# NEWS CLASSIFIER
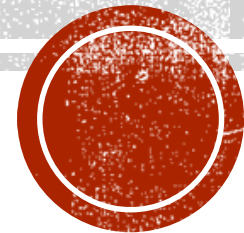
Members:  Mahsa Razavi – Mahsa Anvarian

Student ID: 94523144 – 95521054
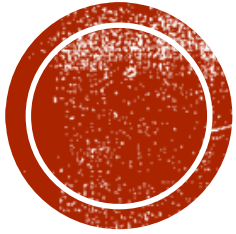
Professor: Dr. Pilevar

Project github link: https://github.com/mahsawz/AI-NewsClassifier

Iran University of Science and Technology
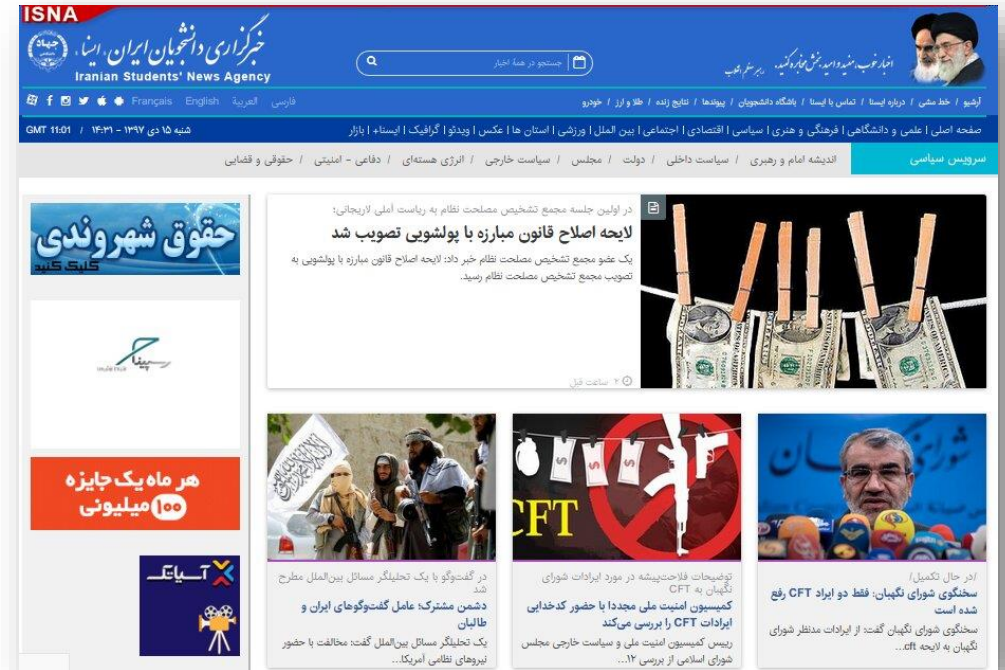
# INTRODUCTION

- **One of the widely used natural language processing task in different business problems is "Text Classification".**

- **Automatically classify the text documents into one or more defined categories.**

