

```
1 !wget https://github.com/hjysam/text_classification/raw/main/data/reuters21578.tar
2 !tar -xvf reuters21578.tar

--2023-12-18 19:47:06-- https://github.com/hjysam/text\_classification/raw/main/data/reuters21578.tar
Resolving github.com (github.com)... 140.82.112.4
Connecting to github.com (github.com)|140.82.112.4|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/hjysam/text\_classification/main/data/reuters21578.tar [following]
--2023-12-18 19:47:06-- https://raw.githubusercontent.com/hjysam/text\_classification/main/data/reuters21578.tar
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.111.133, 185.199.109.133, 185.199.110.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.111.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 28016640 (27M) [application/octet-stream]
Saving to: ‘reuters21578.tar’

reuters21578.tar      100%[=====]  26.72M  ---KB/s    in 0.1s

2023-12-18 19:47:08 (217 MB/s) - ‘reuters21578.tar’ saved [28016640/28016640]

README.txt
all-exchanges-strings.lc.txt
all-orgs-strings.lc.txt
all-people-strings.lc.txt
all-places-strings.lc.txt
all-topics-strings.lc.txt
cat-descriptions_120396.txt
feldman-cia-worldfactbook-data.txt
lewis.dtd
reut2-000.sgm
reut2-001.sgm
reut2-002.sgm
reut2-003.sgm
reut2-004.sgm
reut2-005.sgm
reut2-006.sgm
reut2-007.sgm
reut2-008.sgm
reut2-009.sgm
reut2-010.sgm
```

```
reut2-011.sgm  
reut2-012.sgm  
reut2-013.sgm  
reut2-014.sgm  
reut2-015.sgm  
reut2-016.sgm  
reut2-017.sgm  
reut2-018.sgm  
reut2-019.sgm  
reut2-020.sgm  
reut2-021.sgm
```

```
1 !ls
```

all-exchanges-strings.lc.txt	lewis.dtd	reut2-005.sgm	reut2-012.sgm	reut2-019.sgm
all-orgs-strings.lc.txt	README.txt	reut2-006.sgm	reut2-013.sgm	reut2-020.sgm
all-people-strings.lc.txt	reut2-000.sgm	reut2-007.sgm	reut2-014.sgm	reut2-021.sgm
all-places-strings.lc.txt	reut2-001.sgm	reut2-008.sgm	reut2-015.sgm	reuters21578.tar
all-topics-strings.lc.txt	reut2-002.sgm	reut2-009.sgm	reut2-016.sgm	sample_data
cat-descriptions_120396.txt	reut2-003.sgm	reut2-010.sgm	reut2-017.sgm	
feldman-cia-worldfactbook-data.txt	reut2-004.sgm	reut2-011.sgm	reut2-018.sgm	

```
1 # Data manipulation and analysis
2 !pip install numpy
3 !pip install pandas
4 !pip install datasets
5
6 # Web scraping
7 !pip install beautifulsoup4
8
9 # Text processing and NLP
10 !pip install nltk
11 !pip install gensim
12
13 # Visualization
14 !pip install matplotlib
15 !pip install seaborn
16
17 # Machine learning and scikit-learn
18 !pip install scikit-learn
19 !pip install transformers
20
21 # Model persistence
22 !pip install torch
```



▼ Step 1: Python Module Library

```
1 # Standard libraries
2 import os
3 import re
4 from collections import Counter
5 from itertools import chain
6
7 # Data manipulation and analysis
8 import numpy as np
9 import pandas as pd
10
11 # Web scraping
12 from bs4 import BeautifulSoup
13
14 # Text processing and NLP
15 import nltk
16 from nltk.corpus import stopwords
17 from nltk.stem import PorterStemmer
18 from nltk.tokenize import word_tokenize
19 from gensim.models import Word2Vec
20
21 # Visualization
22 import matplotlib.pyplot as plt
23 import seaborn as sns
24
25 # Machine learning and scikit-learn
26 import sklearn
27 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
28 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer, TfidfVectorizer
29 from sklearn.pipeline import Pipeline, FeatureUnion
30 from sklearn.metrics import classification_report, confusion_matrix, ConfusionMatrixDisplay, multilabel_confusion_matrix
31 from sklearn.preprocessing import MultiLabelBinarizer
32 from sklearn.naive_bayes import MultinomialNB
33 from sklearn.multiclass import OneVsRestClassifier
34 from sklearn.linear_model import LogisticRegression, SGDClassifier
35 from sklearn.ensemble import RandomForestClassifier
36 from sklearn.svm import SVC, LinearSVC
37 from sklearn.neural_network import MLPClassifier
```

```
38 from sklearn.base import BaseEstimator, TransformerMixin
39 from transformers import (
40     AutoModelForSequenceClassification,
41     AutoTokenizer,
42     EvalPrediction,
43     Trainer,
44     TrainingArguments,
45 )
46
47 # Model persistence
48 import pickle
49 import torch
50
51 nltk.download('punkt')
52 nltk.download('stopwords')
53 p = print
54 d = display
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

✓ Step 2: Data Extraction

Downloaded From: <https://archive.ics.uci.edu/dataset/137/reuters+21578+text+categorization+collection>

