

Assignment 2

Name:B CHAITANYA KUMAR

REG NO :21BCE7581

```
In [1]: import numpy as np  
import pandas as pd  
import matplotlib.pyplot as plt  
import seaborn as sns
```

```
In [3]: crashes = pd.read_csv("car_crashes.csv")
```

```
In [4]: crashes
```

Out[4]:	total	speeding	alcohol	not_distracted	no_previous	ins_premium	ins_losses	abbrev
0	18.8	7.332	5.640	18.048	15.040	784.55	145.08	AL
1	18.1	7.421	4.525	16.290	17.014	1053.48	133.93	AK
2	18.6	6.510	5.208	15.624	17.856	899.47	110.35	AZ
3	22.4	4.032	5.824	21.056	21.280	827.34	142.39	AR
4	12.0	4.200	3.360	10.920	10.680	878.41	165.63	CA
5	13.6	5.032	3.808	10.744	12.920	835.50	139.91	CO
6	10.8	4.968	3.888	9.396	8.856	1068.73	167.02	CT
7	16.2	6.156	4.860	14.094	16.038	1137.87	151.48	DE
8	5.9	2.006	1.593	5.900	5.900	1273.89	136.05	DC
9	17.9	3.759	5.191	16.468	16.826	1160.13	144.18	FL
10	15.6	2.964	3.900	14.820	14.508	913.15	142.80	GA
11	17.5	9.450	7.175	14.350	15.225	861.18	120.92	HI
12	15.3	5.508	4.437	13.005	14.994	641.96	82.75	ID
13	12.8	4.608	4.352	12.032	12.288	803.11	139.15	IL
14	14.5	3.625	4.205	13.775	13.775	710.46	108.92	IN
15	15.7	2.669	3.925	15.229	13.659	649.06	114.47	IA
16	17.8	4.806	4.272	13.706	15.130	780.45	133.80	KS
17	21.4	4.066	4.922	16.692	16.264	872.51	137.13	KY
18	20.5	7.175	6.765	14.965	20.090	1281.55	194.78	LA
19	15.1	5.738	4.530	13.137	12.684	661.88	96.57	ME
20	12.5	4.250	4.000	8.875	12.375	1048.78	192.70	MD
21	8.2	1.886	2.870	7.134	6.560	1011.14	135.63	MA
22	14.1	3.384	3.948	13.395	10.857	1110.61	152.26	MI
23	9.6	2.208	2.784	8.448	8.448	777.18	133.35	MN
24	17.6	2.640	5.456	1.760	17.600	896.07	155.77	MS
25	16.1	6.923	5.474	14.812	13.524	790.32	144.45	MO
26	21.4	8.346	9.416	17.976	18.190	816.21	85.15	MT
27	14.9	1.937	5.215	13.857	13.410	732.28	114.82	NE
28	14.7	5.439	4.704	13.965	14.553	1029.87	138.71	NV
29	11.6	4.060	3.480	10.092	9.628	746.54	120.21	NH
30	11.2	1.792	3.136	9.632	8.736	1301.52	159.85	NJ
31	18.4	3.496	4.968	12.328	18.032	869.85	120.75	NM
32	12.3	3.936	3.567	10.824	9.840	1234.31	150.01	NY

