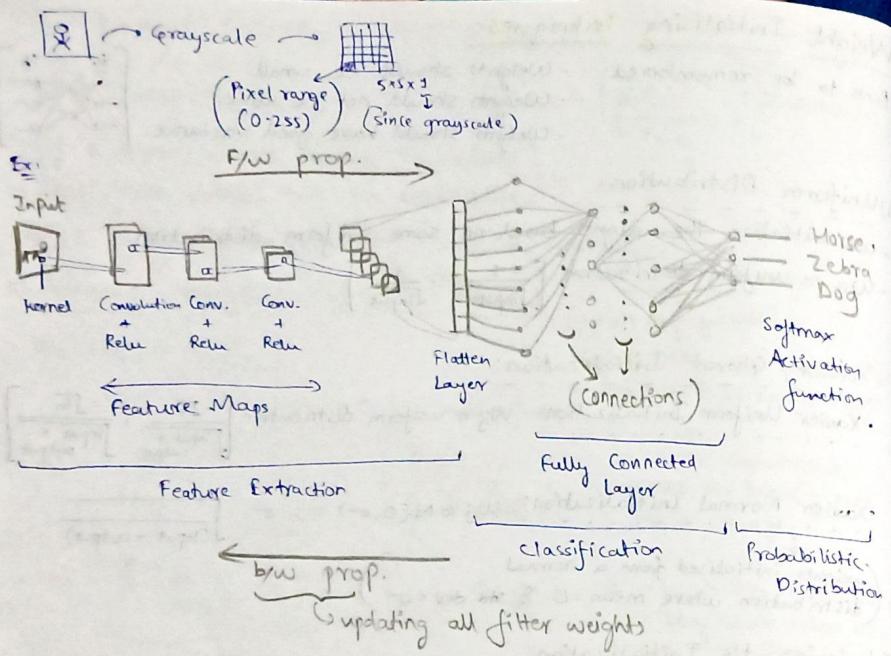
### Dropout Layer:



### Convolutional Neural Network (CNN):

- 1) ANN  $\rightarrow$  supervised learning  $\rightarrow$  Classification
- 2) CNN  $\rightarrow$  ip = images  $\rightarrow$  Regression

- 1) e.g. image classification, object detection.



### Convolution Operation:

- 1) We normalize the pixels (divide by 255)
- 2) We have a 'filter' matrix → which is specialised in a task.  
(e.g. identifying vertical/horizontal edge)
- 3) Apply filter on the image → final matrix

$$\begin{matrix} \text{(let } 3 \times 3) & \downarrow & \text{(let } 6 \times 6) & \downarrow & \text{(let } h-1 = 6-3+1 = 4) \approx 4 \times 4 \\ \text{Info. last (as prev. = } 6 \times 6) & & & \downarrow & \text{Info. last (as prev. = } 4 \times 4) \\ \text{use use padding.} & & & & \end{matrix}$$

→ Padding → Add another layer

↳ Zero padding = (values = 0)

↳ Neighbour padding = (values are equivalent to its neighbour values)

5) Above steps = convolution operation, now we apply 'ReLU'.

We need derivatives in backpropagation to calc. the slopes, which indicate how much wt. in each NN should be adjusted to minimize the loss func. & improve the model's performance.

So that we can find out the derivative & we can update the values in filter because these filters are the ones that are going to learn how we will be able to process this particular image in an efficient way so that we can get the right kind of op.  
(we can have many)

∴ Convolution Layer = Conv. operation + 'ReLU'

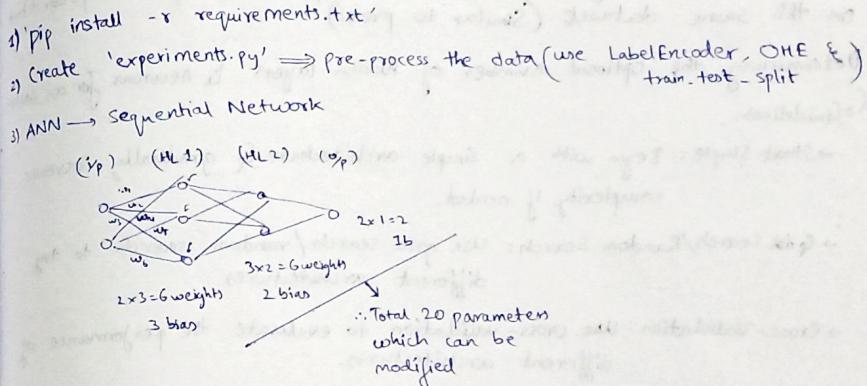
pooling: focuses on summarizing & retaining the most imp. features from the feature maps

types of pooling:

- max pooling → selects max. value in each patch
- min pooling → selects min. of the feature map
- mean pooling → selects avg.

(Refer Pg for more)

### END TO END DEEP LEARNING PROJECT USING ANN



4) When we start creating our ANN, we have to initialize a sequential network.

- Dense - used to create hidden neuron
- We apply activation function to every node (e.g. sigmoid, ReLU, etc)
- Optimizer - useful in back propagation because these are responsible in updating the weights
- Loss func. → we have to reduce this
- Metrics
- Training info. - store the 'log' in a folder, so that I can use Tensorboard (displays the log)

5) Building our ANN model

↳ model = Sequential [  
    → Dense(64, activation = 'relu', input\_shape = (X\_train.shape[1])),  
    → Dense(32, activation = 'relu'),  
    → Dense(1, activation = 'sigmoid') - op layer

])  
model.summary()

6) Compiling & training our model.

7) Setting up Early Stopping & Tensorboard

8) In 'prediction.ipynb':

- Load the ANN trained model, scalar, pickle, onehot
- Take a sample i/p → convert it into suitable format i.e. OHE few columns & scale it.
- Predict on the i/p.

9) Write 'app.py' → deploy on streamlit

10) Practice by predicting a value now i.e. practice regression prob. stat. on this same dataset. (Similar to prev.)

11) Determining the optimal number of hidden layers & neurons for an ANN:

↳ Guidelines:

→ Start Simple: Begin with a simple architecture & gradually increase complexity if needed.

→ Grid Search/Random Search: Use grid search/random search to try different architectures.

→ Cross-Validation: Use cross-validation to evaluate the performance of different architectures.

→ Heuristics & Rules of thumb: → A common practice is to start with 1-2 days.

↓  
The no. of neurons in the hidden layer should be between the size of the i/p layer & the size of o/p layer

→ Refer 'hyperparameter tuning.ipynb'.

// Refer files for above project.

## NLP IN DEEP LEARNING

Simple RNN → LSTM/GRU RNN → Bidirectional RNN → Encoder Decoder

↓  
Transformers ← Self Attention