```
1 def minor_preprocess(file):
2
3     with open(file, 'rb') as f:
4         lines = f.readlines()
5         utf8_safe_lines = [line.decode('utf-8', 'ignore') for line in lines]
6         xml_safe_lines = [re.sub(r'\d*;', '', line) for line in utf8_safe_lines] # Get rid of problematic strings
7         no_newlines = [line.replace('\n', ' ') for line in xml_safe_lines]
8     f.close()
9
10    return ''.join(no_newlines)
11
12 def compile_data(datapath='/content/'):
13     dataset = []
14     for file in os.listdir(datapath):
15         if file.endswith('.sgm'):
16             preprocessed_data = minor_preprocess(datapath + '/' + file)
17             records = [record + '</REUTERS>' for record in preprocessed_data.split('</REUTERS>') if record] # Retain all original form
18
19         dataset.extend(records)
20
21     return dataset
22
23
24 def compile_dictionary(data):
25     data_dict = {
26         'new_id': [],
27         'people': [],
28         'places': [],
29         'orgs': [],
30         'exchanges': [],
31         'title':[],
32         'body': [],
33         'r_topics': [],
34         'lewissplit': [],
35         'cgisplit': [],
36         'topic': []
37     }
```



```
38
39     soup = BeautifulSoup(data, 'xml')
40
41     reuters_tag = soup.find('REUTERS')
42     if reuters_tag:
43         r_topics_value = reuters_tag.get('TOPICS')
44         if r_topics_value:
45             data_dict['r_topics'] = r_topics_value
46
47         new_id_value = reuters_tag.get('NEWID')
48         if new_id_value:
49             data_dict['new_id'] = new_id_value
50
51         lewissplit_value = reuters_tag.get('LEWISSPLIT')
52         if lewissplit_value:
53             data_dict['lewissplit'] = lewissplit_value
54
55         cgisplit_value = reuters_tag.get('CGISPLIT')
56         if cgisplit_value:
57             data_dict['cgisplit'] = cgisplit_value
58
59     people_tag = soup.find('PEOPLE')
60     if people_tag and people_tag.contents:
61         cleaned_people = [tag.text for tag in people_tag.find_all('D')]
62         data_dict['people'] = cleaned_people
63
64     places_tag = soup.find('PLACES')
65     if places_tag and places_tag.contents:
66         cleaned_places = [tag.text for tag in places_tag.find_all('D')]
67         data_dict['places'] = cleaned_places
68
69     orgs_tag = soup.find('ORGS')
70     if orgs_tag and orgs_tag.contents:
71         cleaned_orgs = [tag.text for tag in orgs_tag.find_all('D')]
72         data_dict['orgs'] = cleaned_orgs
73
74     exchanges_tag = soup.find('EXCHANGES')
```

```
75     if exchanges_tag and exchanges_tag.contents:
76         cleaned_exc = [tag.text for tag in exchanges_tag.find_all('D')]
77         data_dict['exchanges'] = cleaned_exc
78
79     body_tag = soup.find('BODY')
80     if body_tag and body_tag.contents:
81         data_dict['body'] = body_tag.contents[0]
82
83     title_tag = soup.find('TITLE')
84     if title_tag and title_tag.contents:
85         data_dict['title'] = title_tag.contents[0]
86
87     topics_tag = soup.find('TOPICS')
88     if topics_tag and topics_tag.contents:
89         cleaned_topics = [tag.text for tag in topics_tag.find_all('D')]
90         data_dict['topic'] = cleaned_topics
91
92     return data_dict
93
```

```
1 # Convert the list into a dictionary with fields of interest
2 reuters_dicts = [compile_dictionary(data) for data in compile_data()]
3
```

```
1 df = pd.DataFrame(reuters_dicts)
2 d(df.head())
3 df.to_csv('raw_data.csv')
```

	new_id	people	places	orgs	exchanges	title	body	r_topics	lewissplit	cgisplit	topic	grid icon	bar chart icon
0	2001	[]	[uk]	[]	[]	JAGUAR SEES STRONG GROWTH IN NEW MODEL SALES	Jaguar Plc > is about to sell its new XJ-6 mod...	YES	TRAIN	TRAINING-SET	[earn]		
1	2002	[]	[]	[]	[]	OCCIDENTAL PETROLEUM COMMON STOCK OFFERING RAI...	[]	NO	TRAIN	TRAINING-SET	[]		
2	2003	[]	[usa, iraq]	[]	[]	CCC ACCEPTS BONUS BID ON WHEAT FLOUR TO IRAQ	The Commodity Credit Corporation, CCC, has acc...	YES	TRAIN	TRAINING-SET	[grain, wheat]		
3	2004	[]	[]	[]	[]	DIAMOND SHAMROCK RAISES CRUDE POSTED PRICES ON...	[]	YES	TRAIN	TRAINING-SET	[crude]		
4	2005	[]	[usa]	[]	[]	NORD RESOURCES CORP > 4TH QTR NET	Shr 19 cts vs 13 cts Net 2,656,000 vs 1,71...	YES	TRAIN	TRAINING-SET	[earn]		

Dataset Overview

```
1 p('Statistical Table')
2 d(df.describe())
3
4 p('\nData Type for each Column')
5 d(df.info())
6
7 # Check for empty lists in each column
8 p('\nNumber of empty list in each column')
9 d(df.applymap(lambda x: isinstance(x, list) and not x.sum()))
10
```

Statistical Table

	new_id	people	places	orgs	exchanges	title	body	r_topics	lewissplit	cgisplit	topic	
count	21600	21600	21600	21600	21600	21600	21600	21600	21600	21600	21600	
unique	21579	218	1097	68	58	20006	18780	4	4	3	656	
top	[]	[]	[usa]	[]	[]	[]	[]	YES	TRAIN	TRAINING-SET	[]	
freq	22	20444	10878	20719	21118	759	2557	13476	14668	20856	10233	

Data Type for each Column

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21600 entries, 0 to 21599
Data columns (total 11 columns):
 #   Column      Non-Null Count  Dtype  
 ---  --  
 0   new_id      21600 non-null   object 
 1   people      21600 non-null   object 
 2   places      21600 non-null   object 
 3   orgs        21600 non-null   object 
 4   exchanges   21600 non-null   object 
 5   title       21600 non-null   object 
 6   body        21600 non-null   object 
 7   r_topics    21600 non-null   int64  
 8   lewissplit  21600 non-null   object 
 9   cgisplit    21600 non-null   object 
 10  topic       21600 non-null   object 
```

Since there is a significant number of empty lists for people, organizations, and exchanges, we have decided to drop these columns.

