

```
In [12]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
#Importing necessary libraries

In [13]: df = pd.read_csv('RTA Dataset.csv')
df.head()
Out[13]:
   Time  Day_of_week  Age_band_of_driver  Sex_of_driver  Educational_level  Vehicle_driver_relation  Driving_experience  Type_of_vehicle  Owner_of_vehicle  Service_year_of_vehicle  ...  Vehicle_movement  Casualty_class  Sex_of_casualty  Age_band_of_casualty  Casualty_severity  Work_o

0  17:02:00      Monday                31-30         Male                Above high school                Employee                1-2yr                Automobile                Owner                Above 10yr  ...  Going straight                na                na                na                na                na

1  17:02:00      Monday                31-30         Male                Junior high school                Employee                Above 10yr                Public (> 45 seats)                Owner                5-10yr  ...  Going straight                na                na                na                na                na

2  17:02:00      Monday                18-30         Male                Junior high school                Employee                1-2yr                Lorry (417 1002)                Owner                NaN  ...  Going straight                Driver or rider                Male                31-50                3

3  1:06:00       Sunday                18-30         Male                Junior high school                Employee                5-10yr                Public (> 45 seats)                Governmental                NaN  ...  Going straight                Pedestrian                Female                18-30                3

4  1:06:00       Sunday                18-30         Male                Junior high school                Employee                2-5yr                NaN                Owner                5-10yr  ...  Going straight                na                na                na                na

5 rows x 32 columns

In [14]: df.shape
Out[14]: (12316, 32)

In [15]: df.columns
Out[15]:
['Time', 'Day_of_week', 'Age_band_of_driver', 'Sex_of_driver', 'Educational_level', 'Vehicle_driver_relation', 'Driving_experience', 'Type_of_vehicle', 'Owner_of_vehicle', 'Service_year_of_vehicle', 'Vehicle_movement', 'Casualty_class', 'Sex_of_casualty', 'Age_band_of_casualty', 'Casualty_severity', 'Work_o']

In [16]: df.isnull().sum()
#Check Missing value
Out[16]:
Time                0
Day_of_week         0
Age_band_of_driver  0
Sex_of_driver       0
Educational_level   741
Vehicle_driver_relation  979
Driving_experience   829
Type_of_vehicle     352
Owner_of_vehicle    482
Service_year_of_vehicle 4427
Defect_of_vehicle   239
Area_accident_occured 385
Lanes_of_Medians    189
Road_alignment      142
Road_surface_type   172
Road_surface_conditions  0
Light_conditions    0
Weather_conditions  0
Type_of_collision   155
Number_of_vehicles_involved  0
Number_of_casualties  0
Vehicle_movement    308
Casualty_class      0
Sex_of_casualty     0
Age_band_of_casualty  0
Casualty_severity   0
Work_of_casualty    3198
Fitness_of_casualty  2435
Pedestrian_movement  0
Cause_of_accident   0
Accident_severity   0
dtype: int64

In [17]: df.duplicated().sum()
#Check for duplicate value
Out[17]: 0

In [18]: df.describe(include="all")
#Descriptive statistics
Out[18]:
   Time  Day_of_week  Age_band_of_driver  Sex_of_driver  Educational_level  Vehicle_driver_relation  Driving_experience  Type_of_vehicle  Owner_of_vehicle  Service_year_of_vehicle  ...  Vehicle_movement  Casualty_class  Sex_of_casualty  Age_band_of_casualty  Casualty_severity  Work_o
count      12316      12316              12316              12316              11575              11737              11487              11366              11834              8388  ...  12008      12316      12316              12316
unique         NaN         NaN                5                3                7                4                7                17                4                6  ...  13                4                3                6
top           NaN         Friday           18-30           Male           Junior high school           Employee           5-10yr           Automobile           Owner           Unknown  ...  Going straight           Driver or rider           Male           na
freq           NaN         2041           4271           11437           7619           9627           3363           3205           10459           2883  ...  8168           4944           5253           4443
mean    2024-08-30  14:17:50.708192040         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
min      2024-08-30  00:01:00         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
25%      2024-08-30  10:31:00         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
50%      2024-08-30  15:10:00         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
