

## Errors

- \* In cmd Import Error i.e., builder.py is not found
  - upgrade pro
- \* Unicode error during path without \n before it
- \* Sparse matrix error - update sklearn version

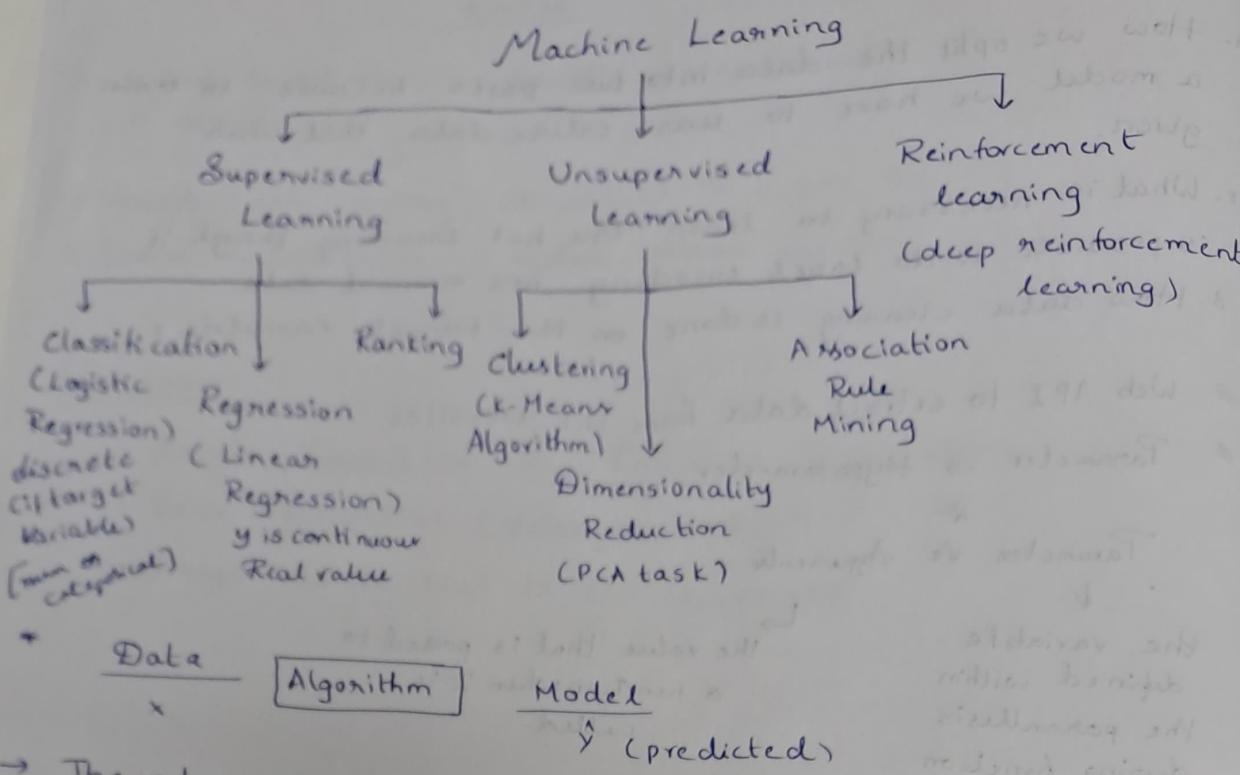
## Use Cases:

- Google Translation i.e., Language Translation
- Voice Recognition (Voice to Text)
- Image inside text to be recognized

## Doubts:

1. How we split the data into two parts because to train a model we have to learn entire data that was given.
  2. What is necessary to learn One hot encoding though it is possible to do label encoding for nominal data.
  3. How data cleaning is done on the target variable?
- = Web API to extract data from the websites
- = Parameter vs Hyperparameter
- Parameter vs Argument
- the variable defined within the parenthesis during function definition
- the value that is passed to a function when it is called.
- = What is the need of Adjusted R square already we have R square?
- +  $R^2$  gives same score even after introducing garbage feature
- ↳ not good metric (It always ↑ when you add new feature).
- \* Adj  $R^2$  is better in this case

- = Jeffrey Hinton - Father of ~~AI~~ ML
- Geoffrey Hinton
- Father of AI - John McCarthy.
- # A sparse vector is a vector that has a large number of zeroes so it takes unwanted space to store these zeroes.
- Sparse      vs      dense vectors also named dense or embeddings.
- A sparse array contains mostly zeroes & few non-zero entries
- A dense array contains mostly non-zeroes.



→ The value which is going to predict is called target variable (tv).

For Example:

In term deposit data the tv is status i.e., Approval or rejected and tv is categorical so we use Classification (logistic regression)

In house pricing data the tv is price i.e., numerical so we use Regression (linear Regression)

Eg: Height

155

180

175

145

185

150

Weight

50

90

85

50

90

55

55

85

Gender

Female

Male

Male

Female

Male

Female

?

?

Queries { Q<sub>1</sub>, Q<sub>2</sub> }

145

180

—

—

Q) How do you train your computer to do this prediction?

Approach No - 1

Step 1 : Store the data

Step 2 : Perform a lookup

Approach No - 2

Male

Female

Height : [175, 185]

Height : [145, 155]

Weight : [85, 90]

Weight : [50, 55]

$$\begin{cases} Th\_H = 160 \\ Th\_W = 65 \end{cases}$$

Code :

```
if Q[Height] <= Th_H and Q[Weight] <= Th_W :
```

    print('Female')

else :

    print('Male')

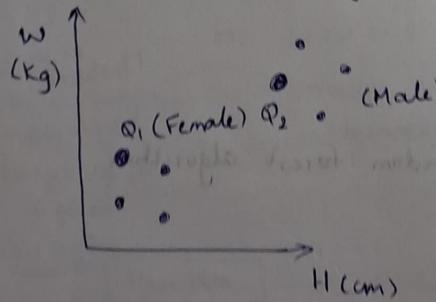
⇒ Here it is human learning / analysis, because machine learning means machine itself creates some conditions & predict the outcomes.

⇒ Features / attributes nothing but columns.

