

# Sprint 9

## Tasca M9 T01

```
In [1]: import pandas as pd
import numpy as np
import nltk
import matplotlib.pyplot as plt
import seaborn as sns
from nltk.corpus import stopwords

from nltk.tokenize import sent_tokenize, word_tokenize, WordPunctTokenizer, regexp_tokenize
from nltk.probability import FreqDist
from nltk.stem import nltk
from nltk.stem.wordnet import WordNetLemmatizer

from nltk.sentiment.vader import SentimentIntensityAnalyzer

import warnings
#warnings.filterwarnings('action=once')
warnings.filterwarnings('ignore')
#nltk.download('all')
```

## Example amazon reviews

```
In [2]: df = pd.read_csv('https://raw.githubusercontent.com/pycaret/pycaret/master/datasets/amazon_reviews_small.csv')
df

Out[2]:
```

	reviewText	Positive
0	This is a one of the best apps according to a b...	1
1	This is a pretty good version of the game for ...	1
2	this is a really cool game, there are a bunch ...	1
3	This is a silly game and can be frustrating, b...	1
4	This is a terrific game on any pad. Hrs of fun...	1
...	...	...
19995	this app is fricken stupid.it froze on the kin...	0
19996	Please add me!!!! I need neighbor! ginger101...	1
19997	love it! this game is awesome, wish it had m...	1
19998	I love love love this app on my side of fashio...	1
19999	This game is a rip off. Here is a list of thi...	0

20000 rows x 2 columns

```
In [3]: # create preprocess text function
def preprocess_text(text):
    # Tokenize the text
    tokens = word_tokenize(text.lower())
    # Remove stop words
    filtered_tokens = [token for token in tokens if token not in stopwords.words('eng...
    # Lemmatize the tokens
    lemmatizer = WordNetLemmatizer()
    lemmatized_tokens = [lemmatizer.lemmatize(token) for token in filtered_tokens]
    # Join the tokens back into a string
    processed_text = ' '.join(lemmatized_tokens)
    return processed_text

# apply the function df
df['reviewText'] = df['reviewText'].apply(preprocess_text)
df

Out[3]:
```

	reviewText	Positive
0	one best apps according bunch people agree bomb...	1
1	pretty good version game free. lot different ...	1
2	really cool game. bunch level find golden egg...	1
3	silly game frustrating, lot fun definitely re...	1
4	terrific game pad. hr fun. grandkids love ...	1
...	...	...
19995	app fricken stupid.it froze kindle wont allow ...	0
19996	please add!!!!!! need neighbor! ginger101...	1
19997	love! game. awesome. wish free stuff house ...	1
19998	love love love app side fashion story fight wo...	1
19999	game rip. list thing make better & bull; fir...	0

20000 rows x 2 columns

```
In [4]: # initialize NLTK sentiment analyzer
analyzer = SentimentIntensityAnalyzer()

# create get_sentiment function
def get_sentiment(text):
    scores = analyzer.polarity_scores(text)
    sentiment = 1 if scores['pos'] > 0 else 0
    return sentiment

# apply get_sentiment function
df['sentiment'] = df['reviewText'].apply(get_sentiment)
df

Out[4]:
```

	reviewText	Positive	sentiment
0	one best apps according bunch people agree bomb...	1	1
1	pretty good version game free. lot different ...	1	1
2	really cool game. bunch level find golden egg...	1	1
3	silly game frustrating, lot fun definitely re...	1	1
4	terrific game pad. hr fun. grandkids love ...	1	1
...	...	...	...
19995	app fricken stupid.it froze kindle wont allow ...	0	0
19996	please add!!!!!! need neighbor! ginger101...	1	1
19997	love! game. awesome. wish free stuff house ...	1	1
19998	love love love app side fashion story fight wo...	1	1
19999	game rip. list thing make better & bull; fir...	0	1

20000 rows x 3 columns

```
In [5]: from sklearn.metrics import confusion_matrix
print(confusion_matrix(df['Positive'], df['sentiment']))

[[ 1131 14636]
 [ 576 3657]]

In [6]: from sklearn.metrics import classification_report
print(classification_report(df['Positive'], df['sentiment']))
```

	precision	recall	f1-score	support
0	0.66	0.24	0.35	4767
1	0.80	0.96	0.87	15233
accuracy	0.79	0.79	0.79	20000
macro avg	0.73	0.60	0.61	20000
weighted avg	0.77	0.79	0.75	20000

## Exercici 1

Agafa un text en anglès que vulguis, i calcula'n la freqüència de les paraules.

