# TEXT CLASSIFICATION

-BUILD A CLASSIFIER MODEL USING NAVIE BAYE'S ALGORITHM

PRESENTED BY:
YUVARAJ.S.S-203052
SREEMANN.S-203017
VASANTH.P-203302

# INTRODUCTION

- A news article discusses current or recent news of either general interest or on a specific topic.

- For example ,political ,trade ,technology's etc......

- Every news websites classifies the news article before when we go to a websites, so visitors can easily click on the type of news of interest.

- This leads to save the visitors time and might helpful for all.
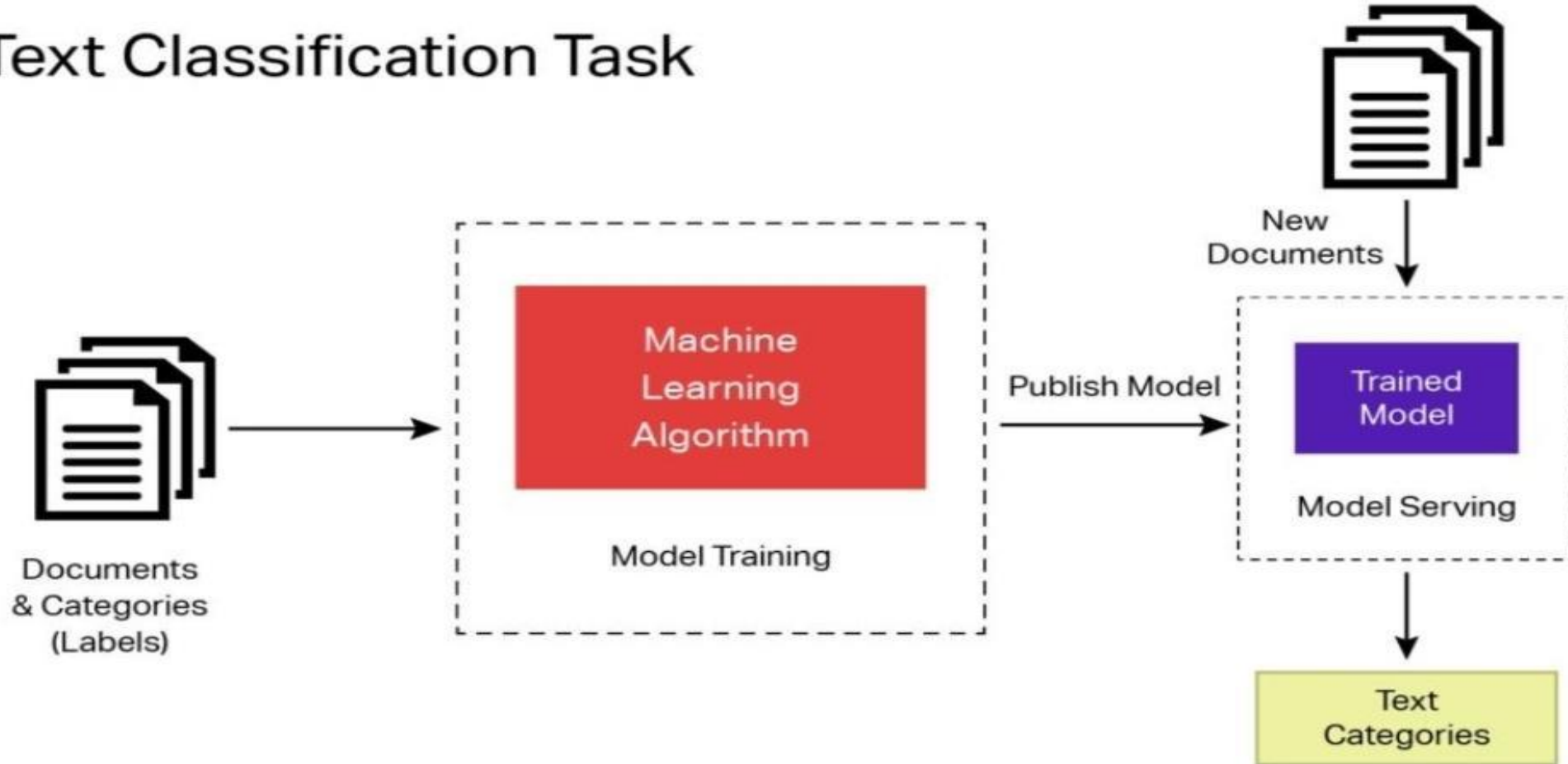
# DATASET-Text classification

- Text classification dataset are used to categorize natural language processing texts according to the content.