Algorithms :

\* Distance based Algorithm

K-NN (K- Nearest Neighbour)

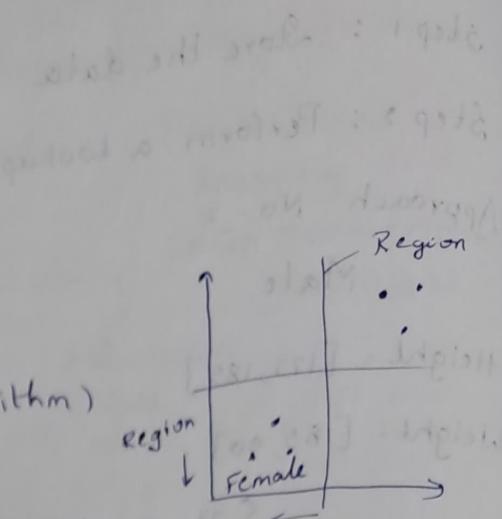
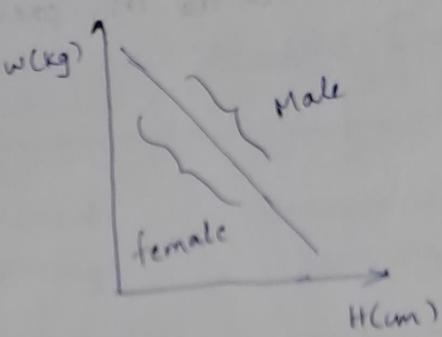


\* Boundary based Algorithm

(Logistic Regression)

It is possible to draw a line that separates male & female.

\* Support vector machines (Non linear separation)



### \* Rule Based Algorithm

(Decision Tree Algorithm)

Approach 2 - Automated

### \* Probabilistic Approach

$$P(F | (He, We))$$

Naive Bayes Algorithm  $P(M | (He, We))$

### \* Ensemble Approach (Collection of Algorithm)

- Random Forest Algorithm = Stacking
- GDA Boost
- GBDT
- xGBoost
- Voting Ensembles
- Cascading

### \* Deep Learning Based Algorithm $\leftarrow$ Artificial Neural Networks

If the data created in non linear structure then it uses deep learning

### 1 In Ensemble Approach

Parallel Models

- + Voting Ensemble
  - + Stacking
  - + BAGging
- $\hookrightarrow$  Random forest algorithm

Sequence  
Models  
(consecutive  
steps)

+ Cascading

+ Boosting

- $\hookrightarrow$  ADA Boost
- $\hookrightarrow$  GBDT
- $\hookrightarrow$  XG Boost

## Algorithms

## Dimensionality Reduction

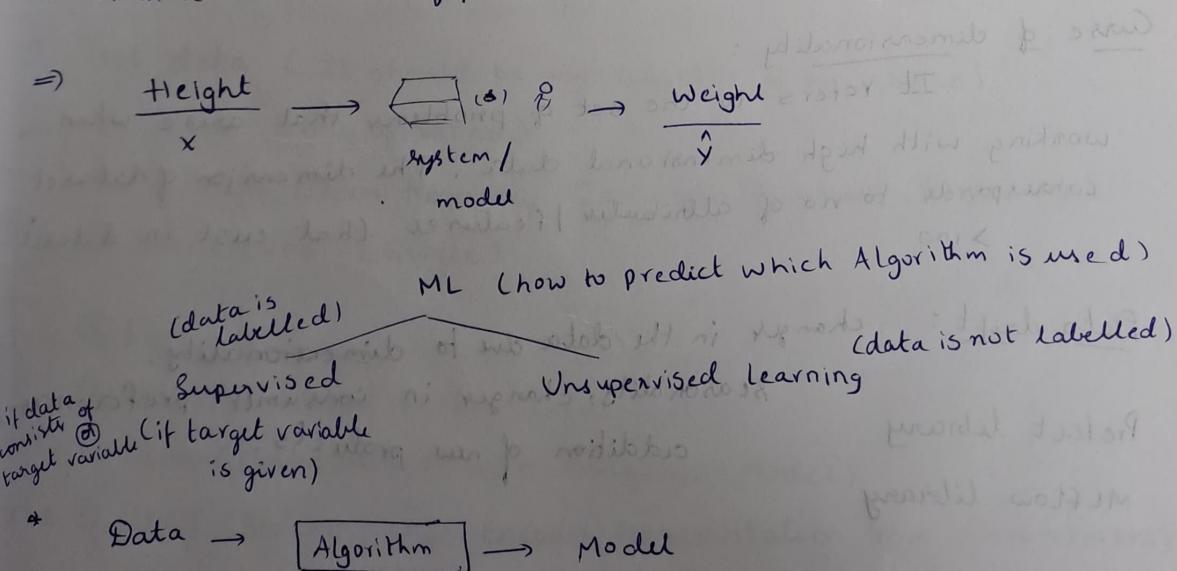
Classification	Regression	Clustering	Reduction
• KNN	• KNN	• K Means	
• Logistic	• linear	Hard clustering	
• System Vector class Support	• SV	Algorithm	
• DT class	• DT	• Hierarchical clustering	
• NB class	• prob of linear Reg (MLE)	• DB Scan	
• Voting Ensembles	• Voting Ensembles	• GMM Algo	
• stacking	• Stacking	• EM Algorithm	
• Random forest (RF)	• RF	Gaussian Mixture Model Expectation Maximization	
• Cascading	• Cascading		
• ADA Boost,	• ADA Boost,		
GBDT, XGBoost	GBDT, XGBoost		
	• ANN Reg		

## MLE - Maximum likelihood Estimation

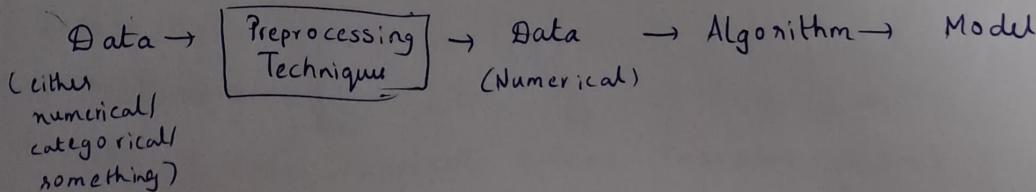
## MAP - Maximum A-Posteriori

## Type of Problems:

- ① Given the height, predict Weight
  - ② Given SL, SW, PL and PW
  - ③ Given an email, op if the email is Spam or not
  - ④ Given a text review of product. Predict the rating



## Actual process



Clustering - grouping similar data or images.

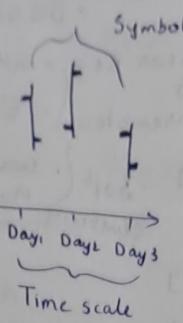
Understanding of Quiz Questions

Curse of dimensionality

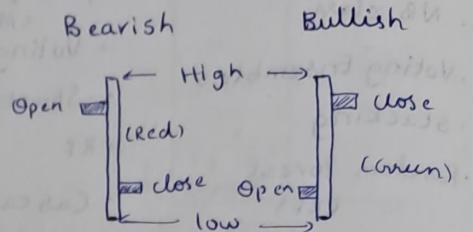
Edit url of linkedin

① Stock marketing : \$40

- Technical Analysis
  - uses historical price of stocks to predict future price
- Qualitative Analysis
  - external factors like scale company's profile, market situation, political & economic factors.



