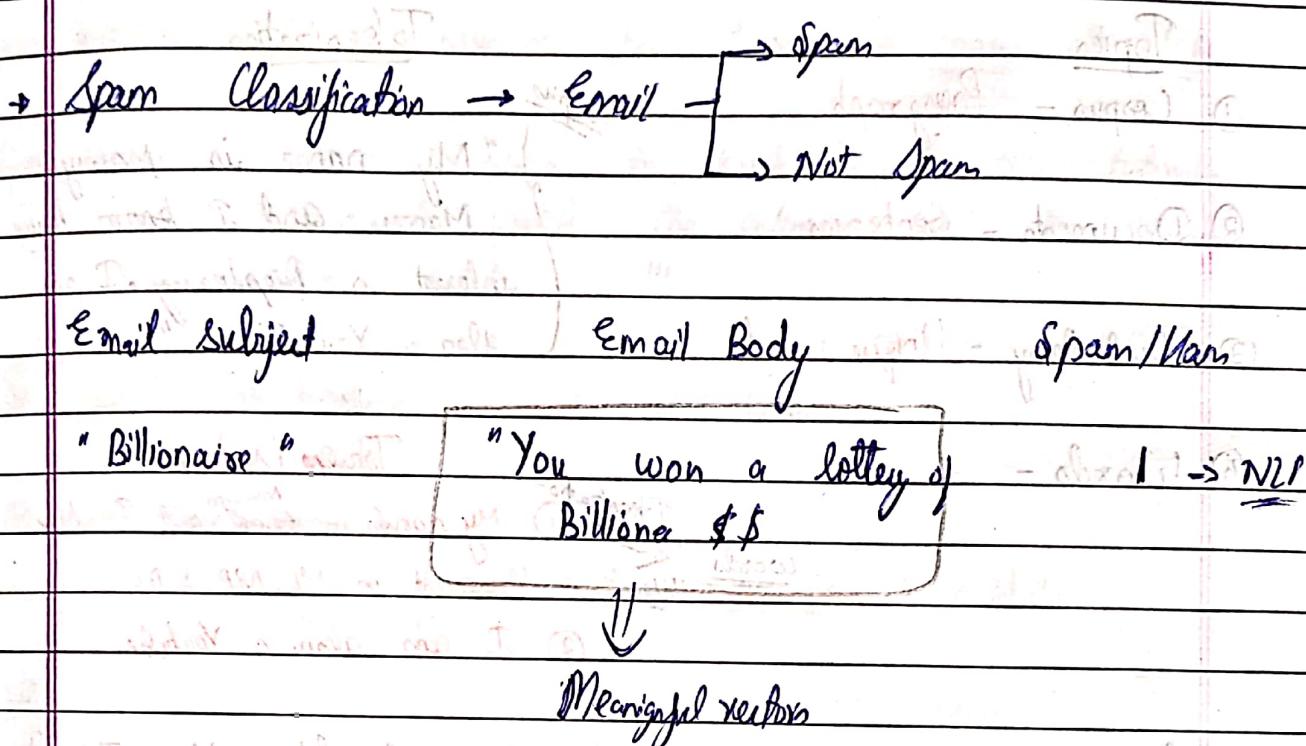


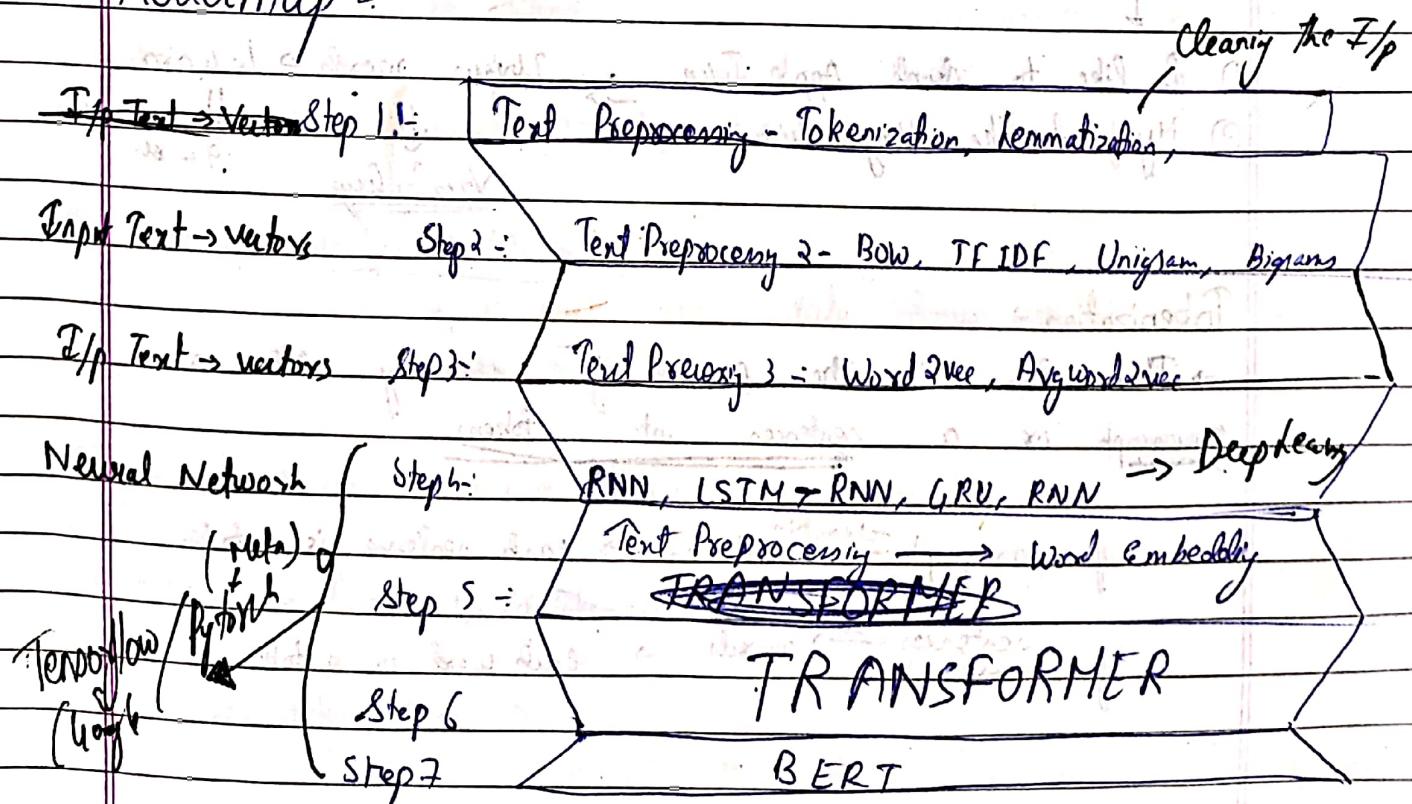
# Roadmap for NLP



e.g. ① Alexa

② Google Home

## Roadmap :-



# Tokenization In NLP

Page No.:  
Date: / /

## Topics

① Corpus - Paragraph

## Tokenization

② Document - Sentence

③ Vocabulary - Unique words

④ Words -

Corpus

"My name is Mawlyn Manzer and I know have interest in DeepLearning. I am also a YouTuber."

↓

## Tokens (sentences)

Tokenization  
Words → Tokens

① My name is <sup>Mawlyn</sup> and I have interest in ML, NLP & DL  
② I am also a YouTuber.

[ I like to drink Apple Juice. My friend likes Mango Juice ]

↓  
Tokenization  
Tokens (sentences)

→ 11 words

- ① I like to drink Apple Juice → Unique words → 10 words  
② My friend likes Mango Juice

Vocabulary  
9 words

## Tokenization :-

→ Tokenization is the process of converting either a paragraph or a sentence into tokens.

if paragraph  $\xrightarrow{\text{tokenize}}$  sentences → each sentence is a token

if sentences  $\xrightarrow{\text{tokenize}}$  words → each word is a token

# Tokenization Practical

Page No.:  
Date: / /

① → Installing nltk = There are two libraries → spacy & nltk

② → corpus = "Hello Welcome, to Krish Naik's NLP tutorial.  
Please do watch the entire course! It's becoming expert in NLP."

③ from nltk.tokenize import sent\_tokenize

④ sent\_tokenize(corpus) → it will make three tokens & print it  
→ type → list

⑤

⑥ from nltk.tokenize import word\_tokenize

⑦ word\_tokenize(corpus)

→ it will form tokens based on words  
it will skip also take 's' & ! & : & ;

⑧ from nltk.tokenize import wordpunct\_tokenize

⑨ wordpunct\_tokenize(corpus)

→ it will also take intentional common words as tokens → !

⑩ from nltk.tokenize import TreeBankWordTokenizer

⑪ tokenizer = TreeBankWordTokenizer()

⑫ tokenizer.tokenize(corpus) → only

→ it will not include first stop (last stop full stop) as a token, & rest divide to word-tokenize



## Text Preprocessing & Stemming using NLTK

Stemming =

Stemming is the process of reducing a word to its word's prefix or affixes. As the roots of words K/o Lemma.

If it is imp in NLP & NLU ( $w \rightarrow \text{wordnet}$ )

eg = # Classification problem

# Comments of products is a +ve review or -ve review.

# Review --> eaty, eat, eaten (go, gone, goes) --> go

### PorterStemmer

① → from nltk.stem import PorterStemmer

② → stemmer = PorterStemmer()

→ for words in words:  
print(word + " --> " + stemmer.stem(word))

O/p → eaty --> eat

eats --> eat

eaten --> eaten

programming --> program

history --> histri

↳ error

Conclusion = PorterStemmer works good for a good number of words, But sometimes it produces error

Page No.:

Date: / /

Page No.:

Date: / /

### RegexpStemmer class

It basically takes a style expression

① uses Regular expression

→ from nltk.stem import RegexpStemmer

→ egg-stemmer = RegexpStemmer('ig\\$|s\\$|e \\$1 ales\\$', min=3)

egg-stemmer.stem('eats')

O/p → eat

reg.stemmer.stem('ingests')

O/p → ingest

### Snowball Stemmer

→ performs better than PorterStemmer

→ from nltk.stem import SnowballStemmer

snowballstemmer = SnowballStemmer('english')

for word in words:

print(word + " --> " + snowballstemmer.stem(word))

O/p → eaty --> eat

eats --> eat

history --> histri

↳ error

great

error

## Lemmatization

Similar to Stemming. The output we get after lemmatization is called Lemma, which is a root word rather than root stem, the output of stemming.

After lemmatization we will be getting a valid word that means the same thing. ~~It uses Morph(1) form~~

### Q/A & Big Notes

→ from nltk.stem import WordNetLemmatizer

→ Lemmatizer = WordNetLemmatizer()

→ Lemmatizer.lemmatize("go")  
o/p → go

POS - Noun - n

Verb - v

adjective - a

adverb - r

## Stopwords

w.r.t different language different tools you can apply different stopwords

→ Stopwords, are common words that are often removed from text before its analyzed because they don't usually carry much meaning on their own.

e.g.: a, the, is, and

Page No.:	1 / 1
Date:	1 / 1

Page No.:	1 / 1
Date:	1 / 1

→ paragraph = "APJ Abdul Kalam Speech"

→ from nltk.corpus import stopwords

→ import nltk

nltk.download("stopwords")

→ ~~from~~ from nltk.stem import SnowballStemmer  
stemmer = SnowballStemmer()

sentence = nltk.word\_tokenize(paragraph)

# Apply Stopwords and filter and then Apply Stemming

for i in range(len(sentence)):

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

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

sentence[i] = ' '.join(words) # connecting all the words into sentence

Now this new sentence file will have the consequential words. like

We could also use Lemmatization

→ from nltk.stem import WordNetLemmatizer  
lemmatizer = WordNetLemmatizer()

# for i in range(len(sentence)):

sentences[i] = sentences[i].lower()

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

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

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



## Parts of Speech Tagging

→ It is the process of assigning a part of speech (such as noun, verb, adjective etc.)