- Some of the dataset providing websites for text classifications are,

1. Kaggle
2. Blogger.com
3. Newsgroup
4. Github.com

# DATA CLEANINNG AND PRE PROCESSING

- Data processing is the process of transforming raw data into an understandable format.

- The quality of data should be checked before applying the machine learning or data mining algorithms.

- It is an important step because, we cannot work on the raw data.

# NAVIE BAYES ALGORITHM

- Navie bayes algorithm is a collection of classification algorithm based on bayes theorem.

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B)}$$

- A,B=events
- P[A|B]=probabilof A given B is true
- P[B|A]=probabilof B given A is true
- P[A],P[B]=the independent probabilities of A aand B

# WORKING

- I can take the dataset from "Github.com" ,which is a news dataset.

- The dataset is "*dsjVoxArticles.tsv* file" and its format is TSV file.

- TSV stands for "TAB SEPERATED VALUES".

- It holds the data in the tablelled manner.

Copy of news classification using navie bayes.ipynb

File Edit View Insert Runtime Tools Help Last saved at 6:04PM

Comment    Share

+ Code   + Text                                                                    Connect ▾

IMPORTING THE REQUIRED LIBRARIES

```python
1  import copy
2  import numpy as np
3  import matplotlib.pyplot as plt
4  import re
5  import nltk
6  nltk.download('stopwords')
7  from sklearn.model_selection import train_test_split
8  from sklearn.feature_extraction.text import CountVectorizer
9  from sklearn.preprocessing import LabelEncoder
10 from sklearn.feature_selection import VarianceThreshold
11 from imblearn.over_sampling import SMOTE
12 from sklearn.dummy import DummyClassifier
13 from sklearn.naive_bayes import MultinomialNB
14 from sklearn.tree import DecisionTreeClassifier
15 from sklearn.neural_network import MLPClassifier
16 from sklearn.ensemble import RandomForestClassifier
17 #from sklearn.metrics import accuracy_score
18 #from sklearn.model_selection import cross_val_score, KFold
19 from sklearn.metrics import confusion_matrix
20 from sklearn.metrics import classification_report
21 import seaborn as sns
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```
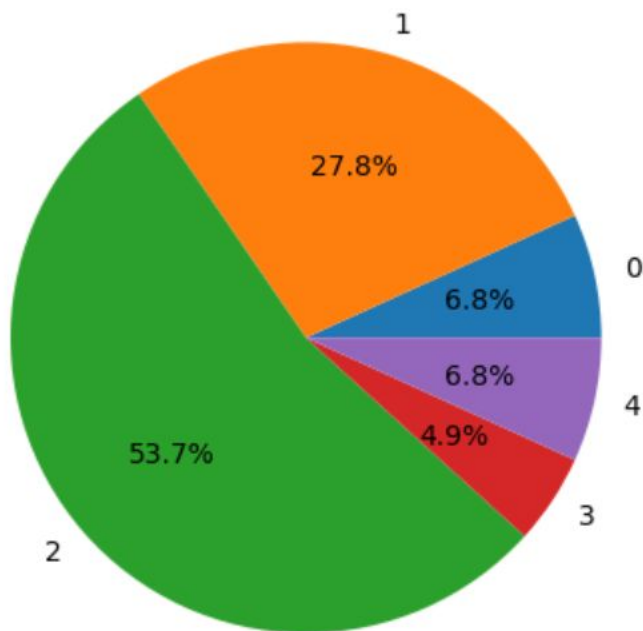
LOADING THE DATA FILE dsjVoxArticles.tsv file

+ Code    + Text                                                      Connect ▾    ⌃

SAMPLING THE DATA

