į	Description: Honour killings in Pakistan are known locally as karo-kari (Urdu: シュー・シュー・(Pakistan has the highest number of documented and estimated honour killings per capita of any country in the world; about one-fifth of the world's honour killings are performed in Pakistan Question 1
	The dataset contains reported cases from 2011 to 2019. Your task is to Visualize the data in best possible manner. For example showing a plot of districts with most number of cases. Not that this is one example, you need to be creative. Answer
	First we converted the file honour-killing-reported-cases-2011-2019(1).xslx into .csv file. We can work with csv files using pandas library. Step 1: Getting the Data
13]:	<pre>import pandas as pd df=pd.read_csv("honour-killing-reported-cases-2011-2019(1).csv") # Reading the csv File => Making sure we got the data df.head()</pre>
13]:	Districts 2011 2012 2013 2014 2015 2016 2017 2018 2019 0 Lahore 10 15 11 24 13 7 9 10 12 1 Sheikhupura 14 13 5 17 11 11 9 10 5
	2 Nankana Sahib 13 10 10 15 14 0 2 1 7 3 Kasur 10 13 10 10 13 13 13 8 7 6 4 Region Total 37 36 25 42 38 24 19 18 18
14]:	Step 2: Exploring the data # Checking the types of data print(df.dtypes) df.info() df.dtypes()
	<pre>df.describe() Districts object 2011 int64 2012 int64 2013 int64 2014 int64</pre>
	2015 int64 2016 int64 2017 int64 2018 int64 2019 int64 dtype: object
	<pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 46 entries, 0 to 45 Data columns (total 10 columns): # Column Non-Null Count Dtype</class></pre>
	2 2012
14]:	9 2019
	mean 23.521739 23.543478 25.065217 25.826087 21.108696 16.021739 11.608696 15.695652 12.586957 std 53.626999 53.757359 57.448490 59.438971 48.782615 36.895341 26.767707 35.767067 29.178817 min 1.000000 0.000000 1.000000 1.000000 0.000000 0.000000 0.000000 0.000000 25% 5.000000 7.000000 6.250000 4.000000 3.250000 2.000000 4.000000 3.000000
	50% 11.000000 11.000000 10.500000 12.000000 9.000000 7.500000 6.500000 6.000000 6.000000 7.500000 19.000000 25.750000 24.000000 13.750000 11.750000 14.000000 10.750000 19.000000 388.000000 404.000000 328.000000 248.000000 181.000000 244.000000 197.00000
	Step 3: Visualizing our data 1. Matplot library is used to make plots from dataset. so, first we will import the required library. 2. Seaborn provides a high-level interface for drawing attractive and informative statistical graphics. So, we will also import the required library.
15]:	<pre>import matplotlib.pyplot as plt import seaborn as sns</pre>
16]:	plt.bar(df["Districts"],height=df['2014'],align="center")
16]:	<pre><barcontainer 46="" artists="" object="" of=""></barcontainer></pre>
	250 - 200 - 150 - 50 - SithaikhapuRasjahiy RatakahikiB ShikabuwahiripiidhasianyiidhasianjiitatiyShiyMothathaanveksahiiwaaNamaskhampitahawahininiah Total
18]:	Trending Plot df.plot()
18]:	<pre><axessubplot:> 400</axessubplot:></pre>
	250 - 2015 - 2016 - 2017 - 2018 - 2019
	Frequency graph
41]:	<pre>import matplotlib.pyplot as plt need_values=total_values[1:] # An "interface" to matplotlib.axes.Axes.hist() method</pre>
	<pre>n, bins, patches = plt.hist(x=need_values) plt.grid(axis='y', alpha=0.75) plt.xlabel('kills') plt.ylabel('Frequency') plt.title('My Very Own Histogram') maxfreq = n.max()</pre>
	My Very Own Histogram 2.00 1.75
	1.50 2.125 1.00 0.75 0.50
	0.50 0.25 0.00 250 300 350 400 kills
19]:	More Exploration of the data total_values = df.max() total_values
19]:	Districts Vehari 2011 364 2012 366 2013 388 2014 404
	2015 328 2016 248 2017 181 2018 244 2019 197 dtype: object
20]:	Plots to see the increase in honour killings over 9 years need_values=total_values[1:] plt.plot(need_values) plt.xlabel("Years")
20]:	plt.ylabel("Number of Cases") plt.title("Trend of Honor killing from 2011-19") Text(0.5, 1.0, 'Trend of Honor killing from 2011-19') Trend of Honor killing from 2011-19
	400 - 350 - 89 300 - 250 -
	200 -
