]: 4]: 8]:	Use pandas to read data as a dataframe called df. import io from google.colab import drive drive.mount('/content/drive') Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
9]:	import pandas as pd path=('/OnlineNewsPopularity.csv') df=pd.read_csv(path) df.head() url timedelta n_tokens_title n_tokens_content n_unique_tokens n_non_stop_words n_non_stop_unique_tokens num_hrefs num_self_hrefs num_ http://mashable.com/2013/01/07/amazon-instant 731 12 219 0.663594 1.0 0.815385 4 2 http://mashable.com/2013/01/07/ap-samsung-spon 731 9 255 0.604743 1.0 0.791946 3 1 http://mashable.com/2013/01/07/apple-40- 731 9 255 0.604743 1.0 0.663866 3 1
⊙]: ⊙]:	ul describe()
]:	50% 339.000000 10.000000 409.000000 0.539226 1.000000 0.690476 8.000000 3.000000 1.000000 0.000000 75% 542.000000 12.000000 716.000000 0.608696 1.000000 0.754630 14.000000 4.000000 4.000000 1.000000 max 731.000000 23.000000 8474.000000 701.000000 1042.000000 650.000000 304.000000 116.000000 128.000000 91.000000 Attribute information: Number of Attributes: 61 including target column shares
	9. mum. just_Nitroto of iningues 11. mum_videor. Number of videos 12. mum_videor. Number of videos 12. mum_videor. Number of videos 13. mum. Jeywords. Number of videos 14. mum_videor. Number of videos 15. mum. Jeywords. Number of videos 16. mum. Jeywords. Number of videos length of the words in the content 16. mum. Jeywords. Number of videovide in the resistation 16. data channel js. centeriamment. Setion Channel Secience 17. data channel js. centeriamment. Setion Channel Secience 17. data channel js. centeriamment. Setion Channel Secience 17. data channel js. section 18 tolas channel Secience 17. data channel js. section 18 tolas channel Secience 17. data channel js. section 18 tolas channel Secience 18. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data channel js. videol js. data channel Secience 19. data videological js. js. data videological js. data videological js. js. data videological js. js. data videological js. js. data videological js. js. data videological
1]: 1]:	<pre>array(['url', 'timedelta', 'n_tokens_title', 'n_tokens_content',</pre>
2]:	timedelta False n_tokens_title False n_tokens_content False n_unique_tokens False n_unique_tokens False title_subjectivity False abs_title_subjectivity False abs_title_sentiment_polarity False abs_title_sentiment_polarity False thength: 61, dtype: bool Q1. Looking at the data above what are your first thoughts about quality of data and modeling?
3]:	df.describe()
5]:	Correlation of Features Lets find the correlation among features (very important for successfull modelling) We will plot correlation matrix using Plotty HeatMap data = [go.Heatmap(z= df.astype(float).corr().values, x=df.columns.values, y=df.columns.values, colorscale='Viridis', reversescale = False, text = True ,
	<pre>text = True,</pre>
	<pre>ValueError</pre>
	<pre>64745</pre>
	<pre>3770</pre>
6]:	ValueError: Invalid value of type 'builtins.bool' received for the 'text' property of heatmap Received value: True The 'text' property is an array that may be specified as a tuple, list, numpy array, or pandas Series Start Of Part 2 Apply various modeling techniques Common imports for modeling from sklearn.linear_model import LinearRegression
7]:	random_state = 101 Function to split the data
9]:	Function to perform Simple Linear Regression def perf_linear_regression(df_data, standscalar=False): X,y,X_train, X_test, y_train, y_test = get_data(df_data) if(standscalar): print ("Regression after applying StandardScaler") X_scaler = StandardScaler() X_train = X_scaler.fit_transform(X_train) X_test = X_scaler.transform(X_test) # y_scaler = StandardScaler() # y_train = y_scaler.fit_transform(y_train[:, None])[:, 0] # y_test = y_scaler.fit_transform(y_test[:, None])[:, 0] # y_train = y_scaler.fit_transform(y_test[:, None])[:, 0] # y_test = y_scaler.fit_transform(y_test[:, None])[:, 0] # y_test = y_scaler.fit_transform(y_test[:, None])[:, 0]
