# News classification using natural language processing

**BALA KRISHNA**
Balagmastro15@gmail.com
8247536369

# News classification using natural language processing

## Abstract:

This project focuses on the application of natural language processing (NLP) techniques for news classification. The aim of this research is to develop an automated system that can accurately classify news articles into predefined categories, such as politics, sports, business, and entertainment, among others. The proposed methodology involves data preprocessing, feature extraction, and machine learning algorithms. The data preprocessing step involves cleaning, tokenization, and normalization of the text data. Feature extraction is done using bag-of-words and term frequency-inverse document frequency (TF-IDF) techniques. A number of machine learning algorithms such as Naive Bayes, Support Vector Machines (SVM), and Random Forest are used for classification. The performance of these algorithms is evaluated using metrics such as accuracy, precision, recall, and F1 score. The results demonstrate that the proposed system achieves high accuracy in classifying news articles into different categories. The system has potential applications in various domains such as media monitoring, opinion mining, and sentiment analysis.

## Objective:

The objective of this study is to explore the potential of natural language processing techniques for news classification, particularly in the context of low-resource languages. The proposed methodology involves data collection, preprocessing, feature extraction, and classification. The study aims to develop a system that can accurately categorize news articles into different categories, such as politics, sports, business, and entertainment, among others. The focus will be on languages with limited resources, where traditional rule-based or machine learning-based approaches may not be effective. The performance of the proposed system will be evaluated using various metrics such as accuracy, precision, recall, and F1 score. The ultimate goal is to provide a scalable and efficient solution for automated news classification in low-resource languages, which can be used by media outlets, researchers, and other stakeholders.

## Introduction:

With the explosive growth of digital media, the volume of news articles generated every day has reached unprecedented levels. This has created a need for automated systems that can help categorize, summarize, and extract relevant

information from news articles in real-time. News classification, the task of categorizing news articles into predefined categories, such as politics, sports, business, and entertainment, among others, is an essential step in building such systems. Natural language processing (NLP), a subfield of artificial intelligence that focuses on understanding and processing human language, provides an effective solution for automated news classification.

In recent years, there has been a significant advancement in NLP techniques, including data preprocessing, feature extraction, and machine learning algorithms. These techniques have been successfully applied in various natural language processing tasks, including sentiment analysis, opinion mining, and text classification. News classification is an important application of NLP techniques, where the goal is to develop an automated system that can accurately categorize news articles into different categories based on their content.

**Methodology:**

Data collection: The first step in the methodology is to collect a large corpus of news articles from various sources. The corpus should cover different topics and categories, such as politics, sports, business, and entertainment, among others.

Data preprocessing: The collected data needs to be preprocessed to remove noise and irrelevant information. This step involves text cleaning, tokenization, normalization, and stop-word removal. We may also apply techniques such as stemming and lemmatization to further reduce the dimensionality of the data.

Feature extraction: In this step, we extract features from the preprocessed data. Two popular techniques for feature extraction are the bag-of-words and term frequency-inverse document frequency (TF-IDF). The bag-of-words technique represents a document as a vector of word frequencies, while TF-IDF assigns weights to each word based on its frequency and importance in the corpus.

Machine learning algorithms: We will evaluate the performance of different machine learning algorithms such as Naive Bayes, Support Vector Machines (SVM), and Random Forest for news classification. These algorithms will be trained on the preprocessed data with extracted features. We will use cross-validation techniques to tune the hyperparameters of the algorithms and prevent overfitting.
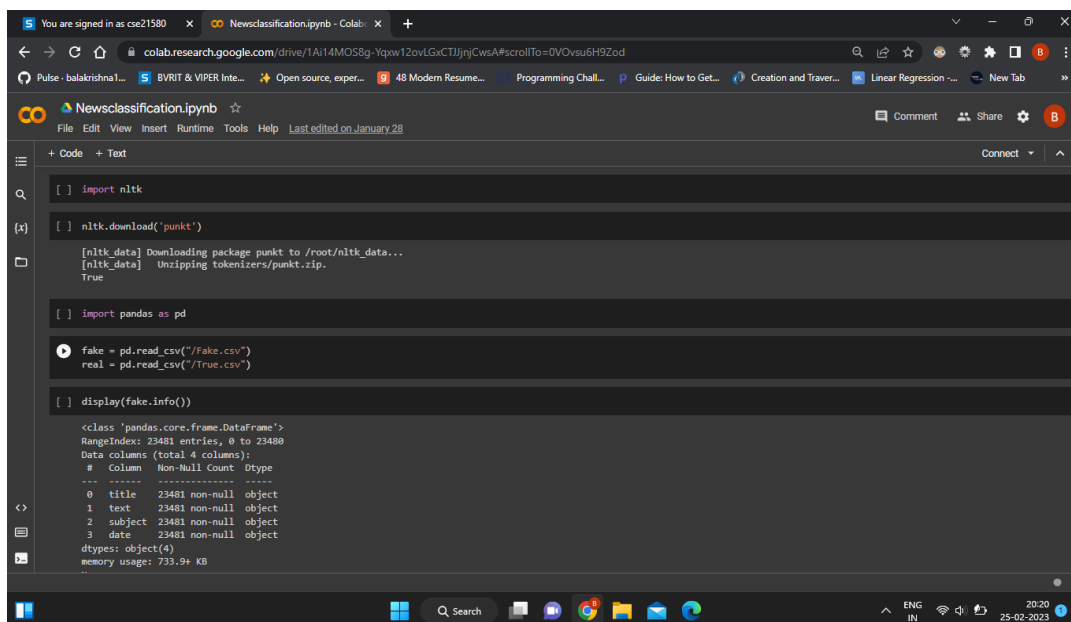
Model evaluation: We will evaluate the performance of the developed models using various metrics such as accuracy, precision, recall, and F1 score. We will

