

Names	IDs
yahia ashraf	20200636
yahia mahmoud	20201222
hamza abdel hamid	20200162
ziad ibrahim	20200193
omar tarek	20200348

```
1 !pip install scikeras
```

```
1 import numpy as np
2 import tensorflow as tf
3 import pandas as pd
4 from sklearn.metrics import accuracy_score
5 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
6 from sklearn.model_selection import GridSearchCV, train_test_split
7 from sklearn.metrics import classification_report
8 from tensorflow.keras import layers, Sequential
9 from scikeras.wrappers import KerasClassifier
10 from matplotlib import pyplot
11 import spacy
12 import matplotlib.pyplot as plt
13 from sklearn.svm import LinearSVC
```

```
1 df = pd.read_csv("sentimentdataset (Project 1).csv")
2 df
```

	Source	ID	Message	Target
0	Yelp	0	Crust is not good.	0
1	Yelp	1	Not tasty and the texture was just nasty.	0
2	Yelp	2	Stopped by during the late May bank holiday of...	1
3	Yelp	3	The selection on the menu was great and so wer...	1
4	Yelp	4	Now I am getting angry and I want my damn pho.	0
...
2740	Amazon	994	The screen does get smudged easily because it ...	0
2741	Amazon	995	What a piece of junk.. I lose more calls on th...	0
2742	Amazon	996	Item Does Not Match Picture.	0
2743	Amazon	997	The only thing that disappoint me is the infra...	0
2744	Amazon	998	You can not answer calls with the unit, never ...	0

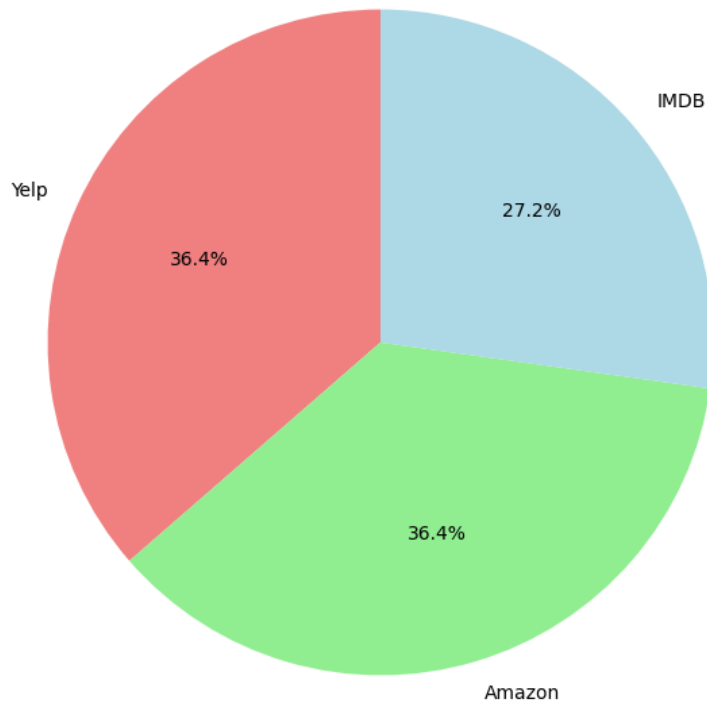
2745 rows × 4 columns

✓ Exploratory Data Analysis

```
1 source_counts = df['Source'].value_counts()
2 target_counts = df['Target'].value_counts()
```

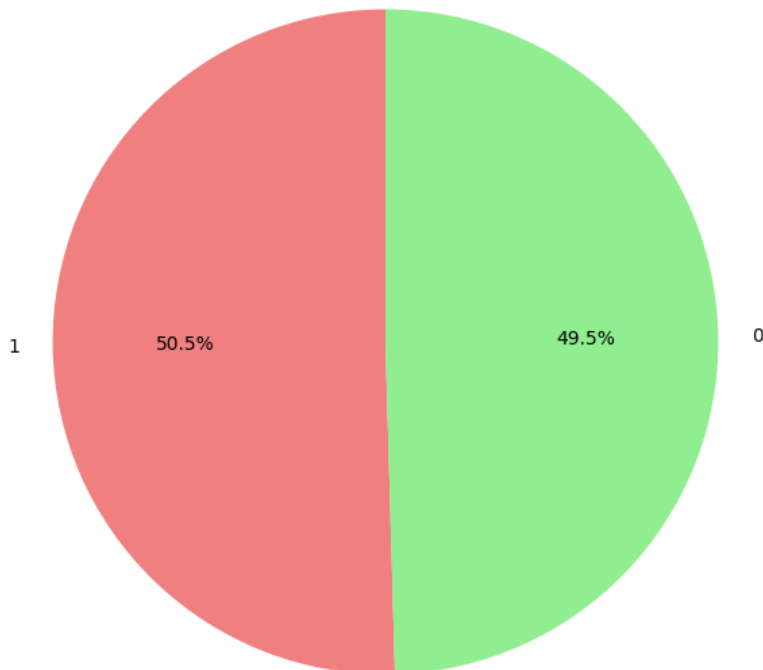
```
1 plt.figure(figsize=(8, 8))
2 plt.pie(source_counts, labels=source_counts.index, autopct='%1.1f%%', startangle=90, colors=|
3 plt.title('Distribution of Categories in the "Source" Column')
4 plt.show()
```