∍]:	<pre># End of If for StandardScaler lm = LinearRegression() lm.fit(X_train,y_train) y_pred = lm.predict(X_test) y_train_pred = lm.predict(X_train) print ("RMSE on Train Data is: {:.2f} :".format(np.sqrt(metrics.mean_squared_error(y_train, y_train_pred)))) print ("RMSE on Test Data is: {:.2f} :".format(np.sqrt(metrics.mean_squared_error(y_test, y_pred)))) return X,y,X_train, X_test, y_train, y_test Approach 1: Linear Regression Q3: What is the RMSE on Test Data for Linear Regression?</pre>
4]:	40000 35000 - 25000 - 20000 - 15000 - 10000 - 0 200000 400000 600000 800000 shares
	2000 1500 1000 500 0 100k 200k 300k 400k 500k 600k 700k 800k Q6: Which data based on shares count should be dropped from dataset ? You can use just vizual cue here , we will have more formal approach in Capstone
6]:	<pre># # Lets get a scatter plot of shares and respective counts y_unique = y.unique() sns.regplot(y_unique, y_unique, data=y, scatter=True)</pre>
7]:	<pre>def filter_threshold_data(dr_copy, threshold, column_name). df_adjusted = df_copy[df_copy[column_name] <= threshold] print (" Original Data Count:", len(df_copy)) print (" After Adjusting Data Count:", len(df_adjusted)) return df_adjusted</pre>
9]: 9]:	RMSE on Train Data is: 5868.35 : RMSE on Test Data is: 5756.43 :
]:	Approach 4: Lets try improve the model by using DecisionTreeRegressor for regression Q8. Why we should use DecisionTree?
1]:	2. Nonlinear relationships between parameters do not affect tree performance 3. It is easy to interpret from sklearn.tree import DecisionTreeRegressor
3]:	<pre>X,y,X_train, X_test, y_train, y_test = get_data(df_adjusted) max_depth_count = len(X.columns) depth_range = [1,2,5,10,15,20,25,30,35,40,max_depth_count] print ("Starting Decision Tree Regression:") for depth_in depth_range:</pre>
	decision_tree_regression(depth,X_train,y_train,X_test,y_test) # End of decision tree testing Original Data Count: 39644 After Adjusting Data Count: 38826 Starting Decision Tree Regression: Tree Max Depth is: 1 RMSE on Train Data is: 2813.15: RMSE on Train Data is: 2720.21: Tree Max Depth is: 2 RMSE on Train Data is: 2793.18: RMSE on Train Data is: 2791.06: Tree Max Depth is: 5 RMSE on Train Data is: 2745.93: RMSE on Train Data is: 2745.93: RMSE on Train Data is: 2698.59: Tree Max Depth is: 10 RMSE on Train Data is: 2443.16:
	RMSE on Test Data is: 2962.95 : Tree Max Depth is: 15 RMSE on Train Data is: 1812.13 : RMSE on Test Data is: 3471.15 : Tree Max Depth is: 20 RMSE on Train Data is: 1031.13 : RMSE on Train Data is: 3844.75 : Tree Max Depth is: 25 RMSE on Train Data is: 486.06 : RMSE on Train Data is: 4014.42 : Tree Max Depth is: 30 RMSE on Train Data is: 4914.42 : Tree Max Depth is: 30 RMSE on Train Data is: 3986.76 : Tree Max Depth is: 35 RMSE on Test Data is: 58.94 : RMSE on Test Data is: 4089.70 : Tree Max Depth is: 40 RMSE on Train Data is: 4089.70 : Tree Max Depth is: 40 RMSE on Train Data is: 11.71 :
4]:	RMSE on Test Data is: 4035.47: Tree Max Depth is: 58 RMSE on Train Data is: 0.00: RMSE on Test Data is: 4059.23: Approach 5: Combining multiple techniques We will use the following 1. Use MinMaxScaler to Scale the Data 2. Use Principal Component Analysis to restrict the no. of dimensions 3. Use Ensemble technique AdaBoost With DecisionTreeRegressor to further improve the model
5]:	The MinMaxScaler is the most famous scaling algorithm, and follows the following formula for each feature: xifmin(x)/(max(x)fmin(x)) It essentially shrinks the range such that the range is now between 0 and 1 (or -1 to 1 if there are negative values). This scaler works better for cases in which the standard scaler might no so well. If the distribution is not Gaussian or the standard deviation is very small, the min-max scaler works better. Applying MinMaxScaler Q12. While applying min max scaling to normalize your features, do you apply min max scaling on the entire dataset before splitting it into training, validation and test data? -Favourite Question df_adjusted = filter_threshold_data(df, 20000, "shares")
]:	Get the PCA for variance upto 95 % Q14. How many features are there after variance is limited to 95% ?