➤ Understanding audience sentiment from social media,

➤ Detection of spam and non-spam emails,

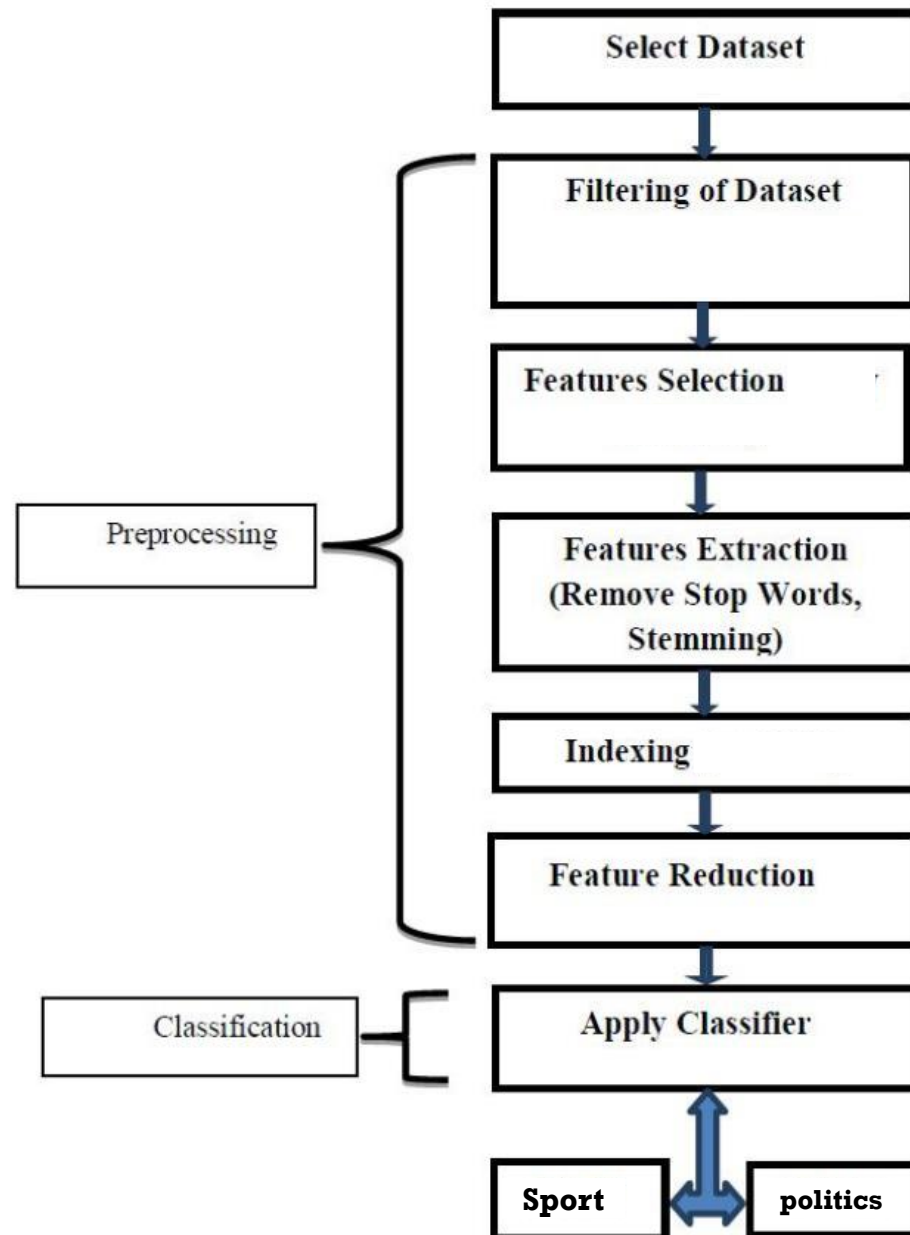➤ Categorization of news articles into defined topics.

# NEWS CLASSIFICATION



sport

**vs**



politics

Select Dataset → Filtering of Dataset → Features Selection → Features Extraction (Remove Stop Words, Stemming) → Indexing → Feature Reduction → Apply Classifier → Sport / politics

Preprocessing

Classification

Main Steps

# 1. DATASET PREPARATION

- **BBC News Dataset**

The dataset consists of text reviews and their labels which can be downloaded at this link.

| Class | Words |
|---|---|
| Sport | 89001 |
| Politics | 99335 |

# 1. DATASET PREPARATION

- **Text Cleaning**

➢ Lower Case
➢ Removing Digits
➢ Removing Punctuation
➢ Removal of Stop Words
➢ Tokenization
➢ Stemming
➢ Lemmatization
➢ …

# 2. NAÏVE BAYES

- Compute the probability of a label in condition of given features
- We need to specify how each feature (word) depends on the class