Distribution of Categories in the "Source" Column



```
1 plt.figure(figsize=(8, 8))
2 plt.pie(target_counts, labels=target_counts.index, autopct='%1.1f%%', startangle=90, colors=|
3 plt.title('Distribution of Categories in the "Target" Column')
4 plt.show()
```

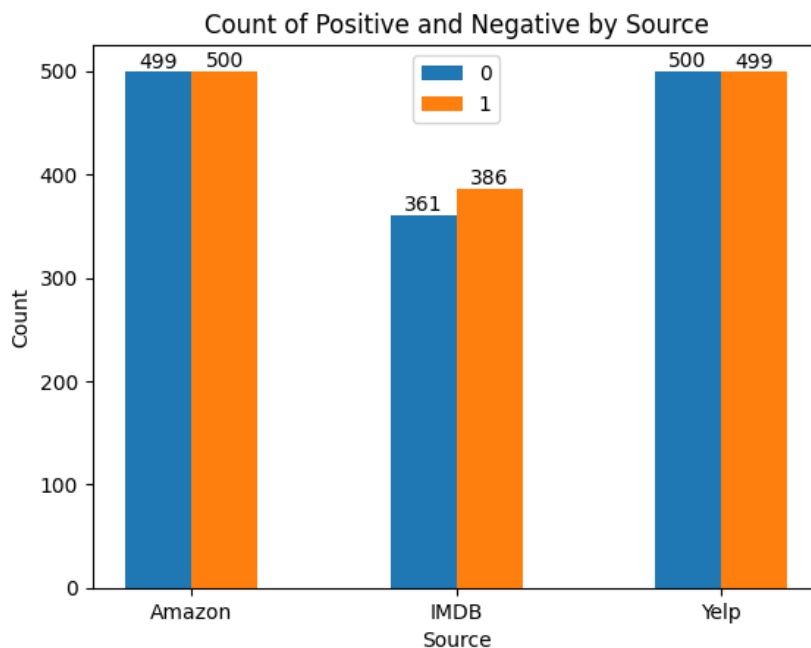
Distribution of Categories in the "Target" Column



```

1 # Group by 'source' and 'target' and get the count
2 grouped_data = df.groupby(['Source', 'Target']).size().unstack()
3
4 # Plotting the bar chart
5 fig, ax = plt.subplots()
6 width = 0.25 # the width of the bars
7
8 sources = grouped_data.index
9 ind = range(len(sources))
10
11 bar_0 = ax.bar(ind, grouped_data[0], width, label='0')
12 bar_1 = ax.bar([i + width for i in ind], grouped_data[1], width, label='1')
13
14 # Adding labels and title
15 ax.set_xlabel('Source')
16 ax.set_ylabel('Count')
17 ax.set_title('Count of Positive and Negative by Source')
18 ax.set_xticks([i + width/2 for i in ind])
19 ax.set_xticklabels(sources)
20 ax.legend()
21
22 for bar in bar_0:
23     yval = bar.get_height()
24     plt.text(bar.get_x() + bar.get_width()/2, yval, round(yval, 2), ha='center', va='bottom')
25
26 for bar in bar_1:
27     yval = bar.get_height()
28     plt.text(bar.get_x() + bar.get_width()/2, yval, round(yval, 2), ha='center', va='bottom')
29
30
31 plt.show()

```



✓ Preprocessing

```

1 nlp = spacy.load('en_core_web_sm')
2 # Add 'not' to spaCy stop words
3 spacy_stop_words = spacy.lang.en.stop_words.STOP_WORDS
4
5 # add others
6 words_to_keep = ["not", "no", "never", "but", "only", "against",
7     "don't", "doesn't", "didn't", "isn't", "aren't", "wasn't", "weren't",
8     "hasn't", "haven't", "hadn't", "won't", "wouldn't", "can't", "cannot",
9     "could've", "should've", "would've", "doesn't", "didn't", "isn't", "ain't"]
10
11 for word in words_to_keep:
12     spacy_stop_words.discard(word)
13
14 # Function for lemmatization and removing stop words
15 def lemmatize_and_remove_stop_words(text):
16     doc = nlp(text)
17     lemmatized_text = [token.lemma_ for token in doc if token.text.lower() not in spacy_stop_
18     return ' '.join(lemmatized_text)
19
20 df = pd.read_csv("sentimentdataset (Project 1).csv")
21
22 df.drop(columns = ['Source', 'ID'], inplace=True)
23
24 # Apply the function to the "Message" column
25 df['Message'] = df['Message'].apply(lemmatize_and_remove_stop_words)
26
27 print(df['Message'])
28
29 # Save the updated DataFrame to a new CSV file
30 df.to_csv('sentimentdataset_stopwords_lemmatized.csv', index=False)