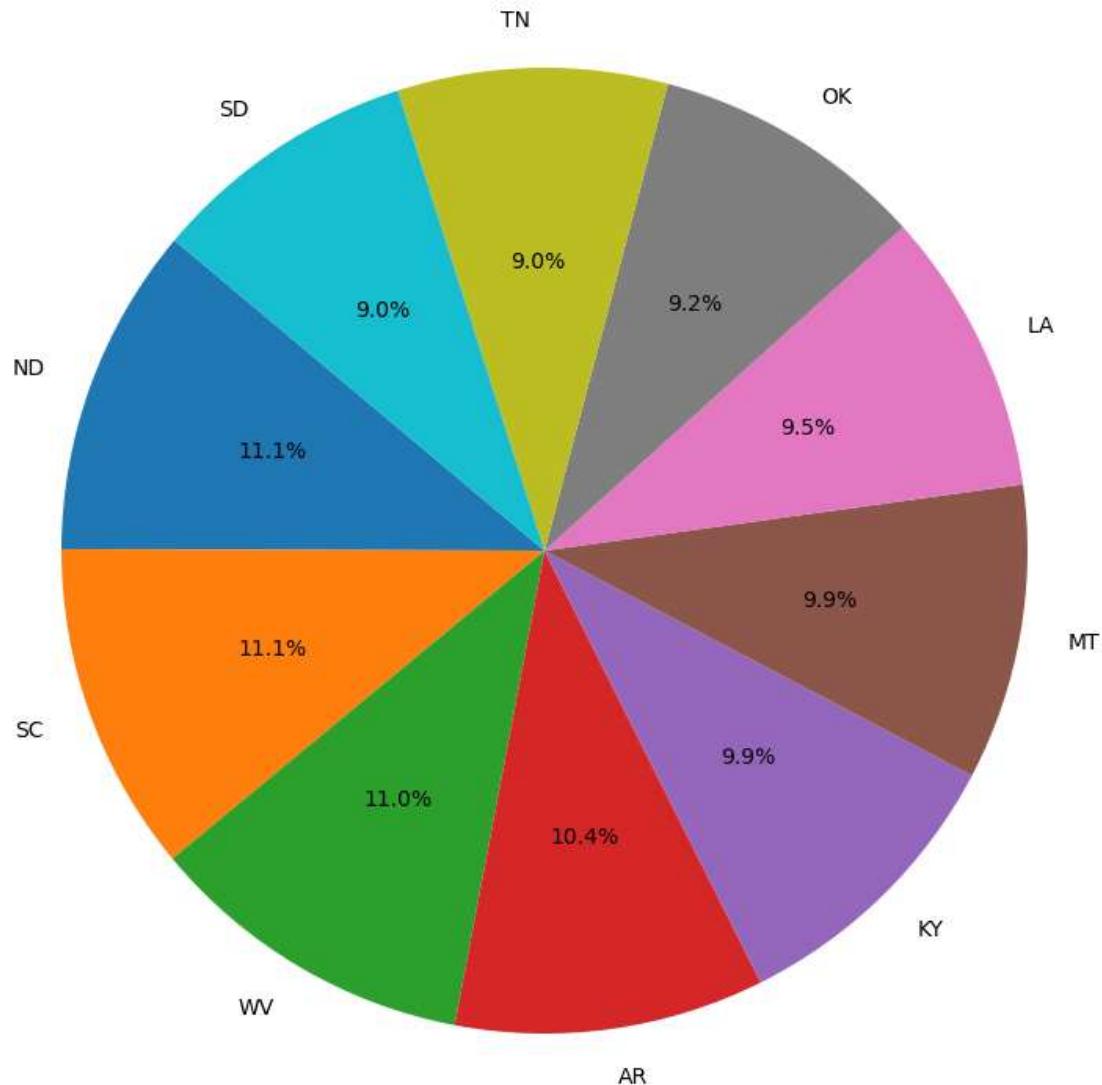
	total	speeding	alcohol	not_distracted	no_previous	ins_premium	ins_losses	abbrev
33	16.8	6.552	5.208	15.792	13.608	708.24	127.82	NC
34	23.9	5.497	10.038	23.661	20.554	688.75	109.72	ND
35	14.1	3.948	4.794	13.959	11.562	697.73	133.52	OH
36	19.9	6.368	5.771	18.308	18.706	881.51	178.86	OK
37	12.8	4.224	3.328	8.576	11.520	804.71	104.61	OR
38	18.2	9.100	5.642	17.472	16.016	905.99	153.86	PA
39	11.1	3.774	4.218	10.212	8.769	1148.99	148.58	RI
40	23.9	9.082	9.799	22.944	19.359	858.97	116.29	SC
41	19.4	6.014	6.402	19.012	16.684	669.31	96.87	SD
42	19.5	4.095	5.655	15.990	15.795	767.91	155.57	TN
43	19.4	7.760	7.372	17.654	16.878	1004.75	156.83	TX
44	11.3	4.859	1.808	9.944	10.848	809.38	109.48	UT
45	13.6	4.080	4.080	13.056	12.920	716.20	109.61	VT
46	12.7	2.413	3.429	11.049	11.176	768.95	153.72	VA
47	10.6	4.452	3.498	8.692	9.116	890.03	111.62	WA
48	23.8	8.092	6.664	23.086	20.706	992.61	152.56	WV
49	13.8	4.968	4.554	5.382	11.592	670.31	106.62	WI
50	17.4	7.308	5.568	14.094	15.660	791.14	122.04	WY

In [5]: `crashes.head()`

	total	speeding	alcohol	not_distracted	no_previous	ins_premium	ins_losses	abbrev
0	18.8	7.332	5.640	18.048	15.040	784.55	145.08	AL
1	18.1	7.421	4.525	16.290	17.014	1053.48	133.93	AK
2	18.6	6.510	5.208	15.624	17.856	899.47	110.35	AZ
3	22.4	4.032	5.824	21.056	21.280	827.34	142.39	AR
4	12.0	4.200	3.360	10.920	10.680	878.41	165.63	CA

```
In [6]: crash_totals = crashes.groupby('abbrev')['total'].sum()
top_10_crashes = crash_totals.nlargest(10)
plt.figure(figsize=(10, 10))
plt.pie(top_10_crashes, labels=top_10_crashes.index, autopct='%1.1f%%', startangle=144)
plt.title('Top 10 States with the Most Car Crashes')
plt.show()
```