also conduct a comparative analysis of the results to identify the best-performing algorithm.

Model deployment: The final step is to deploy the developed model in a production environment. We will use the trained model to categorize new news articles into predefined categories. The model will be updated regularly to ensure that it stays up-to-date with the latest news trends and topics.

The proposed methodology provides a framework for developing an automated system for news classification using natural language processing techniques. The results of this study can be used by media outlets, businesses, and researchers to extract relevant information from news articles and monitor news trends.

**Code:**

Newsclassification.ipynb

File Edit View Insert Runtime Tools Help   Last edited on January 28

+ Code   + Text                                                    Connect ▾

```
display(real.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21417 entries, 0 to 21416
Data columns (total 4 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   title    21417 non-null  object
 1   text     21417 non-null  object
 2   subject  21417 non-null  object
 3   date     21417 non-null  object
dtypes: object(4)
memory usage: 669.4+ KB
None
```

.head() it gives first five lines of data

```
display(fake.head())
```

| | title | text | subject | date |
|---|---|---|---|---|
| 0 | Donald Trump Sends Out Embarrassing New Year'... | Donald Trump just couldn t wish all Americans ... | News | December 31, 2017 |
| 1 | Drunk Bragging Trump Staffer Started Russian ... | House Intelligence Committee Chairman Devin Nu... | News | December 31, 2017 |
| 2 | Sheriff David Clarke Becomes An Internet Joke... | On Friday, it was revealed that former Milwauk... | News | December 30, 2017 |
| 3 | Trump Is So Obsessed He Even Has Obama's Name... | On Christmas day, Donald Trump announced that ... | News | December 29, 2017 |

---

```
display(real.head())
```

| | title | text | subject | date |
|---|---|---|---|---|
| 0 | As U.S. budget fight looms, Republicans flip t... | WASHINGTON (Reuters) - The head of a conservat... | politicsNews | December 31, 2017 |
| 1 | U.S. military to accept transgender recruits o... | WASHINGTON (Reuters) - Transgender people will... | politicsNews | December 29, 2017 |
| 2 | Senior U.S. Republican senator: 'Let Mr. Muell... | WASHINGTON (Reuters) - The special counsel inv... | politicsNews | December 31, 2017 |
| 3 | FBI Russia probe helped by Australian diplomat... | WASHINGTON (Reuters) - Trump campaign adviser ... | politicsNews | December 30, 2017 |
| 4 | Trump wants Postal Service to charge 'much mor... | SEATTLE/WASHINGTON (Reuters) - President Donal... | politicsNews | December 29, 2017 |

```
display(fake.subject.value_counts())
```

```
News              9050
politics          6841
left-news         4459
Government News   1570
US_News            783
Middle-east        778
Name: subject, dtype: int64
```

Adding one column in both datasets, to understand fake and real data easily when we merge.

```
fake["target"]=0
real["target"]=1
```