- He decidit agafar un text sobre els beneficis de la IA, generat per ChatGPT-3 i també unes reviews positives i negatives de Twitter per fer un anàlisis posterior.

```
In [7]: text_ID = """ Artificial intelligence (AI) has the potential to revolutionize the way we work. Firstly, AI has the potential to improve efficiency and productivity in numerous industries. Secondly, AI has the potential to create new industries and job opportunities. As AI advances, it will help solve some of the world's most pressing challenges. Finally, AI has the potential to improve safety and security in various contexts. For instance, self-driving cars equipped with AI can help prevent accidents caused by human error. AI-powered security systems can also help detect and prevent crime in public places, making cities safer for residents. In conclusion, AI has enormous potential to improve our lives in numerous ways. While there are certainly challenges associated with the use of AI, the benefits outweigh the risks, and it is important that we continue to invest in this technology to unlock its full potential.

Artificial intelligence (AI) has the potential to revolutionize the way we live and work, and many experts believe that it will be one of the most significant technological advancements of the 21st century. There are several reasons why AI is viewed as a positive force for humanity. Firstly, AI has the potential to improve efficiency and productivity in numerous industries. By automating tasks that were previously done manually, businesses can save time and resources, allowing them to focus on more important tasks. For example, in the healthcare industry, AI-powered systems can help doctors and nurses make more accurate diagnoses, leading to better patient outcomes. Secondly, AI has the potential to create new industries and job opportunities. As AI technology continues to advance, new products and services will emerge, creating jobs for developers, engineers, and other professionals. In addition, the use of AI can free up workers to focus on more creative and higher-level tasks that require human skills and judgment. Thirdly, AI has the potential to improve safety and security in various contexts. For instance, self-driving cars equipped with AI can help prevent accidents caused by human error. AI-powered security systems can also help detect and prevent crime in public places, making cities safer for residents. Finally, AI has the potential to help solve some of the world's most pressing challenges, such as climate change and disease outbreaks. By analyzing vast amounts of data and identifying patterns, AI can help researchers develop new treatments, predict the spread of diseases, and improve environmental sustainability. In conclusion, AI has enormous potential to improve our lives in numerous ways. While there are certainly challenges associated with the use of AI, the benefits outweigh the risks, and it is important that we continue to invest in this technology to unlock its full potential.

tokens = nltk.word_tokenize(text_ID.lower())
print('tokens:', len(tokens))

tagged = nltk.pos_tag(tokens)
print('tagged:', len(tagged))

fdist = FreqDist(tokens)
print('fdist:', len(fdist))

tokens = ['artificial', 'intelligence', 'the', 'ai', 'has', 'the', 'potential', 'to', 'revolutionize']

tagged_tokens = [('artificial', 'JJ'), ('intelligence', 'NN'), ('the', 'the'), ('ai', 'NN'), ('has', 'VBZ'), ('the', 'the'), ('potential', 'NN'), ('to', 'to'), ('to', 'to'), ('revolutionize', 'VB')]

Paraules més comuns al text:
('the', 24),
('and', 16),
('is', 15),
('a', 15),
('to', 12),
('of', 12),
('in', 9),
('potential', 7),
('that', 7),
('has', 6),
('can', 6),
('help', 5),
('that', 4),
('improve', 4),
('for', 4)
```