# 3. EVALUATING RESULTS

- Word Clouds

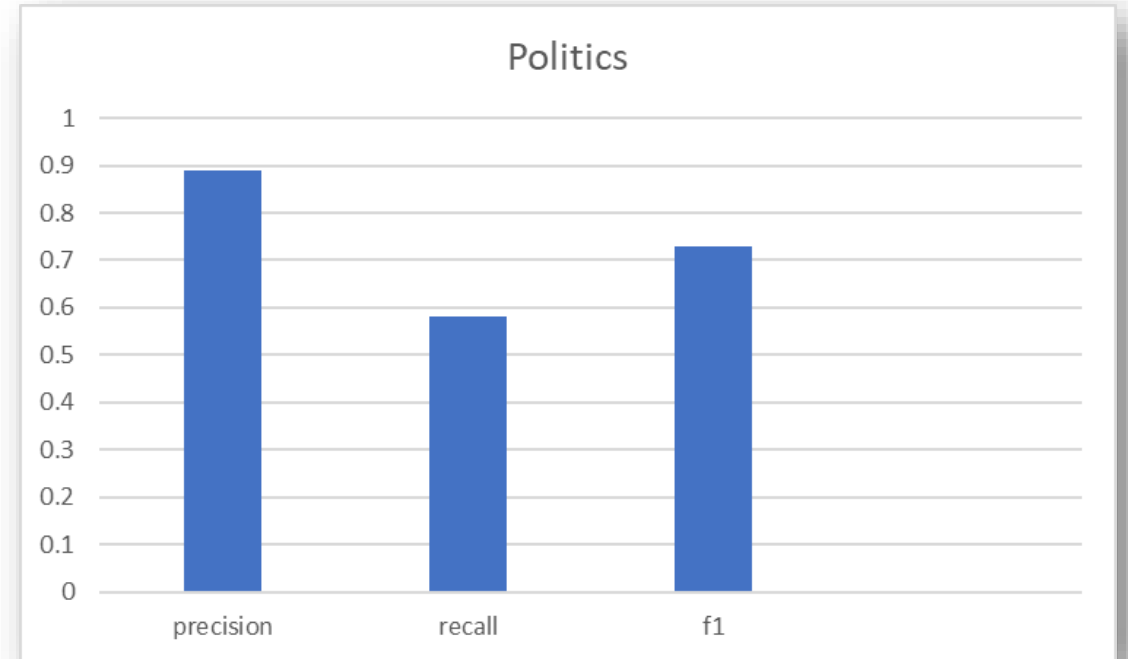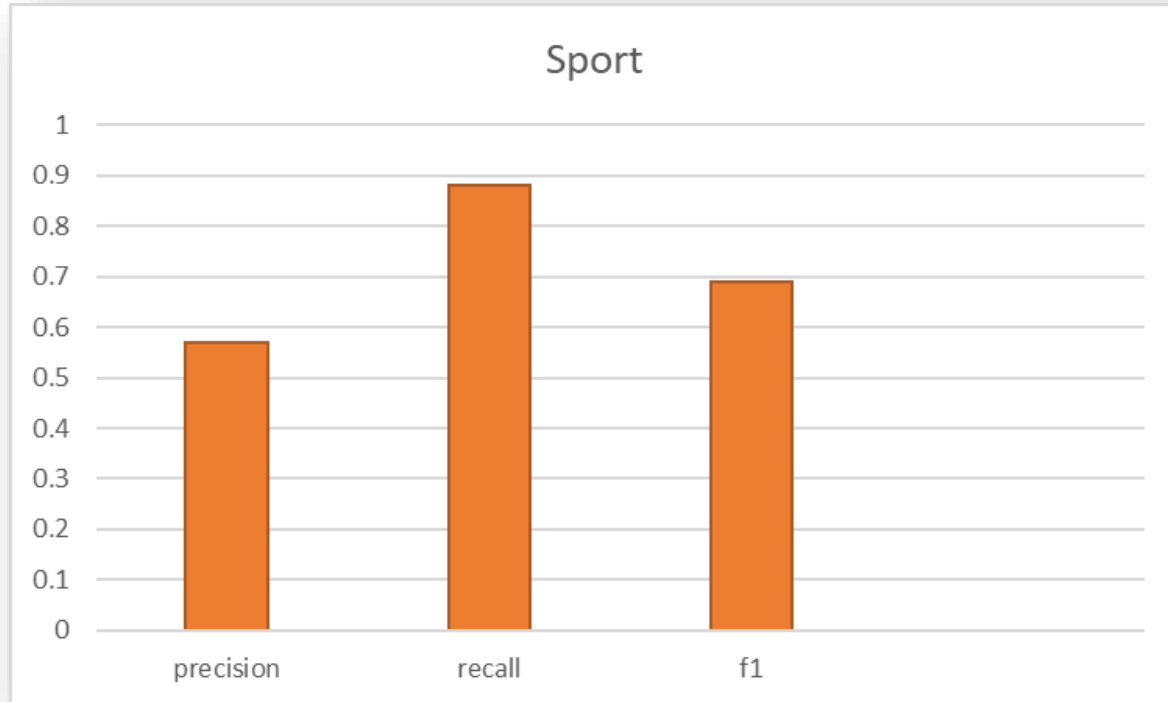# 3. Evaluating Results

