In [1]:	PROBLEM DEFINITION # The research questions we would work upon are: # 1.Did PERCEPTION about covid vaccines played any role in cases worldwide?
	# 1.Did PERCEPTION about Covid Vaccines played any fole in cases worldwide? # 2.Which continent's population was more inclined towards the perception of not taking a vaccine? # The data we intend to use is based on the perception of the people about Covid Vaccines. # There has always been a conflict in the population about the vaccines. # Some had positive perceptions and the others, negative. We have taken into consideration the following scenario: # The risks of COVID-19 disease are greater than the risks of the vaccine. # The people those who are in favour of this statement took the vaccines, and those who aren't, did the opposite.
In [2]:	DATA DESCRIPTION #To get to the bottom of the situation above, my approach includes usage of two datasets: #(1) owid-covid-data.csv (renamed to Covid Data New).
In [3]:	# It contains the data about Covid-19 and its effects on the world (deaths, cases, etc) #(2) Raw_Data_Vaccine_hesitancy_globally.csv (renamed to Perception) # This one contains the data about the public survey held in June of 2021. This specifically contains how people were perceptive about Covid vaccines. Some people had positive pe ## To import the libraries
	<pre>import seaborn as sns import pandas as pd import numpy as np import matplotlib.pyplot as plt %matplotlib inline from collections import Counter from imblearn.over_sampling import RandomOverSampler</pre>
	from sklearn.linear_model import LogisticRegression from sklearn.model_selection import train_test_split from sklearn.metrics import classification_report from sklearn.metrics import confusion_matrix from sklearn.metrics import accuracy_score
In [4]:	<pre>from sklearn.metrics import precision_score from sklearn.metrics import recall_score from sklearn.metrics import f1_score import warnings warnings.filterwarnings('ignore')</pre>
In [5]:	#Importing the datasets. Perceptions = pd.read_csv('D:\Group Projects 4070\Group Project New\Perception.csv') Covid=pd.read_csv('Covid Data New.csv') ##Renaming the column for ease of use Perceptions.rename(columns={'Perception/The risks of Covid-19 disease are greater then the risks of vaccine': 'Perception'}, inplace=True)
In [6]: Out[6]:	Perceptions.head() Country Perception India Somewhat agree
	 India Somewhat agree India Somewhat agree India Somewhat agree
In [7]: Out[7]:	Covid head Fazi Covid head Fazi Country date Years total_cases total_deaths new_cases_smoothed total_deaths.1 new_deaths n
	2 BRA South America Brazil 28/02/2020 2020 1.0 NaN NaN NaN NaN 2.2 75.88 0.765 215313504 3 BRA South America Brazil 29/02/2020 2020 2.0 NaN NaN NaN NaN 2.2 75.88 0.765 215313504 4 BRA South America Brazil 01/03/2020 2020 2.0 NaN NaN NaN NaN 2.2 75.88 0.765 215313504
	5 rows × 70 columns FEATURE ENGINEERING
In [8]:	## The perception column has 5 Categories # Strongly agree, # Somewhat agree, # Unsure/no opinion, # Somewhat disagree, # Strongly disagree.
In [11]: In [12]: In [13]:	## We've removed the 'unsure/nop opinion', as we don't know whether they took the vaccines or not, and their number was less. Perceptions=Perceptions['Perception']!='Unsure/no opinion'] ## We've further divided the groups- # Strongly Agree + Somewhat Agree = TOOK VACCINES. # Strongly Disagree + Somewhat Disagree = DID NOT TAKE VACCINES
In [14]: In [15]:	Perceptions['Perception']=Perceptions['Perception'].replace({'Somewhat disagree':'Strongly disagree'}) Perceptions['Perception']=Perceptions['Perception'].replace({'Somewhat agree':'Strongly agree'}) ## Keeping only the necessary columns Covid_columns=Covid[['continent','Country','date','Years','total_cases','total_deaths','new_cases']]
	<pre>## Merging the two datasets ## Combining both datasets on common columns. ## Note that They can be joined basis one column 'COUNTRY'. Perceptions_Covid = pd.merge(Perceptions, Covid_columns, on='Country') Perceptions_Covid</pre>