Additionally, places are not expected to contribute significantly to the topic classification task, therefore, this column will also be removed.

As a result:

- the selected features are 'title' and 'body'
- the label is 'topic'.

```
dtypes: object(11)
```

Next, check if the data size is extracted correctly based on different split methods.

```
1 p('The Modified Apte ("ModApte") Split')
2 p('train size: ' + str(len(df[(df.lewissplit=='TRAIN') & (df.r_topics=='YES')])))
3 p('test size: ' + str(len(df[(df.lewissplit=='TEST') & (df.r_topics=='YES')])))
4 p('not-used size: ' + str(len(df[df.lewissplit=='NOT-USED'])))
5 p()
6
7 p('The Modified Hayes ("ModHayes") Split')
8 p('train size: ' + str(len(df[df.cgisplit=='TRAINING-SET'])))
9 p('test size: ' + str(len(df[df.cgisplit=='PUBLISHED-TESTSET'])))
10 p('not-used size: ' + str(len(df[df.cgisplit=='NOT-USED'])))
```

```
The Modified Apte ("ModApte") Split
train size: 9603
test size: 3299
not-used size: 722
```

```
The Modified Hayes ("ModHayes") Split
train size: 20856
test size: 722
not-used size: 0
```

The size mentioned above aligns with the information provided in the dataset readme file, confirming that there are no mistakes in the data extraction. Please refer to:

- Line 459: VIII.B. The Modified Apte ('ModApte') Split
- Line 510: VIII.C. The Modified Hayes ('ModHayes') Split

For simplicity, we will choose ModApte over ModHayes, as it offers a more balanced distribution of training (9603) and testing (3299) data sizes, equivalent to a ratio of 65% versus 35%.

▼ Step 3: Data Preprocessing

Combine title and body together as text column

```
1 p("Combine title and body together")
2 df['text'] = df['title'].apply(lambda x: ''.join(x)) + ' ' + df['body'].apply(lambda x: ''.join(x))
3 df.insert(df.columns.get_loc('body') + 1, 'text', df.pop('text'))
4 df = df.drop(['body', 'title'], axis=1)
```

Combine title and body together

Remove special characters, newlines and number

```
1 df['text'] = df['text'].replace('[^a-zA-Z0-9\s]', '', regex=True)      # Special Characters
2 df['text'] = df['text'].replace('\n', ' ', regex=True)                  # Remove newlines
3 df['text'] = df['text'].replace('\d', ' ', regex=True)                  # Remove numbers
4
```

Removing stopwords through tokenization

```
1 # Define function to remove stopwords
2 def remove_stopwords(text):
3     stop_words = set(stopwords.words('english'))
4     words = word_tokenize(text)
5     filtered_words = [word for word in words if word.lower() not in stop_words]
6     return ' '.join(filtered_words)
7
8 # Apply the function to the DataFrame
9 df['text'] = df['text'].apply(remove_stopwords)
10
```

Stemming through tokenizing

```
1 # Initialize Porter Stemmer
2 ps = PorterStemmer()
3
4 # Define function to apply stemming
5 def apply_stemming(text):
6     words = word_tokenize(text)
7     stemmed_words = [ps.stem(word) for word in words]
8     return ' '.join(stemmed_words)
9
10 # Apply the function to the DataFrame
11 df['text'] = df['text'].apply(apply_stemming)
12
```

Result of the text columns after preprocessing

```
1 df.text.head()

0    jaguar see strong growth new model sale jaguar...
1    occident petroleum common stock offer rais mln...
2      ccc accept bonu bid wheat flour iraq commod cr...
3    diamond shamrock rais crude post price one dlr...
4    nord resourc corp th qtr net shr ct vs ct net ...
Name: text, dtype: object
```

Recreate a new data dictionary based on ModApte split; for more information, refer to step 2 for the reason.

- Train
- Test
- Unused

```

1 #Create a new function to cater for ModApte split for data dict
2 def create_data_dict(df, lewissplit, r_topics):
3     data_dict = {
4         'text': df[(df.lewissplit == lewissplit) & (df.r_topics == r_topics)].text.to_list(),
5         'topic': df[(df.lewissplit == lewissplit) & (df.r_topics == r_topics)].topic.to_list(),
6     }
7     return data_dict
8
9 # 'test' dict
10 test_dict = create_data_dict(df, 'TEST', 'YES')
11
12 # 'train' dict
13 train_dict = create_data_dict(df, 'TRAIN', 'YES')
14
15 # 'unused' dict (Will not be ignored)
16 unused_dict = create_data_dict(df, 'NOT-USED', '')

```

▼ Step 4: Data Exploration / Cleaning

Word Frequency Analysis for Text

- Identify the most frequent words.
- Identify the words that appear for one time.
- Count the number of unique words

```

1 from nltk.tokenize import word_tokenize
2 from collections import Counter
3
4 # List of sentences
5 sentences = train_dict['text']
6
7 # Tokenize the text
8 all_tokens = [word_tokenize(sentence) for sentence in sentences]
^