Test data : 10%  of dataset devoted for testing

- Precision
- Recall
- F-score
- Accuracy

# 3. Evaluating results



Sport

| | precision | recall | f1 |
|---|---|---|---|

Politics

| | precision | recall | f1 |
|---|---|---|---|

# Naïve Bayes Summary

## Advantages:

- Fast to train (single scan through data)
- Fast to classify

## Disadvantages:
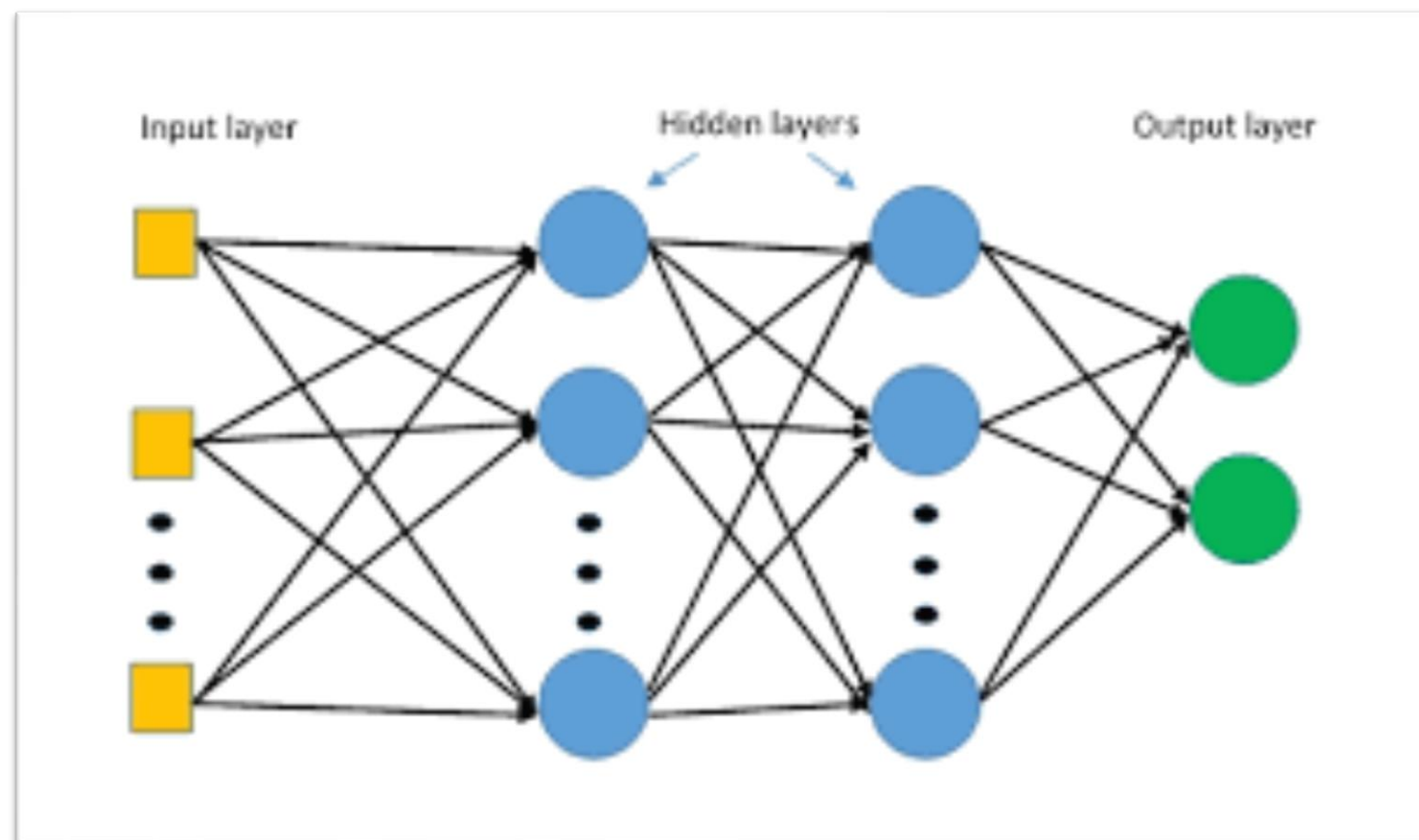
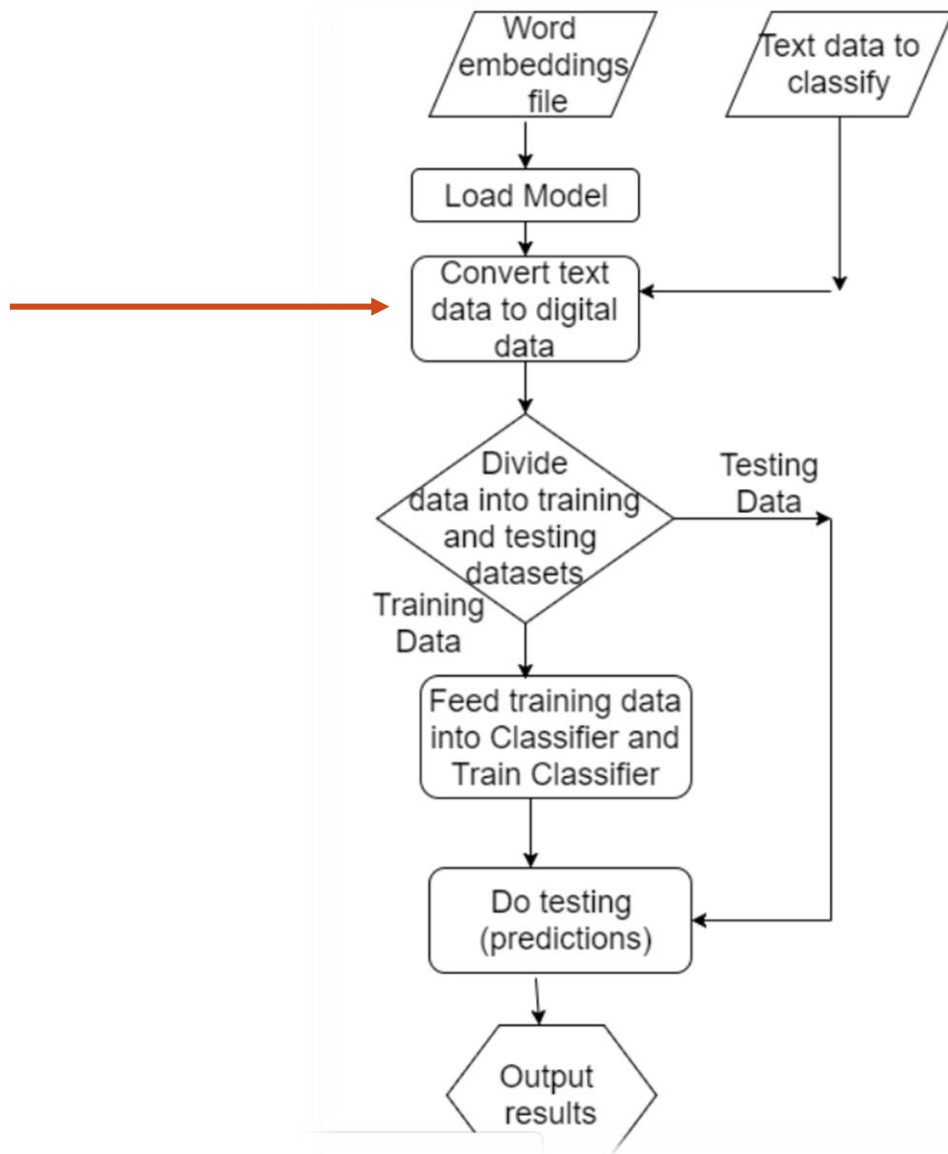- Assumes independence of features