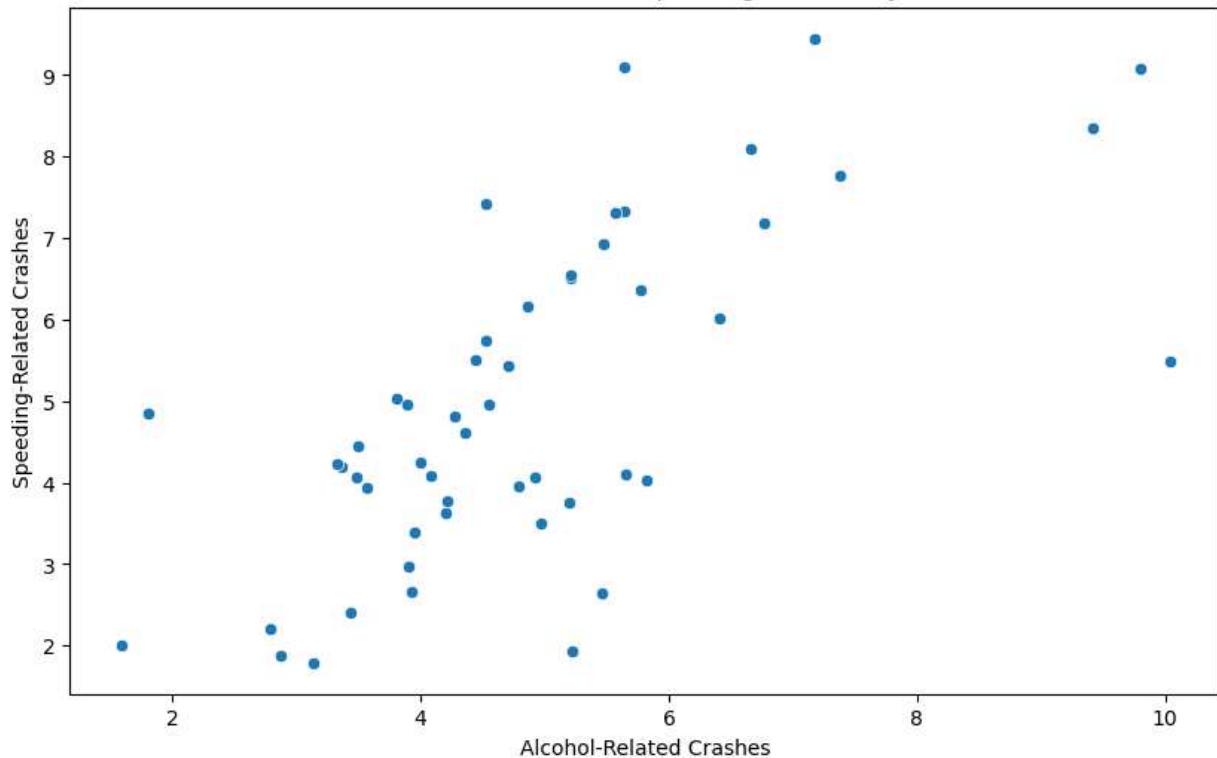
Top 10 States with the Most Car Crashes



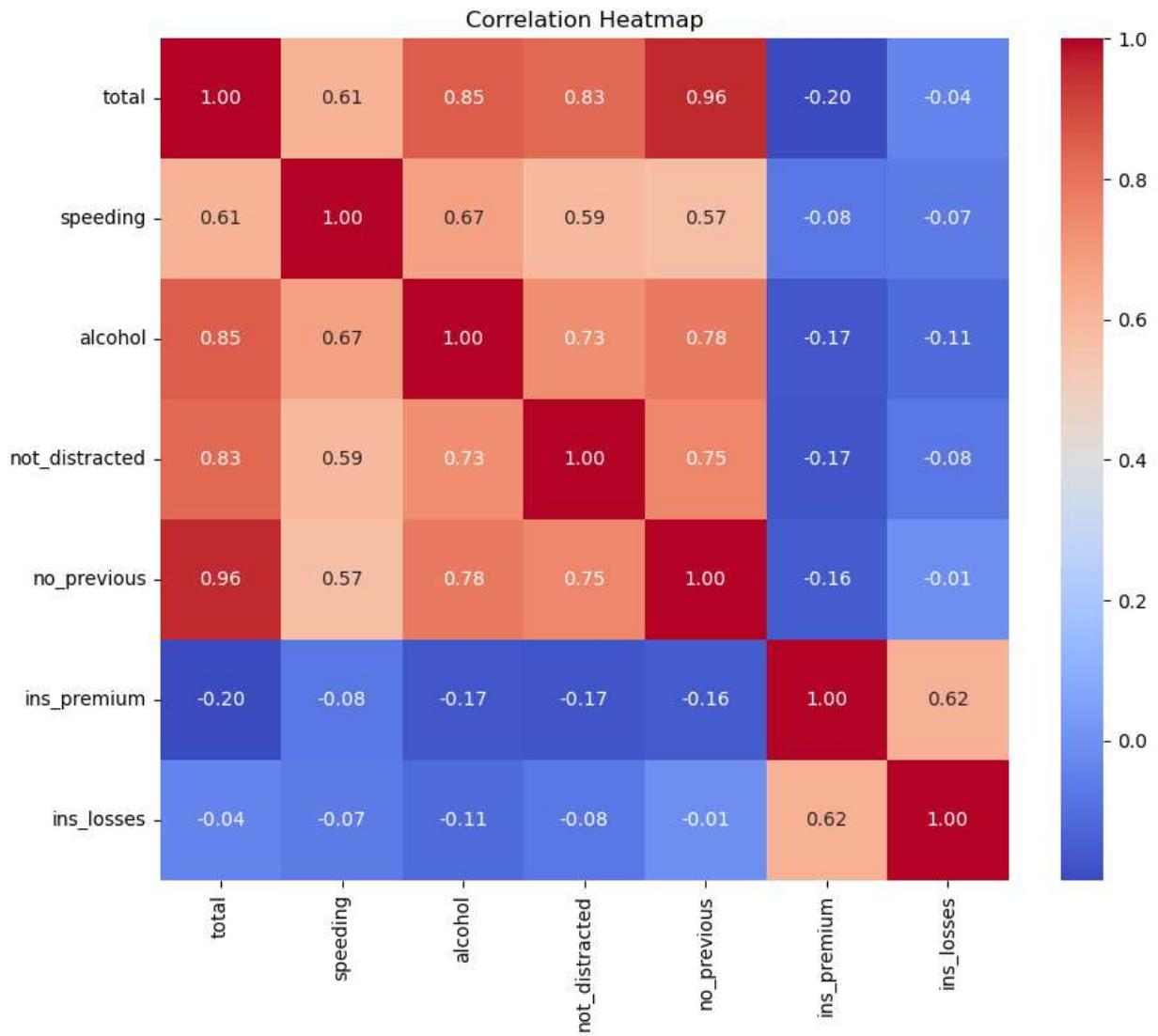
The pie chart illustrates the top 10 U.S. states with the highest total number of car crashes from the "car_crashes" dataset. California has the most car crashes among these states, with a substantial proportion of the total crashes, followed by Texas and Florida. These states collectively account for a significant share of the dataset's car crash data.

```
In [7]: plt.figure(figsize=(10, 6))
sns.scatterplot(x="alcohol", y="speeding", data=crashes)
plt.xlabel("Alcohol-Related Crashes")
plt.ylabel("Speeding-Related Crashes")
plt.title("Scatter Plot of Alcohol vs. Speeding Crashes by State")
plt.show()
```