```

```
9
10 # Flatten the list of tokens
11 all_tokens_flat = [token for sublist in all_tokens for token in sublist]
12
13 # Create a Counter to count word frequencies
14 word_freq = Counter(all_tokens_flat)
15
16 # Identify the most frequent words (top N)
17 most_frequent_words = word_freq.most_common(10)
18
19 # Find words that appear only once
20 words_appear_one_time = [word for word, freq in word_freq.items() if freq == 1]
21
22 # Count the number of unique words
23 num_unique_words = len(word_freq)
24
25 # Display the results
26 p("Most Frequent Words:")
27 p(most_frequent_words)
28 p("\nNumber of Words Appear for One Time:")
29 p(len(words_appear_one_time))
30 p('\nExample of Words Appear for One Time')
31 p(words_appear_one_time)
32 p("\nNumber of Unique Words:", num_unique_words)
33
```

Most Frequent Words:

```
[('said', 23727), ('mln', 15324), ('dlr', 12669), ('pct', 9805), ('reuter', 9193), ('vs', 9187), ('year', 7138), ('bank', 6485)]
```

Number of Words Appear for One Time:

```
7469
```

Example of Words Appear for One Time

```
['init', 'barcelo', 'gloeilampenfabriek', 'kph', 'anatolia', 'ataturk', 'seeker', 'tb', 'withdrfaw', 'taiwain', 'resourec', 'in
```

Number of Unique Words: 20133



- As stopwords have applied on previous step, so there isn't any more words like the, will and is.
- There are 7231 words that appear for one time as compared to a total of 19737 unique words

Label Analysis

- Topic is considered as the label for this text classification

```
1 # Topic list
2 my_list = train_dict['topic']
3
4 p('Average token for each element')
5 average_token_per_element = sum([len(element) for element in my_list])/len(my_list)
6 p(average_token_per_element)
7
8 # Flatten the list of lists
9 flattened_list = [item for sublist in my_list for item in sublist]
10
11 # Use Counter to count occurrences
12 element_counts = Counter(flattened_list)
13 p('\nNumber of unique topic')
14 p(len(element_counts))
15
16 p('\nDensity')
17 p(average_token_per_element/len(element_counts))
18
19 # Sort the items by frequency in descending order
20 sorted_counts = sorted(element_counts.items(), key=lambda x: x[1], reverse=True)
21
22 p('\nTopic that Occurence for only one time')
23 filtered_counts = [item for item in sorted_counts if item[1] == 1]
24 p(len([item[0] for item in filtered_counts ]))
25
26 p("\ntopic occurences for one time")
27 # Filter tuples where the value is equal to 1
28 occurrence_only_one = [element[0] for element in filtered_counts]
29 p(occurrence_only_one)
30
31 # Extract & plot top 10
32 top_10 = sorted_counts[:10]
33 plt.figure(figsize=(10, 5))
34 bars = plt.barh(range(len(top_10)), [count for _, count in top_10], align='center')
35
36 # Add labels to the bars
37 for bar, (item, count) in zip(bars, top_10):
```

```
38     plt.text(bar.get_width(), bar.get_y() + bar.get_height() / 2, f'{count}', ha='left', va='center')
39
40 plt.yticks(range(len(top_10)), [item for item, _ in top_10])
41 plt.title('Top 10 Topic Occurrences')
42 plt.xlabel('Frequency')
43 plt.ylabel('Element')
44 plt.show()
45
```

Average token for each element

1.0047901697386235

Number of unique topic

115

Density

0.008737305823814117

Check if we can remove some of the one-time occurance by comparing to test set

20

```
1 # Initialize a list to store samples with one time occurance labels
2 topic_occurrence_one_test = []
3
4 for sample in test_dict['topic']:
5     if any(label in occurrence_only_one for label in sample):
6         topic_occurrence_one_test.append(sample)
7
8 p('Number of topic occurence one time in test')
9 p(len(topic_occurrence_one_test), 'out of', len(test_dict['text']))
10 p('\nShare in percentage')
11 p(str(round(len(topic_occurrence_one_test)/len(test_dict['text'])*100,2)) + '%')
```

Number of topic occurence one time in test

7 out of 3299

Share in percentage

0.21%



Safe to drop these one-time occurrence for train and test to reduce the label size.



```

1 def remove_unwanted_topic(dict):
2     df1 = pd.DataFrame(dict)
3     df1 = df1[~df1['topic'].apply(lambda x: any(topic in x for topic in occurrence_only_one))]
4     return df1.to_dict(orient='list')
5
6 train_dict = remove_unwanted_topic(train_dict)
7 test_dict = remove_unwanted_topic(test_dict)
8
9 p('new number of train data', len(train_dict['text']) )
10 p('new number of test data',len(test_dict['text'])) )
11

new number of train data 9588
new number of test data 3292

```

▼ Step 5: Prediction based on Classifier Model

This code prepares training data by extracting text (x_train) and corresponding topics (y_train). It then transforms the text into numerical vectors using TF-IDF, and the topics into a binary matrix, facilitating the training of machine learning models.

```

1 # Custom transformer to convert text to Word2Vec embeddings
2 class Word2VecTransformer(BaseEstimator, TransformerMixin):
3     def __init__(self, model):
4         self.model = model
5
6     def fit(self, X, y=None):
7         return self
8
9     def transform(self, X):
10        return [self.model.wv[x] if x in self.model.wv else [0] * self.model.vector_size for x in X]