## Exercici 2

Treu les stopwords i realitza stemming al teu conjunt de dades.

```
In [10]: # Regex Tokens
stokr = regexp_tokenize(text_ID.lower(), "[\w']+")
print('tokens:', len(stokr))
display(stokr[15])

filtered_tokens = [token for token in stokr if token not in stopwords.words('english')]
print('filtered tokens:', len(filtered_tokens))

lemmatizer = WordNetLemmatizer()
lemmatized_tokens = [lemmatizer.lemmatize(token) for token in filtered_tokens]
display(lemmatized_tokens[15])

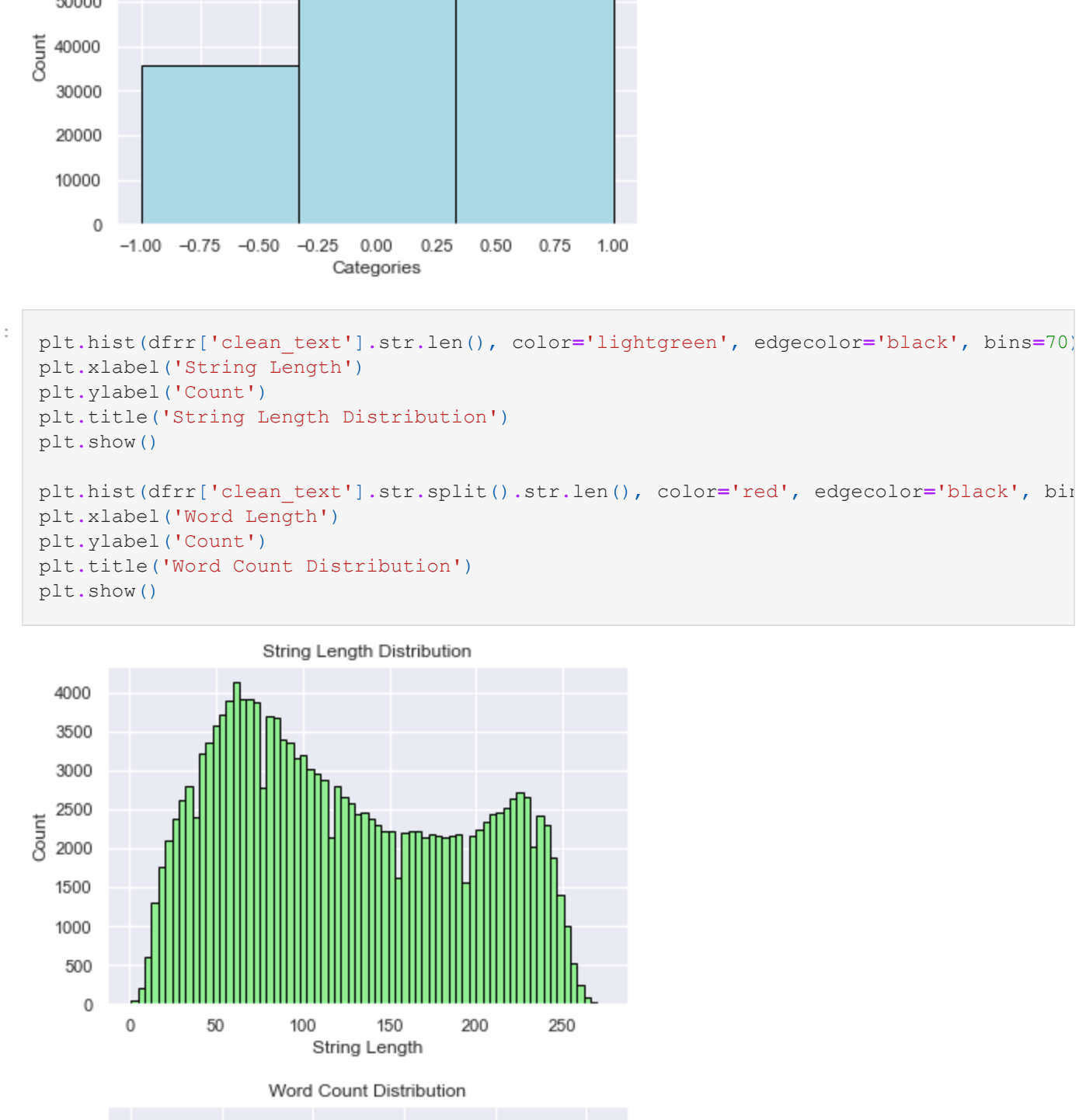
# Regex Tokens
['artificial',
'intelligence',
'ai',
'has',
'the',
'potential',
'to',
'revolutionize',
'the',
'way',
'we',
'live',
'and',
'work',
'and']

Without Stopwords:
['artificial',
'intelligence',
'ai',
'potential',
'revolutionize',
'way',
'live',
'work',
'many',
'experts',
'believe',
'one',
'significant',
'technological',
'advancements']

Lemmatized:
['artificial',
'intelligence',
'ai',
'potential',
'revolutionize',
'way',
'live',
'work',
'many',
'expert',
'believe',
'one',
'significant',
'technological',
'advancement']

In [11]: fdist = FreqDist(lemmatized_tokens)

plt.subplots(figsize=(20,10))
plt.title('Word frequency after stopwords and lemmatization', fontsize=20)
fdist.plot(30, cumulative=False)
plt.show()
```