CONCLUSION

# MLP (MULTILAYER PERCEPTRON)
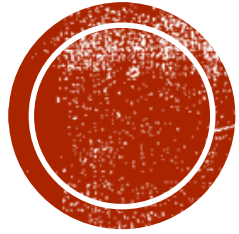
**Second Approach**

Total Accuracy 0.799

# MLP VS NB ?

- Acciracy
- Precision
- Recall
- Time taken to build model

| Approach | Accuracy | Precision | Recall | Time Taken To Build Model |
|----------|----------|-----------|--------|---------------------------|
| MLP | 93 | 93.2 | 93 | 10.94 Sec |
| NB | 88 | 88.2 | 88 | 0.14 Sec |

Performance measurement of both classifiers

# REFERENCES

- https://pdfs.semanticscholar.org/d945/29ae612b76287c7dcdccd5bf09f3c5e772af.pdf
- https://www.analyticsvidhya.com/blog/2018/02/the-different-methods-deal-text-data-predictive-python/
- https://www.saedsayad.com/naive_bayesian.htm
- https://www.analyticsvidhya.com/blog/2018/04/a-comprehensive-guide-to-understand-and-implement-text-classification-in-python/
- https://medium.com/@minbaekim/text-mining-preprocess-and-naive-bayes-classifier-da0000f633b2
- https://www.python-course.eu/text_classification_python.php
- https://www.datacamp.com/community/tutorials/wordcloud-python
- http://ai.intelligentonlinetools.com/ml/fasttext-word-embeddings-text-classification-python-mlp/