```

```

1 x_train, y_train = train_dict['text'], train_dict['topic']
2 x_test, y_test = test_dict['text'], test_dict['topic']
3
4 # TF-IDF Vectorization
5 vectorizer = TfidfVectorizer(stop_words='english', max_features=5000, ngram_range=(1,2))
6 vectorizer.fit(x_train)
7
8 # Word2Vec Model
9 word2vec_model = Word2Vec(sentences=x_train, vector_size=100, window=5, min_count=1, workers=4)
10
11 # Create a pipeline with both transformations
12 combined_features = FeatureUnion([
13     ('tfidf', vectorizer),
14     ('word2vec', Word2VecTransformer(word2vec_model))
15 ])
16
17 # Combined_features as pickle
18 with open('combined_features.sav', 'wb') as combined_features_file:
19     pickle.dump(combined_features, combined_features_file)
20
21 # Transform the data
22 X_train_combined = combined_features.fit_transform(x_train)
23 X_test_combined = combined_features.transform(x_test)
24
25
26 # Transform Labels to Binary Matrix
27 mlb = MultiLabelBinarizer()
28 y_train_bin = mlb.fit_transform(y_train)
29 y_test_bin = mlb.transform(y_test)

```

WARNING:gensim.models.word2vec:Each 'sentences' item should be a list of words (usually unicode strings). First item here is in /usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_label.py:895: UserWarning: unknown class(es) ['cottonseed', 'f-c
warnings.warn(

Transforms test text data (test_dict['text']) into numerical vectors using the pre-trained TF-IDF vectorizer. It also converts test topics (test_dict['topic']) into a binary matrix using the MultiLabelBinarizer. Finally, it prints the classes (unique topics) present in the binary matrix,

providing insight into the labeled topics used in testing.

```
1 p(mlb.classes_)

['acq' 'alum' 'austdlr' 'barley' 'bop' 'can' 'carcass' 'cocoa' 'coconut'
 'coconut-oil' 'coffee' 'copper' 'copra-cake' 'corn' 'cornglutenfeed'
 'cotton' 'cpi' 'cpu' 'crude' 'dfl' 'dlr' 'dmk' 'earn' 'fishmeal' 'fuel'
 'gas' 'gnp' 'gold' 'grain' 'groundnut' 'heat' 'hog' 'housing' 'income'
 'instal-debt' 'interest' 'inventories' 'ipi' 'iron-steel' 'jet' 'jobs'
 'l-cattle' 'lead' 'lei' 'linseed' 'livestock' 'lumber' 'meal-feed'
 'money-fx' 'money-supply' 'naphtha' 'nat-gas' 'nickel' 'nzdlr' 'oat'
 'oilseed' 'orange' 'palladium' 'palm-oil' 'palmkernel' 'pet-chem'
 'platinum' 'plywood' 'pork-belly' 'potato' 'propane' 'rand' 'rape-oil'
 'rapeseed' 'reserves' 'retail' 'rice' 'rubber' 'saudriyal' 'ship'
 'silver' 'sorghum' 'soy-meal' 'soy-oil' 'soybean' 'stg' 'strategic-metal'
 'sugar' 'sun-oil' 'sunseed' 'tapioca' 'tea' 'tin' 'trade' 'veg-oil'
 'wheat' 'wool' 'wpi' 'yen' 'zinc']
```

Training multiple classifiers is advantageous for diverse insights and enhanced performance which provides robustness, balances bias and variance trade-offs, and allows for exploratory analysis, ultimately contributing to a more versatile and reliable machine learning system. Due to limited time, only 3 classifiers are run to give the quick result.

```
1 #Wait for more than 1min to run
2
3 classifiers = [
4     OneVsRestClassifier(LogisticRegression(random_state=0)),
5     OneVsRestClassifier(SGDClassifier(alpha=0.0001, random_state=0)),
6     OneVsRestClassifier(MultinomialNB(alpha=0.1)),
7     #OneVsRestClassifier(RandomForestClassifier(n_estimators=50, max_depth=3, random_state=0)),
8     #OneVsRestClassifier(SVC(kernel='linear', C=1.0, random_state=0)),
9     #OneVsRestClassifier(MLPClassifier(hidden_layer_sizes=(50,), max_iter=200, alpha=0.0001, solver='adam', random_state=0))
10
11 ]
12
13 entries = []
14 cv_df = pd.DataFrame()
15 for clf in classifiers:
16     print(str(clf))
17     scores = cross_val_score(clf, X_train_combined, y_train_bin, cv=5, scoring='accuracy')
18     for fold_idx, accuracy in enumerate(scores):
19         entries.append((clf.estimator, fold_idx, accuracy))
20
21 cv_df = pd.DataFrame(entries, columns=['classifier', 'fold_idx', 'accuracy'])

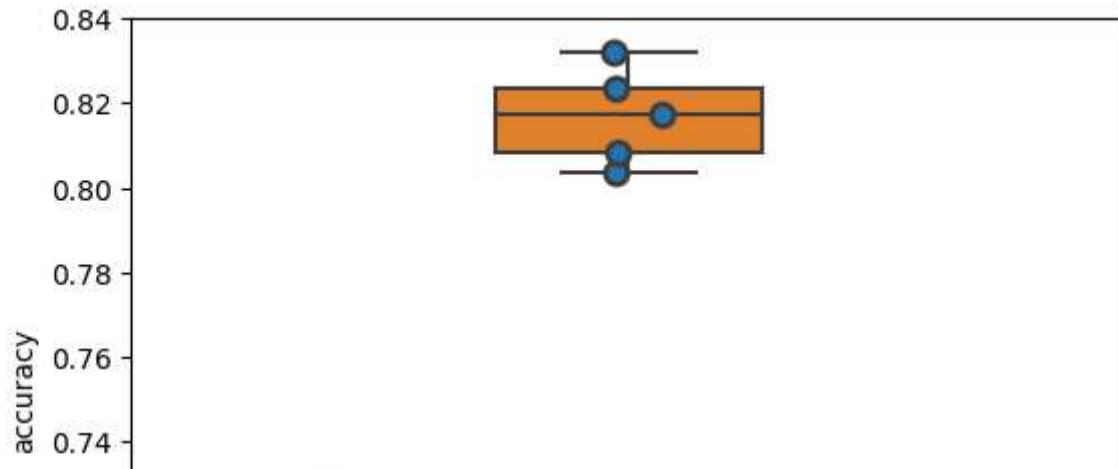
    warnings.warn(
        "/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 34 is present in all training exemplar"
    )
    warnings.warn(
        "/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 57 is present in all training exemplar"
    )
    warnings.warn(
        "/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 17 is present in all training exemplar"
    )
    warnings.warn(
        "/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 44 is present in all training exemplar"
    )
    warnings.warn(
        "/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 63 is present in all training exemplar"
    )
    warnings.warn(
```

```
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 14 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 23 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 34 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 57 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 17 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 44 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 63 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 12 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 53 is present in all training exempl
warnings.warn(
OneVsRestClassifier(estimator=MultinomialNB(alpha=0.1))
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 5 is present in all training example
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 14 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 23 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 34 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 57 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 17 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 44 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 63 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 12 is present in all training exempl
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/multiclass.py:77: UserWarning: Label not 53 is present in all training exempl
warnings.warn(
```