75%      2024-08-30  18:10:00         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
max      2024-08-30  23:59:00         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN
std           NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN         NaN  ...  NaN         NaN         NaN         NaN

11 rows x 32 columns

In [19]: df.groupby('Accident_severity')['size']
#Display the column size groups
Out[19]:
Accident_severity
Fatal injury      158
Serious injury   1783
Slight injury    10415
dtype: int64

In [20]: df['Number_of_casualties'].value_counts()
#Print count for each unique value in column
Out[20]:
Number_of_casualties
1      8397
2      2290
3       909
4       294
5       207
6        89
7         22
Name: count, dtype: int64

In [21]: plt.figure(figsize=(10,7))
sns.barplot(data=df, y='Number_of_vehicles_involved', x='Number_of_casualties')
#Display of bar chart
Out[21]:
Figure with 2 Axes
Number_of_vehicles_involved
Number_of_casualties
1 2 3 4 5 6 7 8
1 2 3 4 5 6 7 8
Number of vehicles involved
Number of casualties

In [22]: sns.barplot(data=df, y='Number_of_casualties')
plt.show()
Out[22]:
Figure with 1 Axes
Number_of_casualties
1 2 3 4 5 6 7 8
1 2 3 4 5 6 7 8
Number of casualties
Number of vehicles involved

In [23]: sns.barplot(data=df, y='Number_of_vehicles_involved')
plt.show()
Out[23]:
Figure with 1 Axes
Number_of_vehicles_involved
1 2 3 4 5 6 7 8
1 2 3 4 5 6 7 8
Number of vehicles involved
Number of casualties

In [24]: sns.scatterplot(x=df['Number_of_vehicles_involved'], y=df['Number_of_casualties'])
plt.show()
#Scatter plot to determine relationships in columns
Out[24]:
Figure with 1 Axes
Number_of_casualties
Number_of_vehicles_involved
1 2 3 4 5 6 7
1 2 3 4 5 6 7
Number of casualties
Number of vehicles involved

In [25]: sns.pairplot(df[['Number_of_vehicles_involved', 'Number_of_casualties']])
plt.show()
#Pairwise relationships in dataset
Out[25]:
Figure with 4 Axes
Number_of_vehicles_involved
Number_of_casualties
1 2 3 4 5 6 7 8
1 2 3 4 5 6 7 8
Number of vehicles involved
Number of casualties

In [26]: correlation_matrix = df[['Number_of_vehicles_involved', 'Number_of_casualties']].corr()
sns.heatmap(correlation_matrix, annot=True)
plt.show()
#Display of heatmap having correlation of columns as attribute
Out[26]:
Figure with 1 Axes
Number_of_vehicles_involved
Number_of_casualties
1 0.21
1 0.21
Number of vehicles involved
Number of casualties

In [27]: plt.figure(figsize=(10,7))
plt.pie(df['Accident_severity'].value_counts().values,
       labels=df['Accident_severity'].value_counts().index,
       autopct='%2.2f%%')
plt.show()
#Display of pie chart
Out[27]:
Figure with 1 Axes
Slight injury 84.56%
Fatal injury 1.28%
Serious injury 14.15%

In [28]: grid = sns.FacetGrid(data=df, col='Accident_severity', height=4, aspect=1, sharex=False)
grid.map(sns.countplot, 'Number_of_vehicles_involved', palette='black', 'brown', 'orange')
plt.show()
#Users[Vagpet: Kaushikarandad]Lib[matplotlib-packages:seaborn:Oldcore.py:119: FutureWarning: Using the countplot function without specifying 'order' is likely to produce an incorrect plot.
warnings.warn(warnings)
Out[28]:
Figure with 3 Axes
Accident_severity = Slight injury
Accident_severity = Serious Injury
Accident_severity = Fatal injury
Count
Number of vehicles involved
1 2 3
1 2 3
Number of vehicles involved

In [29]: lists=['Vehicle_driver_relation', 'Work_of_casualty', 'Fitness_of_casualty', 'Day_of_week',
             'Casualty_severity', 'Time', 'Sex_of_driver', 'Educational_level', 'Defect_of_vehicle', 'Owner_of_vehicle', 'Service_year_of_vehicle',
             'Road_surface_type', 'Sex_of_casualty']
df.drop(labels = lists, inplace=True)
#Drop the particular columns
Out[29]:
Figure with 1 Axes
In [30]: df.shape
Out[30]: (12316, 19)

In [31]: df.columns
Out[31]:
['Time', 'Day_of_week', 'Age_band_of_driver', 'Driving_experience', 'Type_of_vehicle', 'Area_accident_occured', 'Lanes_of_Medians', 'Road_alignment', 'Types_of_collision', 'Weather_conditions', 'Light_conditions', 'Type_of_collision', 'Number_of_vehicles_involved', 'Number_of_casualties', 'Vehicle_movement', 'Casualty_class', 'Age_band_of_casualty', 'Pedestrian_movement', 'Cause_of_accident', 'Accident_severity']

In [32]: df.isnull().sum()
#Checking the missing values
Out[32]:
Age_band_of_driver 0
Driving_experience 0
Type_of_vehicle 0
Area_accident_occured 0
Lanes_of_Medians 0
Road_alignment 0
Types_of_collision 0
Weather_conditions 0
Type_of_collision 0
Number_of_vehicles_involved 0
Number_of_casualties 0
Vehicle_movement 0
Casualty_class 0
Age_band_of_casualty 0
Pedestrian_movement 0
Cause_of_accident 0
Accident_severity 0
dtype: int64

In [33]: target_count = df['Accident_severity'].value_counts()
print('Class 0:', target_count[0])
print('Class 1:', target_count[1])
print('Proportion:', round(target_count[0] / target_count[0] + target_count[1], 2), '%')
target_count.plot(kind='bar', title='Count (target)')
#Display of counts for each categorical bin using bars.