```
1 labels = list(set(Ytr))
2 counts = []
3 for label in labels:
4     counts.append(np.count_nonzero(Ytr == label))
5 plt.pie(counts, labels=labels, autopct='%1.1f%%')
6 plt.show()
```
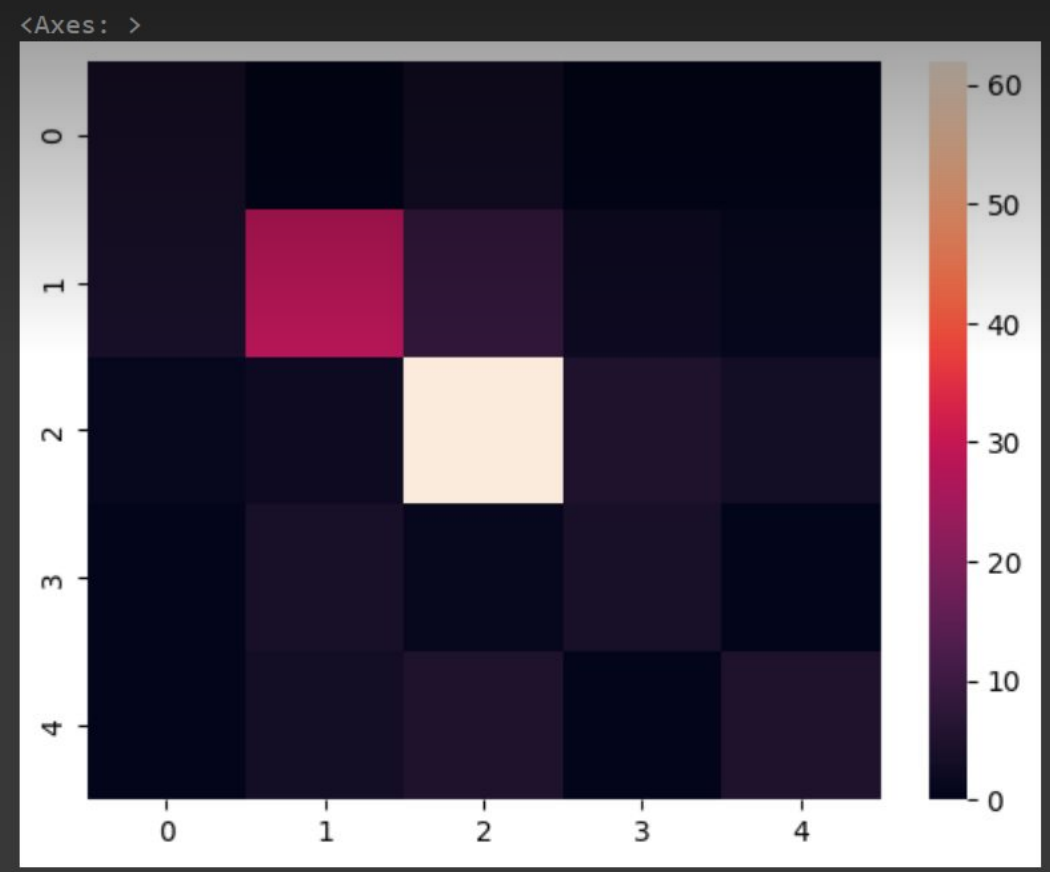
+ Code    + Text    Connect ▾    ⌃

| | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Politics & Policy | | 0.36 | 0.44 | 0.40 | 9 |
| Science & Health | | 0.56 | 0.38 | 0.45 | 13 |
| | | | | | |
| accuracy | | | | 0.71 | 145 |
| macro avg | | 0.58 | 0.58 | 0.57 | 145 |
| weighted avg | | 0.71 | 0.71 | 0.71 | 145 |

<Axes: >

THANK YOU