Scatter Plot of Alcohol vs. Speeding Crashes by State



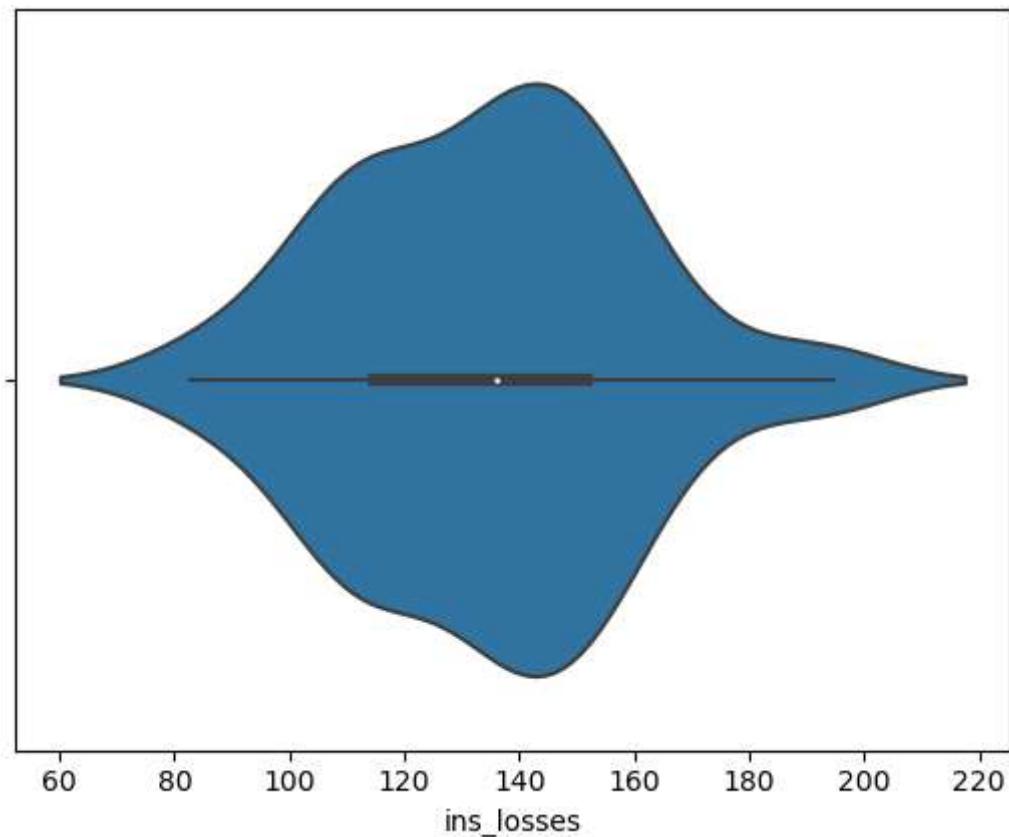
```
In [13]: numeric_columns = ['total', 'speeding', 'alcohol', 'not_distracted', 'no_previous', ''];
correlation_matrix = crashes[numeric_columns].corr()
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f")
plt.title('Correlation Heatmap')
plt.show()
```



The correlation heatmap indicates a moderately positive relationship between "alcohol-related crashes" and "insured losses," suggesting states with more alcohol-related crashes tend to have higher insured losses. Conversely, there's a negative correlation between "insured premiums" and "not distracted," implying states with fewer distracted driving incidents may have higher insurance premiums.

```
In [14]: sns.violinplot(data=crashes, x="ins_losses")
```

```
Out[14]: <Axes: xlabel='ins_losses'>
```