7]:	<pre>X_train_normreduced = pd.DataFrame(pca.transform(X_train_norm)) X_train_normreduced = X_train_normreduced.loc[:,pca.explained_variance_ratiocumsum()<0.95] print (X_train_normreduced.shape) (27178, 22) We got 22 features for around 95 % variance. Lets use this info to do modeling further</pre>
	<pre>X_test_norm = min_max_norm.transform(X_test) X_train_norm = pca.fit_transform(X_train_norm) X_test_norm = pca.transform(X_test_norm) X_train_norm = pd.DataFrame(X_train_norm) X_test_norm = pd.DataFrame(X_test_norm) Original Data Count: 39644 After Adjusting Data Count: 38826 # Transforming y . We can choose different log base also such as 2, 10 etc y_train_log = np.log(y_train) y_test_log = np.log(y_test) Q15. Why use AdaBoostRegressor?</pre>
9]:	<pre>Generic function to use AdaBoostRegression for different dataset and estimator values def ada_boost_regression(X_train_norm,y_train_log,X_test_norm,y_test_log,n_estimators): AdaDecision = AdaBoostRegressor(DecisionTreeRegressor(max_depth=4, min_samples_leaf= 5, min_samples_split= 5),</pre>
	Q16. For which n_estimators values do you see lowest RMSE for Test Data? Try changing the values in n_estimators_list and see what RMSE you get n_estimators_list = [10,20,30,40,50,100,125,150] for n_estimator in n_estimators_list: ada_boost_regression(X_train_norm,y_train_log,X_test_norm,y_test_log,n_estimator) Estimator Count in AdaBoostRegressor: 10 RMSE on Train Data is: 2870.27 : RMSE on Test Data is: 2777.38 : Estimator Count in AdaBoostRegressor: 20 RMSE on Train Data is: 2882.41 : RMSE on Test Data is: 2793.28 : Estimator Count in AdaBoostRegressor: 30
∍]:	RMSE on Train Data is: 2890.05 : RMSE on Test Data is: 2801.35 : Estimator Count in AdaBoostRegressor: 40 RMSE on Train Data is: 2905.45 : RMSE on Test Data is: 2815.61 : Estimator Count in AdaBoostRegressor: 50 RMSE on Train Data is: 2927.80 :
∍]:	RMSE on Test Data is: 2836.73 : Estimator Count in AdaBoostRegressor: 100 RMSE on Train Data is: 3019.68 : RMSE on Test Data is: 2932.70 : Estimator Count in AdaBoostRegressor: 125 RMSE on Train Data is: 3001.36 : RMSE on Test Data is: 2911.01 : Estimator Count in AdaBoostRegressor: 150 RMSE on Train Data is: 3051.36 : RMSE on Train Data is: 2962.80 :
9]:	Estimator Count in AdaBoostRegressor: 190 RMSE on Train Data is: 3919.68: RMSE on Test Data is: 2932.70: Estimator Count in AdaBoostRegressor: 125 RMSE on Train Data is: 3001.36: RMSE on Test Data is: 2001.36: RMSE on Test Data is: 2001.3