Out[15]:	Country Perception continent date Years total_cases total_deaths new_cases 0 India Strongly agree Asia 30/01/2020 2020 1.0 NaN 1.0 1 India Strongly agree Asia 01/02/2020 2020 1.0 NaN 0.0 2 India Strongly agree Asia 02/02/2020 2020 1.0 NaN 0.0
	3 India Strongly agree Asia 02/02/2020 2020 2.0 NaN 1.0 4 India Strongly agree Asia 03/02/2020 2020 3.0 NaN 1.0 23191073 Brazil Strongly agree South America 04/02/2023 2023 36867401.0 697365.0 1118.0 23191074 Brazil Strongly agree South America 05/02/2023 2023 36867401.0 697365.0 0.0
	23191075 Brazil Strongly agree South America 06/02/2023 2023 36878774.0 697439.0 11373.0 23191076 Brazil Strongly agree South America 07/02/2023 2023 36887991.0 697533.0 9217.0 23191077 Brazil Strongly agree South America 08/02/2023 2023 36897683.0 697626.0 9692.0 23191078 rows × 8 columns
In [16]:	<pre>## Data Mapping - Mapping the data such that # STRONGLY AGREE = 1 # STRONGLY DISAGREE = 0 Perceptions_Covid['Perception']=Perceptions_Covid['Perception'].map({'Strongly agree':1,'Strongly disagree':0}) Perceptions_Covid</pre>
Out[16]:	Country Perception continent date Years total_cases total_deaths new_cases 0 India 1 Asia 30/01/2020 2020 1.0 NaN 1.0 1 India 1 Asia 31/01/2020 2020 1.0 NaN 0.0 2 India 1 Asia 01/02/2020 2020 1.0 NaN 0.0 3 India 1 Asia 02/02/2020 2020 2.0 NaN 1.0
	4 India 1 Asia 03/02/2020 2020 3.0 NaN 1.0
	23191076 Brazil 1 South America 07/02/2023 2023 36887991.0 697533.0 9217.0 23191077 Brazil 1 South America 08/02/2023 2023 36897683.0 697626.0 9692.0 23191078 rows × 8 columns ## Taking care of the NULL VALUES
<pre>In [17]: Out[17]:</pre>	<pre>## Dropping the Null Values Perceptions_Covid.isnull().sum().sum() no_null=Perceptions_Covid.dropna() no_null.isnull().sum().sum()</pre>
In [18]: Out[18]:	<pre>sns.countplot(x='Perception', data=no_null) <axes: ,="" xlabel="Perception" ylabel="count"></axes:></pre>
	2.00 - 1.75 - 1.50 -
	1.25 - 1.00 - 0.75 -
	0.50 - 0.25 - 0.00 1
In [20]:	# As the above the graph depicts that the majority of people were in favour of TAKING THE VACCINES # and very less population were against the idea. sns.countplot(x='continent', hue='Perception', data=no_null)
Out[20]:	<pre><axes: ,="" xlabel="continent" ylabel="count"> le6 6- Perception</axes:></pre>
	5 - 4 - ting 3 -
In [21]:	Asia Africa South America Europe North America continent # The above graph shows that mostly the group of population not in favour of taking the vaccines is from
In [22]: Out[22]:	<pre># Asia and Europe. sns.stripplot(x='Perception', y='new_cases', data=no_null) <axes: ,="" xlabel="Perception" ylabel="new_cases"> le6</axes:></pre>
	1.2 - 1.0 -
	0.8 - 0.6 - 0.4 -
In [23]: In [24]:	# With the above graph, it is clear that the people who weren't willing to take vaccines were also leading to the increase ## in newly prevailed cases of Covid. Thereby increasing the cases count
In [25]: Out[25]:	DATA MODELLING no_null['Perception'].value_counts() 1 20388174 0 2136000
In [26]: In [59]:	Name: Perception, dtype: int64 df=no_null # Calculating the correlation amongs the features # We can get corelation between features based on values ranging (-1, 1). Near to -1 or 1 shows strong corelation # in OPPOSITE directions. # Near to zero means lesser corelation.
In [28]:	<pre># Visualization to check the features corelation and multicollinearity sns.heatmap(correlation, annot=True, cmap='coolwarm') plt.title('Feature Correlation Heatmap') plt.show()</pre>
	Perception - 1 0.001 0.0013 0.07 -0.013 - 0.8
	total_cases - 0.0013
	total_deaths - 0.07
	Perception - Years - total_deaths - total_deaths -
In [55]:	## DATA BALANCING # Our data needs to be balanced because against the hughe values fo 1's, we have less values of 0's. # We would do the Over Sampling technique. The reason behind this is because Under Sampling technique reduces the # size of the sampled data as per the lower value. I did both the sampling technique. UNDER & OVER.
