n [3]:	df_true = pd.read_csv("C:/Users/Collins PC/Downloads/Fake News Detection Project/True.csv")
[4]: t[4]:	title text subject date Donald Trump Sends Out Embarrassing New Year' Donald Trump just couldn t wish all Americans News December 31, 2017 Drunk Bragging Trump Staffer Started Russian House Intelligence Committee Chairman Devin Nu News December 31, 2017 Sheriff David Clarke Becomes An Internet Joke On Friday, it was revealed that former Milwauk News December 30, 2017
[5]: t[5]:	
	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 The code below adds a new column named "class" to both DataFrames, df_fake and df_true, and assigns a value of 0 to all entries in the "class" column of df_fake and a value of 1 to all entries in the "class" column of df_true. This step is crucial supervised machine learning tasks where we need to define a target variable (in this case, the authenticity of the news articles) that the models will learn to predict based on the features of the data. By labeling the data in this manner, we establish the ground truth for classification, with "0" representing fake news and "1" representing true news.
	The code below displays concise summary information about the DataFrames df_true and df_fake using the info() method in pandas. The summary includes the total number of entries in each column, the data type of each column, and the number non-null values in each column. This summary helps in understanding the structure of the dataset, identifying any missing values, and assessing the data types of the features, which is crucial for data preprocessing and analysis. df_true.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 21417 entries, 0 to 21416 Data columns (total 5 columns): # Column Non-Null Count Dtype </class>
[9]:	2 subject 21417 non-null object 3 date 21417 non-null object 4 class 21417 non-null int64 dtypes: int64(1), object(4) memory usage: 836.7+ KB df_fake.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 23481 entries, 0 to 23480 Data columns (total 5 columns): # Column Non-Null Count Dtype</class>
	0 title 23481 non-null object 1 text 23481 non-null object 2 subject 23481 non-null object 3 date 23481 non-null object 4 class 23481 non-null int64 dtypes: int64(1), object(4) memory usage: 917.4+ KB The code snippet below removes the last 10 rows from the DataFrames df_fake and df_true for manual testing purposes. It creates new DataFrames named df_fake_manual_testing and df_true_manual_testing containing these last 10 rows, which be used for manual evaluation of the trained models. The removal is done by iterating over a range of index values corresponding to the rows to be removed, and using the drop() method with the axis=0 parameter to drop rows based on their index. This process ensures that the original DataFrame modified accordingly, removing the specified rows from the dataset.
[10]:	However, it's worth noting that directly modifying a DataFrame within a loop can sometimes lead to unexpected behavior due to indexing changes. It's advisable to use caution when modifying DataFrames in place, especially within loops. Alternat you can create a copy of the DataFrame and modify the copy to avoid modifying the original data. # Removing last 10 rows for manual testing df_fake_manual_testing = df_fake.tail(10) for i in range(23480,23470,-1): df_fake.drop([i], axis = 0, inplace = True) df_true_manual_testing = df_true.tail(10) for i in range(21416,21406,-1):
[11]: [11]:	The code below labels the manual testing data by assigning the value "0" to the "class" column of the DataFrame df_fake_manual_testing and the value "1" to the "class" column of the DataFrame df_true_manual_testing. This process is essential establishing the ground truth, indicating whether each news article in the manual testing datasets is fake or true. By using the .loc[] accessor to select all rows and the "class" column, followed by the assignment operation, we ensure that each new article is correctly labeled according to its authenticity. These class labels serve as reference points for evaluating the performance of machine learning models during manual testing, providing valuable insights into the models' ability to distinguish between fake and true news articles.