Accuracy table for different classifiers

Best classifier is SGD

```
1 # Your existing code for creating the boxplot
2 cv_df['classifier'] = cv_df['classifier'].astype(str)
3 sns.boxplot(x='classifier', y='accuracy', data=cv_df)
4 sns.stripplot(x='classifier', y='accuracy', data=cv_df,
5                 size=8, jitter=True, edgecolor="gray", linewidth=2)
6
7 # Shorten x-axis labels
8 plt.xticks(rotation=45, ha='right') # Adjust the rotation angle as needed
9
10 plt.show()
```



```
1 # Train Classifier
2 clf = OneVsRestClassifier(SGDClassifier(alpha=0.0001, random_state=0))
3 clf.fit(X_train_combined, y_train_bin)
4 # Predictions and Evaluation
5 y_pred = clf.predict(X_test_combined)

1 print("Classification Report:")
2 print(classification_report(y_test_bin, y_pred, target_names=mlb.classes_, zero_division=1))
```



Classification Report Analysis

The classification report highlights the model's proficiency in specific topics, as evident from high precision, recall, and F1-scores. For example, the 'acq' topic boasts a precision of 0.98, recall of 0.93, and an F1-score of 0.95. However, the model's performance is unevenly distributed across various topics.

The macro average for recall and F1-score, standing at 0.49 and 0.55 respectively, suggests an opportunity to enhance the model's performance on minority classes.

Conversely, the weighted average, signaling high precision (0.95), recall (0.80), and F1-score (0.84), indicates effective performance on common topics. Yet, this may not accurately represent the model's efficacy on less frequent topics.

The samples average, portraying elevated precision (0.97), recall (0.87), and F1-score (0.86), demonstrates commendable performance on individual instances. However, it might not fully encapsulate the model's effectiveness for each specific topic.

In summary, while the model excels in specific areas, there exists room for improvement, particularly in achieving a more balanced performance across all topics.

Save the classifier model as pickle file

```
1 # Save the model to a file
2 with open('classifier.sav', 'wb') as file:
3     pickle.dump(clf, file)
```