Symbol Anatomy :



Here Open means opening time price  
Close - closing time price

② Spam data in emails are found. If it is labelled email then we use Supervised learning. So : the tasks we use are

- Logistic Regression
- SV machines
- Random forest classifiers.

③ forecast water utilization for different input data

④ Pricing of house based on features like bedroom, sqft area, location (latitude, longitude), floor, condition of house.  
For all these data at first we must know the appropriate selection of features for our outcomes.

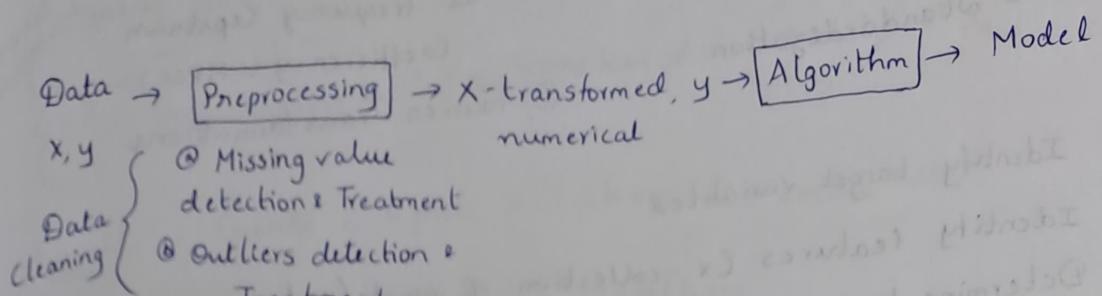
Curse of dimensionality :

It refers to the set of problems that arise when working with high dimensional data. The dimension of dataset corresponds to no of attributes / features that exist in dataset. > 100

Data drift : change in the data due to dimensionality.

Prefect library  
MLflow library  
seasonality, change in consumer preference, addition of new products.

# Data Preprocessing and [Process] Lifecycle:



→ Categorical Column (predefined & certain categories)  
Eg: Crop type, Soil type, Gender, Nationality column

- ① One hot encoding | Dummy Encoding } help to convert categorical col to numerical col
- ② Label encoding }

categorical

Nominal  
(No order)

Eg: Gender

M  
F  
M  
F  
F

Ordinal  
(Order)

Eg: Grade

Ex  
Avg  
Good  
Ex  
Avg

S-300	I-300
L	L
P	E
O	O
H	S

\* Use One hot encoding

\* Label Encoding

→ Text data (It should be any length & of our choice)

- ML { ① Bow (Bag of Word)  
 ② TF-IDF (Term frequency - Inverse Document frequency)

DL { ③ Word2Vec (Google)  
 (2013)

④ FastText

⑤ Glove (Global vector Encoding)

2017 ⑥ ELMo (Embeddings from language Model)

2018 ⑦ BERT (Bidirectional Encoder Representation from Transformers)

SORT

→ Image Data

⑧ Flattening } ML

⑨ SIFT (Scale Invariant feature Transformation)

⑩ CNN [extract features]

- Numerical Data
  - ④ Normalization
  - ④ Standardization
  - Audio Data
  - ④ Mel-Scaled filter bank
  - ④ Mel-frequency Cepstrum Coefficient (MFCC)
  - ④ Fourier Transformations
- ④ Identify target variables
- ④ Identify features ( $x$ -collection of input features)
- ④ Determine type (Supervised / Unsupervised)
- ④ Determine task →  $y \in \{ \text{'spam'}, \text{'ham'} \}$  → classification  
 $\uparrow$   
 belongs to  $\nwarrow$  Discrete values
- $y \in \mathbb{R}$  (Real no. value → Regression  
 in 1-D)

Example:

$x$	$y$
Col-1	Col-2
1	1
3	9
0	0
2	4
5	25

label  
 (spam or ham)  
 (spam or ham)

clear vs label  
 (spam or ham)  
 in a column  
 spam & ham  
 are diff  
 clearr,  
 2 or more

quadratic mean square of one variable

$$\text{Col-1} = 5 \rightarrow$$

$$\text{Model} \rightarrow \text{Col-2} = \text{Col-1} * \text{Col-1}$$

Multi label classification: (multi output classes)

Movie Summary →

Comedy	Horror	Drama
T	F	T

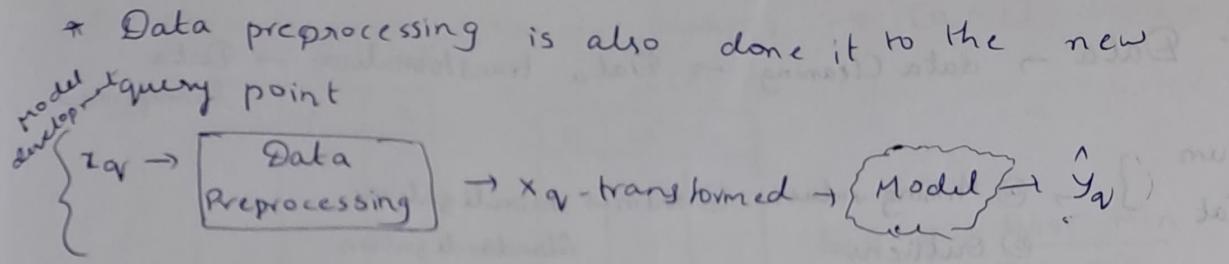
\* Multiple target variables. If  $y$  has more than 1 target

Binary Classification: (either 0 or 1) Variable.

→  $y \in \{ \text{'spam'}, \text{'ham'} \}$  If  $y$  has only two categories

Multiclass Classification:

→  $y \in \{ \text{'setosa'}, \text{'versicolor'}, \text{'virginica'} \}$  If  $y$  has more than 2 categories.