→ POS Tagging helps to understand the grammatical structure of a sentence and the definition of each word within it.

import nltk

from nltk.stem import

nltk.download('averaged\_perceptron\_tagger')

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)

CC -

Tag

Description

NN -

NNP -

NNPS -

PRP -

## Named Entity Recognition

Person eg: Krish C Naik

Place or location: eg - India

Date Eg: September, 26-09-1989

Time Eg: 4:30pm

Money Eg: 1 million dollars

Organization Eg: Microsoft Private Limited

Code

① import nltk

words = nltk.word\_tokenize(sentence)

② tag\_element = nltk.pos\_tag(words)

③ nltk.download('maxent\_ne\_chunker') → nltk.download('words')

④ nltk.ne\_chunk(tag\_element).draw()

(S) O/P → S

person

Gustav NNP Eg: NNP



## What we have learnt?

Text

O/P

Input text for model

① The food is good

The food is bad

Dataset

Text Preprocessing-1

Text Preprocessing-2

① Tokenization

① Stemming

② Lowercase the words

② Lemmatization

③ Regular Expression → ③ Stopwords

Text → Vectors

ML Algorithms

① One Hot Encoding

② Bag of Words (BOW)

③ Tf-IDF

④ Word2Vec

⑤ Avg Word2Vec

## One Hot Encoding

Text O/P

Vocabulary & unique words

D1 The food is good

D2 The food is bad

D3 Pizza is Amazing

D1  $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$ , D2  $\begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$ , D3  $\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$

$4 \times 7$

each time when it runs it will check  
all the unique words, i.e. 7 so in  
this way it runs 4 times & whenever  
it will check all the words

### Advantages

① Easy to implement with Python

[sklearn, ONE, pd.get\_dummies()]

### Disadvantages

① Sparse matrix (lot of zeros)

leads to overfitting

② for ml algorithm → we need

for sized input

③ am not able to understand what the most

important word, how this word is related to  
another.

③ No semantic (similar) meaning is

getting captured



(OOV) → So other will not work if the word is not in Vocabulary

Out of Vocabulary (OOV)② Bags of Words

Text      O/P  
 We are a good boy      1

She is a good girl      1 →  
 ↓ overall the word words  
 cap & Stopword

Boy & girl are good [1 1] → 1      Food is not good [0 1]  
 1 0 0 0 1 0 1      1 0 0 0 0 1 0 1  
 1 0 0 0 1 0 1      1 0 0 0 0 1 0 1  
 1 0 0 0 1 0 1      1 0 0 0 0 1 0 1  
 S1 → good, boy      Food is not good

S2 → good, girl

S3 → Boy, girl, good

## Vocabulary frequency

good	3	good boy girl O/P
boy	2	→ S1      1 0 → 1
girl	2	→ S2      1 0 1 → 1
		S3      1 1 1 → 1

Now, there are two kinds of BOW

Binary BOW

S1 and O/P

BOWCount will get updated  
based on frequencyAdvantages

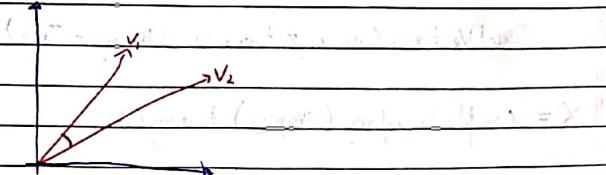
Simple &amp; intuitive

Disadvantages① Sparse matrix or array  
leads to Overfitting2 Fixed size input → Helps in ML algorithms  
iterating② Order of word is getting  
changed③ Out of vocabulary is still a  
problem

which is the  
most freq word  
④ Semantic meaning is still not  
getting captured

→ Q: ① The food is good → [1 1 1 0 1] → v,

② The food is not good → [1 1 1 1 1 1] → v,



## Spam Classifier Practical

```
① import pandas as pd
messages = pd.read_csv('sms clas.csv', sep='|', names=['label',
                                                       'message'])
```

### Data Cleaning & Preprocessing

```
import re
import nltk
```

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

corpus = []

```
for i in range(0, len(messages)):
    review = re.sub('[^a-zA-Z]', ' ', messages['message'][i])
    review = review.lower()
    review = review.split()
    review = [ps.stem(word) for word in review if not
              word in stopwords.words('english')]
    review = ' '.join(review)
    corpus.append(review)
```

### Create the Bag of Words model

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
cv = CountVectorizer(max_features=100, binary=True)
```

```
X = cv.fit_transform(corpus).toarray()
```

### N-grams (BOW)

e.g.: bigrams, trigrams

s1 → The food is not good.

s2 → The food is not good.

Bigram: [ food not good ] food good food not not good?

s1 → [ food not good ] food good food not not good?

s2 → [ food not good ] food good food not not good?

Only one word per vocabulary

sklearn → n-gram = (1,1) → Unigrams

= (1,2) → Unigram, bigram

⇒ (1, 3) → Unigram, bigram, trigram

⇒ (2, 3) → Bigram, trigram

### Create the Bag of words model with ngram

```
from sklearn.feature_extraction.text import CountVectorizer
```

# for Binary BOW enable binary=True

```
cv = CountVectorizer(max_features=100, binary=True, ngram_range=(2,2))
```

```
X = cv.fit_transform(corpus).toarray()
```

cv

# Words → Vectors

Page No.:  
Date: / /

## 4) TF-IDF [Term Frequency - Inverse Document Frequency]

- S1 → good boy : Term Freq (TF) = No. of rep of words  
 S2 → good girl in sentence  
 S3 → boy, girl, god No. of words in sentence

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

### Term Frequency

	S1	S2	S3
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

### IDF

Words	IDF
good	$\log_e (3/2) \approx 0$
boy	$\log_e (3/2) \approx 0$
girl	$\log_e (3/2) \approx 0$

### Final → TF-IDF

good boy girl

sent.1 0  $1/2 \cdot \log_e (3/2)$  0

sent.2 0 0  $\frac{1}{2} \log_e (3/2)$

sent.3 0  $\frac{1}{2} \log_e (3/2)$   $\frac{1}{3} \log_e (3/2)$

### Advantages

1) Intuitive

2) Fixed Size → Vocab size

3) Word importance is getting captured

### Disadvantages

1) Sparsity will exist

2) OOV

## TF-IDF Practicals

```
from nltk.stem import WordNetLemmatizer
lemmatize = WordNetLemmatizer()
corpus = [' '.join([lemmatize.lemmatize(word) for word in sent]) for sent in corpus]
for i in range(0, len(messages)):
    review = re.sub('[^a-zA-Z]', ' ', messages['message'][i])
    review = review.lower()
    review = review.split()
    review = [word for word in review if not word in stopwords.words('english')]
    review = ' '.join(review)
    corpus.append(review)
```

### Create TF-IDF and NGram

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer(max_features=100) → tfidf = TfidfVectorizer(max_features=100,
                                                                    ngram_range=(2,2))
X = tfidf.fit_transform(corpus).toarray()
```

import numpy as np

```
np.set_printoptions(edgeitems=30, linewidth=10000, formatter=lambda x: "% .3g" % x))
```

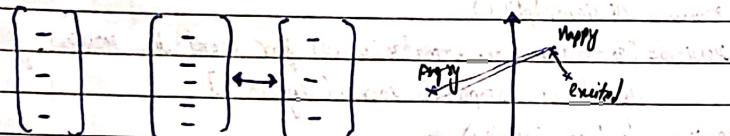
tfidf.vocabulary\_



## Word Embedding

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

→ Angry Happy excited



## Word Embedding

Count or frequency

Deep learning Train model

① ONE

② BOW

③ TF-IDF

CROW

Skipgrams

(Continuous Bag of WORDS)

## Word2Vec → Feature Representation

It uses a neural network model to learn word associations from a large corpus of text.

Once trained, such a model can predict absent synonym words or suggest additional words for a particular sentence.

## WORD2VEC

each distinct word → <sup>comes</sup> list of numbers called vector

Vocabulary → Unique Words → Corpus (paragraph)

	Boy	Girl	King	Queen	Apple	Mango
Gender	-1	1	-0.92	+0.93	0.61	0.03
Royal	0.01	0.02	0.95	0.96	-0.02	0.02
Age	0.03	0.02	0.75	0.68	0.95	0.96
Food	-	-	-	-	0.91	0.92
nm	-	-	-	-	-	-

with this we can get conclusion like King - Boy + Queen = Girl

g → King [0.95, 0.96]

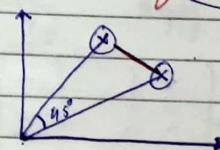
Queen [-0.96, 0.95]

Man [0.95, 0.98]

Woman [-0.94, -0.96]

$$\boxed{\text{King} - \text{Man} + \text{Queen} = \text{Woman}}$$

Cosine Similarity



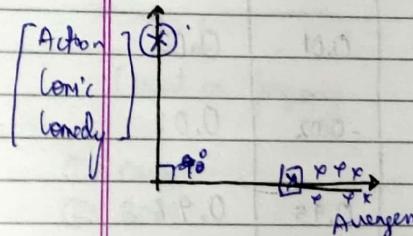
$$\text{Distance} = 1 - \text{Cosine Similarity}$$

$$\text{Cosine Similarity} = \cos \theta$$

$$\text{for e.g. } \theta = 45^\circ, \cos 45^\circ = 0.7071$$

$$\text{Distance} = 1 - 0.7071$$

$$\Rightarrow 0.29$$



$$\text{Distance} = 1 - 0 = 1$$

if Distance is closer to zero, this means words are similar

if Distance is equal to one → There is no similarity

(+)

Word2Vec

→ CBOW [Continuous Bag of Words]