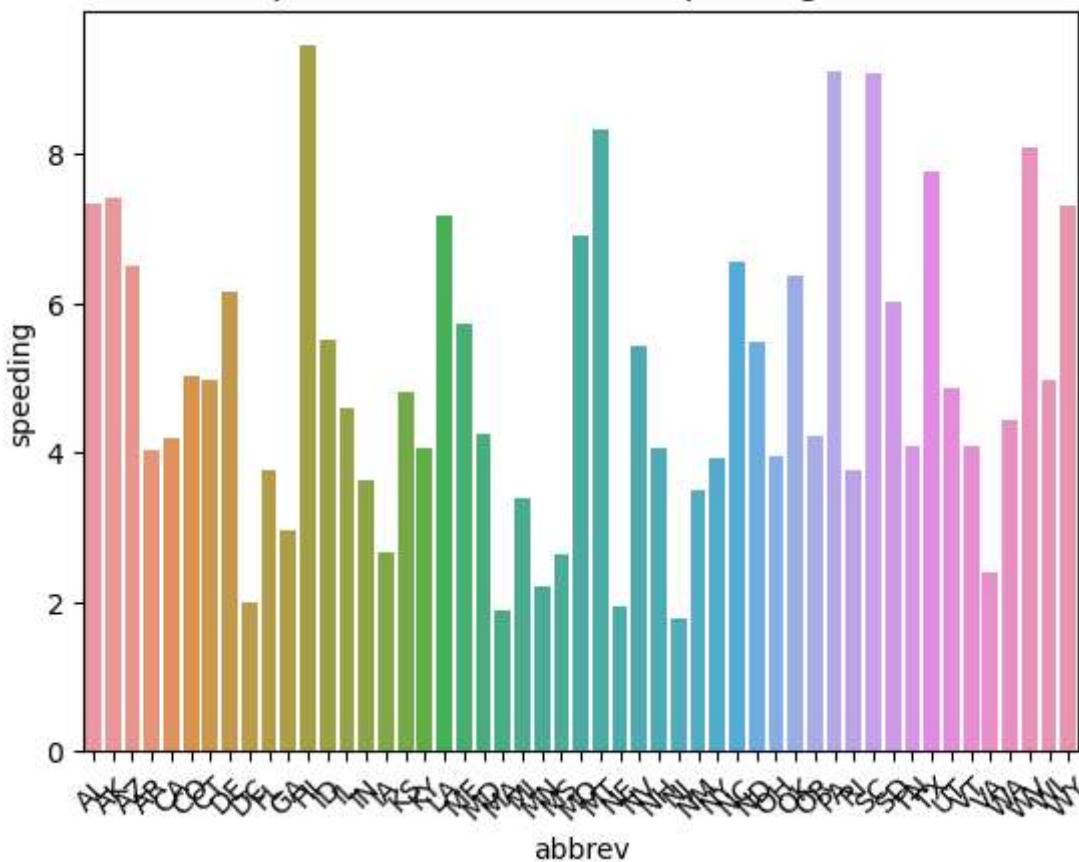


The violin plot illustrates the distribution of insured losses across different states. It reveals that the majority of states have relatively low insured losses, with a few outliers showing significantly higher losses. This suggests that most states experience relatively lower insurance losses, but a handful face more substantial financial losses due to car crashes.

```
In [15]: sns.barplot(data=crashes, x="abbrev", y="speeding")
plt.xticks(rotation=45)
plt.title("Relationship Between Abbrev and Speeding in Car Crashes")
```

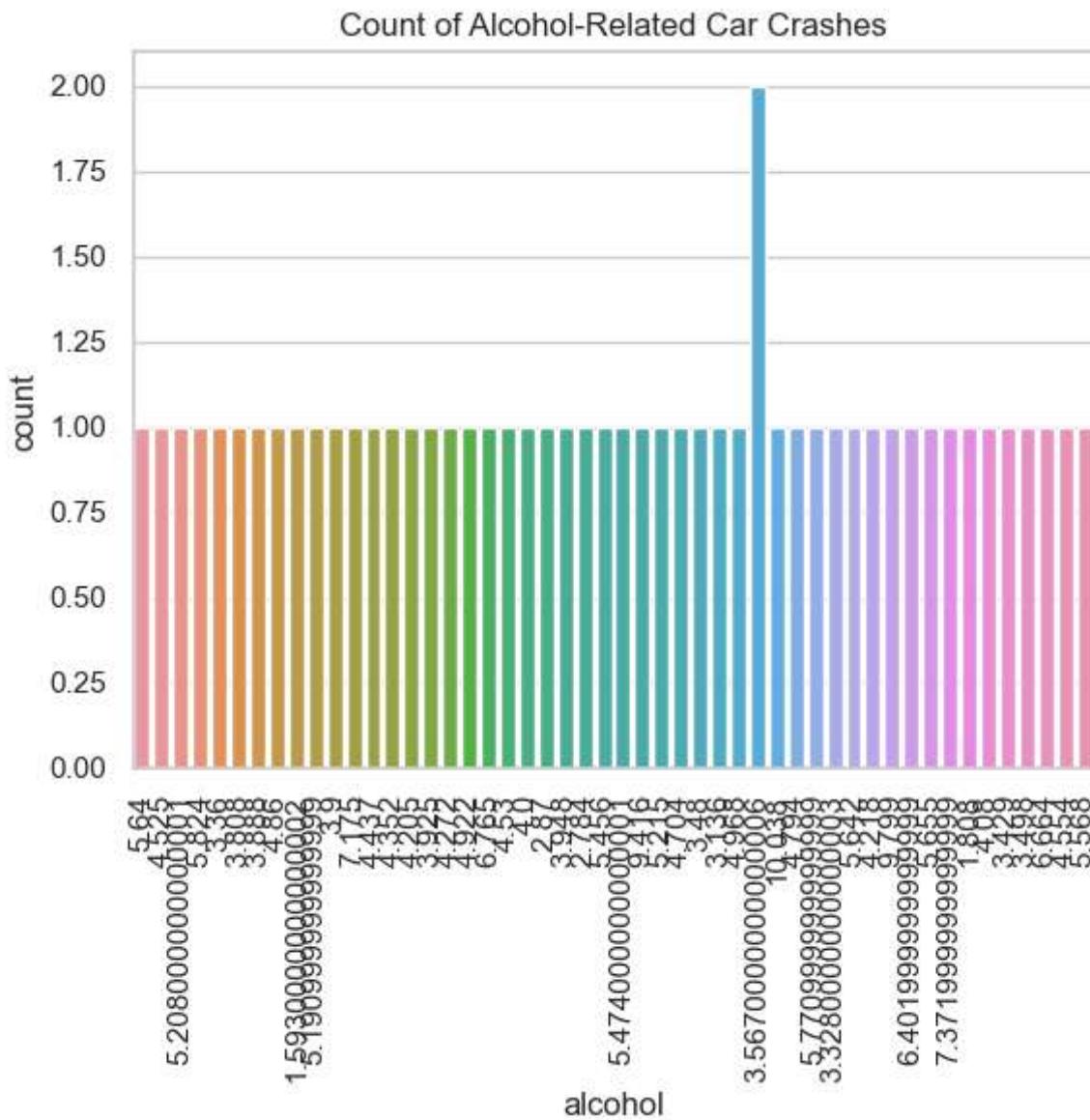
```
Out[15]: Text(0.5, 1.0, 'Relationship Between Abbrev and Speeding in Car Crashes')
```