[13]: [13]:	
	Hillary Clinton: 'Israel First' (and no peace Robert Fantina CounterpunchAlthough the United Middle-east January 18, 2016 0 McPain: John McCain Furious That Iran Treated 21st Century Wire says As 21WIRE reported earl Middle-east January 16, 2016 0 JUSTICE? Yahoo Settles E-mail Privacy Class-ac 21st Century Wire says It s a familiar theme Middle-east January 16, 2016 0 Sunnistan: US and Allied 'Safe Zone' Plan to T Patrick Henningsen 21st Century WireRemember Middle-east January 15, 2016 0 How to Blow \$700 Million: Al Jazeera America F 21st Century Wire says Al Jazeera America will Middle-east January 14, 2016 0 10 U.S. Navy Sailors Held by Iranian Military 21st Century Wire says As 21WIRE predicted in Middle-east January 12, 2016 0
[14]:	
	21412 'Fully committed' NATO backs new U.S. approach BRUSSELS (Reuters) - NATO allies on Tuesday we worldnews August 22, 2017 1 21413 LexisNexis withdrew two products from Chinese LONDON (Reuters) - LexisNexis, a provider of I worldnews August 22, 2017 1 21414 Minsk cultural hub becomes haven from authorities MINSK (Reuters) - In the shadow of disused Sov worldnews August 22, 2017 1 21415 Vatican upbeat on possibility of Pope Francis MOSCOW (Reuters) - Vatican Secretary of State worldnews August 22, 2017 1 21416 Indonesia to buy \$1.14 billion worth of Russia JAKARTA (Reuters) - Indonesia will buy 11 Sukh worldnews August 22, 2017 1 The code below concatenates the DataFrames df_fake_manual_testing and df_true_manual_testing along the rows (axis=0) to create a single DataFrame named df_manual_testing. This combined DataFrame contains the manual testing data for
[15]:	fake and true news articles. Subsequently, the to_csv() method is used to export the df_manual_testing DataFrame to a CSV file named "manual_testing.csv". This file will contain the manual testing data, which can be used for further analysis or evaluation purposes. By saving the manual testing data to a CSV file, it becomes easily accessible for future reference and can be seamlessly integrated into the overall workflow of the project.
[16]: [16]:	Concatenating the two DataFrames allows for a unified dataset that can be used for further analysis, feature engineering, or model training in the fake news detection project. df_merge = pd.concat([df_fake, df_true], axis =0) title text subject date class Donald Trump Sends Out Embarrassing New Year' Donald Trump just couldn t wish all Americans News December 31, 2017 0 Drunk Bragging Trump Staffer Started Russian House Intelligence Committee Chairman Devin Nu News December 31, 2017 0
	Sheriff David Clarke Becomes An Internet Joke On Friday, it was revealed that former Milwauk News December 30, 2017 0 Trump Is So Obsessed He Even Has Obama's Name On Christmas day, Donald Trump announced that News December 29, 2017 0 Pope Francis Just Called Out Donald Trump Dur Pope Francis used his annual Christmas Day mes News December 25, 2017 0 Racist Alabama Cops Brutalize Black Boy While The number of cases of cops brutalizing and ki News December 25, 2017 0 Fresh Off The Golf Course, Trump Lashes Out A Donald Trump spent a good portion of his day a News December 23, 2017 0 Trump Said Some INSANELY Racist Stuff Inside In the wake of yet another court decision that News December 23, 2017 0 Becember 23, 2017 0 Many people have raised the alarm regarding th News December 22, 2017 0
[17]: [17]: [18]:	9 WATCH: Brand-New Pro-Trump Ad Features So Muc Just when you might have thought we d get a br News December 21, 2017 0 df_merge.columns Index(['title', 'text', 'subject', 'date', 'class'], dtype='object') df = df_merge.drop(["title", "subject", "date"], axis = 1) The code below calculates the sum of missing values (NaN) for each column in the DataFrame df using the isnull() method, which returns a DataFrame of Boolean values indicating whether each element is missing or not, and then the sum() met which calculates the sum of True values (i.e., missing values) along each column.
[19]: [19]:	text 0 class 0 dtype: int64 The code below randomly shuffles the rows of the DataFrame df using the sample() method with the frac parameter set to 1. By setting frac to 1, it means that the entire DataFrame will be sampled (i.e., all rows will be included), and since no replacement is specified, each row will be selected exactly once. Shuffling the rows of the DataFrame can be useful for various purposes, such as randomizing the order of data for training machine learning models or ensuring that the data is not biased by any particular order. It helps in creating more robust and the class of the data is not biased by any particular order. It helps in creating more robust and the class of the data is not biased by any particular order.
	reliable models by preventing any patterns or biases that may exist due to the original order of the data. df = df.sample(frac = 1) df.head()
	Leftist agitators showed up tonight at the ope 0 21835 The laws don't apply to the Clinton's they re 0 The code below resets the index of the DataFrame df to the default integer index and drops the original index column. The reset_index() method resets the index of the DataFrame, and the inplace=True parameter ensures that the changes are applied to the DataFrame df directly. Then, the drop() method is used to remove the original index column, specified by label "index", along the columns (axis=1). Again, inplace=True ensures that the changes are applied to the DataFrame df directly. This process ensures that the DataFrame df has a clean, sequential integer index without retaining any previous incinformation, which can be helpful for further analysis or modeling tasks.