→ Skipgram

Pretrained Model

Train A

model for sentences from

(i) CBOW →

Corpus is our Dataset

[iNeuron Company Is Related To DATA SCIENCE] → Variables

300 dimension ← [Window Size = 5] ⇒ Feature Representations

[-----] OHE

I/P

O/P

→ [iNeuron, Company, Related, To] IS

→ [Company, Is, To, DATA] Related

→ [Is, Related, DATA, SCIENCE] To

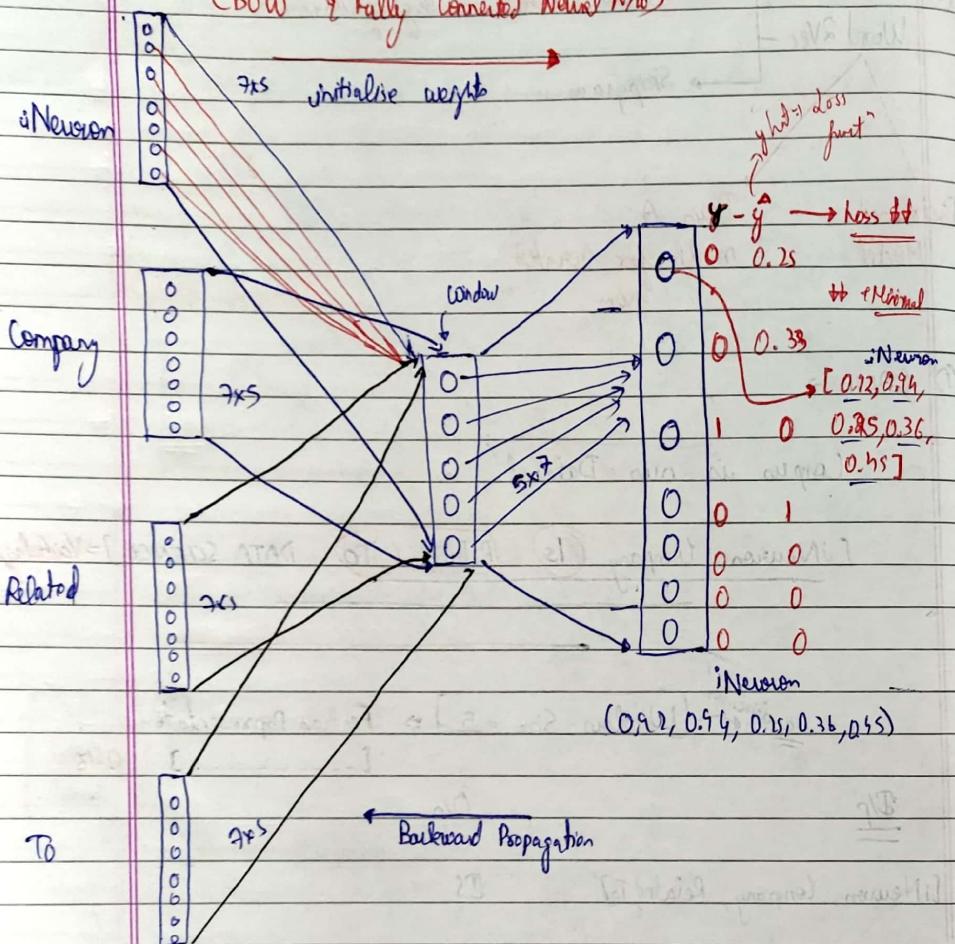
OHE

iNeuron	1	0	0	0	0	0	1
Company	0	1	0	0	0	0	0
Related	0	0	1	0	0	0	1
To	0	0	0	1	0	0	1



Scanned with OKEN Scanner

(CBOW & Fully Connected Neural Net)



Skip Exam → opp. of (B/W)

Dataset: iNeuron Company is related to Data Science

Window size = 5

O/P

→ [Newsman, Company, Related, To]

I/p

ES

Refoto, J

Te

→ Company, Is, To, Date

→ PIs, Related, Data Survey

→ PIs, Related, Data Survey

~~bottom~~ → here Forward propagation happens

Q When should we apply (BOW or skipgram

small Dataset = CBOW  
Large Dataset  $\rightarrow$  Skipgram

Q How to improve CBOW or skipgram

## 1) Preparing the Training Data

2) Increase the Window size & vector Dimension is also increasing as a consequence.

Digitized by Google Word2Vec

3 Billion words → Google News # 300 columns hoga ek line will be like.

features representation of 300 dimension vectors?  
[vector]  $\rightarrow$  [-----]

## Word2Vec Practical Implementation (Gensim)

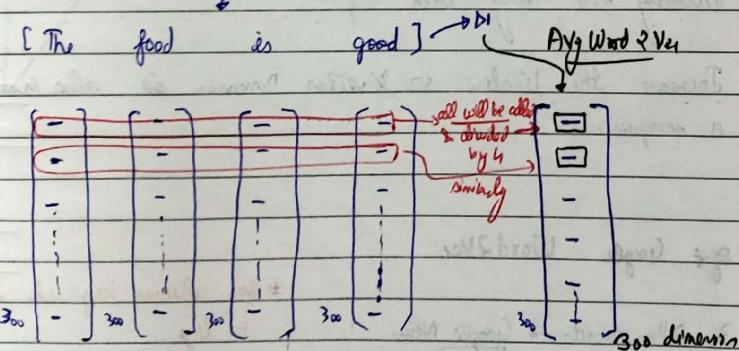
### Advantages of Word2Vec

- \* Sparse Matrix  $\rightarrow$  Dense Matrix
- \* Semantic Info is getting captured [Honest, good]
- \* Vocabulary size  $\rightarrow$  Fixed set of dimension  
 ↳ Google Word2Vec [300 dimensions]
- \* COV is also solved.

300 → Avg Word2Vec

	Text	O/p	300 Avg Word2Vec	O/p
D1	The food is good.	1	[ - - - ]	1
D2	The food is bad.	0	[ - - - - ]	
D3	Pizza is Amazing.	1	[ - - - ]	

Google pretrained Word2Vec Model



So, first of all Now after this for every sentence this Avg Word2Vec will be stored & based on it we will get the output.

- 1) ! pip install gensim  
import gensim
- 2) from gensim.models import Word2Vec, KeyedVectors
- 3) import gensim.downloader as api  
wv = api.load('word2vec-google-news-300')  
wv['king'] ↳ 1.6GB file downloaded

4) wv['cricket'] ↳ wv.most\_similar('cricket')  
 ↳ 53% similarity  
 ↳ 'gig': 0.83  
 'Cricketess', 0.81

Test\_cricket, 0.80  
Cricket, 0.75

5) wv.most\_similar('happy')  
 ↳ 'glad', 0.760  
 'pleased', 0.66  
 'eustatic', 0.662

6) wv.similarity('hockey', 'sports')  
 ↳ 0.535

7) vec = wv['king'] - wv['man'] + wv['woman']

8) wv.most\_similar('vec')  
 ↳ 'king', 0.85  
 'queen', 0.73  
 ↳ 84% simili



~~Imp~~

## Best Practice for ML Algorithms

- ① Preprocessing and Cleaning
- ② BOW And TF-IDF (Sentences → Vectors)
- ③ Train Test Split
- ④ Trained Our Models

~~Wrong method~~

- (1) Preprocessing and Cleaning
- (2) Training Test split
- (3) BOW & TF-IDF (Sentences → vectors) {Priority Data design}
- (4) Trained Our Models

## Spam Ham Classification Project Using BOW

1) import pandas as pd  
 messages = pd.read\_csv('spamcollection.csv', sep='|',  
 names=['label', 'message'])

### Data Cleaning & Preprocessing

2) ~~import re~~  
~~import nltk~~  
~~nltk.download('stopwords')~~  
~~from nltk.corpus import stopwords~~  
~~from nltk.stem.porter import PorterStemmer~~  
~~ps = PorterStemmer()~~

3) Corpus = []

```
for i in range(0, len(messages)):
    review = re.sub('[^a-zA-Z]', ' ', messages['message'][i])
    review = review.lower()
    review = review.split()
    review = [ps.stem(word) for word in review if not word in
              stopwords.words('english')]
    review = ' '.join(review)
    corpus.append(review)
```

### 3) Create Bag of Words

#### # Output Feature

```
y = pd.get_dummies(messages['label'])
y = y.iloc[:, 0].values
```

#### # Train Test split

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(corpus, y, test_size=0.20)
```

#### # Create the Bag of Words model

```
from sklearn.feature_extraction.text import CountVectorizer
```

#### # for Binary BOW enable Binary=True

```
CV = CountVectorizer(max_features=2500, ngram_range=(1, 2))
```

#### # Independent Features

```
X_train = CV.fit_transform(X_train).toarray()
X_test = CV.transform(X_test).toarray()
```

## ⑤ ML Model

→ from sklearn.naive\_bayes import MultinomialNB  
spam\_detect\_model = MultinomialNB().fit(X\_train, y\_train)  
y\_pred = spam\_detect\_model.predict(X\_test)

## ⑥ Performance Metrics

from sklearn.metrics import accuracy\_score, classification\_report  
accuracy\_score(y\_test, y\_pred)  
print(classification\_report(y\_test, y\_pred))

## ⑦ Similarly, the TF-IDF Model will be implemented

→ ~~import gensim~~  
gensim.models.Word2Vec()

→ AvgWord2Vec → It ~~don't~~ gives you the entire sentence in vector form

Page No.:  
Date: / /

## Open Ham Projects Using Word2Vec, AvgWord2Vec

Simple.preprocess() → Convert a document into a list of lowercase tokens, ignoring tokens that are too short or too long.

from gensim.utils import simple\_preprocess

```
words = []
for sent in corpus:
    sent_tokens = sent_tokenize(sent)
    for sent in sent_tokens:
        words.append(simple_preprocess(sent))
```

words list → ['go', 'until', 'jury', 'point', 'ok', 'by', ]

# Let's train Word2Vec from scratch

import gensim

model = gensim.models.Word2Vec(words)

# To get all the Vocabulary  
model.wv.index\_to\_key

model.corpus\_count → 5569

model.epochs → 5

epochs & better model value



Avg-word 2 Vec  
Let's make

def avg\_word2vec(doc):

- # remove out-of-vocabulary words
- # sent = [word for word in doc if word in model.wv.index\_to\_key]
- # print(sent)

return np.mean([model.wv[word] for word in doc if word in model.wv.index\_to\_key])

→ So, our previous code only made a single sentence into vector, to convert entire sentences into vector

① in corpus

We will, from tf\_idf import TfidfVectorizer

# apply for the entire sentences

import numpy as np

X = []

for i in tf\_idf (range(len(words))):  
X.append(avg\_word2vec(words[i]))

# making independent features

X\_new = np.array(X)

X\_new.shape → 5569 , X\_new[0].shape → (100, )

# Dependent Features, # Output Features

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

y = pd.get\_dummies(y['label'])

y = y.iloc[:, 0].values

→ y.shape → (5569)

→ X[0].reshape(1, -1).shape → (1, 100)

# This is the 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() → Their 3 dimensions

0 1 2

(apple Banana)	0.35	0.25	0.65
(apple Chiku)	0.45	0.75	0.75

# lets say apple → (0.3, 0.2, 0.7)  
Banana → (0.4, 0.3, 0.0)

Applying Avg Word2Vec

$(0.3+0.4)/2 \rightarrow 0.35$

$0.2+0.3/2 \rightarrow 0.25$

$0.7+0.6/2 \rightarrow 0.65$

Continue Code

Now for sentence "apple Banana",  
vec → [0.35, 0.25, 0.65]

→ df.dropna(inplace=True)

→ X.isnull().sum()

→ y = df['Output']

# Independent feature

X = df

# Train Test split

from sklearn.model\_selection import train\_test\_split

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.20)

# Model Training

from sklearn.ensemble import RandomForestClassifier

classifier = RandomForestClassifier()

classifier.fit(X\_train, y\_train)

y\_pred = classifier.predict(X\_test)