Email: Create a spam detection for email data

email-text  $\rightarrow$  **R**  $\rightarrow$  spam/ham

Step:

① target  $\rightarrow$  Email-label  $\in \{\text{spam}, \text{ham}\}$

Data  $\rightarrow$

email	label
...	spam
...	ham

Type  $\rightarrow$  Super learning algo can be applied

Task  $\rightarrow$  Classification

Dataset  $\rightarrow$   $x$ ,  $y$

text label  $\rightarrow$  **data preprocessing**  $\rightarrow$

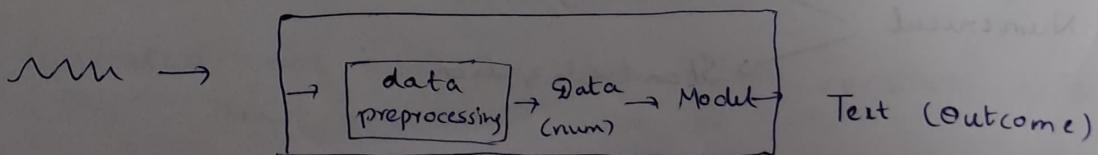
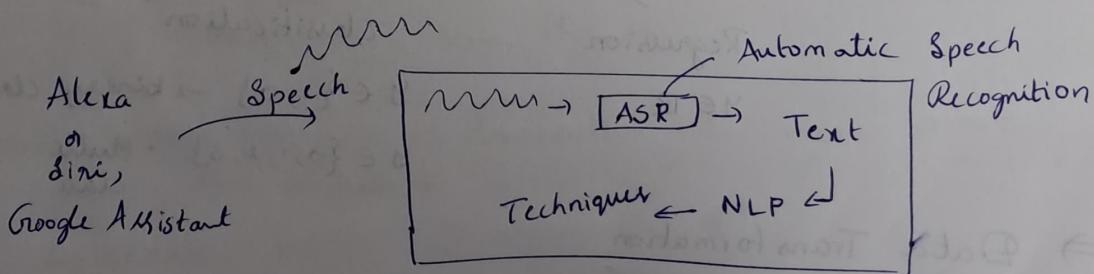
Q) How do you use this model?  $x$ -transformed by  $y$  (num)

email  $\rightarrow$  Model  $\rightarrow$  Label

email  $x_q \rightarrow$  Data  
(text) Preprocessing  $\rightarrow$   $x_q\text{-transformed}$  (num)  
also using  $\leftarrow$  Model

\* SME - <sup>object</sup> matter expert

\* feature Engineering - deriving new features from existed



Alexa, YouTube subtitles.

\* Data  $\rightarrow$  data Cleaning  $\rightarrow$  Data transformation  $\rightarrow$  Data transformed

Num

Cat

$\rightarrow$  ① Missing values ② Num  $\rightarrow$  Normalization/ Standardization  
③ Outliers

④ Remove noise & ⑤ Cate  $\rightarrow$  OHE  
abnormalities from LE

data

Non numerical to  
numerical data

Text

$\rightarrow$  ① Special char Removal  
② lower case Conversion  
③ Stopword Removal  
④ Stemming /  
Lemmatization

① BOW ② ELMO  
③ TF-IDF ④ BERT  
⑤ W2V ⑥ fasttext  
⑦ Glove

All are  
required

Vectorization means  
Convert Non numerical to  
numerical

Why these techniques? (for categorical features)

To convert categorical into numerical features

Why? (for numerical columns)

Rescaling because for each feature it has some own scale  
which is different for everyone / every feature

Eg: Height, Soft, Weight, age . So that all numerical  
features follows same scale.

$x$  - transformed,  $y$   $\rightarrow$

Algorithm

$\rightarrow$  Model

Regression

Classification

$y \in \mathbb{R}$

$y \in \{1, 0\} \rightarrow$  binary classification

$y \in \{0, 1, 2, 3\}$  - Multi class

## Data Transformation

- Numerical

$\rightarrow$  Normalization

$\rightarrow$  Standardization

① Normalization (Min-max scaling)  $\rightarrow [0 - 1]$

Cg

[MS. 190]

H	w'	G
-	-	
-	-	
-	-	
-	-	

(0-1)

$$\text{for } H(\text{cd}) : \frac{\text{Obs}i - \text{min}}{\text{max} - \text{min}}$$

$$\text{for } w(\text{cd}): \frac{\text{Obs}_{i(w)} - \min}{\max - \min}$$

$$\text{for H (sol)}: \frac{140 - 140}{190 - 140} = 0 \text{ (min)}$$

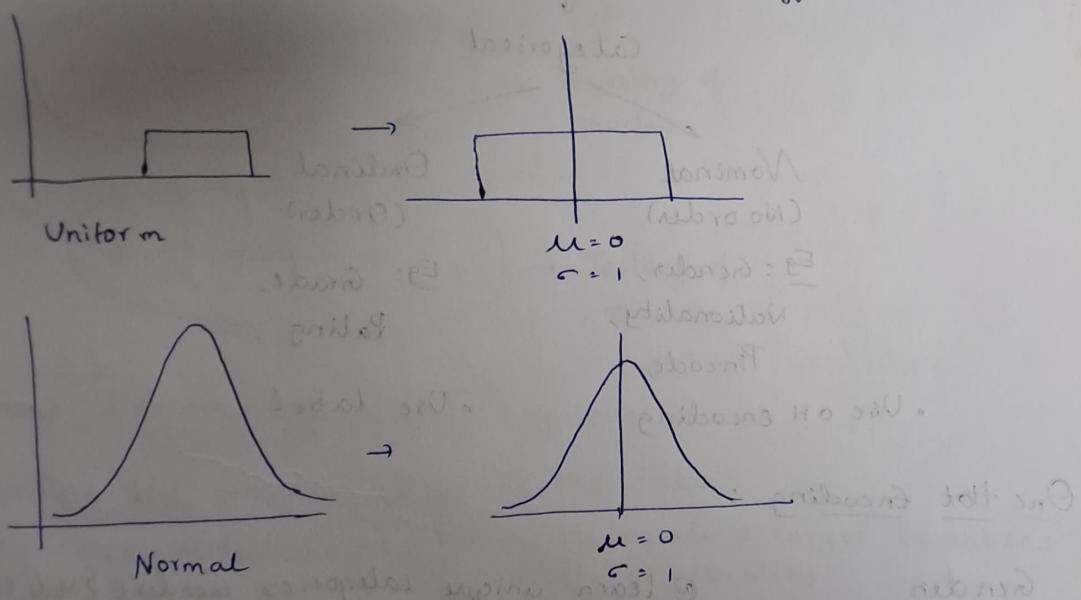
$$= \frac{190 - 140}{190 - 140} = 1 \text{ (max)}$$

- \* It brings all the values to scale  
[0 - 1]

## ② Standardization (z-transformation)

$$\begin{cases} \rightarrow \mu = 0 \\ \rightarrow \sigma = 1 \end{cases}$$

Robert Scaler



$Z$  score  $\rightarrow$  how many standard deviation away is the observation from the  $\mu$