## Exercici 3

Realitza sentiment analysis al teu conjunt de dades.

```
In [13]: # create preprocess text function
def preprocess_text(text):
    # Tokenize the text
    tokens = word_tokenize(text.lower())
    # Remove stop words
    filtered_tokens = [token for token in tokens if token not in stopwords.words('eng...
    # Lemmatize the tokens
    lemmatizer = WordNetLemmatizer()
    lemmatized_tokens = [lemmatizer.lemmatize(token) for token in filtered_tokens]
    # Join the tokens back into a string
    processed_text = ' '.join(lemmatized_tokens)
    return processed_text

# create get_sentiment function
def get_sentiment(text):
    scores = analyzer.polarity_scores(text)
    sentiment = 1 if scores['pos'] > scores['neg'] else -1
    return sentiment

In [14]: # initialize NLTK sentiment analyzer
analyzer = SentimentIntensityAnalyzer()

pt = ' '.join(lemmatized_tokens)
scores = analyzer.polarity_scores(pt)
print('scores:', scores)

Polarity scores: ('neg': 0.048, 'neu': 0.611, 'pos': 0.341, 'compound': 0.9956)

In [15]: # initialize NLTK sentiment analyzer
analyzer = SentimentIntensityAnalyzer()

pt=preprocess_text(text_ID)
print(pt)

scores = analyzer.polarity_scores(pt)
s=get_sentiment(pt)

if s==1:
    print('Positive text sentiment: 1', '\n33[0m')
else:
    print('Negative text sentiment: -1', '\n33[0m')
print('scores:', scores)
```

artificial intelligence (ai) potential revolutionize way live, many expert believe one significant technological advancement 21st century several reason ai viewed positive light firstly, ai potential improve efficiency productivity numerous industries, automating task previously done manually, business save time resource, allowing focus important task example, healthcare industry, ai-powered system help doctor nurse make accurate diagnosis, leading better patient outcome secondly, ai potential create new industry job opportunity ai technology continues advance, new product service emerge, creating job developer, engineer, professional addition, use ai free worker focus creative higher-level task require human skill judgment thirdly, ai potential improve safety security various context instance, self-driving car equipped ai help prevent accident caused human error ai-powered security system also help prevent crime public place, making city safer resident finally, ai potential help solve world's pressing challenge, climate change disease outbreak analyzing vast amount data identifying pattern, ai help researcher develop new treatment, predict spread disease, improve environmental sustainability conclusion, ai benefit potential improve life numerous ways certainly challenge associated use ai, benefit outweigh risk, important continue invest technology unlock full potential

Positive text sentiment: 1

Polarity scores: ('neg': 0.048, 'neu': 0.606, 'pos': 0.345, 'compound': 0.9956)

- Volem que entre fer el word\_tokenize o el regexp\_tokenize no hi ha una diferència massa significativa en aquest cas.

Reviews extretes de: [https://www.kaggle.com/datasets/cosmos98/twitter-and-reddit-sentimental-analysis-dataset?resource=download&select=Twitter\\_Data.csv](https://www.kaggle.com/datasets/cosmos98/twitter-and-reddit-sentimental-analysis-dataset?resource=download&select=Twitter_Data.csv)