```
from sklearn.metrics import accuracy_score, classification_report
print(accuracy_score(y-test, y-pred))
```

NLP →	Amazon Kindle Review Sentiment Analysis	ML Algorithm
① Train Test Data preparation	② Train Test split	③ Bent → Ver
④ Checking Null values,	Simple train test split	Apply BOW, TFIDF
⑤ Checking unique values	Put your X & Y	Word2Vec,
⑥ Encoding the output to 0 or 1 (if its a classification problem)	according to your independent & dependant factors	X-train-Wv = fit-train-to-vector() Y-test-Wv = transform
⑦ Lower all the text		↓
⑧ Removing special characters, stopwords		
⑨ Applying Lemmatization		

# Deep Learning

- part of ML used to mimic human brain
- called "deep" because it uses multiple layers of neural networks to help the system understand & interpret data.

## Deep Learning

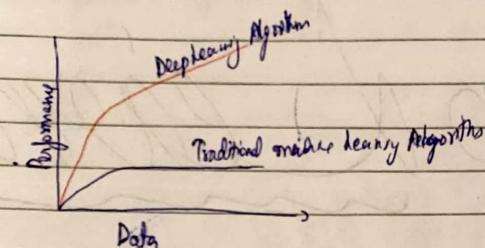
- ① ANN → Artificial Neural N/w → Classification  
Regression
- ② CNN → Convolutional Neural N/w → I/p & Images Video → RCNN, Masked RCNN, Detection, YOLO, Vg, Vf
- ③ RNN → Recurrent Neural N/w → NLP → Time ~~Series~~, NLP series  
I/p → Text, Time series

## Framework

↳ Tensorflow, Pytorch

Word Embedding, LSTM RNN, GRU RNN  
Bidirectional LSTM RNN,  
Encoder Decoder, Transformers, BERT

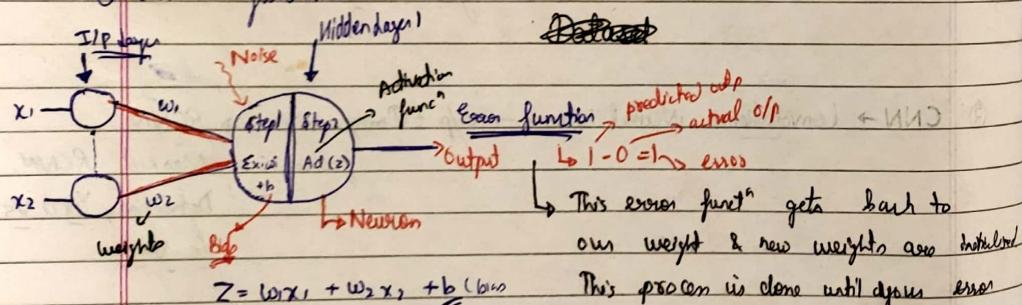
Why Deep Learning is becoming Popular?



## P Perceptron [Artificial Neuron or Neural Network Unit]

used in solving Binary Classifiers Dataset

① Input layer ✓	IQ	No. of study hours	O/p Pass/Fail
② Hidden layer ✓		95	3
③ Weight ✓		4	0 → FAIL
④ Activation function ✓	→ 110	1	1 → PASS



$$Z = \sum_{i=1}^n x_i w_i + b$$

Bias → It makes your data input data float.

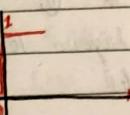
Activation function

Transform the output b/w

- ① 0 to 1
- ② -1 to 1

Two types of Activation func<sup>n</sup>

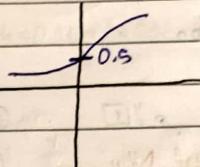
Step function



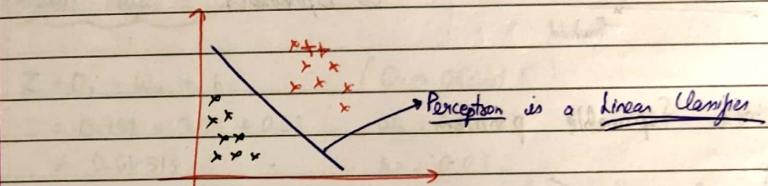
Threshold val = 0

$$\begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

Sigmoid function



$$\begin{cases} 1 & z > 0.5 \\ 0 & z \leq 0.5 \end{cases}$$



Step

$$Z = \sum_{i=1}^n w_i x_i + b$$

$$Z = b + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

$$y = m x + c$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

Linear problem statement

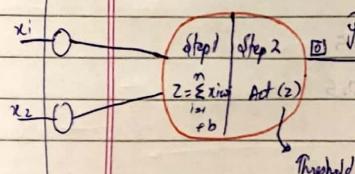


## Perception Models [ANN]

Single Layer Perception Model

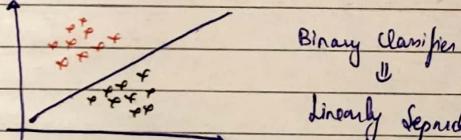
② Multi layered Perception Model

Feed Forward Neural Net



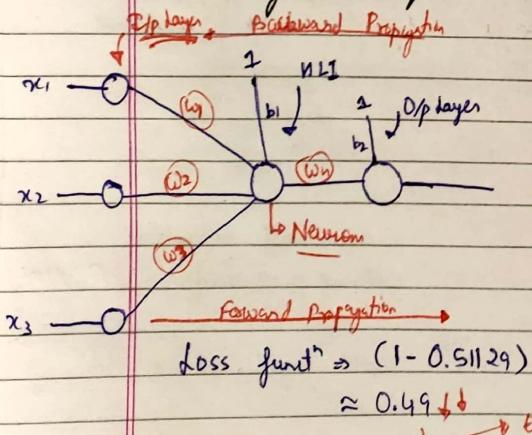
- ① Forward propagation
- ② Backward propagation
- ③ Loss function
- ④ Activation functs
- ⑤ Optimizers

linear separable problem



→ The single layer perception is only good for linearly separable dataset.

## Multi layered Perception Model [ANN]



	$n_1$	$n_2$	$n_3$	Pass/Fail
IQ	95	4	4	1
Study hrs	100	5	2	1
Play hrs	95	2	7	0
O/P				

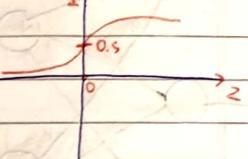
Loss funcn  $\Rightarrow (1 - 0.51129)$

$$\approx 0.49 \downarrow \rightarrow \text{Error}$$

## Hidden layer 1

$$\begin{aligned} \text{Step 1: } Z &= 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01 \\ &= 1.151 \end{aligned}$$

$$\text{Sigmoid} = \frac{1}{1+e^{-Z}}$$



Step 2: Activation ( $Z$ )

$$f(z) = \frac{1}{1+e^{-1.151}} \Rightarrow 0.759$$

## Hidden layer 2

$$\begin{aligned} \text{Step 1: } Z &= 0_1 + w_4 + b_2 & (0_1 \Rightarrow \text{Output 1}) \\ &= 0.759 \times 0.02 + 0.03 \\ &= 0.04518 \end{aligned}$$

$$w_4 = 0.02$$

$$b_2 = 0.03$$

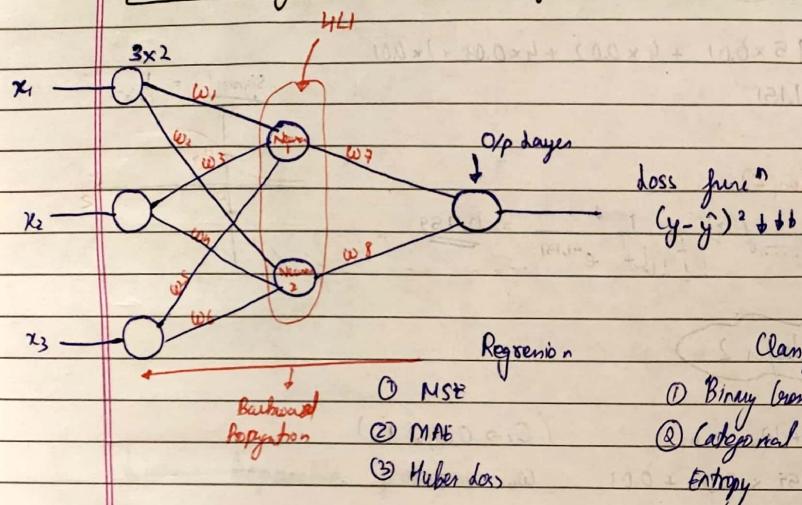
Step 2: Activation ( $Z$ )

$$\frac{1}{1+e^{-(0.04518)}} \Rightarrow 0.51129$$

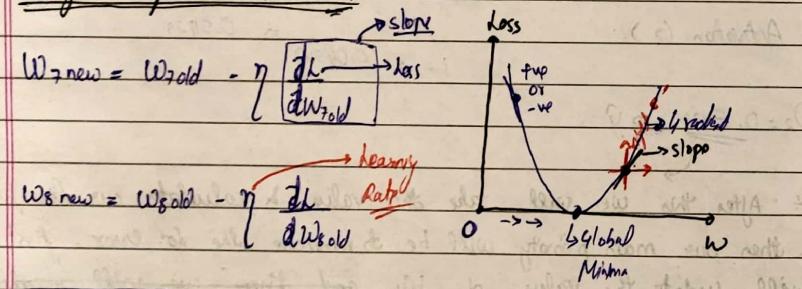
$$0_2 = 0.51129 \Rightarrow \hat{y}$$

Note: After this we will take this value & calculate our loss funcn & then our main priority will be to minimize the loss error. first we will update the value of  $w_4$  and then we will recalculate if loss funcn not minimized then we will go further back & update values of  $w_1, w_2, w_3$

## Back Propagation And Weight Updation Formulas



## Weight Updation Formulas



$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dh}{dW_{\text{old}}} \Rightarrow \text{Weight Updation formula}$$

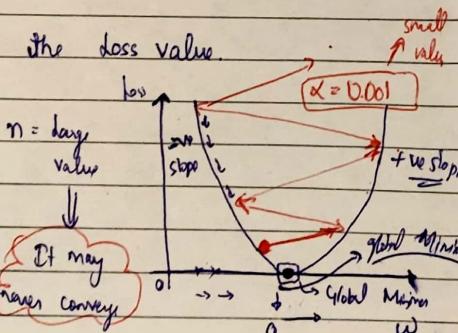
Optimizers: To reduce the loss value.

$$\begin{aligned} W_{\text{new}} &= W_{\text{old}} - \eta (-ve) \\ &= W_{\text{old}} + \eta (+ve) \end{aligned}$$

$\eta$  = learning value

$\rightarrow$  learning rate

$(W_{\text{new}} > W_{\text{old}})$



$$\frac{d(\eta y)}{dt} \stackrel{\text{Chain Rule}}{=} \frac{dy}{dt} = \frac{dy}{dx} \times \frac{dx}{dt}$$

$$W_{\text{new}} = W_{\text{old}} - \eta (-ve)$$

$$= W_{\text{old}} - \eta (+ve)$$

$$(W_{\text{new}} < W_{\text{old}})$$

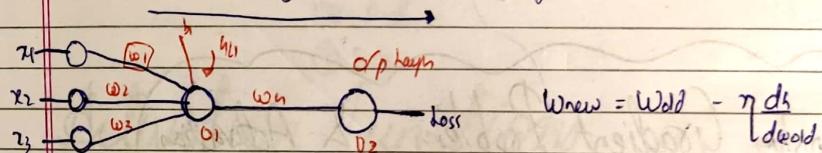
When  $W_{\text{old}}$  reaches global Minima

∴

$$(W_{\text{new}} = W_{\text{old}})$$

## Chain Rule of Derivative

- In Backward propagation, when you go backwards.



$$W_{\text{new}} = W_{\text{old}} - \eta \frac{dh}{dW_{\text{old}}}$$

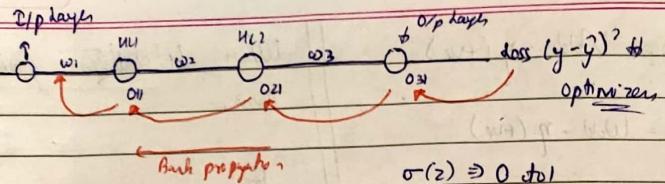
$$\frac{dh}{dW_{\text{old}}} = \frac{dh}{dL} * \frac{dL}{dO_2} * \frac{dO_2}{dW_{\text{old}}} \rightsquigarrow \text{Chain Rule of Derivatives}$$

$$W_{\text{new}} = \boxed{\frac{dh}{dW_{\text{old}}} = \frac{dl}{dO_2} * \frac{dO_2}{dO_1} * \frac{dO_1}{dW_{\text{old}}}}$$

### Front propagation

Page No.:

Date: / /



# The weight update is dependent on  $\frac{dL}{dW_{old}}$  & this is dependent on sigmoid activation func.

$$\rightarrow W_{new} = W_{old} - \eta \frac{dL}{dW_{old}} \quad \text{small value}$$

### Vanishing Gradient Problem & Activation funct.

$$\frac{dL}{dW_{old}} = \frac{dL}{dO_{31}} * \frac{dO_{31}}{dO_{21}} * \frac{dO_{21}}{dO_{11}} * \frac{dO_{11}}{dW_{old}}$$

$$O_{31} = \sigma(w_3 + O_{21} + b_3) \quad z = w_3 + O_{21} + b_3$$

$$O_{31} = \sigma(z)$$