၁၂

how far away is observation from mean  
std deviation

$$Z\text{-score} = \frac{\text{Obsi} - \mu}{\sigma}$$

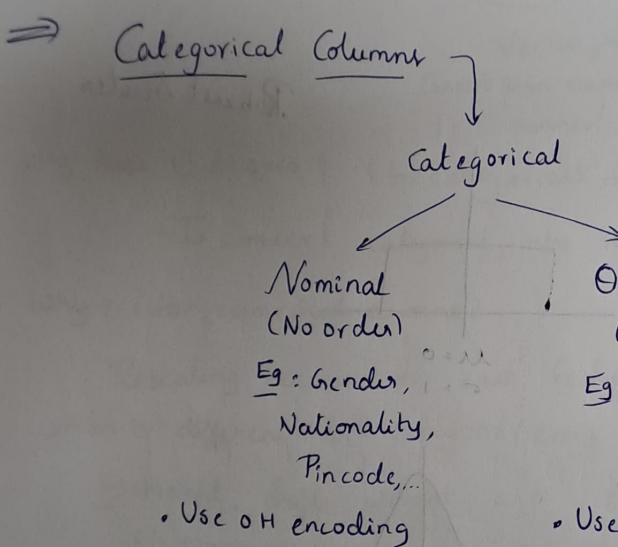
\* deprecated - not applicable for future versions just throw errors.

- \* DNS - Domain name server
- \* Spread - how far data points away from mean
  - σ (for sample)
  - σ (for population)

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

- $H \xrightarrow{\text{Standardization}} H' = \frac{obs_i - \mu_H}{\sigma_H}$
- Here, outliers are going to corrupt / effect the data
- Conduct same technique either standardization / normalization for all features in a dataset



One Hot Encoding :-

① Gender

Male
Male
Female
Male
Female



② learn unique categories  $\text{Gender} \in \{\text{Male, Female}\}$

③ Apply Onehot Encoding

(drop the first column) → curse of dimensionality

because if we do  
Gender - M then

Encoding for male - 1

for female - 0

(unique right)

④ Multi  
collinearity  
btw features

Gender - M	Gender - F
1	0
1	0
0	1
1	0
0	1

Q: What is encoding for  
Gender = Male?

Ans

## Gender

Male
Female
Male
Male
Female

① Identify unique categories

$$\text{Gender} \in \{\text{'Male'}, \text{'Female'}\}$$

→ ② Label categories

$$\begin{cases} \text{'Male'} : 1 \\ \text{'Female'} : 2 \end{cases}$$

③ Apply label encoding

1
2
1
1
2

## Grades

Excellent
Good
Avg
Good
Avg
Excellent

① Unique

$$\text{Grades} \in \{\text{'Excellent'}, \text{'Good'}, \text{'Avg'}\}$$

②  $\begin{cases} \text{'Excellent'} : 3, \\ \text{'Avg'} : 1, \\ \text{'Good'} : 2 \end{cases}$

③ Apply

3
2
1
2
1
3

Algorithm → Model

→ Training of model

$x_q \rightarrow \text{model} \rightarrow \hat{y}_q$

→ Testing of model

## Machine Learning Project Life Cycle:

① Define the problem & Check the data

i.e., Price prediction → identify all inputs & target variables  
 $x \rightarrow \text{sqft, age, locality}$   
 $y \rightarrow \text{Price}$

Spam prediction

$x \rightarrow \text{email}$

$y \rightarrow \text{label (spam, ham)}$

Task

spam
ham

Regression

Classification

Algorithm

KNN

linear

SV

DT

Evaluation matrix

1. Mean squared error

2. " absolute "

3. Root mean squared error

4. R<sup>2</sup> (coefficient of determination)

5. Adjusted R<sup>2</sup>

Algorithm

KNN class

logistic

SV

DT

RF

GBM

Evaluation Matrix

1. Accuracy

2. Confusion matrix

3. Precision & Recall

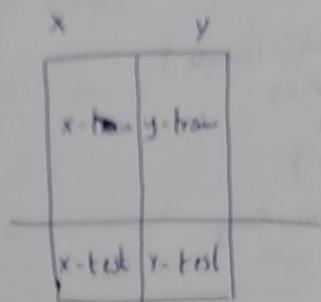
4. F1-score

5. ROC-AUC Curve

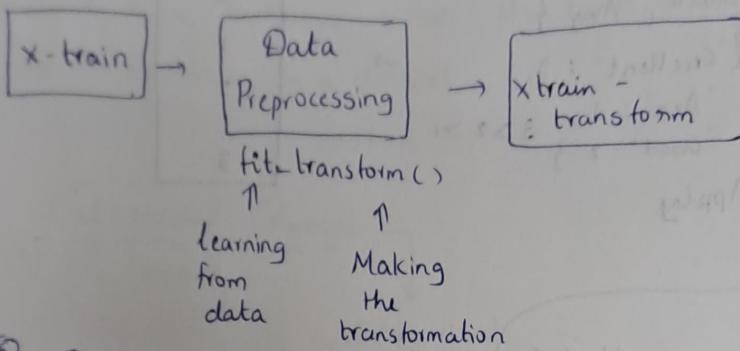
6. Log loss

② Identify input & Output

③ Split the data for training & testing



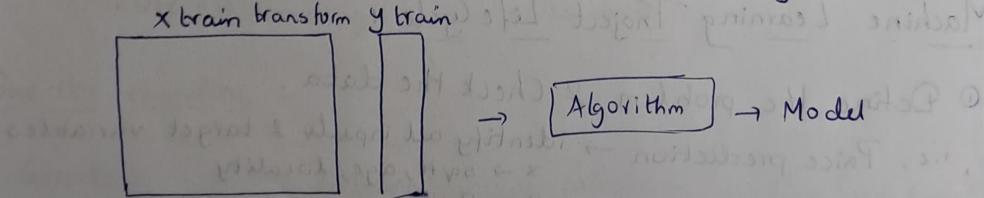
④ Data Preprocessing on x-axis train



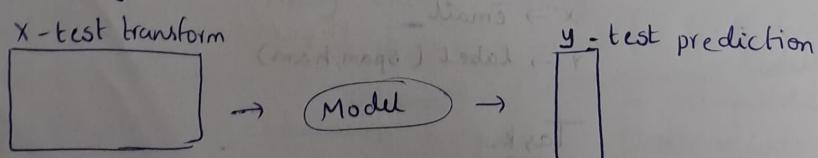
⑤ Data preprocessing on x-test



⑥ Apply ML Algorithm to learn from training data



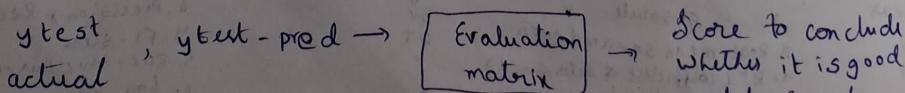
⑦ Testing your model on unseen data



Here we already have y-test so just compare it with model y-test predictions.