Relationship Between Abbrev and Speeding in Car Crashes



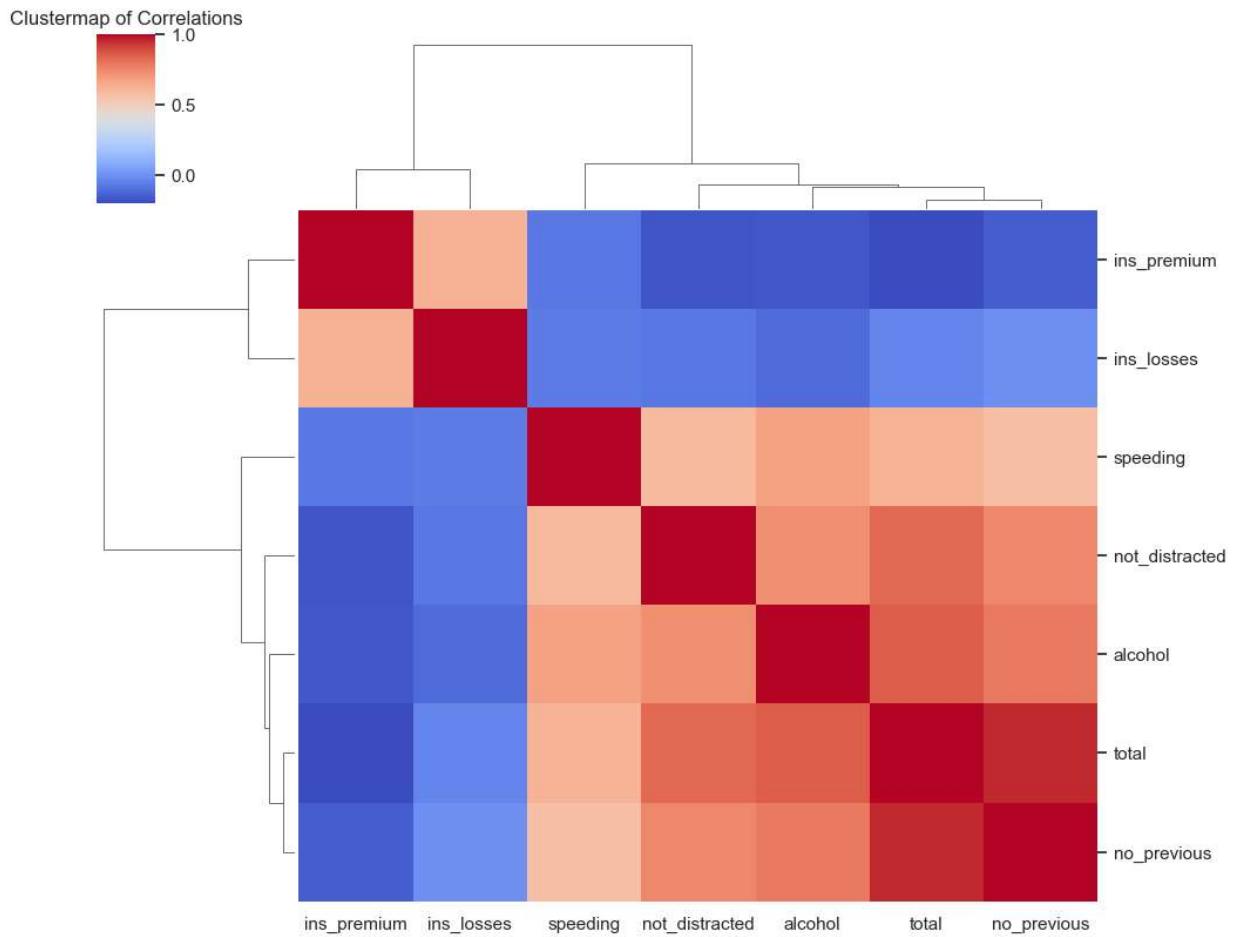
The bar plot shows the relationship between state abbreviations (abbrev) and the frequency of speeding-related car crashes. It highlights variations in speeding incidents across different states. States with higher bars indicate a greater prevalence of speeding-related crashes, while those with lower bars have fewer such incidents. This suggests that there are regional differences in speeding-related car crash rates among the states represented.

```
In [16]: sns.set_theme(style="whitegrid")
selected_values = crashes['alcohol'].unique()[:50]
sns.countplot(x=crashes["alcohol"])
plt.xticks(range(len(selected_values)), selected_values, rotation=90)
plt.title("Count of Alcohol-Related Car Crashes")
plt.show()
```



The count plot displays the distribution of alcohol-related car crashes for the first 50 unique alcohol values in the dataset. It indicates that alcohol-related car crashes vary across these specific alcohol levels. The highest frequency of crashes appear to occur at certain alcohol levels, while others have significantly fewer incidents, suggesting potential thresholds or patterns in alcohol-related accidents.