$$\frac{dO_{31}}{dO_{21}} = \frac{d(\sigma(z))}{d(z)} * \frac{dz}{dO_{21}} \quad \{ \text{Chain rule} \}$$

$$0 \leq \sigma(z) \leq 0.25 \Rightarrow \sigma(w_3 + O_{21} + b_3)$$

$$\sigma(O_{21})$$

Derivative of  
Sigmoid

$$0 - 0.025 * W_{3, old}$$

$\Rightarrow$  Now,  $W_{new} \approx W_{old}$

(1)

So, There is no change in gradient descent, & the value of loss function is stagnant

(2)

To fix this problem, Researchers started exploring other Activation funct.

- ① Tanh
- ② ReLU
- ③ Prelu
- ④ Swish

### Activation Functions

#### ① Sigmoid Activation function $\therefore [0 \text{ to } 1]$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Advantages

① Suitable for Binary Classification.

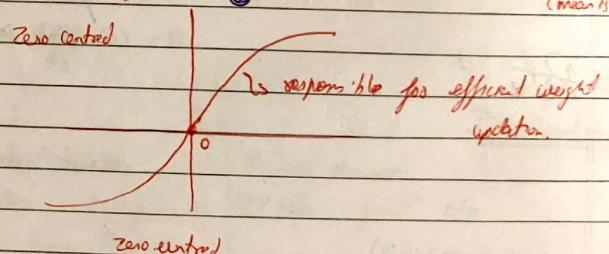
② Can prediction lie very close to 1 or 0

Disadvantages

③ mathematical operat<sup>n</sup> is relatively Time consuming

④ Prone to Vanishing Gradient Problems.

⑤ Function output is not Zero Centred  
(mean is  $0.5e^{-z}$ )

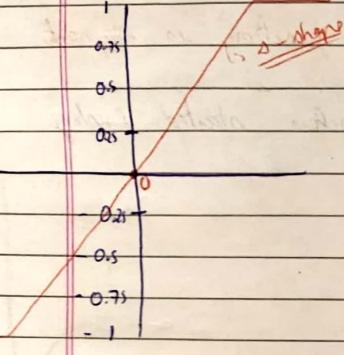


Scanned with OKEN Scanner

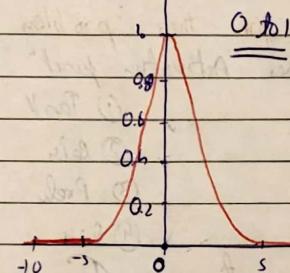
## Activation

### ② Tanh function

Tanh(x)



$d \text{Tanh}(x) / dx$



$0 \text{ to } 1$

Range from  $-1 \text{ to } 1$

Derivative = range from  $0 \text{ to } 1$

#### Advantages

① Zero Centric ✓

② weight updation is effec

#### Disadvantages

① prone to vanishing gradient

Problem

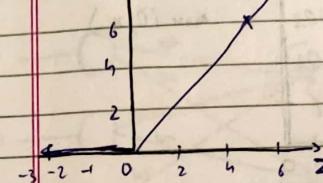
② tanh derivative (matemati  
func' is relatively  
time consuming

### ③ ReLU

## Rectified Linear Unit

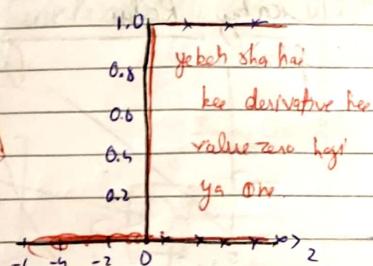
### ③ ReLU Activation Funct<sup>n</sup>

ReLU(z)



$d(\text{ReLU}(z)) / dz$

$$Z = (O_2 * w_3) + b_3$$



Range from  $0 \text{ to } +\infty (x) \rightarrow$  if  $x$  is negative,  $O/p = 0$ , if  $x$  is  $+\infty$ , whatever is the value of  $x$  you will get the  $O/p \Rightarrow$  o/p value of  $x$ .

# 8) Derivative of ReLU output as  $\boxed{II} \rightarrow$  weight updation will happen

$$\# 9) \rightarrow W_{\text{new}} = W_{\text{old}} - \eta \frac{dW}{dW_{\text{old}}} \Rightarrow 0$$

↓

if derivative of  $\text{ReLU}(z)$  is 0  $\rightarrow W_{\text{new}} \approx W_{\text{old}} \Rightarrow$  Dead Neuron  
(A neuron which does not function as it)

$$\# 10) \text{ if } z = +\infty \rightarrow \frac{d \text{ReLU}(z)}{dz} = 1$$

$$\# 11) \text{ if } z = -\infty \rightarrow \frac{d \text{ReLU}(z)}{dz} = 0$$

#### Advantages

① Solving Vanishing Gradient Problem.

②  $\text{Max}(0, x) \rightarrow$  Calculation is superfast. The ReLU function has a linear relationship

③ It is much faster than sigmoid or Tanh.

Tanh, ReLU  
are nonlinear

#### Disadvantages

① Dead Neuron

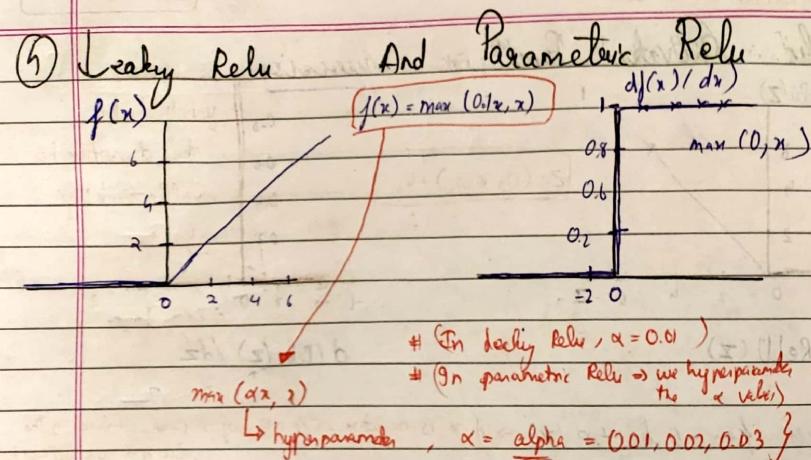
② ReLU funct<sup>n</sup> o/p

$(0, x) \rightarrow 0 \text{ or } +\infty$  numbers

↓

It is not zero centric.  
(as mean is not zero)





Relu  $\rightarrow$  Dead Neuron  $\rightarrow$  Dead Relu Problem

Advantages

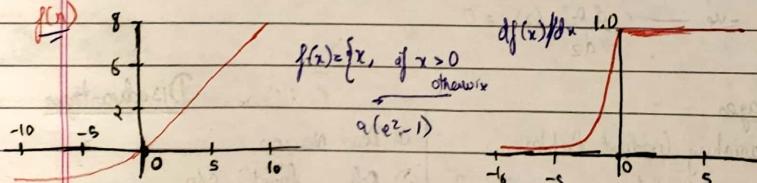
- ① Leaky Relu Relu has all the advantage of Relu

Disadvantages

- ② It is not zero centered
- ③ It removes the Dead Relu

Problem (zero value discourage, very small value due to underflow)

## ⑥ ELU (Exponential Linear Units)



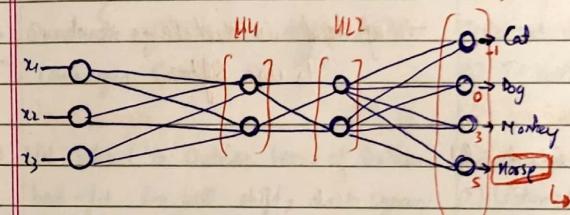
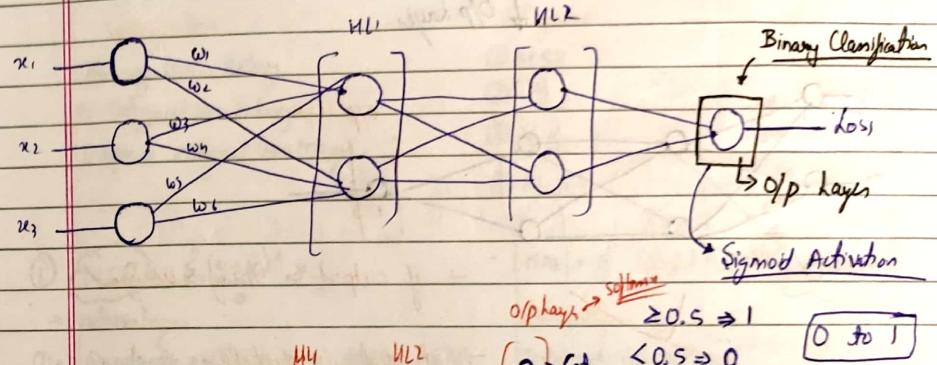
Advantages

- ① No Dead Relu Issues
- ② Zero Centred

Disadvantages

- ① Slightly more computationally intensive.

## ⑥ Softmax Activation function [Multiclass Classification Problem]



↳ softmax Activation

$$\text{Softmax} = \frac{e^{y_i}}{\sum_{k=0}^K e^{y_k}}$$

Softmax

$$\rightarrow \text{Cat} = \frac{e^1}{e^{-1+0+3+5}} = 0.00033$$

↳ predicted output      transformed output

$$\rightarrow \text{Dog} = \frac{e^0}{e^{-1+0+3+5}} = 0.0024$$

$$\rightarrow \text{Monkey} = \frac{e^3}{e^{-1+0+3+5}} = 0.0183$$

$$\rightarrow \text{Horse} = \frac{e^5}{e^{-1+0+3+5}} = 0.1353$$

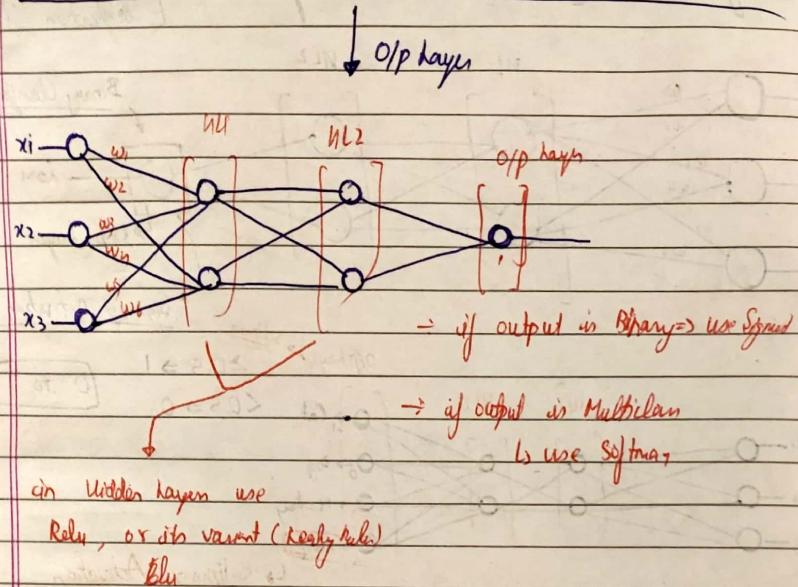
$$\text{Pr(Horse)} = 0.1353$$

$$0.00033 + 0.0024 + 0.0183$$

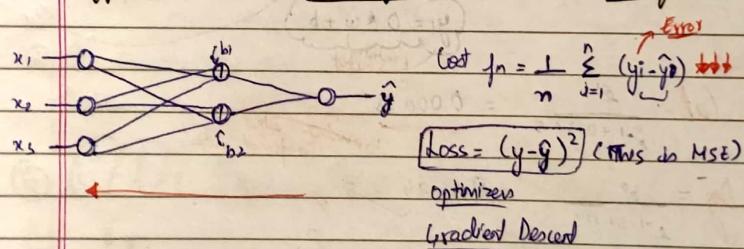
$$+ 0.1353$$

$$\approx 86\%$$

## ⑦ Which Activation Function To Use When?



## Difference b/w Loss function And Cost function



$x_1$	$x_2$	$x_3$	O/p
-	-	-	0
-	-	-	1
-	-	-	0

loss function

$$\rightarrow \text{MSE} = (y - \hat{y})^2$$

Cost function

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

# calculate for one datapoint at a time

# It calculates the error for all data points & then ...

ANN

Classification

- ① Binary Cross Entropy
- ② Categorical Cross Entropy
- ③ Sparse Categorical Cross Entropy

Regression

- ① MSE
- ② MAE
- ③ Huber Loss
- ④ RMSE

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

→ advantages

- ① Quadratic eq  $\Rightarrow$  so it is differentiable
- ② Converges Faster

Disadvantages

- ① Not robust to outliers, becoz of outliers best-fit line will shift, due to square

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

MAE

Advantages

- ① Robust to outliers
- ② If will be in same const

Disadvantages

- ① Convergence usually takes more time
- ② Optimization as a complex task

③ Huber loss =

③ combines the advantage of MSE & MAE properly balance b/w the two by being quadratic for smaller errors & linear for larger errors

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

MSE  
MAE

$K_3 b \Rightarrow |y - \hat{y}| \leq \delta$ , means no outliers are present.

④ RMSE

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

Advantages

- ① Differentiable

Disadvantages

- ① Not robust to outliers



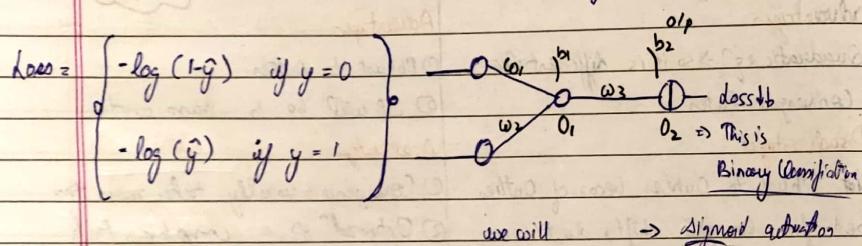
Classification  $\rightarrow$  Cross Entropy

- $\rightarrow$  Binary {Binary}
- $\rightarrow$  Categorical {Multiclass}
- $\rightarrow$  Sparse Categorical {Multilabels}

## ① Binary Cross Entropy

$$\text{loss} = -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