⑧ Predict on unseen data

⑨ Evaluation of your prediction



\* Version should be 1.1 or above  
Coding part :  $\Rightarrow \text{print}(\text{sklearn}.\text{--version})$  → class

» from sklearn.tree import DecisionTreeClassifier  
model = DecisionTreeClassifier() [object creation]  
model.fit(x-train-transformed, y-train)

y-test-pred:  $\underbrace{\text{model}.predict(\text{x-test transformed})}_{\text{predictions}}$

from sklearn.metrics import accuracy\_score  
 $\text{print}(\text{accuracy-score}(y-test, y-test-pred))$

④ for splitting the data into train and test

» import sklearn ( $! \text{pip install -U scikit-learn}$ ) → here data is randomly shuffled.  
from sklearn.model\_selection import train\_test\_split  
x-train, x-test, y-train, y-test = train\_test\_split (x-data-name,  
y-data-name, train\_size = 0.8/0.7, random\_state = 100)

⑤ Separating Categorical and Numerical Columns  
categorical:

» x-train-cat = x-train.select\_dtypes (include = ['object'])  
numerical:

» x-train-num = x-train.select\_dtypes (include = ['int64', 'float64'])

⑥ Scaling the numerical features

» from sklearn.preprocessing import StandardScaler  
scalen = StandardScaler()  
x-train-num-rescaled = pd.DataFrame (scalen.fit\_transform  
(x-train-num), columns = x-train-num.columns,  
index = x-train-num.index)  
x-train-num-rescaled

≡ Based on the Mean Absolute error which value is less means best model i.e., Random forest Regression  
≡ Maximum time for training took by Random forest Regression Algorithm.

Another technique to rescale the numerical features i.e.,  
④ Normalization

» from sklearn.preprocessing import Normalizer

» from sklearn.preprocessing import MinMaxScaler

fit\_transform means here fit represents learner from  
the data i.e., u & s for every feature for transforming  
using Standardization

① Applying OneHotEncoding for Categorical Columns

» from sklearn.preprocessing import OneHotEncoder  
encoder = OneHotEncoder (drop='first', min\_frequency=3000,  
sparse=False)

# Image Preprocessing and Computer vision:

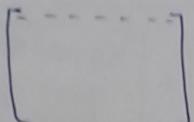
import pil → Python imaging library



Here along  $w \rightarrow 720$  pixels

along H → 540 pixels

1 pixel = 1 Numerical value



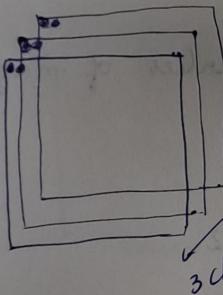
$540 \times 720$

[ R G B ]

White - R G B  
255 255 255

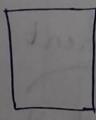
Black - R G B  
O O O

intensities



→ Every colour is formed with the intensities of R G B

colours (13) diamonds, pink



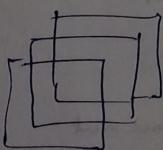
→

Black & white  
(Grey scale)

→ No need of multiple channels

$$\rightarrow \text{No. of pixels} = 64$$

Aspect ratio =



→ Coloured image  
(RGB)

$$\rightarrow \text{No. of pixels} = 64 * 3$$

» from matplotlib import image  
(read the image & displays the opf. of image)  
\* The image is stored in python as numpy array

## Image Manipulation with PIL:

① Rotate

② Crop

③ Flipping

④ Resize

flip left to right

flip top to bottom

flip transpose (also done with rotate)

\* from PIL import Image

from PIL import ImageDraw

from PIL import ImageFont

→ Resizing is not automatic but thumbnail itself takes the size or per aspect ratio (W/H)

→ Image Enhancer (To enhance / ↑ the properties of image)

enhance(0) - black & white

enhance(1) - original

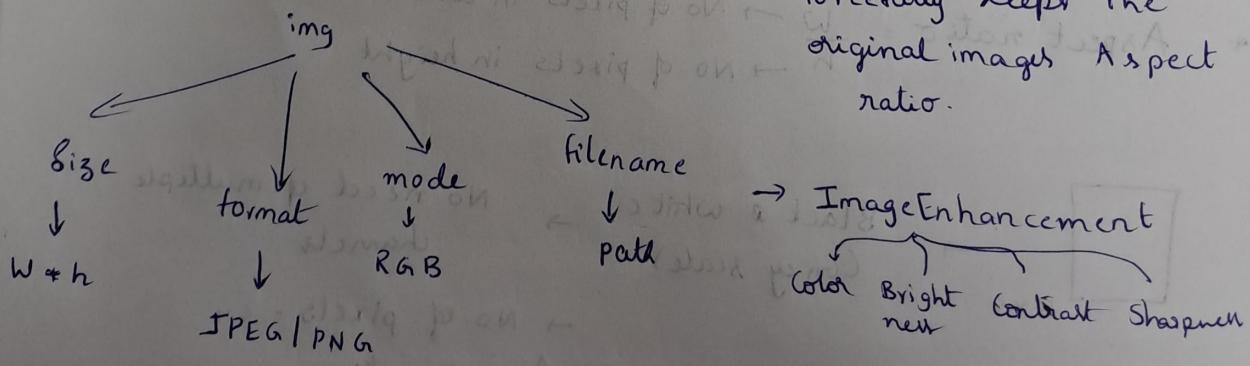
enhance(2) - ↑ intensity of colour by 2

① Color

② Contrast

\* img.resize((W, H))

## Properties of image



\* Aspect ratio = img.size[0]  
img.size[1]