_	<pre>df.drop(["index"], axis = 1, inplace = True) df.columns Index(['text', 'class'], dtype='object') df.head()</pre>
	1 The results of this town s gun ownership manda 0 2 Proving that the biggest obstacle to his presi 0 3 Leftist agitators showed up tonight at the ope 0 4 The laws don't apply to the Clinton's they re 0 The wordopt function preprocesses text data by applying various cleaning operations. It first converts all text to lowercase for consistency. Then, it removes content within square brackets, non-word characters, URLs, HTML tags, punctuation, new characters, and words containing numbers. This comprehensive cleaning process aims to standardize the text data and remove any irrelevant or noisy elements, preparing it for further analysis or modeling in natural language processing tasks.
[25]:	text = text.lower() text = re.sub('\[.*?\]', '', text) text = re.sub("\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\
	preprocessing step ensures that the text data is cleaned and standardized according to the operations defined in the wordopt function, making it suitable for downstream analysis or modeling tasks in natural language processing. If "text"] = df["text"] apply(wordopt) The code below sets the X and y features in preparation for further processing using machine learning models x = df["text"]
[28]: [29]:	x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.25) The provided code below is part of a machine learning project focused on text data preprocessing and transformation using TF-IDF (Term Frequency-Inverse Document Frequency) vectorization. The goal is to convert textual data into a numerical format suitable for machine learning models.
[30]:	In this segment of the machine learning project, a Logistic Regression model from the scikit-learn library is being utilized to perform supervised learning. The LogisticRegression object, denoted as LR, is first instantiated, representing the model to trained. The fit method is then called with the previously transformed training data (xv_train) and the corresponding target labels (y_train). This step involves training the Logistic Regression model by finding the optimal weights that best map the TIDF features of the training data to their respective target labels, thus enabling the model to make predictions on new, unseen data. from sklearn.linear_model import LogisticRegression LR = LogisticRegression() LR.fit(xv_train, y_train) v LogisticRegression() LogisticRegression()
	pred_lr=LR.predict(xv_test) The score method is used to evaluate the performance of the trained Logistic Regression model by calculating its accuracy on the test data (xv_test) against the true labels (y_test). LR.score(xv_test, y_test)
[33]:	print(classification_report(y_test, pred_lr)) precision recall f1-score support 0 0.99 0.99 0.99 5891 1 0.98 0.99 0.99 5329 accuracy 0.99 11220 macro avg 0.99 0.99 0.99 11220 weighted avg 0.99 0.99 0.99 11220 The code below calculates and visualizes the Receiver Operating Characteristic (ROC) curve, a crucial metric for evaluating the performance of a binary classification model. By computing the false positive rate (FPR) and true positive rate (TPR)
[34]:	the roc_curve function, and the area under the curve (AUC) with the auc function, the code assesses the model's ability to distinguish between the two classes. The ROC curve plot, generated with matplotlib, shows the trade-off between the FPR TPR at various threshold settings, with the AUC value providing a single scalar to summarize the model's overall performance. An AUC closer to 1 indicates a better performing model. ###################################
	plt.plot ([0, 1], [0, 1], color='navy', lw=lw, linestyle='') plt.xlim([0.0, 1.0]) plt.ylim([0.0, 1.05]) plt.xlabel('False Positive Rate') plt.ylabel('True Positive Rate') plt.title('Receiver Operating Characteristic (ROC)') plt.legend(loc = "lower right") plt.show() Receiver Operating Characteristic (ROC) 1.0
	Lace Bositive Bate 0.8 - 0.8 - 0.8 - 0.8 - 0.8 - 0.8 - 0.8 - 0.9
	In this part of the machine learning project, a Decision Tree Classifier from the scikit-learn library is being employed for supervised learning. A DecisionTreeClassifier object, denoted as DT, is instantiated and subsequently trained using the fit met with the TF-IDF transformed training data (xv_train) and the corresponding target labels (y_train). This process involves building a decision tree model that learns to map the input features to the target labels by recursively splitting the data based feature values, thus enabling the model to make predictions on new data.