$y$  = Actual value  
 $\hat{y}$  = Predicted value



$$\hat{y} = o_1 = \sigma(w_1 y_1 + w_2 y_2 + b_1)$$

Review again next!

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

## ② Categorical Cross Entropy (Multiclass Classification)

$f_1$	$f_2$	$f_3$	O/P	Good	Bad	Neutral	C = No. of Classes
→ 2	3	4	Good	[ 1 <sup>y1</sup> 0 <sup>y2</sup> 0 <sup>y3</sup> ]			
→ 5	6	7	Bad	[ 0 <sup>y1</sup> 1 <sup>y2</sup> 0 <sup>y3</sup> ]			
→ 8	9	10	Neutral	[ 0 <sup>y1</sup> 0 <sup>y2</sup> 1 <sup>y3</sup> ]			

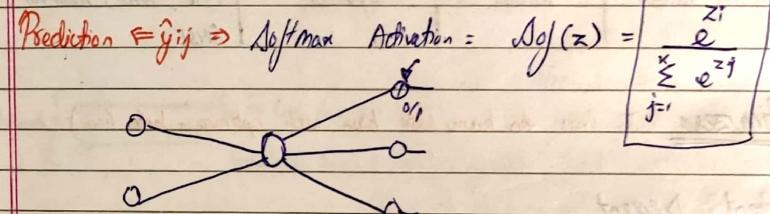
$$\text{loss}(x, y) = -\sum_{j=1}^C y_{ij} * \ln(\hat{y}_{ij}) \quad i = 1 \text{ to } n$$

Actual value  $\Rightarrow y_{ij} = [y_{11}, y_{12}, y_{13}, \dots, y_{1C}] \quad C = \text{No. of categories}$   
 $[y_{21}, y_{22}, y_{23}, \dots, y_{2C}]$

Actual value  $\Rightarrow y_{ij} = [y_{11}, y_{12}, y_{13}, \dots, y_{1C}] \quad C = \text{No. of categories}$

$$[y_{21}, y_{22}, y_{23}, \dots, y_{2C}]$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in the class} \\ 0 & \text{otherwise} \end{cases}$$



O/p  $\Rightarrow$   $\hat{y}$  The output will be in the form of Probabilities

$$\text{O/p} \Rightarrow \hat{y}_{ij} = \text{Probabilities}$$

$$\text{Categorical} \Rightarrow [0.2, 0.3, 0.5]$$

## Cross Entropy

[This also gives the probability of other categories]

## ③ Sparse Categorical Cross Entropy

$$[0.2, 0.3, 0.5] \rightarrow \text{category}$$

What will be output happen if there are equal probabilities?

Instead of giving a list of probabilities it will just return the index with the max probability

Find Index  $\Rightarrow$  O/p Here it was neutral

## Disadvantages

- ① Loses info about the probab. of other categories



Which loss function to use when?

### # Combinations

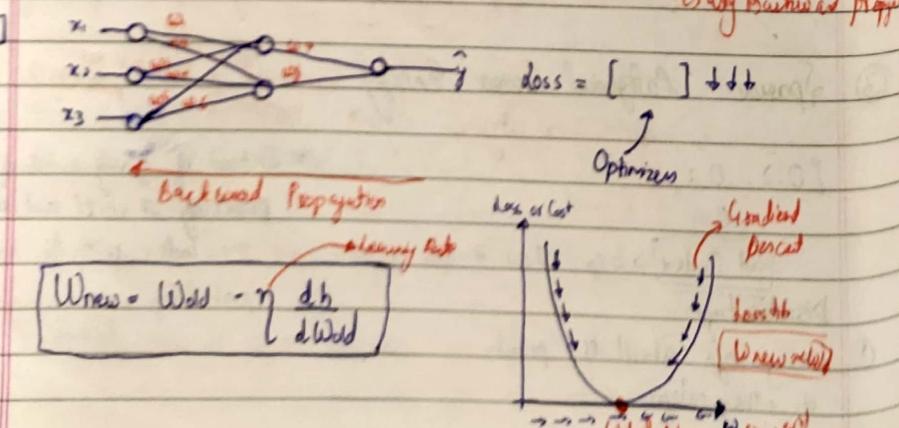
Hidden layers	Activation & func applied here	Problem Statement	Loss function
① Rule or to Yarns	Sigmoid	Binary (Yes/No)	Binary Cross Entropy
② Rule or No Yarns	Sigmoid	MultiClass	Categorical cross entropy
③ Rule or No Yarns	Linear	Regression	MSE, MAE, Huber loss, RMSE

→ Optimizers (In loss to learn how to train with optimizers help from)

- ① Gradient Descent
- ② Stochastic Gradient Descent (SGD)
- ③ MiniBatch SGD
- ④ SGD with Momentum
- ⑤ Adagrad and RMSProp
- ⑥ Adam Optimizers

(1) Gradient Descent Optimizer: [weight update formula]

↳ Why Backward pass



MSE

$$\text{loss } f_n = (y - \hat{y})^2$$

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

Dataset = 1000 Data points

Epochs, Iteration

1 Epoch { 1000 datapoint  
weight will be updated

1 Epoch = 1 Iteration

$\hat{y}_i \Rightarrow \text{Cost function} \rightarrow$

Gradient

Descent optimizer [weight update]

2nd Epoch {  $\hat{y}_i \Rightarrow \text{Cost function} \rightarrow$

Gradient Descent [weight update formula]

1000 batch {  $\hat{y}_i \Rightarrow \text{Cost function} \rightarrow$

1 million

Advantages

① Convergence will happen.

Disadvantages

② Huge Resources RAM, GPU

③

Resource Intensive

Epoch & Iteration are used to describe training process of a model.

Epoch refers to one complete pass through the entire training dataset. During one epoch, the model sees every example in training dataset once.

Iteration refers to one update of the model's parameters (weights & bias) using a single batch of data. No. of iterations depends on the batch size.



### ② Stochastic Gradient Descent (SGD)

here, w.r.t. to every iteration will be only one data point and updating the weight.

Epoch 1  
Iteration 1 → 1 datapoint → loss function → weight update → optimizes

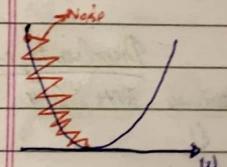
Iteration 2 → 2nd datapoint →

Iteration 1000 → 1000 datapoint

Epoch 2  
Iteration 1 →  
Iteration 1000 →

#### Advantages

① Solves Resource Issue



① Time Complexity ↑ (since in one iteration we do one point do one point)

② Convergence will also take more time

③ Noise gets introduced

#### Disadvantages

Noise has got reduced.  
→ Global Minima

### ③ Mini Batch SGD

Epoch, Iteration, Batch-size

No. of Iterations =  $\frac{1000}{100} \Rightarrow 10$   
1000 iterations

$$\text{No. of iterations} = \frac{\text{Total data points}}{\text{Batch size}} = \frac{1000}{100} = 10$$

Epoch 1

Iteration 1 →

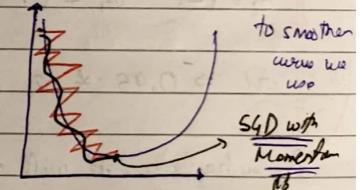
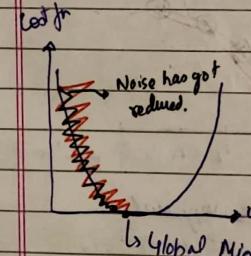
$$\text{cost } f_{\theta} = \sum_{i=1}^{100} (y_i - \hat{y}_i)^2 \quad \text{↑}$$

← change the weight  
Loss will not be used  
Optimizer → Mini Batch SGD

Iteration 2 →

Iteration 3 →

Iteration 100 →



#### Advantages

① Convergence speed will increase

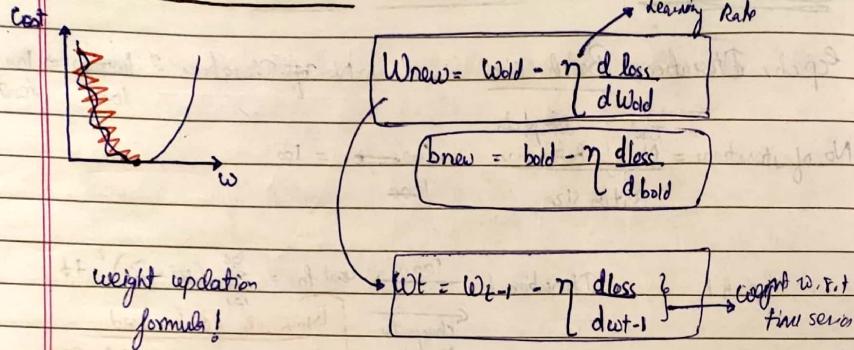
② Noise will be less when compared to SGD.

③ Efficient Resource usage (RAM)

#### Disadvantages

① Noise still exists.

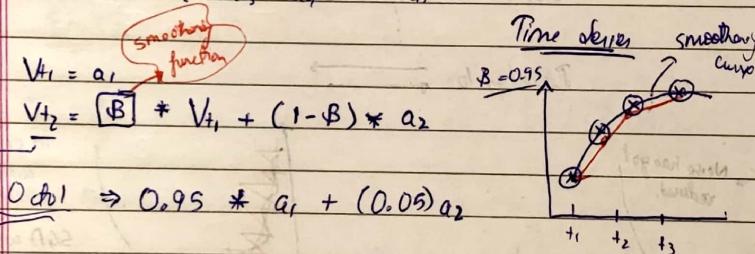
### 3) SGD with Momentum ( $\omega$ )



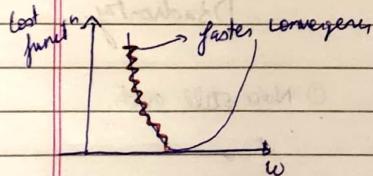
Exponential Weight Average { smoothing }  $\Rightarrow$  ADAM, SARIMAX

Time =  $t_1, t_2, t_3, t_4, \dots, t_n$

Values =  $a_1, a_2, a_3, a_4, \dots, a_n$



so basically, it will make the next value very much related to previous value, resulting in decrease of difference b/w them & as a result smooth curve & faster convergence.



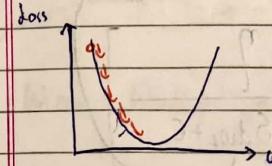
Advantages

- ① Reduces the noise
- ② Quicker convergence

### 5) Adagrad $\Rightarrow$ Adaptive Gradient Descent

$$w_t = w_{t-1} - \eta \frac{d\text{loss}}{d w_{t-1}}$$

learning rate =  $\frac{0.001}{d w_{t-1}}$



As the convergence happens  $\rightarrow$  when it can not reach the local minima  $\rightarrow$  learning rate should be high  $\rightarrow$  when it has reached near global minima ( $\eta$  should be reduced)

$$w_t = w_{t-1} - \eta' \frac{d\text{loss}}{d w_{t-1}}$$

$$\eta' = \frac{\eta}{\sqrt{\alpha t + E}}$$