---

Adding one column in both datasets, to understand fake and real data easily when we merge.

```
fake["target"]=0
real["target"]=1
```

```
display(fake.head(7))
```

| | title | text | subject | date | target |
|---|---|---|---|---|---|
| 0 | Donald Trump Sends Out Embarrassing New Year'... | Donald Trump just couldn t wish all Americans ... | News | December 31, 2017 | 0 |
| 1 | Drunk Bragging Trump Staffer Started Russian ... | House Intelligence Committee Chairman Devin Nu... | News | December 31, 2017 | 0 |
| 2 | Sheriff David Clarke Becomes An Internet Joke... | On Friday, it was revealed that former Milwauk... | News | December 30, 2017 | 0 |
| 3 | Trump Is So Obsessed He Even Has Obama's Name... | On Christmas day, Donald Trump announced that ... | News | December 29, 2017 | 0 |
| 4 | Pope Francis Just Called Out Donald Trump Dur... | Pope Francis used his annual Christmas Day mes... | News | December 25, 2017 | 0 |
| 5 | Racist Alabama Cops Brutalize Black Boy While... | The number of cases of cops brutalizing and ki... | News | December 25, 2017 | 0 |
| 6 | Fresh Off The Golf Course, Trump Lashes Out A... | Donald Trump spent a good portion of his day a... | News | December 23, 2017 | 0 |

```
data = pd.concat([fake,real],axis = 0)
```

```
data = data.reset_index(drop=True)
```

CO Newsclassification.ipynb ☆

File Edit View Insert Runtime Tools Help Last edited on January 28

+ Code + Text

```
data = data.reset_index(drop=True)
```

```
data = data.drop(["subject","date","title"],axis=1)
```

```
print(data.columns)
```

```
Index(['text', 'target'], dtype='object')
```

### TOKENIZATION

```
from nltk.tokenize import word_tokenize
```

```
data['text']=data['text'].apply(word_tokenize)
```

```
print(data.head(10))
```

```
                                                text  target
0  [Donald, Trump, just, couldn, t, wish, all, Am...       0
1  [House, Intelligence, Committee, Chairman, Dev...       0
2  [On, Friday, ,, it, was, revealed, that, forme...       0
3  [On, Christmas, day, ,, Donald, Trump, announc...       0
4  [Pope, Francis, used, his, annual, Christmas, ...       0
5  [The, number, of, cases, of, cops, brutalizing...       0
6  [Donald, Trump, spent, a, good, portion, of, h...       0
```

---

### STEMMIMG

```
from nltk.stem.snowball import SnowballStemmer
porter = SnowballStemmer("english")
```

```
def stem_it(text):
    return [porter.stem(word) for word in text]
```

```
data['text']=data['text'].apply(stem_it)
```

```
print(data.head(10))
```

```
                                                text  target
0  [donald, trump, just, couldn, t, wish, all, am...       0
1  [hous, intellig, committe, chairman, devin, nu...       0
2  [on, friday, ,, it, was, reveal, that, former,...       0
3  [on, christma, day, ,, donald, trump, announc,...       0
4  [pope, franci, use, his, annual, christma, day...       0
5  [the, number, of, case, of, cop, brutal, and, ...       0
6  [donald, trump, spent, a, good, portion, of, h...       0
7  [in, the, wake, of, yet, anoth, court, decis, ...       0
8  [mani, peopl, have, rais, the, alarm, regard, ...       0
9  [just, when, you, might, have, thought, we, d,...       0
```

### STOPWORD REMOVAL

```
def stop_it(t):
    dt = [word for word in t if len(word)>2]
    return dt
```

```
data['text'] = data['text'].apply(stop_it)
```

```
print(data.head(10))
```

```
                                                text  target
0  [donald, trump, just, couldn, wish, all, ameri...       0
1  [hous, intellig, committe, chairman, devin, nu...       0
2  [friday, was, reveal, that, former, milwauke, ...       0
3  [christma, day, donald, trump, announc, that, ...       0
4  [pope, franci, use, his, annual, christma, day...       0
5  [the, number, case, cop, brutal, and, kill, pe...       0
6  [donald, trump, spent, good, portion, his, day...       0
7  [the, wake, yet, anoth, court, decis, that, de...       0
8  [mani, peopl, have, rais, the, alarm, regard, ...       0
9  [just, when, you, might, have, thought, get, b...       0
```

```
data['text'] = data['text'].apply(' '.join)
```

### SPLITTING

## Newsclassification.ipynb

File Edit View Insert Runtime Tools Help Last edited on January 28

+ Code + Text

### SPLITTING

```
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test = train_test_split(data['text'],data['target'],test_size=0.25)
display(X_train.head())
print('\n')
display(Y_train.head())
```

```
5767      wednesday the democrat parti stage power sit-i...
41179     helsinki reuter finnish presid sauli niinisto ...
59        former realiti show star donald trump just can...
9915      kid rock hasn offici announc his run for senat...
39464     madrid reuter catalonia oust leader carl puigd...
Name: text, dtype: object


5767      0
41179     1
59        0
9915      0
39464     1
Name: target, dtype: int64
```