	I'm only showing the output with Undersampling because doing that also gets us lot of data to utilize. The oversampling also lead to approximately the same accuracy of the model.
In [31]: In [32]:	<pre>X=df #Independent features y=df['Perception'] # Dependent Features # Under Sampling technique for balancing the Data. from imblearn.under_sampling import RandomUnderSampler rus = RandomUnderSampler(random_state=0)</pre>
In [33]: In [34]:	<pre>X_resampled, y_resampled = rus.fit_resample(X,y) print(sorted(Counter(y_resampled).items()),y_resampled.shape) [(0, 2136000), (1, 2136000)] (4272000,) A=X_resampled A.drop(['Country','continent','date','Years'], axis=1, inplace=True)</pre>
out[34]:	Perception total_cases total_deaths new_cases
	3 0 359.0 4.0 167.0 4 0 670.0 9.0 311.0
	4271997
In [35]: Out[35]:	B=y_resampled B 0 0 1 0 2 0 3 0
	3 0 4 0 4271995 1 4271996 1 4271997 1 4271998 1 4271999 1 Name: Perception, Length: 4272000, dtype: int64
In [56]:	# We would apply LOGISTIC REGRESSION because our Depedent Feature is Categorical. # We would divide our dataset into 70:30 Ratio. # 70% - Train Data # 30%-Test Data # Random State = 1 means that whenever the code is run, our data would be split into small but same datasets.
In [37]: In [38]:	<pre>from sklearn.linear_model import LogisticRegression classifier=LogisticRegression() # Splitting the data into test and train. from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(A, B, test_size=0.3, random_state=1)</pre>
In [39]: Out[39]:	<pre>classifier.fit(X_train,y_train) v LogisticRegression LogisticRegression()</pre>
<pre>In [40]: In [41]: Out[41]:</pre>	<pre>predictions=classifier.predict(X_test) from sklearn.metrics import classification_report classification_report(y_test, predictions) '</pre>
In [45]:	<pre>class_report = classification_report(y_test, predictions) print("Classification Report:") print(class_report) Classification Report:</pre>
Ŧ	0 0.66 0.58 0.62 640245 1 0.62 0.71 0.66 641355 accuracy 0.64 1281600 macro avg 0.64 0.64 0.64 1281600 weighted avg 0.64 0.64 0.64 1281600 from sklearn.metrics import confusion_matrix
In [42]: In [43]: In [44]:	<pre>from sklearn.metrics import confusion_matrix from sklearn.metrics import accuracy_score #Calculating precision, recall and f1 score from sklearn.metrics import precision_score from sklearn.metrics import recall_score from sklearn.metrics import recall_score from sklearn.metrics import f1 score</pre>
In [49]:	<pre>from sklearn.metrics import recall_score # ACCURACY # PRECISION # RECALL # F1-SCORE accuracy = accuracy_score(y_test, predictions) print('Precision:', 100* precision_score(y_test, predictions)) print('Recall:', 100* recall_score(y_test, predictions))</pre>
In [50]: Out[50]:	<pre>from sklearn import metrics cm = metrics.confusion_matrix(y_test, predictions) cm array([[368603, 271642],</pre>
In [51]:	<pre>plt.figure(figsize=(2,2)) sns.heatmap(cm, annot=True, fmt = '1.0f', linewidth = 1.0,</pre>
	sales o - 368603 271642
	The state of the s
In [52]:	# We hereby conclude that the accuracy of the model is 64%. This accuracy is not that bad, yet not so good. # As our data is purely based on the available data. # The Sampling technique has been used here because the accuracy without using Random Sampler was 86%.
In ۲۹٦	# Which shows that our model was overfitting. Thus to avoid the situation, we took the path of Sampling. # Further, the values of No's was much less then values of Yes's, which lead to incorrect assumption. CONCLUSION ## ANSWER 1
. [1]:	## ANSWER 1 ## Continents Asia and Europe were the ones with people who weren't willing to take the vaccines, as per the data. ## ANSWER 2 ## the majority of people were in favour of vaccines and very less were against. ## BONUS, ## The group of population who weren not willing to take vaccines, were not much, but they were also playing a role ## in increasing the number of cases worldwide.
In [53]:	
In [54]:	Cons ## Our dataset is huge with many missing values, which hamper the predictions at the first place. ## Also there is so much imbalance in our predicting feature. ## By doing Undersampling and Oversampling, we got Accuracy of 64.08% and 64.10%. This shows that our data isn't perfect,
In []:	## we wouldn't get better predictions if we try either method.