$\eta' \uparrow \uparrow \quad x_t \uparrow \uparrow$

$\alpha t + E \rightarrow$  Epsilon = small value

$$\begin{array}{l} t=1 \quad t=2 \quad t=3 \\ \eta' = 0.01 \quad \eta' = 0.005 \quad \eta' = 0.003 \end{array}$$

### Disadvantage

①  $\eta' \rightarrow$  Possibility to become a very small value  $\approx 0$

### 6) Adadelta And RMS PROP

Exponential Weighted Avg.

$$\eta' = \frac{\eta}{\sqrt{S d w_t + E}}$$

$$\beta = 0.95$$

$$\begin{aligned} S d w_t &= \text{initialized with zero} \\ S d w_t &= \beta * S d w_{t-1} + (1-\beta) \left( \frac{d\text{loss}}{d w_t} \right)^2 \end{aligned}$$

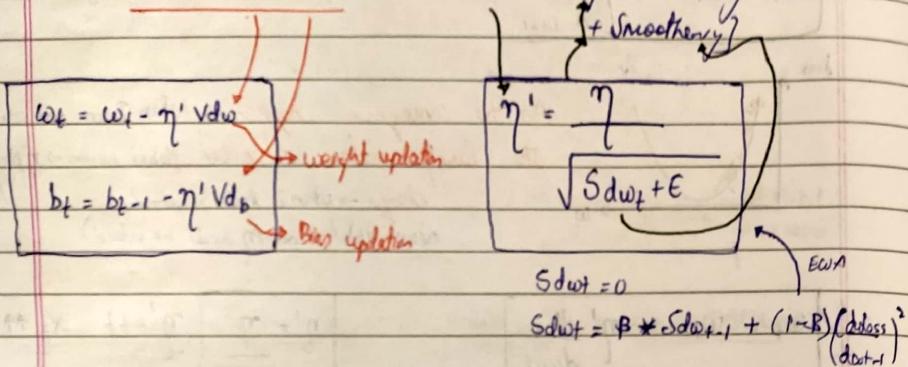
$$w_t = w_{t-1} - \eta' \frac{d\text{loss}}{d w_{t-1}}$$

[Dynamic LR + Smoothing EWMA]

S = Smoothing EWMA

## ⑦ Adam Optimizer

SGD with Momentum + RMSProp [Dynamic Learning Rate]

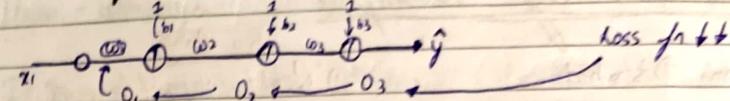
EWA (Exponential Weighted Avg)

$$Vdw_t = \beta * Vdw_{t-1} + (1-\beta) \frac{dloss}{dw_{t-1}}$$

$$Vdb_t = \beta * Vdb_{t-1} + (1-\beta) \frac{dloss}{db_{t-1}}$$

⇒ Momentum  
Smoothness

## Exploding Gradient Problem



$$z = (O_2 + w_3 + b_3) \rightarrow \sigma(z)$$

$$W_{\text{new}} = W_{\text{old}} - \eta \left[ \frac{dloss}{dW_{\text{old}}} \right]$$

two possibilities could come

if  $\eta \frac{dloss}{dW_{\text{old}}}$  is a very high +ve value

then  $[W_{\text{new}} \ggg W_{\text{old}}] - ①$

if  $\eta \frac{dloss}{W_{\text{old}}}$  is a very high -ve value

then  $[W_{\text{new}} \lll W_{\text{old}}]$

$$\frac{dloss}{dW_{\text{old}}} = \frac{dloss}{dO_3} \times \frac{dO_3}{dO_2} \times \frac{dO_2}{dO_1} \times \frac{dO_1}{dW_{\text{old}}}$$

$$big \times big \times big \times big = [big value]$$

weight initialization  $\Rightarrow$  very big value

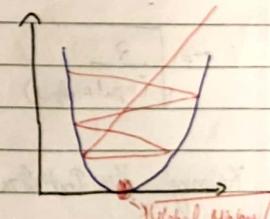
$$\frac{dO_3}{dO_2} = \left| \frac{d\sigma(z)}{dz} \right| \times dz$$

$$= [0 - 0.25] * d(O_2 \times w_3 + b_3)$$

zero  
zero  
0.25

$$dO_2$$

$$\Rightarrow [0 - 0.25] * w_3$$



To solve exploding gradient problem these techniques, we ~~for~~ invented

Page No.:

Date:

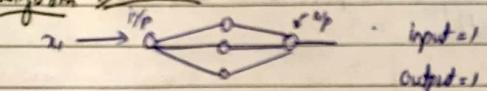
Page No.:

Date:

## Weight Initialization Techniques

- 1) Uniform Distribution
- 2) Xavier / Glorot Initialization
- 3) Kaiming or He Initialization

### ① Uniform Distribution



$$\text{weight} = \left[ \frac{-1}{\sqrt{\text{input}}}, \frac{1}{\sqrt{\text{input}}} \right]$$

$$\left[ \frac{-1}{\sqrt{1}}, \frac{1}{\sqrt{1}} \right]$$

### ② Xavier / Glorot Initialization

#### ① Xavier Normal Init

$$W_{ij} \approx N(0, \sigma^2 \text{std deviation})$$

$$\sigma = \sqrt{\frac{2}{(\text{input} + \text{output})}}$$

#### ② Xavier Uniform

$$W_{ij} \approx \text{Uniform Distribution} \left[ \frac{-\sqrt{6}}{\sqrt{\text{input} + \text{output}}}, \frac{\sqrt{6}}{\sqrt{\text{input} + \text{output}}} \right]$$

### ③ Kaiming He Initialization

#### ① He Normal

$$W_{ij} \approx N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{\text{Input}}}$$

#### ② He uniform

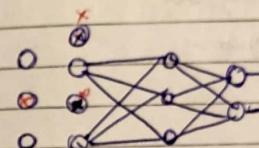
$$W_{ij} \approx \text{Uniform Distribution}$$

$$\left[ \frac{-\sqrt{6}}{\sqrt{\text{Input}}}, \frac{\sqrt{6}}{\sqrt{\text{Input}}} \right]$$

## Dropout

### Drop - Out Layer

It is a regularization technique used to prevent overfitting. In neural network, dropout randomly drops a fraction of neurons during training, forcing the network to learn more robust features.



The features are only dropped in training data not in test data.

$$0 \leq p \leq 1$$

→ selected by hyperparameters

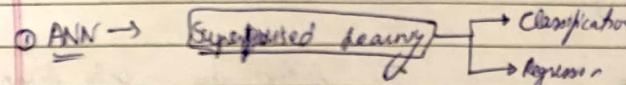


Scanned with OKEN Scanner

# Convolutional Neural Network (CNN)

DNN → Supervised learning →

CNN is a type of Deep learning algorithm primarily designed for Image recognition and Classification.



Defn: Input feature O/P

② CNN → I/p → Images e.g.: Image classification, Object Detection, Segmentation

Brain: Two main parts w.r.t CNN

• Visual Area

→ ① Visual Area: Eye, using sensory, visual recognition  
Motor

→ ② Motor Area: Eye movement and orientation

V<sub>1</sub>, Imp  
Kraissell's area  
Visual Cortex (V<sub>1</sub>-V<sub>5</sub>)

→ The region of brain that receives, integrates and processes visual information relayed from the eyes.

→ V<sub>1</sub> → Primary Visual Cortex [Orientation of edges and lines]

V<sub>2</sub> [Differences in cells, complex patterns, object motion]

V<sub>3</sub> V<sub>4</sub> V<sub>5</sub> [Object Recognition]

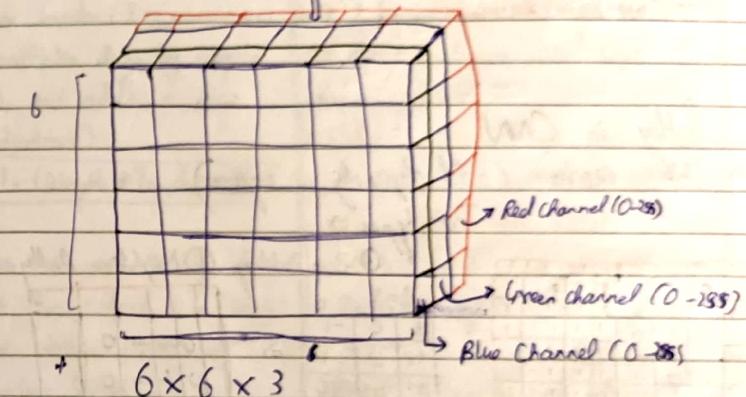
after all visualize the image

# All about Images

③

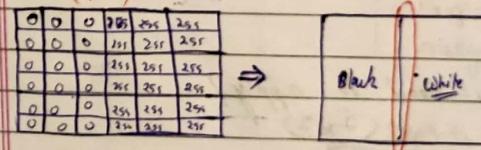
## RGB Image and Grayscale Image

As Image  
RGB → Grayscale  
Red → Green → Let's Visualise



## Convolution Operation In CNN

→ (0, 1) → black color  
white color



Step 1  
① Normalization  
Divide by 255

n=6 f=3  
Convolution operation

0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1
0	0	0	1	1	1

+1	0	-1
+2	0	-2
+3	0	-1

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

6x6x1

3x3 (filter)

Vertical edge filter

4x4  
n-f+1=6-3

[O/P size = (I/p size - filter size) + 1]

Step 2 → De.

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

⇒

white	black	white
-------	-------	-------

(here we have denoising  
the)

Padding in CNN

padding cushion (add dependency on softmax)

Two types ⇒

① Zero Padding ② Neighbor Padding

$f=3$

$$\begin{matrix} +1 & 0 & -1 \\ +2 & 0 & -2 \\ +1 & 0 & -1 \end{matrix}$$

$$n = 6$$

$8 \times 8$

$$n - f + 2p + 1 = 9p$$

$$6 - 3 + 2p + 1 = 6$$

$$3 + 2p + 1 = 6$$

$p = 1$

$$\begin{matrix} 0 & -1 & 0 \\ 0 & -1 & 0 \\ 0 & -1 & 0 \end{matrix}$$

$6 \times 6$

Q How much padding you need to apply?

(7x2)

(3x3)

(7x7)

(3x3)

⇒ 1 layer of padding

## 6) Operation of CNN Vs ANN

basic brain model for computer

ANN

CNN