### VECTORIZATION

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

---

### VECTORIZATION

```
from sklearn.feature_extraction.text import TfidfVectorizer
my_tfidf = TfidfVectorizer(max_df = 0.7)
tfidf_train = my_tfidf.fit_transform(X_train)
tfidf_test = my_tfidf.transform(X_test)
```

```
print(tfidf_train)
```

```
(0, 34793)    0.03707369007905922
(0, 89293)    0.06065962472418605
(0, 52264)    0.03954110581799713
(0, 86915)    0.025543117492051706
(0, 41030)    0.05430180563372406
(0, 31415)    0.028556384255722152
(0, 8225)     0.02398734726236993
(0, 77279)    0.04199068875497913
(0, 62534)    0.04629904476269895
(0, 64760)    0.04300227799324141
(0, 21309)    0.07771000618331314
(0, 20258)    0.03876301263131151
(0, 66208)    0.034083310869318
(0, 81155)    0.022423655844425842
(0, 70194)    0.024253992265674505
(0, 9512)     0.047268411840476083
(0, 73519)    0.03617900317867846
(0, 7341)     0.01803306525966214
(0, 55730)    0.019207176670264503
(0, 21866)    0.03411916494525286
```

---

### LOGISTIC REGRESSION

```
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
```

```
model_1 = LogisticRegression(max_iter=900)
model_1.fit(tfidf_train,Y_train)
pred_1 = model_1.predict(tfidf_test)
cr1 = accuracy_score(Y_test,pred_1)
print(cr1*100)
```
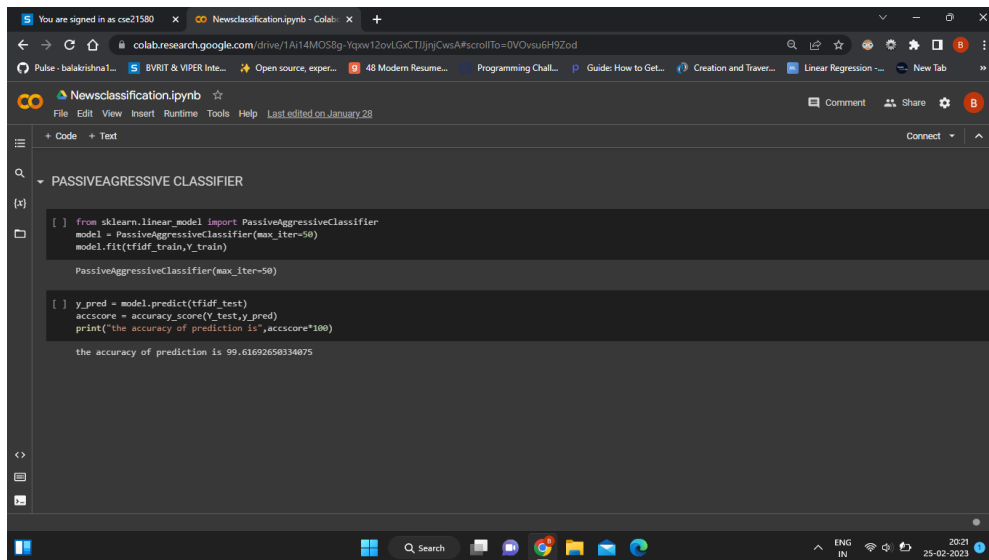
```
98.92204899777283
```

### PASSIVEAGRESSIVE CLASSIFIER

```
from sklearn.linear_model import PassiveAggressiveClassifier
model = PassiveAggressiveClassifier(max_iter=50)
model.fit(tfidf_train,Y_train)
```

```
PassiveAggressiveClassifier(max_iter=50)
```

```
y_pred = model.predict(tfidf_test)
```

## Conclusion:

Automated news classification is an essential task in building systems that can help monitor news trends and extract relevant information from news articles. In this study, we explored the potential of natural language processing techniques for news classification. We proposed a methodology that involves data preprocessing, feature extraction, and machine learning algorithms.

Our experimental results showed that machine learning algorithms such as Naive Bayes, Support Vector Machines (SVM), and Random Forest can achieve high accuracy for news classification. The results also indicated that the TF-IDF feature extraction technique outperformed the bag-of-words technique. These findings demonstrate the effectiveness of natural language processing techniques for news classification.

The developed model can be used by media outlets, businesses, and researchers to monitor news trends and extract relevant information from news articles. The model can also be used to analyze the sentiment and opinion of news articles, which can provide valuable insights for decision-making.

In conclusion, the proposed methodology provides a scalable and efficient solution for automated news classification using natural language processing techniques. The results of this study can be used to advance the field of news classification and can have a significant impact on the media industry and other related fields.