Classifier Prediction

```
1 #When handling the list for new text input
2 def preprocess_text_advanced(text_input):
3     # If text_input is a list, process each element separately
4     if isinstance(text_input, list):
5         processed_texts = []
6         for single_text in text_input:
7             # Process each element using the single-text processing function
8             processed_text = preprocess_text_advanced_single(single_text)
9             processed_texts.append(processed_text)
10    return processed_texts
11 else:
12     # If text_input is a single string, process it directly
13     return preprocess_text_advanced_single(text_input)
14
15 # Download NLTK resources
16 nltk.download('stopwords')
17 nltk.download('punkt')
18
19 #When handling the single for new text input
20 def preprocess_text_advanced_single(text_input):
21     # Lowercasing
22     text_input = text_input.lower()
23
24     # Removing HTML tags (if any)
25     text_input = re.sub(r'<.*?>', '', text_input)
26
27     # Removing special characters, numbers, and newlines
28     text_input = re.sub(r'[^a-zA-Z\s]', '', text_input)
29     text_input = re.sub(r'\d+', '', text_input)
30     text_input = text_input.replace('\n', ' ')
31
32     # Tokenization
33     tokens = word_tokenize(text_input)
34
35     # Removing stop words
36     stop_words = set(stopwords.words('english'))
37     tokens = [word for word in tokens if word not in stop_words]
```



```
38
39     # Stemming
40     stemmer = PorterStemmer()
41     stemmed_words = [stemmer.stem(word) for word in tokens]
42
43     # Joining tokens back into a string
44     preprocessed_text = ' '.join(stemmed_words)
45
46     return preprocessed_text
47
48 def classifier_topic_prediction(text_input):
49
50     # Preprocess the new text
51     processed_text = preprocess_text_advanced(text_input)
52
53     # Load the TF-IDF vectorizer (already fitted on your training data)
54     with open('combined_features.sav', 'rb') as combined_features_file:
55         combined_features = pickle.load(combined_features_file)
56
57     # Transform the data
58     X_test_combined = combined_features.transform(processed_text)
59
60     # Load the classifier model
61     with open('classifier.sav', 'rb') as classifier_model:
62         classifier = pickle.load(classifier_model)
63
64     # Predict the topic
65     one_hot_array = classifier.predict(X_test_combined)
66     position = np.argmax(one_hot_array)
67     p('Predicted Topic: ', mnb.classes_[position])
68
69
70 def classifier_topic_prediction_array(text_input):
71
72     if len(text_input) > 1:
73
74         for i in text_input:
```

```
75     classifier_topic_prediction([i])
76
77 else:
78     classifier_topic_prediction(text_input)
79
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

Prepared 3 kinds of Texts:

- train_text (Extracted from train data)
- new_text (Copy from internet)
- test_text_array (Extracted from test data)

1 train_text = ["""The U.S. Agriculture Department
2 reported the farmer-owned reserve national five-day average
3 price through February 25 as follows (Dlrs/Bu-Sorghum Cwt) -
4 Nat'l Loan Release Call
5 Avge Rate-X Level Price Price
6 Wheat 2.55 2.40 IV 4.65 --
7 V 4.65 --
8 VI 4.45 --
9 Corn 1.35 1.92 IV 3.15 3.15
10 V 3.25 --
11 X - 1986 Rates.
12
13 Nat'l Loan Release Call
14 Avge Rate-X Level Price Price
15 Oats 1.24 0.99 V 1.65 --
16 Barley n.a. 1.56 IV 2.55 2.55
17 V 2.65 --
18 Sorghum 2.34 3.25-Y IV 5.36 5.36
19 V 5.54 --
20 Reserves I, II and III have matured. Level IV reflects
21 grain entered after Oct 6, 1981 for feedgrain and after July
22 23, 1981 for wheat. Level V wheat/barley after 5/14/82,
23 corn/sorghum after 7/1/82. Level VI covers wheat entered after
24 January 19, 1984. X-1986 rates. Y-dlrs per CWT (100 lbs).
25 n.a.-not available.""]
26
27 new_text = ["""TBEIJING: China is set to showcase the strides made towards self-reliance at an inaugural global trade fair this
28
29 The participation of Western businesses at the China International Supply Chain Expo will also be in focus as Sino-US ties warm
30
31 China has promoted the five-day event - which begins in Beijing on Nov 28 - as the world's first supply chain expo at the nation
32
33 It will be the latest in a series of global trade shows China has held in the second half of the year. At least two have already
34
35 Chinese state media have framed these events as a sign of Beijing's contributions to the global economy and its commitment to in
36
37 A quarter of the 515 companies and institutions participating in the supply chain expo are foreign, hailing from 55 countries an

38
39 US firms make up around 20 per cent of the foreign registrations, with the likes of Amazon, Apple, Tesla, Intel, HP and Qualcomm
40
41 The number of US exhibitors is "much higher than expected", according to the expo organiser.
42
43 We hope that the participating US companies actively engage in the expo, yielding fruitful outcomes," said CCPIT Vice-Chairman M
44
45 We will continue to vigorously promote high-level opening up and better protect the rights and interests of foreign investors pe
46
47 test_text_array = ['indonesian commod exchang may expand indonesian commod exchang like start trade least one new commod possib
48
49

```
1 p('Classifier')
2 p()
3 classifier_topic_prediction_array(train_text)
4 p()
5 classifier_topic_prediction_array(new_text)
6 p()
7 classifier_topic_prediction_array(test_text_array)
```

Classifier

Predicted Topic: barley

Predicted Topic: trade

Predicted Topic: rubber

Predicted Topic: grain

Predicted Topic: gold

Predicted Topic: acq

✓ Step 6: Prediction based on BERT Transformer

To execute Step 6, make sure to open the IPython notebook in Google Colab since a GPU is necessary. Afterward, click on "Runtime," navigate to "Change Runtime Type," and choose "T4 GPU" from the options.

This series of functions and operations help in creating and preprocessing data dictionaries for the ModApte split with the removal of unwanted topics

```
1 #Create a new function to cater for ModApte split for data dict
2 def create_data_dict(df, lewissplit, r_topics):
3     data_dict = {
4         'text': df[(df.lewissplit == lewissplit) & (df.r_topics == r_topics)].text.to_list(),
5         'topic': df[(df.lewissplit == lewissplit) & (df.r_topics == r_topics)].topic.to_list(),
6     }
7     return data_dict
8
9 # 'test' dict
10 test_dict = create_data_dict(df, 'TEST', 'YES')
11
12 # 'train' dict
13 train_dict = create_data_dict(df, 'TRAIN', 'YES')
14
15 # 'unused' dict (Will not be ignored)
16 unused_dict = create_data_dict(df, 'NOT-USED', '')
17
18 def remove_unwanted_topic(dict):
19     df1 = pd.DataFrame(dict)
20     df1 = df1[~df1['topic'].apply(lambda x: any(topic in x for topic in occurrence_only_one))]
21     return df1.to_dict(orient='list')
22
23 train_dict = remove_unwanted_topic(train_dict)
24 test_dict = remove_unwanted_topic(test_dict)
```

New dataset Dictionary

```
1 from datasets import DatasetDict, Dataset
2
3 train_dataset = Dataset.from_dict(train_dict)
4 test_dataset = Dataset.from_dict(test_dict)
5 unused_dataset = Dataset.from_dict(unused_dict)
6
7 dataset = DatasetDict({
8     'train': train_dataset,
9     'test': test_dataset,
10    'unused': unused_dataset,
11})
12
```