→ Image Filters

→ Blur

→ Smooth

→ Edge Enhancement

→ Emboss

→

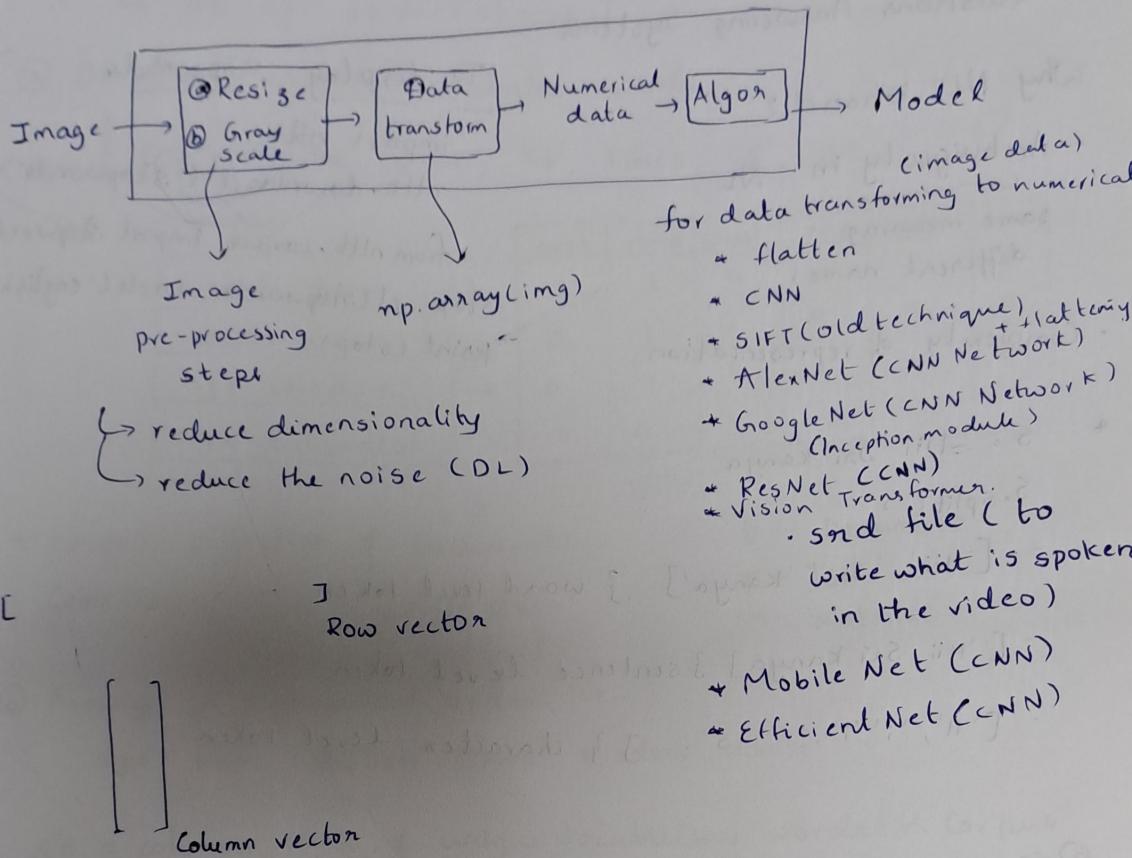
## Reducing the Dimensionality :

1. Reducing the Channel → Convert to Grayscale

plt.imshow(img.convert('L'))

gray-image = image.convert('L')

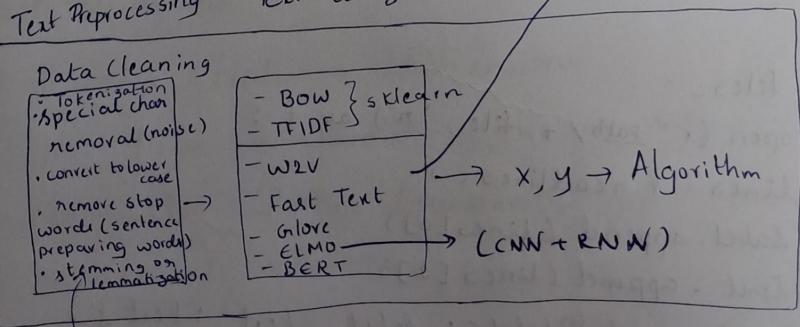
2. Resize of the image



→ To clean the outliers we make use of log transformation

Text data :  $\xrightarrow{\text{Techniques to transform text to numerical}}$  RNN | LSTM | GRU  
Transformer (SOTA)

DL techniques for text vectors  
Text Embedding i.e., dense vectors in lower dimensional space.



nltk / spacy (Natural language tool kit)  
for W2V {Gensim  
GloVe }  
BERT {sbert}

1. Spell checker, keyword search etc
2. Sentiment analysis, Spam classification
3. Machine Translation → language translation
4. Chatbots / Dialog Systems
5. Question Answering Systems

Why NLP is hard?

- Ambiguity in NL  
same meaning  
different names
- Complexity of representation

To display stop words

```
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stop = set(stopwords.words('english'))
print(stop)
```

- \* S = 'Hii Sri Kanya'
- S.split()
- ['Hii', 'Sri', 'Kanya'] } word level token
- ['Hii Sri Kanya'] } sentence level token
- ['H', 'i', 'i', ' ', 'S', ...] } character level token

### Quiz :

```
> import os
      ↗ path
      print(os.listdir('text'))
```

files\_ = os.listdir(r'path')

label = []

text = []

for file in files\_:

    with open(r'path/ + file, 'r') as f:

        lines = f.readlines()

        label.append(lines[0])

        text.append(lines[1])

df = pd.DataFrame([{'label': label, 'text': text}])

## Vectorization:

### Bag of Word:

[ 'it was the Best of Times \$',  
 'It Was The worst of times.',  
 'IT war g the age of wisdom',  
 'it war THE age of foolishness' ]

'best times',  
 'worst times',  
 'age wisdom',  
 'age foolishness' ]

① Learn the unique vocabulary words

{ 'it', 'was', 'the', 'Best', 'of', 'Times', '\$', 'It', ... }

best times
worst times
age wisdom
age foolishness

BOW

best	times	worst	age	wisdom	foolish
1	1	0	0	0	0
0	1	1	0	0	0
0	0	0	1	1	0
0	0	0	1	0	1

\* corpus = collection of Documents

D - Document

T - Term

M - Matrix

Each Document = A single data

② Arrange in alphabetical order i.e.,

age best foolish times wisdom worst

No. of columns = No. of unique vocabulary words in Corpus

Most common value = 0

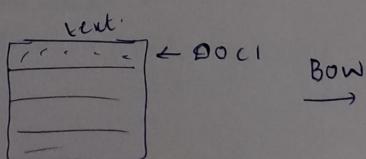
During the process of BOW:

③ Learn all the unique words from the Corpus

④ Sort the unique words & form the features

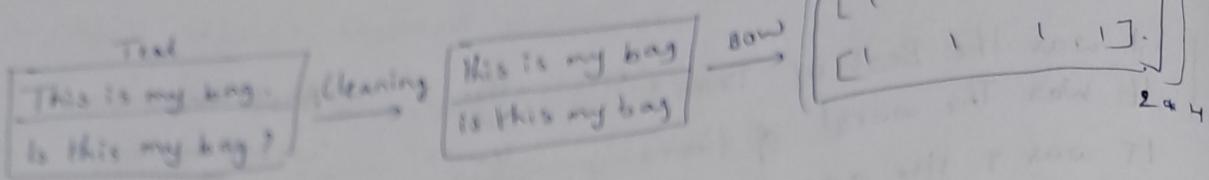
\* No. of col's / features / dimensions  
 - No. of unique words in the Corpus.