```

```

0          crust not good .
1      not tasty texture nasty .
2  stop late bank holiday Rick Steve recommendati...
3      selection menu great price .
4      get angry want damn pho .
   ...
2740  screen smudge easily touch ear face .
2741      piece junk .. lose call phone .
2742          item not match Picture .
2743  only thing disappoint infra red port ( irda ) .
2744      not answer call unit , never work !
Name: Message, Length: 2745, dtype: object

```

```

1 # read the preprocessed file
2 df = pd.read_csv("sentimentdataset_stopwords_lemmatized.csv")
3 df

```

```

1 X = df['Message']
2 y = df['Target']

# Create a TfIdfVectorizer
3
4 # Transform the training and testing data
5 X_tfidf = vectorizer.fit_transform(X)
6
7 # Split the data into training and testing sets
8 X_train, X_test, y_train, y_test = train_test_split(X_tfidf, y, test_size=0.2, random_state=42)

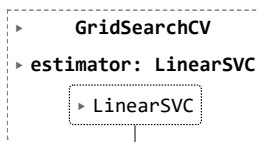
```

✓ SVC Model

```

# Create a DataFrame
1 # to convert it to a data frame
2 df_tfidf = pd.DataFrame(X_train_tfidf.toarray(), columns=vectorizer.get_feature_names_out())
3
4 svc = LinearSVC(max_iter=100000, dual=False)
5
6 # Define the parameter grid for grid search
7 param_grid = {
8     'C': [0.001, 0.01, 0.1, 1, 10, 100],
9     'penalty': ['l1', 'l2'],
10    'loss': ['squared_hinge'],
11    'tol': [1e-3, 1e-2, 1e-1, 1],
12    'class_weight': [None, 'balanced']
13 }
14
15 # Perform grid search with 5-fold cross-validation
16 grid_search = GridSearchCV(svc, param_grid, scoring='accuracy', error_score='raise', cv=5)
17 grid_search.fit(X_train, y_train)

```



```

1 # Print the best parameters found by grid search
2 print("Best Parameters: ", grid_search.best_params_)

Best Parameters: {'C': 0.1, 'class_weight': 'balanced', 'loss': 'squared_hinge', 'penalty': 'l2', 'tol': 0.01}

1 # Predict on the testing set with the best model
2 y_pred = grid_search.predict(X_test)
3
4 # Evaluate the model
5 accuracy = accuracy_score(y_test, y_pred)
6 print("Accuracy on Testing Set: {:.2f}%".format(accuracy * 100))

```

Accuracy on Testing Set: 83.79%

```
1 # Display additional metrics
2 print("\nClassification Report:")
3 print(classification_report(y_test, y_pred))
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.81      0.86      0.84       265
     1       0.87      0.81      0.84       284

 accuracy          0.84
 macro avg         0.84
 weighted avg      0.84
```

✓ ANN Model

```
1 def build_model(neurons=16, learning_rate=0.01):
2     model = Sequential()
3     model.add(layers.Dense(units=neurons, activation='relu', input_dim=X_train.shape[1]))
4     model.add(layers.Dense(units=neurons, activation='relu'))
5     model.add(layers.Dense(units=1, activation='sigmoid'))
6
7     model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=learning_rate),
8                   loss='binary_crossentropy',
9                   metrics=['accuracy'])
10    return model

1 model = KerasClassifier(build_fn=build_model,neurons = [16, 32, 64, 128], learning_rate = [0.
2
3 print(model.get_params().keys())
4 param_grid = {
5     'neurons': [16, 32, 64, 128],
6     'learning_rate': [0.001, 0.01, 0.1],
7     'batch_size': [16, 32, 64, 128]
8 }
9
dict_keys(['model', 'build_fn', 'warm_start', 'random_state', 'optimizer', 'loss', 'metrics', 'batch_size', 'validation
10
11 grid = GridSearchCV(estimator=model, param_grid=param_grid, cv=3)
12 grid_result = grid.fit(X_train, y_train)

1 print("Best: ANN model is %f using %s" % (grid_result.best_score_, grid_result.best_params_))

Best: ANN model is 0.791439 using {'batch_size': 16, 'learning_rate': 0.01, 'neurons': 16}

1 best_model = grid_result.best_estimator_
2 best_model.fit(X_train, y_train)
3 prediction = best_model.predict(X_test)
4 acc = accuracy_score(y_pred=prediction, y_true=y_test)
5 print(f'Test Accuracy: {acc}')
```

/usr/local/lib/python3.10/dist-packages/scikeras/wrappers.py:915: UserWarning: ``build_fn`` will be renamed to ``model``
X, y = self._initialize(X, y)
Test Accuracy: 0.8415300546448088

```
1 report = classification_report(y_test, prediction)
2 print("Classification Report:\n", report)
```

```
Classification Report:
              precision    recall  f1-score   support

     0       0.84         0.83         0.83         265
     1       0.84         0.85         0.85         284

 accuracy          0.84          0.84          0.84          549
 macro avg         0.84          0.84          0.84          549
 weighted avg         0.84          0.84          0.84          549
```