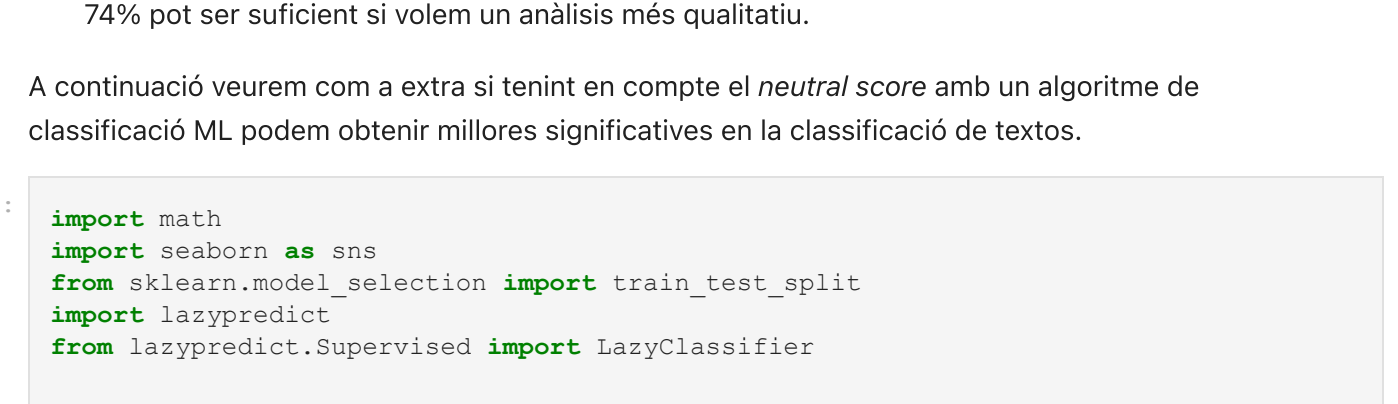
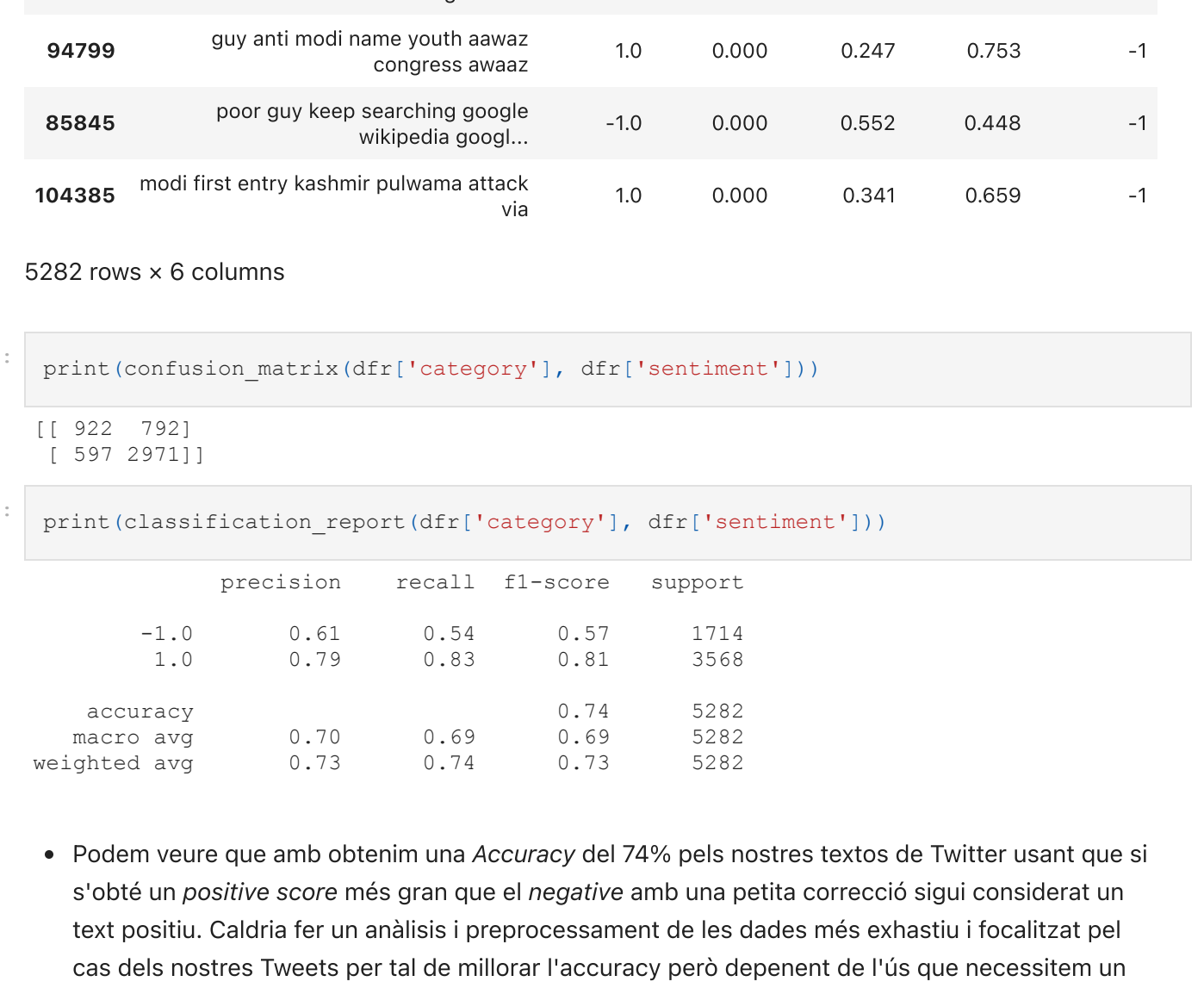
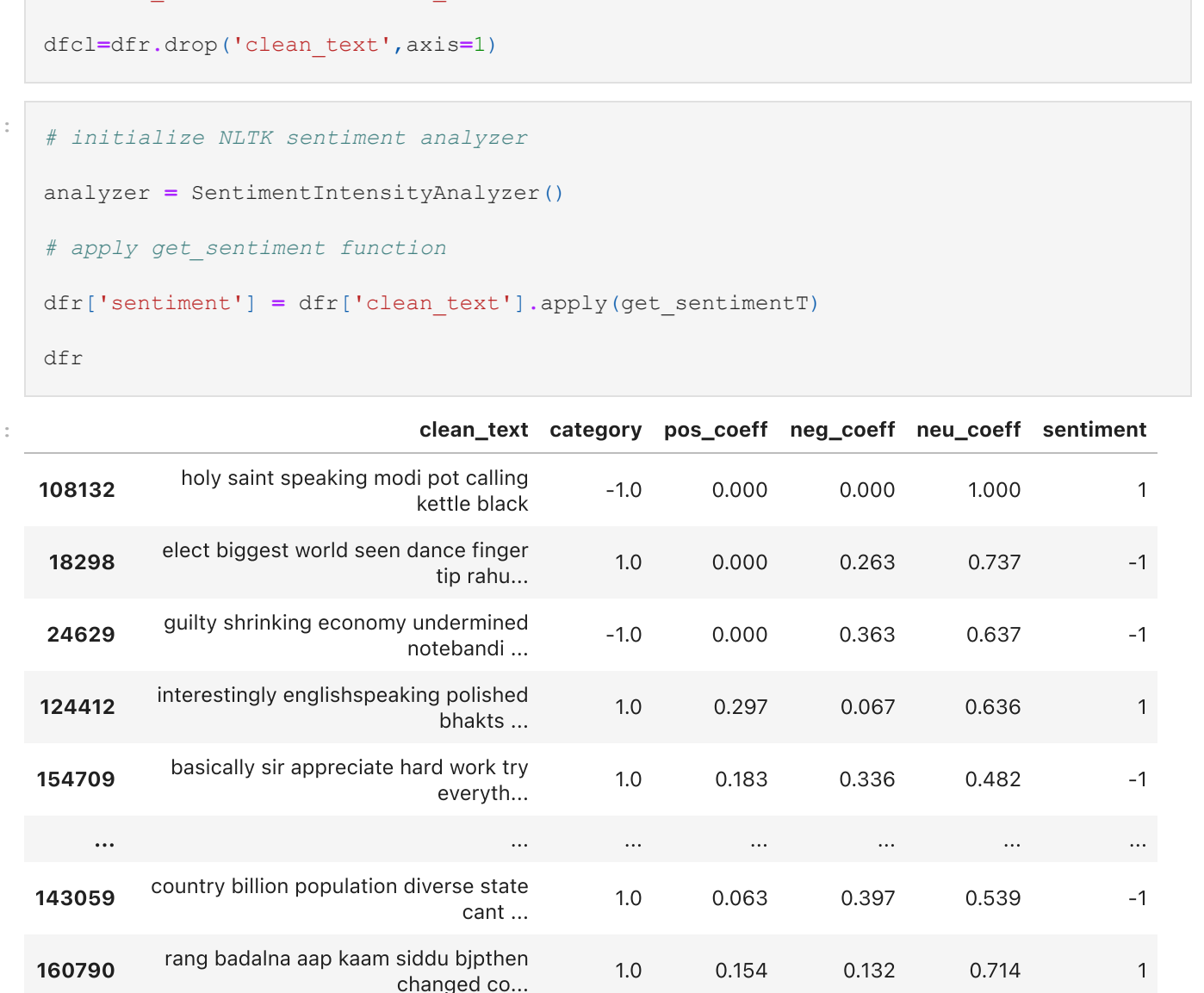
Content: These tweets and Comments Were Made on Narendra Modi and Other Leaders as well as Peoples Opinion Towards the Next Prime Minister of The Nation ( In Context with General Elections Held In India – 2019).

```
In [16]: dfrr=pd.read_csv('Twitter_Data.csv', delimiter=',')
dfrr=dfrr.dropna()
dfrr

Out[16]:
```