[35]: [35]:	<pre>pred_dt = DT.predict(xv_test)</pre> from sklearn.tree import DecisionTreeClassifier DT = DecisionTreeClassifier() DT.fit(xv_train, y_train) v DecisionTreeClassifier() pred_dt = DT.predict(xv_test)
[38]:	0.995632798573975 print(classification_report(y_test, pred_dt)) precision recall f1-score support 0 0.99 1.00 1.00 5891 1 1.00 0.99 1.00 5329 accuracy
	This code segment is used to compute and visualize the Receiver Operating Characteristic (ROC) curve for the predictions made by a Decision Tree Classifier. The false positive rate (FPR) and true positive rate (TPR) are calculated using the roc_curve function based on the true labels (y_test) and the predicted labels (pred_dt). The area under the curve (AUC) is then computed with the auc function, providing a summary metric of the classifier's performance. The ROC curve is plotted using matplotlib, with the FPR on the x-axis and the TPR on the y-axis, illustrating the trade-offs between sensitivity and specificity at various threshold levels. The AUC value is included in the plot's legend to indicate the overall discriminative abilithe classifier, with values closer to 1 indicating better performance. #Compute ROC curve and ROC area fpr, tpr, _ = roc_curve(y_test, pred_dt) roc_auc = auc(fpr, tpr)
	<pre># Plot the ROC Curve plt.figure() lw = 2 plt.plot (fpr, tpr, color= "darkorange",</pre>
	Receiver Operating Characteristic (ROC) 0.8 -
	0.2 - ROC curve (area = 1.00) 0.0 0.2 0.4 0.6 0.8 1.0 False Positive Rate
[39]: t[39]:	<pre>GBC = GradientBoostingClassifier(random_state=0) GBC.fit(xv_train, y_train)</pre>
	GBC.score(xv_test, y_test) 0.9957219251336898 print(classification_report(y_test, pred_gbc)) precision recall f1-score support
	0 1.00 0.99 1.00 5891 1 0.99 1.00 1.00 5329 accuracy 1.00 1.220 macro avg 1.00 1.00 1.00 11220 weighted avg 1.00 1.00 1.00 1.00 11220 This code segment is designed to compute and visualize the Receiver Operating Characteristic (ROC) curve for the predictions made by a Gradient Boosting Classifier. The roc_curve function is used to calculate the false positive rate (FPR) and to positive rate (TPR) based on the true labels (y_test) and the predicted labels (pred_gbc). The area under the curve (AUC) is then computed using the auc function, providing a single value summary of the model's overall performance, where a high AUC indicates better discriminative ability. The ROC curve is plotted using matplotlib, with the FPR on the x-axis and the TPR on the y-axis, showing the trade-off between the true positive rate and false positive rate at various threshold levels. The plot includes a diagonal line representing a random classifier, and the ROC curve is annotated with the AUC value in the legend, helping to visually assess the model's performance.
[49]:	<pre>#Compute ROC curve and ROC area fpr, tpr, _ = roc_curve(y_test,pred_gbc) roc_auc = auc(fpr, tpr) # Plot the ROC Curve plt.figure() lw = 2 plt.plot (fpr, tpr, color= "darkorange",</pre>
	plt.xlabel('False Positive Rate') plt.title('Receiver Operating Characteristic (ROC)') plt.legend(loc = "lower right") plt.show() Receiver Operating Characteristic (ROC) 1.0
	0.6 - O.2 - O.2 - O.2 - O.2 - O.3 - O.4 - O.2 - O.3 -
	ROC curve (area = 1.00) 0.0 0.2 0.4 0.6 0.8 1.0 False Positive Rate In this part of the machine learning project, a Random Forest Classifier from the scikit-learn library is employed for supervised learning. The RandomForestClassifier object, denoted as RFC, is instantiated with a specified random state to ensure reproducibility. The fit method is called with the TF-IDF transformed training data (xv_train) and the corresponding target labels (y_train). This process involves training the Random Forest model, which is an ensemble method that builds multiple decision trees and merges their results to improve predictive accuracy and control overfitting. Each tree is trained on a random subset of the data and features, and the final prediction is made by averaging the outputs (regression) or majority voti (classification) from all the individual trees.