⑤ Transformation  $\rightarrow$  All the documents gets transformed to numerical vectors.



Count of the corresponding word  
 n + d  
 no. of unique words

Issue → Position is not Preserved



Try to somehow include the positional information

n-gram approach:

1-gram - { 'this', 'is', 'my', 'bag' }

2-gram - { 'this', 'is', 'my', 'bag', 'this', 'is', 'is', 'my', 'my', 'bag', 'is', 'this', 'this', 'my' }

# Bag of words

from sklearn.feature\_extraction.text import CountVectorizer

Step:

① Identify target variable

Classification      Regression

CSN  
↓  
Compressed sparse matrix  
format

② Split the data into train & test

③ Apply text cleaning & transform on the train data

④ Apply cleaning on test data.

⑤ Build the model using train data

$x_{\text{train-dtm}}$  → Model →

⑥ Use the model for pred on unseen data

$x_{\text{test-dtm}}$  → Model →  $y_{\text{test-pred}}$

⑦ Evaluate the performance of Model

## In Chatbots (CRASA)

① Intent Classification is a 'Classification' problem. [intention]

$\text{intent} \in \{\text{'greet'}, \text{'good bye'}, \dots\}$ , check\_bank\_balance  
check\_course\_detail  
check\_weather, call\_sombody

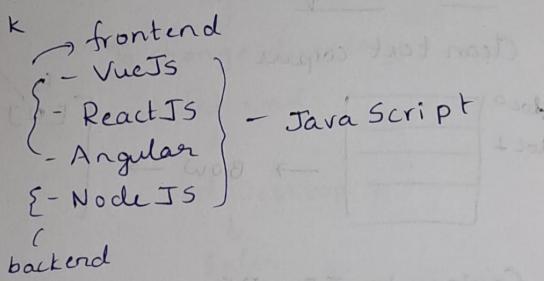
For creating backend application, frameworks used are

- Flask
  - Django
  - Fast API
- } to create backend server.

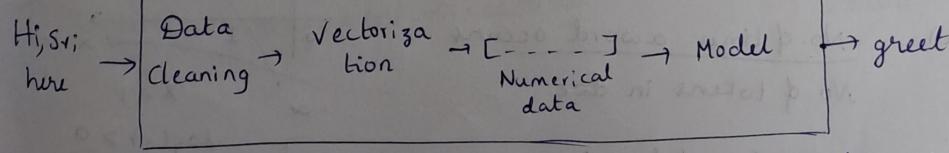
Server responder to the request given by any person

CDN - Content delivery network

- Spring
  - Spring boot
  - .Net
- } Java framework



## ② Entity Extraction



Intent classification

vs

Intent classification +  
Entity extraction

I/p : Hi, Sri here

I/p : Hi, Sri here

O/p : Hi user

O/p : Hi Sri  
How are you?

How are you ?

## Spam/ham detection:

### 1. Data Visualization [Countplot]

To get count of values in %, we use `Normalize = True`.

- `tqdm` is a library used for creating progress bars.
- `wordcloud` to represent the words in a visual pattern. Larger words represent the high frequently occurred words.

Check 31 cell in this python file i.e., (AEDA) ~~related at~~

→ `spam_df = x_train[x_train['y_train'] == 'spam']`

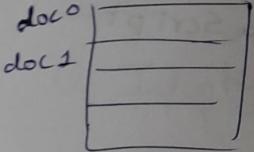
## Preprocessing Technique:

### \* TF-IDF (Term Frequency - Inverse Document frequency)

- Tell us the importance of words.

#### BOW Technique

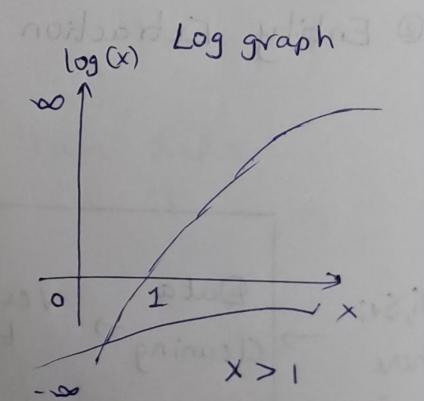
Clean text corpus



Coming to TF-IDF

- how important the word is on doc level

(TF) → how important this word is on corpus level



$$* \text{TF}(w, \text{doc}) = \frac{\text{No. of times a word occurs}}{\text{No. of tokens in doc}}$$

$$\text{IDF}(w, \text{corpus}) = \log_n \left( \frac{\text{No. of docs in corpus}}{\text{No. of doc which contains words}} \right)$$

$$\log(x) > 0$$

- (Python code)
1. python installed
  2. Virtual Environment using python  
cmd → python → in venv name-of-env
  3. Activate the environment  
cmd → name-of-env\Scripts\activate
  4. Install Rasa { Rasa is a framework that helps in chatbot development }

The version of python should be '3.9'.  
environment  
→ for linux / mac we bin scripts activate.

- Good server gives data in less time
5. Create a Rasa Project

cmd → rasa init  
↓  
command prompt

Some rasa Commands :

- Initialize the project

• rasa init

- To interact with the chatbot

• rasa shell

- To train the chatbot

• rasa train

To jump back to previous directory

- cd ..

To jump into current directory

- cd .

\* To Know python version

python -v

python --version

data/inlu.yml

Contains intents

Eg:

data/stories.yml

contains flow of dialogues

domain.yml

- define the intents

- contains utterances

Create new folder



click cmd : cd m



To change directory

click cmd



cd Desktop



Enter : addresses

To check whether the rasa is installed

↓ pip list

To stop the server

- /stop

- .ipynb → interactive python notebook
  - .ipynb → Only open on jupyter notebook
  - .py → Open on notepad i.e., text file
  - .yml → "
  - intent → 'greet'
  - action → 'utter-greet'
  - intent → 'mood-unhappy'
  - action → 'utter-cheer-up'
- Story  
↳ defines dialogue flow

To train a chatbot :

- Identify the intent

### data/inlu.yml :

- nlu:
  - intent : greet
    - examples:
      - hey
      - hello
- intent : goodbye
  - examples:
    - cu
    - goodbye
    - ceeyou later
    - bye
    - Goodnight

### data/stories.yml :

- stories:
- story : happy path
  - steps:
    - intent : greet
    - action : utter greet
    - intent : mood great
    - action : utter happy

### domain.yml

- For Eg:
- intents:

  - greet
  - goodbye

- responses:
- utter greet :

  - text: 'Hey ! How are you ?'

- utter goodbye :

  - text: 'Bye'

- utter happy :

  - text: 'Great , carry on'

- Single Shot detection
- Zero Shot detection.