Out[33]:
Figure with 1 Axes
Class 0: 1743
Class 1: 1743
Proportion: 5.98 %
Count (target)
Slight injury
Serious injury
Fatal injury
Count
Number of vehicles involved
1 2 3
1 2 3
Number of vehicles involved

In [34]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df=df.apply(le.fit_transform)
#LabelEncoder is used to convert categorical columns into numerical ones
Out[34]:
Figure with 1 Axes
In [35]: df.corr()
#Display pairwise correlation of columns, excluding null values.
Out[35]:
   Time  Day_of_week  Age_band_of_driver  Driving_experience  Type_of_vehicle  Area_accident_occured  Lanes_of_Medians  Road_alignment  Types_of_collision  Types_of_collision  Weather_conditions  Light_conditions  Type_of_collision  Number_of_vehicles_involved  Number_of_casualties  Vehicle_movement  Casualty_class  Age_band_of_casualty  Pedestrian_movement  Cause_of_accident  Accident_severity
Time                1.000000              -0.013407              0.004630              0.011472              0.008178              0.000486              -0.030205              -0.007766              0.002482              0.000496              0.001690              -0.006222
Day_of_week         -0.013407              1.000000              0.014333              -0.002090              -0.010444              -0.015602              0.009245              -0.015577              -0.000094              0.012574              0.007008              -0.006222
Age_band_of_driver   0.004630              0.014333              1.000000              -0.007631              -0.020156              -0.020910              -0.005553              0.019018              0.003739              0.000903              0.025991              0.001140
Driving_experience    0.011472              -0.002090              -0.007631              1.000000              0.022780              -0.025859              0.014580              -0.002790              -0.013597              0.003127              -0.003940              -0.016498
Type_of_vehicle      0.008178              -0.010444              -0.020156              0.022780              1.000000              0.042815              0.011774              0.004624              0.006828              0.007008              -0.001199
Owner_of_vehicle     0.000486              -0.015602              -0.020910              -0.025859              0.042815              1.000000              0.007188              -0.000918              0.000625              0.016304              0.020418              -0.001536
Service_year_of_vehicle 0.024825              0.009245              -0.005553              0.014580              0.002125              0.007188              1.000000              0.004724              0.037063              0.001816              0.007516              -0.020214
Vehicle_movement     -0.020214              -0.003696              0.019018              -0.002790              0.011774              -0.000918              0.004624              1.000000              0.040485              0.000456              0.000456              -0.020963
Casualty_class        0.007766              -0.015577              -0.003739              0.003127              0.005625              0.037063              0.019847              0.000466              1.000000              0.007367              0.000000              -0.013347
Sex_of_casualty      0.025482              -0.000094              0.000903              0.003127              0.005625              0.016304              -0.001816              0.000485              0.007367              1.000000              -0.013347              0.037733
Age_band_of_casualty 0.000496              0.012574              0.025991              -0.003940              0.007008              0.020418              0.007516              0.000456              0.014883              0.013347              1.000000              0.007304
Work_of_casualty     0.001690              -0.006222              0.007008              -0.016498              -0.001199              0.007516              0.000456              -0.020214              -0.020963              0.014138              0.020418              1.000000
Fitness_of_casualty  -0.006222              0.004474              0.003022              -0.004690              -0.002188              -0.014134              -0.004369              0.027828              0.056116              0.014138              -0.020694              0.220531
Vehicle_movement     -0.006662              0.004794              0.001084              0.004691              -0.002089              -0.004267              -0.017671              -0.004667              0.002237              0.007648              0.000640              -0.004480
Casualty_class        0.005386              -0.002783              -0.006116              0.011325              0.000634              -0.003317              0.002423              -0.002213              0.012499              0.003313              0.012333              0.006281
Age_band_of_casualty 0.000043              -0.010791              -0.008304              0.008001              0.005640              0.007946              0.006815              -0.011720              0.014429              0.006157              0.014229              0.000981
Pedestrian_movement  0.000372              0.018190              0.010500              0.000949              0.005648              0.005219              0.008585              0.006157              0.014229              0.016292              -0.000801              0.021196
Cause_of_accident    -0.000703              -0.007728              -0.004200              -0.009698              0.016304              0.002287              -0.005861              0.003911              0.004210              0.011320              0.017157              -0.021815
Accident_severity    0.013185              0.000509              -0.003709              -0.016374              -0.011844              -0.005967              -0.009861              0.004579              0.025741              0.010172              0.025967              0.100205

In [36]: plt.figure(figsize=(25,15))
sns.heatmap(df.corr(),annot=True)
Out[36]:
Figure with 1 Axes
In [37]:
```