[48]: [48]:	RandomForestClassifier(random_state=0)
	RFC.score(xv_test, y_test) 0.9888591800356507 print(classification_report(y_test, pred_rfc)) precision recall f1-score support 0 0.99 0.99 0.99 5891 1 0.99 0.98 0.99 5329 accuracy
	macro avg 0.99 0.99 0.99 11220 This code segment is aimed at computing and visualizing the Receiver Operating Characteristic (ROC) curve for the predictions made by a Random Forest Classifier. It utilizes the roc_curve function to calculate the false positive rate (FPR) and trepositive rate (TPR) based on the true labels (y_test) and the predicted labels (pred_rfc). The area under the curve (AUC) is then computed using the auc function, serving as a summary metric for the classifier's overall performance. The ROC curve plotted using matplotlib, with the FPR on the x-axis and the TPR on the y-axis, illustrating the trade-off between sensitivity and specificity at various threshold levels. The plot includes a diagonal line representing a random classifier, and the ROC is annotated with the AUC value in the legend, aiding in the visual assessment of the model's performance ##Compute ROC curve and ROC area fpr, tpr, _ = roc_curve(y_test, pred_rfc) roc_auc = auc(fpr, tpr) ##Plot the ROC Curve
	<pre># Flot the Rot Curve plt.figure() lw = 2 plt.plot (fpr, tpr, color= "darkorange",</pre>
	Receiver Operating Characteristic (ROC) 1.0 - 0.8 -
	0.2 - ROC curve (area = 0.99) 0.0 0.2 0.4 0.6 0.8 1.0 False Positive Rate
	This function, manual_testing(news), facilitates manual testing of news articles by predicting their authenticity using trained machine learning models. It takes a news article as input, preprocesses it, and then uses the trained models (Logistic Regression, Decision Tree, Gradient Boosting, and Random Forest) to predict whether the news is fake or not. The function returns the predictions made by each model for the provided news article. The output_lable(n) function maps the model predictions (n) to meaningful labels indicating whether the news is classified as fake or not. def output_lable(n): if n == 0: return "This News is Fake" elif n == 1: return "NOT Fake News"
[54]:	<pre>def manual_testing(news): testing_news = {"text":[news]} new_def_test = pd.DataFrame(testing_news) new_def_test["text"] = new_def_test["text"].apply(wordopt) new_x_test = new_def_test["text"] new_xv_test = vectorization.transform(new_x_test) pred_LR = LR.predict(new_xv_test) pred_DT = DT.predict(new_xv_test) pred_GBC = GBC.predict(new_xv_test) pred_GBC = GBC.predict(new_xv_test) pred_RFC = RFC.predict(new_xv_test) return print("\n\nLR Prediction: {} \nGBC Prediction: {} \nRFC Prediction: {} \nFC Prediction: {} \nGBC Pred_LR[0]),</pre>
[54]:	
[55]:	Feel free to input a news article, and I'll run it through the trained models to predict its authenticity. news = str(input()) manual_testing(news) LR Prediction: NOT Fake News DT Prediction: NOT Fake News GBC Prediction: NOT Fake News RFC Prediction: NOT Fake News RFC Prediction: NOT Fake News
[55]:	Feel free to input a news article, and I'll run it through the trained models to predict its authenticity. news = str(input()) manual_testing(news) LR Prediction: NOT Fake News DT Prediction: NOT Fake News GBC Prediction: NOT Fake News
[55]: [57]:	Feel free to input a news article, and I'll run it through the trained models to predict its authenticity. Inews = str(input()) manual_testing(news) LR Prediction: NOT Fake News OT Prediction: NOT Fake News GEC Prediction: NOT Fake News RFC Prediction: NOT Fake News RFC Prediction: NOT Fake News Inews = str(input()) manual_testing(news) LR Prediction: This News is Fake GEC Prediction: This News is Fake RFC Prediction: This News is Fake