21]:	# Maximum kills during 2011-19 and the corresponding Year dict_of_values=dict(need_values) years=list(dict_of_values.keys())
	<pre># print(years) values=list(dict_of_values.values()) # print(values) # print(list_of_values) print(f"Maximum Honor killings = {need_values.max()}, during the year {years[values.index(404)]}")</pre>
	Maximum Honor killings = 404, during the year 2014 Overal Data Grapgh
22]:	Creating different Models and thier Visualization # import sys # !{sys.executable} -m pip install statsmodel
	<pre># imports import pandas as pd import seaborn as sns # import statsmodels.formula.api as smf from sklearn.linear_model import LinearRegression from sklearn import metrics</pre>
23]:	<pre>import numpy as np # allow plots to appear directly in the notebook %matplotlib inline</pre>
23]:	data=pd.read_csv("CleanData.csv",index_col=0) data.head() 2012 2013 2014 2015 2016 2017 2018 2019 2020 Total kills 2011
	10 15 11 24 13 7 9 10 12 3 114 14 13 5 17 11 11 9 10 5 8 103 13 10 10 15 14 0 2 1 7 8 80 10 13 10 10 13 13 8 7 6 5 95
24]: 24]:	11 9 7 7 9 8 5 15 13 13 97 data.shape (36, 10)
25]:	<pre>sns.pairplot(data, x_vars=["2017","2018","2019","2020"], y_vars='Total kills', size=4, aspect=0.9) /home/jawad/.local/lib/python3.8/site-packages/seaborn/axisgrid.py:1969: UserWarning: The `size` parameter has been renamed to `height`; please up e your code.</pre>
25]:	warnings.warn(msg, UserWarning) <seaborn.axisgrid.pairgrid 0x7f1fb3660400="" at=""> 300 -</seaborn.axisgrid.pairgrid>
	250
	Question 2 Get reported cases from 2019 to 2021 (from any source, this can be only 10 rows), handle it as a regression problem and use the model built on previous data. The task
	should be to predict the district wise total number of cases. Compare your machine learning models performance on unseen data. Step 1: Obtaining Dataset
	We use data from Pujab police data from 2011-2020. There was no data for this Year Step 2: Data Processing for Multiple Regression Algorithm
26]:	We pre-processed the data and formed the file CleanData.csv which we will use for training our Multiregression algorithm and test its performance in the following section. We are using MultiRegression Algorithm because the data is for several years i.e 2011-2020. #import the libraries
a f	<pre>#import the libraries import numpy as py import matplotlib.pyplot as plt import pandas as pd #Import the data set from Desktop dataset = pd.read_csv('CleanData.csv')</pre>
	<pre>dataset.head() X=dataset.iloc[:,:-1].values y=dataset.iloc[:,10].values #Training and Testing Data (divide the data into two part)</pre>
	<pre>from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test =train_test_split(X,y,test_size=0.30, random_state=40) # We are using 30 % for testing and 70% for training our Algo. #regression from sklearn.linear_model import LinearRegression reg = LinearRegression()</pre>
	<pre>reg = LinearRegression() reg.fit(X_train,y_train) #for predict the test values y_prdict=reg.predict(X_test) print(y_prdict); print('\nIntercept: ',reg.intercept_)</pre>
	[80. 112. 26. 107. 97. 43. 74. 20. 23. 131. 94.] Intercept: -1.4210854715202004e-14 Step 3: Algorithm Evluation
	Step 3: Algorithm Eviuation There are many ways to evaluate model performance but in classical statistics, the performance of linear regression models is evaluated with R ² — which gives a value between 0 and 1 and the higher the R ² the better the model. from sklearn.metrics import r2_score
1.	<pre>from sklearn.metrics import * import numpy as np score=r2_score(y_test,y_prdict) print(score) # Our Algorithm gives 100% accurate results</pre>
	<pre># Evaluating the model using MAE Evaluation Metric print(f"mean_squared_error is {mean_absolute_error(y_test, y_prdict)}") # Evaluating the model using RMSE Evaluation Metric print(f"root_mean_squared_error is {np.sqrt(mean_absolute_error(y_test, y_prdict))}")</pre> 1.0
	mean_squared_error is 2.5514943693203597e-14 root_mean_squared_error is 1.5973397789200517e-07 The End
[]:	