	clean_text	category
0	when modi promised "minimum government maximum...	-1.0
1	talk all the nonsense and continue all the dra...	0.0
2	what did just say vote for modi welcome bjp...	1.0
3	asking his supporters prefix chowkidar their n...	1.0
4	answer who among these the most powerful mod...	1.0
...	...	...
162975	why these 456 crores paid nearav modi not reco...	-1.0
162976	dear rrs terrorist payal gawar what about mod...	-1.0
162977	did you cover her interaction forum where she ...	0.0
162978	there big project came into india modi dream p...	0.0
162979	have you ever listen about india gukral where ...	1.0

162969 rows x 2 columns





```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

```
In [29]: data = dfc1
X = data.iloc[:,1:]
y = data['category']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.2, random_state=1)

# Normalize the features
scaler = StandardScaler()
X_trainS = scaler.fit_transform(X_train)
X_testS = scaler.transform(X_test)

clf = LazyClassifier(verbose=0, ignore_warnings=True, custom_metric=None)
models, predictions = clf.fit(X_train, X_test, y_train, y_test)
models
```

100% ██████████ 29/29 [00:05<00:00, 5.42it/s]

Out [29]:

	Accuracy	Balanced Accuracy	ROC AUC	F1 Score	Time Taken
Model					
NearestCentroid	0.70		0.71	0.71	0.02
BernoulliNB	0.75		0.69	0.69	0.74
GaussianNB	0.75		0.67	0.67	0.73
XGBClassifier	0.74		0.66	0.66	0.73
LogisticRegression	0.75		0.66	0.66	0.73
CalibratedClassifierCV	0.75		0.66	0.66	0.73
LinearDiscriminantAnalysis	0.75		0.66	0.66	0.73
SVC	0.75		0.66	0.66	0.73
LinearSVC	0.75		0.65	0.65	0.73
RidgeClassifier	0.75		0.65	0.65	0.73
RidgeClassifierCV	0.75		0.65	0.65	0.73
LGBMClassifier	0.74		0.65	0.65	0.72
AdaBoostClassifier	0.74		0.65	0.65	0.72
SGDClassifier	0.75		0.65	0.65	0.72
RandomForestClassifier	0.72		0.65	0.65	0.71
BaggingClassifier	0.72		0.65	0.65	0.71
QuadraticDiscriminantAnalysis	0.73		0.64	0.64	0.71
KNeighborsClassifier	0.72		0.64	0.64	0.71
ExtraTreesClassifier	0.70		0.63	0.63	0.69
DecisionTreeClassifier	0.69		0.62	0.62	0.68
ExtraTreeClassifier	0.68		0.62	0.62	0.68
PassiveAggressiveClassifier	0.73		0.60	0.60	0.69
NuSVC	0.62		0.56	0.56	0.62
LabelSpreading	0.68		0.50	0.50	0.56
LabelPropagation	0.68		0.50	0.50	0.56
DummyClassifier	0.68		0.50	0.50	0.56
Perceptron	0.58		0.46	0.46	0.55

```
In [30]: # Train the logistic regression model
lr2 = LogisticRegression()
lr2.fit(X_trainS, y_train)
lr_pred2 = lr2.predict(X_testS)
lr_acc2 = accuracy_score(y_test, lr_pred2)
print("Logistic regression accuracy:", lr_acc2)
```

Logistic regression accuracy: 0.7549668874172185

- Amb el model Logistic no sembla que poguem trobar un mètode de classificació significativament més precís a partir dels coeficients que obtenim del *Polarity score*.

```
In [ ] :
```