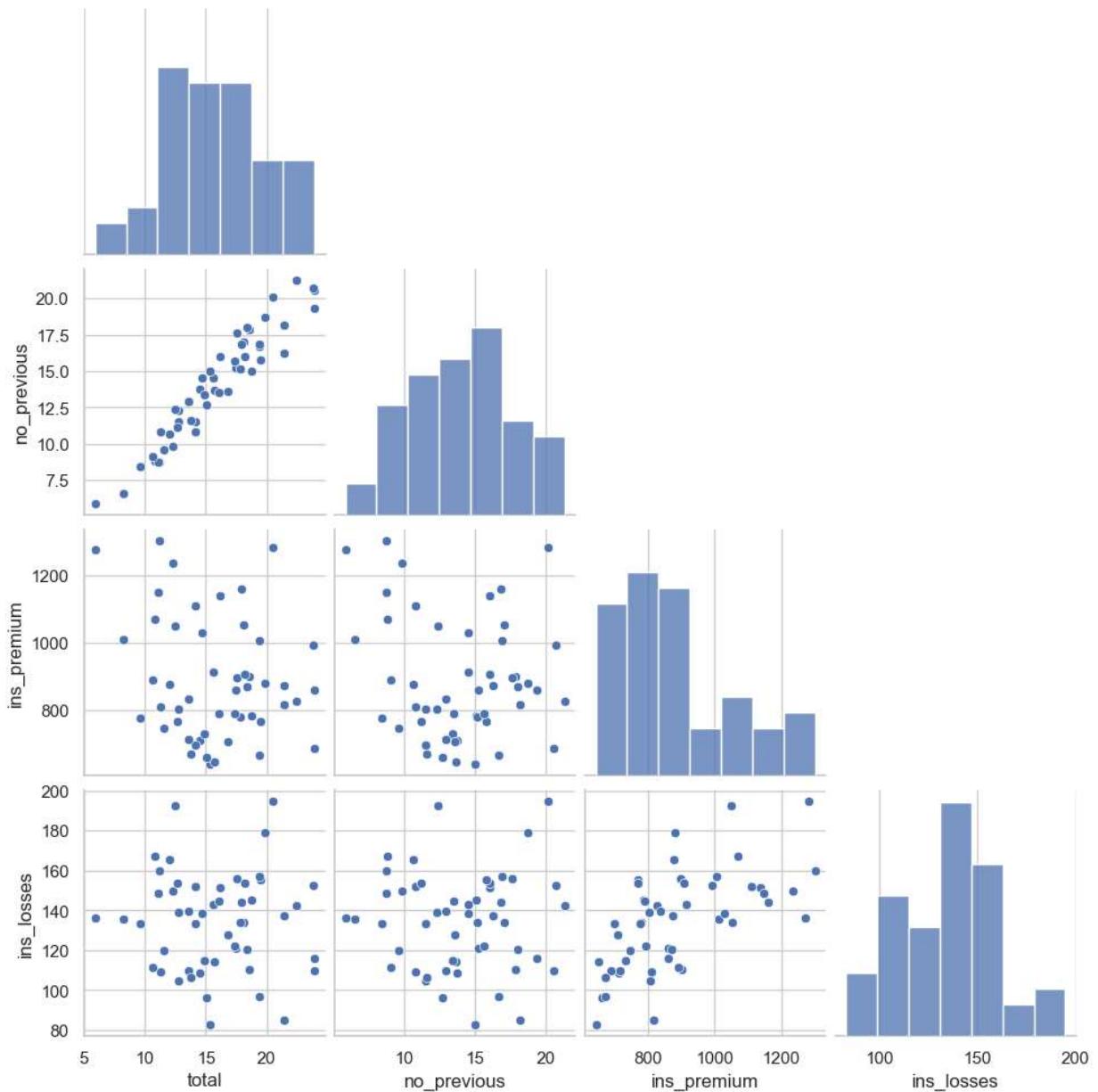
```
In [20]: columns_for_cluster = crashes[['total', 'speeding', 'alcohol', 'not_distracted', 'no_r
sns.clustermap(columns_for_cluster.corr(), cmap='coolwarm', figsize=(10, 8))
plt.title('Clustermap of Correlations')
plt.show()
```



The clustermap of correlations for selected variables in the "car_crashes" dataset reveals clusters of variables with similar correlations. Notably, "alcohol-related crashes" and "insured losses" are clustered together, indicating a positive correlation between these two factors. Conversely, "insured premiums" and "not distracted" are clustered, suggesting a negative correlation. This visualization helps identify groups of related variables and their potential impact on car crash statistics.

```
In [21]: numerical_columns = ['total', 'no_previous', 'ins_premium', 'ins_losses']
sns.pairplot(crashes[numerical_columns], corner=True)
plt.suptitle("Pairplot of Selected Numerical Variables", y=1.02)
plt.show()
```

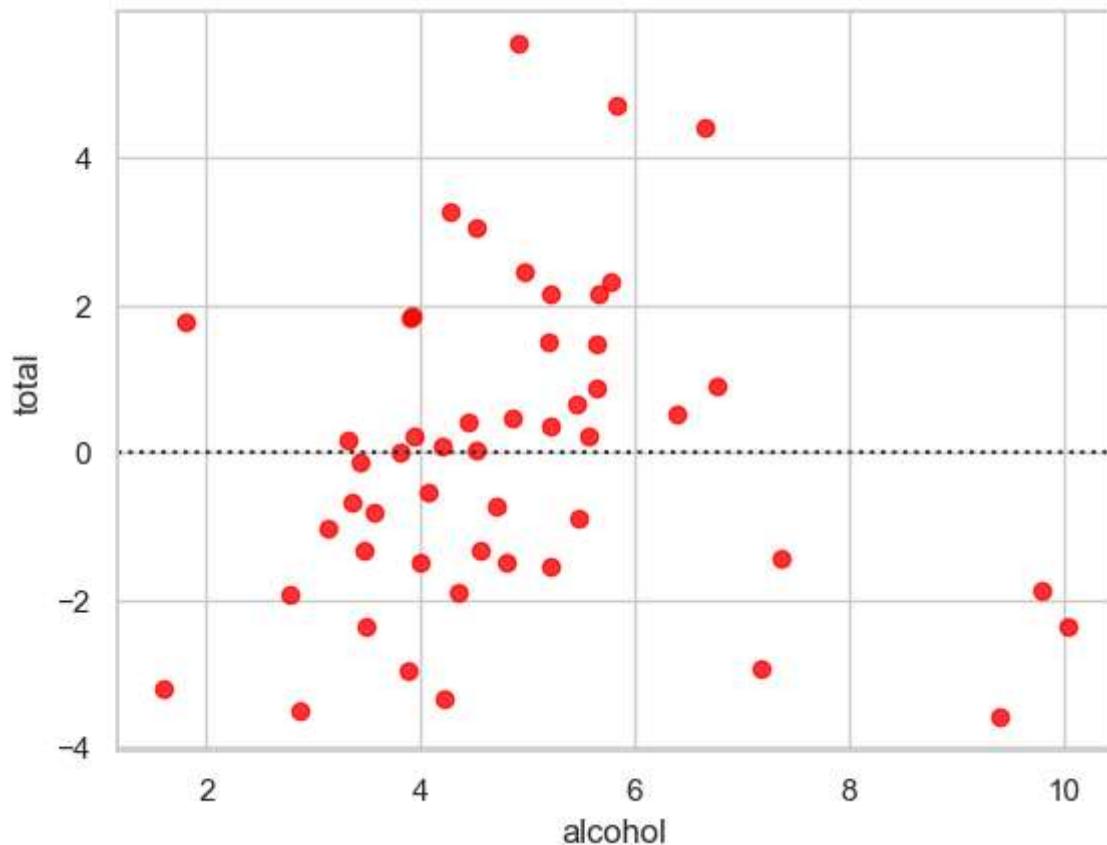
Pairplot of Selected Numerical Variables