```
1 dataset
2
3 DatasetDict({
4     'train': Dataset({
5         'features': ['text', 'topic'],
6         'num_rows': 9588
7     })
8     'test': Dataset({
9         'features': ['text', 'topic'],
10        'num_rows': 3292
11    })
12     'unused': Dataset({
13         'features': ['text', 'topic'],
14         'num_rows': 0
15     })
16 })
```

Prepares the **labels** for a multi-label classification task by transforming them into a multi-hot encoding format and creating dictionaries for label mappings. The assertion check ensures the correctness of the **label** processing

```

1 # Check number of unique labels
2 unique_labels = set(chain.from_iterable(dataset['train'][["topic"]]))
3 print(f"We have {len(unique_labels)} unique labels:\n{unique_labels}")
4
5 # Transform topics into multi-hot encoding format
6 mlb = MultiLabelBinarizer()
7 mlb.fit(dataset['train'][['topic']])
8 dataset = dataset.map(
9     lambda x: {"labels": torch.from_numpy(mlb.transform(x[["topic"]])).float()}, batched=True)
10
11 labels = mlb.classes_
12 id2label = {idx:label for idx, label in enumerate(labels)}
13 label2id = {label:idx for idx, label in enumerate(labels)}
14 num_labels = len(id2label)
15
16 assert num_labels == len(unique_labels)

```

```

We have 95 unique labels:
{'yen', 'lumber', 'oilseed', 'propane', 'zinc', 'dfl', 'livestock', 'copra-cake', 'inventories', 'tea', 'sugar', 'pet-chem', 'c
Map: 100%                                         9588/9588 [00:00<00:00, 46533.43 examples/s]

Map: 100%                                         3292/3292 [00:00<00:00, 38430.69 examples/s]
/usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_label.py:895: UserWarning: unknown class(es) ['cottonseed', 'f-c
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/preprocessing/_label.py:895: UserWarning: unknown class(es) ['sfr'] will be ign
  warnings.warn(

```

Useful to quickly inspect a subset of label mappings and ensure that the indices and labels are correctly mapped in the id2label dictionary. It can help catch any unexpected issues with the label processing

```
1 # sanity check:  
2 for idx, label in id2label.items():  
3     if idx>=10:  
4         break  
5  
6     print(f"{idx}: {label}")
```

```
0: acq  
1: alum  
2: austdlr  
3: barley  
4: bop  
5: can  
6: carcass  
7: cocoa  
8: coconut  
9: coconut-oil
```

Prepares the training dataset for model training by tokenizing the text and removing unnecessary columns

```
1 tokenizer = AutoTokenizer.from_pretrained("distilbert-base-cased")  
2  
3 # Tokenize and remove unwanted columns  
4 def tokenize_function(example):  
5     return tokenizer(example["text"], padding="max_length", truncation=True, max_length=512)  
6  
7 columns = dataset["train"].column_names  
8 columns.remove("text")  
9 columns.remove("labels")  
10  
11 #The dataset.map function is used to apply the tokenize_function to each example in the training dataset.  
12 #The batched=True parameter indicates that tokenization should be applied in batches for efficiency,  
13 #and the specified columns are removed from the resulting tokenized dataset  
14 tokenized_dataset = dataset.map(tokenize_function, batched=True, remove_columns=columns)
```

```
tokenizer_config.json: 100%                                29.0/29.0 [00:00<00:00, 2.31kB/s]

config.json: 100%                                         465/465 [00:00<00:00, 33.0kB/s]

vocab.txt: 100%                                         213k/213k [00:00<00:00, 1.61MB/s]

tokenizer.json: 100%                                     436k/436k [00:00<00:00, 6.12MB/s]

Map: 100%                                              9588/9588 [00:04<00:00, 2438.35 examples/s]

1 example = tokenized_dataset['train'][0]
2 print(example.keys())

dict_keys(['text', 'labels', 'input_ids', 'attention_mask'])

1 [id2label[idx] for idx, label in enumerate(example['labels']) if label == 1.0]

['earn']
```

To convert the tokenized dataset to the PyTorch format.

```
1 tokenized dataset.set format("torch")
```

This code initializes a sequence classification model using the `AutoModelForSequenceClassification` class from the `transformers` library. Let's break down the code:

```
problem_type="multi_label_classification", # Problem type (multi-label classification in this case)
id2label=id2label,                      # Mapping from label indices to label names
label2id=label2id,                      # Mapping from label names to label indices
)
```

Explanation:

- 1. Model Initialization:** The `AutoModelForSequenceClassification.from_pretrained` function initializes a pre-trained sequence classification model. In this case, it uses the "distilbert-base-cased" model.
- 2. Number of Labels:** The `num_labels` parameter is set to the number of unique labels in the dataset. This is necessary for initializing the classification head of the model.
- 3. Problem Type:** The `problem_type` parameter is set to "multi_label_classification" to indicate that it's a multi-label classification problem.
- 4. Label Mapping:** The `id2label` and `label2id` parameters are used to provide mappings between label indices and label names. These mappings are helpful during model evaluation to convert between indices and human-readable labels.

Now, the `model` variable holds an instance of a pre-trained sequence classification model ready for fine-tuning on the specific multi-label classification task.

```
1 model = AutoModelForSequenceClassification.from_pretrained(
2     "distilbert-base-cased",
3     num_labels=num_labels,
4     problem_type="multi_label_classification",
5     id2label=id2label,
6     label2id=label2id
7 )
```

model.safetensors: 100%

263M/263M [00:01<00:00, 241MB/s]

Some weights of DistilBertForSequenceClassification were not initialized from the model checkpoint at distilbert-base-cased and You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
1 # test forward pass
2 tokenized_dataset['train'][0]['labels'].type()
3
4 'torch.FloatTensor'
```

Generate classification report using pipeline class