[Step 1]: It takes inputs (like numbers or data), multiplies them by weights (importance values), and adds a bias (small adjustment)

$$\text{Formula} = \text{Output} = (\text{Weights} * \text{Inputs}) + \text{Bias}$$

[Step 1]: It takes an image and applies filter to detect patterns like edges, colors or shapes.

▪ Filters are like small feature detectors. For  $g = 1$ , one filter might detect vertical edges, while another detects horizontal edges.

▪ Initially, filters are random.

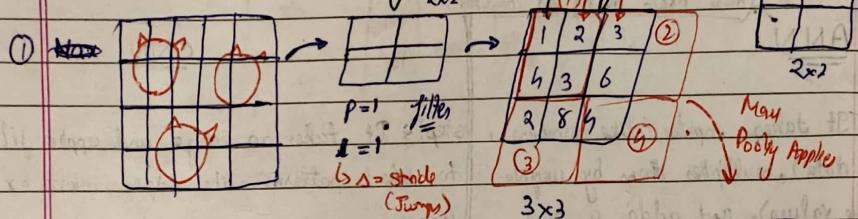
[Step 2]: After applying filters, it performs convolution operation (multiply & add values) to create feature maps.

[Step 3]: It applies the ReLU activation to feature maps. This helps in making network more efficient.

[Step 4]: After ReLU, it performs Max Pooling.



## Max Pooling, Min Pooling, Mean Pooling



## Pooling in CNN

- ① Max Pooling: Takes the highest value in a region. Keeps strong features.
- ② Min Pooling: Takes the lowest value in a region. Rarely used.
- ③ Mean Pooling: Takes the average of all values in a region. Smoothes features.

## Purpose of Pooling

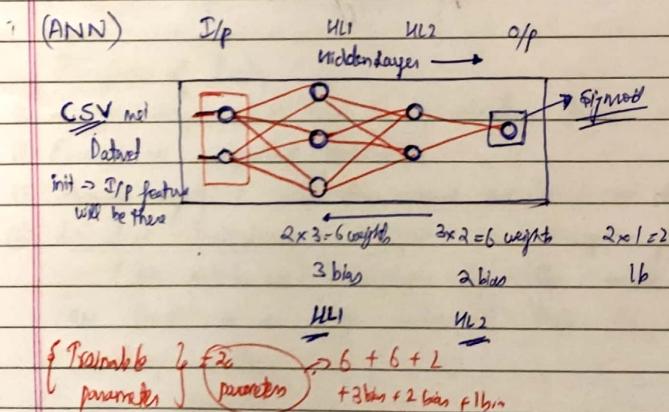
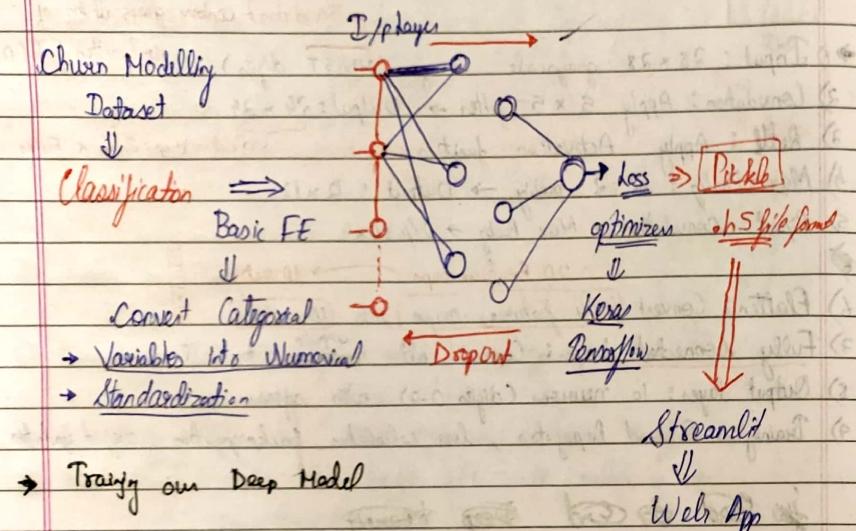
- ① Reduces image size: Makes the img smaller, networks run faster.
- ② Focuses on important features.
- ③ Prevents overfitting: Helps the network generalise to new data.
- ④ Location Invariance: A CNN can recognise features (like objects or patterns) no matter where they are located in an image. Achieved through pooling.

How CNN Works in the context of Handwritten digit dataset.

- This dataset contains grayscale img of handwritten digits (0-9).
- 1) Input:  $28 \times 28$  grayscale image (MNIST digits).
  - 2) Convolution: Apply  $5 \times 5$  filter → Output:  $24 \times 24$ . ( $\text{Feature Map} = \text{Image} \times \text{Filter}$ )
  - 3) ReLU: Apply Activation function.
  - 4) Max Pooling:  $2 \times 2$  pooling → Output:  $12 \times 12$
  - 5) Repeat: Convolution + Max Pooling → Output:  $4 \times 4$
  - 6) Flatten: Convert  $4 \times 4$  feature map into 1D vector.
  - 7) Fully Connected Layer: Connect all inputs to all neurons.
  - 8) Output Layer: 10 neurons (Digits 0-9) with softmax.
  - 9) Training: Forward Propagation, loss calculation, backpropagation, weight updates

for end-to-end Deep learning

## End to End Deep Learning Project: ANN



- ① Sequential N/W → This goes in sequence
- ② Dense → To set no. of neurons in hidden layer
- ③ Activation func → sigmoid, tanh, ReLU, Leaky ReLU
- ④ Optimizer → Back Propagation → Update the weight
- ⑤ Loss func →
- ⑥ Metrics → [accuracy], [mse, mae]
- ⑦ Training → logs → folders → Tensorflow → Visualize

## ANN Implementation

- ① Import tensorflow as tf  
`import tensorflow as tf  
import tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Dense  
from tensorflow.keras.callbacks import EarlyStopping, TensorBoard  
import Datetime`
- ② (X.train.shape [17, 12])  
`# Build our ANN model`
- ③ model = Sequential([  
 `Dense(64, activation = "relu", input_shape = (X.train.shape [17,]), # H1 connected with input layer`  
 `Dense(32, activation = "relu"), # H2`  
 `Dense(1, activation = "sigmoid") # O/p layer`  
`])`

④ → model.summary()

⑤ import tensorflow

opt = tensorflow.keras.optimizers.Adam(learning\_rate = 0.01)

loss = tensorflow.keras.losses.BinaryCrossentropy()

# Compile the Model

model.compile(optimizer = opt, loss = "binary\_crossentropy", metrics = ["accuracy"])

# Set up the Tensorboard

from tensorflow.keras.callbacks import EarlyStopping, TensorBoard

log\_dir = "logdir" + datetime.datetime.now().strftime("%Y%m%d - %H-%M-%S")

tensorboard\_callback = TensorBoard(log\_dir = log\_dir, histogram\_freq = 1)



## Experiment.ipynb

### (5) # Set up Early Stopping

early\_stopping\_callback = EarlyStopping(monitor='val\_loss', patience=6, restore\_best\_weights=True)

### (6) # Train the model

history = model.fit(X\_train, y\_train, validation\_data=(X\_test, y\_test), epochs=100, callbacks=[tensorflow\_callback, early\_stopping\_callback])

(7) model.save('model.h5')

### (8) # Load Tensorflow extension

\* load\_ext tensorflow

(9) for visualization: % tensorboard --logdir logdir/ logitfit

(1) import pickle, scalar, encoder & others & train-test-split

(2) Read the data (read\_csv)

(3) Preprocess the data

(a) Drop unrel. columns (b) Encode categorical & OHE (pd.get\_dummies)  
b[ data[['4eo']] ] . binary

(c) convert to df (OHE) (d) combine of the original data (cate.)

(e) Save label encoder & OHE (using pickle dump)

(f) Do scaling & save it

(g) \* ANN -> import tensorflow, sequential, Dense, EarlyStopping, tensorflow, date\_time

(h) Using (X\_train.shape[1]) we get no. of i/p of first hidden layer

(i) # Build our ANN model

model = Sequential([Dense(64, activation='relu'), input\_shape=(X\_train.shape[1]),  
] UL2 [ UL3 ]) )

# Optimizer

(j) import tensorflow -> select optimizer & loss

opt = tensorflow.keras.optimizers.Adam(learning\_rate=0.01)

loss = f

(k) # Compile the model -> model.compile(optimizer=opt, loss="binary\_crossentropy", metrics=['accuracy'])

(l) # Setup Early Stopping / Train Model -> early\_stopping (monitor='val\_loss', patience=10, restore\_best\_weights=True)  
history = model.fit(X\_train, y\_train, validation\_data=(X\_test, y\_test), epochs=100  
callbacks=[tensorflow\_callback, early\_stopping\_callback])

(m) model.save(model.h5)

## Prediction.ipynb

1) Importing ~~as tensorflow.keras.models~~ ~~load\_model~~, ~~pickle~~, ~~many~~ ~~openpyxl~~

2) ~~model~~ ~~load~~ the trained model  
model = ~~load~~ ~~model~~("model.h5")

3) Load other encoders & scales

using pickle.load(file)

# Example w/p data

a) input\_data = { " " :  
Convert input into Df}

④ Out: 'Geography'

⑤ Encode categorical variables

⑥ Concatenate & scale

⑦ Predict class