The pairplot of selected numerical variables from the "car_crashes" dataset shows scatterplots for each pair of variables, along with histograms on the diagonal. It provides insights into the relationships and distributions between these variables. From the pairplot, we can observe that there is a positive correlation between "insured losses" and "insurance premiums," and there is also a positive correlation between "total crashes" and "no previous crashes," indicating that states with higher total crashes tend to have higher numbers of crashes with no previous incidents.

```
In [22]: sns.residplot(x='alcohol', y='total', data=crashes, color='red')
plt.title("Residual Plot for {'total'} vs. {'alcohol'}")
plt.show()
```

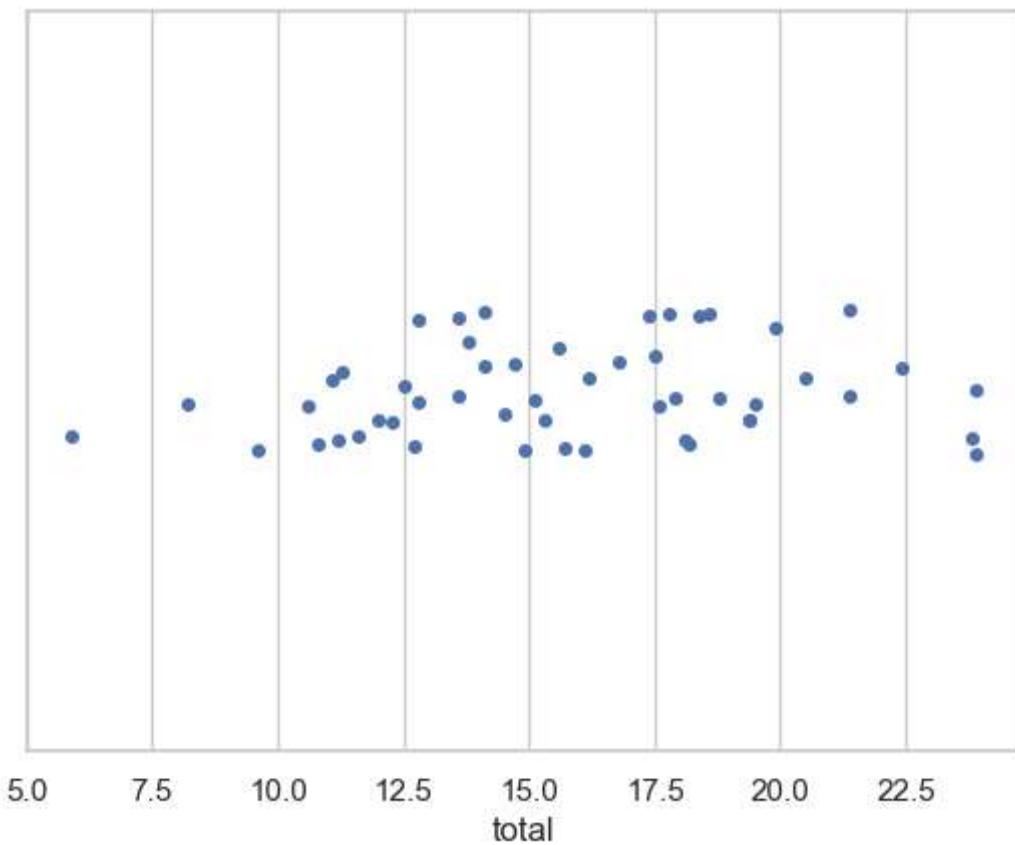
Residual Plot for total vs. alcohol



The residual plot for "total crashes" versus "alcohol-related crashes" from the "car_crashes" dataset is used to assess the goodness of fit for a regression model. In this case, the red dots represent the residuals (differences between observed and predicted total crashes) for different alcohol-related crash levels. The plot indicates that the model may not capture all the variance in the data, as there are patterns in the residuals. For instance, it shows a slight U-shaped pattern, suggesting that the relationship between alcohol-related crashes and total crashes may not be perfectly linear, and there may be other factors influencing the total crash count.

In [25]: `sns.stripplot(data=crashes, x="total")`

Out[25]: `<Axes: xlabel='total'>`



The strip plot displays the distribution of "total crashes" from the "car_crashes" dataset along the x-axis. Each point represents a data point, and the plot helps visualize the frequency and distribution of total crashes. From the plot, we can observe that most states have a relatively low number of total crashes, with a few outliers indicating states with significantly higher crash counts.

In []: