

Exploratory Data Analysis on US accidents

Steps 1) Selecting large real world dataset 2>Data preparation & cleaning using pandas 3)Exploratory analysis & visualization using matplotlibb & seaborn 4)Ask and answer questions 5)Summarize and write conclusions 6)Document, Publish & Present the notebook online

Data Preparation & Cleaning:

load file using pandas

```
In [2]: import pandas as pd

data = pd.read_csv("C:\\Users\\Goudise Bharani\\Downloads\\archive (1)\\US_Accidents_msc21_updated.csv")
len(data.columns)

Out[2]: 47

look at some info about data & columns

In [3]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2845341 entries, 0 to 2845341
Data columns (total 47 columns):
 #   Column      Dtype
---  --
 0   ID           object
 1   Severity    int64
 2   Start_Time  object
 3   End_Time    object
 4   Start_Lat   float64
 5   Start_Lng   float64
 6   End_Lat     float64
 7   End_Lng     float64
 8   Distance(mi) float64
 9   Description  object
10   Number      float64
11   Street       object
12   Side        object
13   City        object
14   County      object
15   State       object
16   Zipcode     object
17   Country     object
18   Timezone    object
19   Airport_Code object
20   Weather_Timestamp object
21   Temperature(F) float64
22   Wind_Chill(F) float64
23   Humidity(%) float64
24   Pressure(in) float64
25   Visibility(mi) float64
26   Wind_Direction object
27   Wind_Speed(mph) float64
28   Precipitation(in) float64
29   Weather_Condition object
30   Amenity     bool
31   Bump        bool
32   Crossing    bool
33   Give_Way   bool
34   Junction   bool
35   No_Exit     bool
36   Railway    bool
37   Roundabout bool
38   Station    bool
39   Stop        bool
40   Traffic_Calming bool
41   Traffic_Signal bool
42   Turning_Loop bool
43   Sunrise_Sunset object
44   Civil_WhiteLight object
45   Nautilal_WhiteLight object
46   Astronomical_WhiteLight object
dtypes: bool(13), float64(13), int64(1), object(12)
memory usage: 773.4+ MB

In [4]: data.describe()

Out[4]:
```

	Severity	Start_Lat	Start_Lng	End_Lat	End_Lng	Distance(mi)	Number	Temperature(F)	Wind_Chill(F)	Humidity(%)	Pressure(in)	Visibility(mi)	Wind_Speed(mph)	Precipitation(in)
count	2845341.00	2845341.00	2845341.00	2845341.00	2845341.00	181021.00	2.770000e+06	2.770000e+06	2.772200e+06	2.787100e+06	2.777300e+06	2.887300e+06	2.285800e+06	7.291600e+00
mean	4.877721e+01	32.65024e+01	-87.71762e+01	32.65024e+01	-87.71762e+01	1.020779e+02	1.899900e+03	6.778556e+01	6.465556e+01	6.465556e+01	9.999999e+00	3.995466e+00	9.348871e+00	0.000000e+00
std	1.478717e+01	1.935379e+01	1.837183e+01	1.935379e+01	1.837183e+01	1.960361e+01	1.830000e+04	1.862353e+01	2.174070e+01	2.387451e+01	1.065286e+00	2.375466e+00	5.537454e+00	9.348871e+00
min	1.000000e+00	2.456602e+01	-1.256481e+02	2.456602e+01	-1.256481e+02	0.000000e+00	0.000000e+00	4.900000e+01	4.900000e+01	1.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00	0.000000e+00
25%	2.000000e+00	3.344071e+01	-1.180313e+02	3.344071e+01	-1.180313e+02	5.200000e+02	1.270000e+03	5.000000e+01	4.900000e+01	4.800000e+01	1.000000e+00	1.500000e+00	0.000000e+00	0.000000e+00
50%	2.000000e+00	3.629816e+01	-8.242180e+01	3.629799e+01	-8.241772e+01	2.440000e+01	4.007000e+03	6.400000e+01	6.300000e+01	6.700000e+01	1.000000e+00	2.000000e+00	0.000000e+00	0.000000e+00
75%	2.000000e+00	4.191604e+01	-8.037243e+01	4.016105e+01	-8.037338e+01	7.640000e+01	8.567000e+03	7.600000e+01	7.600000e+01	8.300000e+01	1.000000e+00	1.000000e+00	0.000000e+00	0.000000e+00
max	4.000000e+00	4.900050e+01	-6.711317e+01	4.907050e+01	-6.710914e+01	1.551860e+02	9.999971e+06	1.960000e+02	1.960000e+02	1.000000e+02	5.890000e+01	1.400000e+02	1.887000e+03	2.400000e+00

```
In [5]: numeric = ['float64','int64']
new = data.select_dtypes(include=numeric)
len(new.columns)

Out[5]: 14

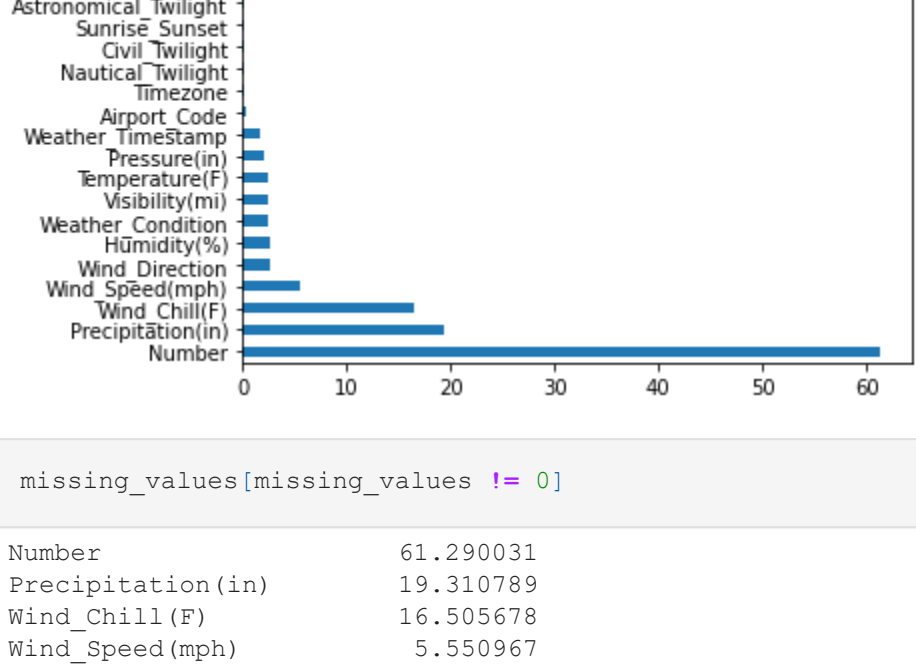
Fix any missing values or incorrect values
```

```
In [6]: missing_values = data.isna().sum().sort_values(ascending=False)/len(data)*100
missing_value

Out[6]:
Number                61.290031
Precipitation(in)     19.310789
Wind_Chill(F)         16.300678
Wind_Speed(mph)       5.550987
Wind_Direction        2.559284
Humidity(%)           2.568830
Weather_Condition     2.482514
Visibility(mi)         2.478250
Temperature(F)        2.436466
Pressure(in)          2.400553
Weather_Timestamp     1.783125
Airport_Code          0.350601
Timezone              0.128956
Nautilal_WhiteLight  0.100761
Civil_WhiteLight      0.100761
Sunrise_Sunset        0.100761
Astronomical_WhiteLight 0.100761
Zipcode              0.046356
City                 0.048815
Street               0.000070
Country              0.000000
Junction              0.000000
Chave_Time            0.000000
End_Time              0.000000
Start_Lat             0.000000
Turning_Loop          0.000000
Traffic_Signal        0.000000
Traffic_Calming       0.000000
Stop                  0.000000
Station               0.000000
Roundabout           0.000000
Railway              0.000000
No_Exit               0.000000
Give_Way              0.000000
Amenity               0.000000
Junction              0.000000
Start_Lng             0.000000
End_Lng               0.000000
End_Lat               0.000000
Distance(mi)          0.000000
Description           0.000000
Severity              0.000000
Side                  0.000000
City                  0.000000
State                 0.000000
ID                    0.000000
dtype: float64

In [7]: missing_values[missing_values != 0].plot(kind='bar')

Out[7]: <AxesSubplot>
```



```
In [8]: missing_values[missing_values != 0]

Out[8]:
Number                61.290031
Precipitation(in)     19.310789
Wind_Chill(F)         16.300678
Wind_Speed(mph)       5.550987
Wind_Direction        2.559284
Humidity(%)           2.568830
Weather_Condition     2.482514
Visibility(mi)         2.478250
Temperature(F)        2.436466
Pressure(in)          2.400553
Weather_Timestamp     1.783125
Airport_Code          0.350601
Timezone              0.128956
Nautilal_WhiteLight  0.100761
Civil_WhiteLight      0.100761
Sunrise_Sunset        0.100761
Astronomical_WhiteLight 0.100761
Zipcode              0.046356
City                 0.048815
Street               0.000070
dtype: float64

Exploratory Analysis & Visualization:

Pick up columns for analysis:

1)City 2)Start time 3)Start latitude 4)Start longitude 5)Temperature

City:

In [9]: len(data['City'].unique())

Out[9]: 11482

In [10]: cities_by_accidents = data['City'].value_counts()
cities_by_accidents

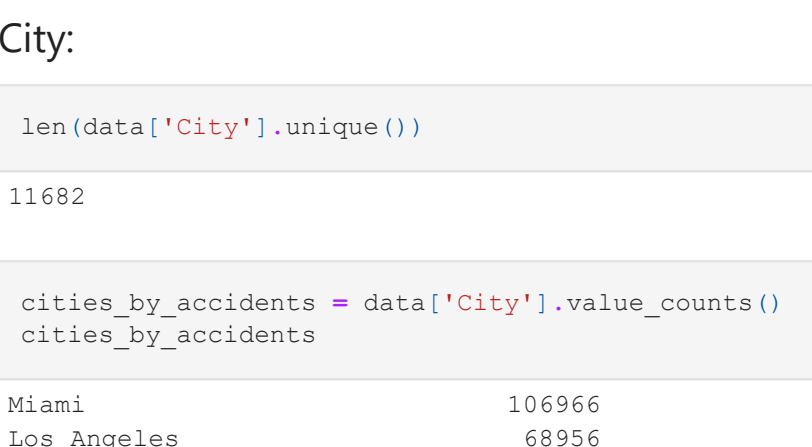
Out[10]:
Miami                106966
Los Angeles          69956
Orlando              54691
Dallas               41979
Houston              39448
Ridgdale              ...
Seki                  1
Westridge             1
Bullock               1
American_Park_Pleasant_Grove 1
Name: City, Length: 11482, dtype: int64

In [11]: cities_by_accidents[[10]]

Out[11]:
Miami                106966
Los Angeles          69956
Orlando              54691
Dallas               41979
Houston              39448
Charlotte            33152
Sacramento           32559
San_Diego             26827
Raleigh              22860
Minneapolis           22768
Name: City, dtype: int64

In [12]: cities_by_accidents[[10]].plot(kind='bar')

Out[12]: <AxesSubplot>
```

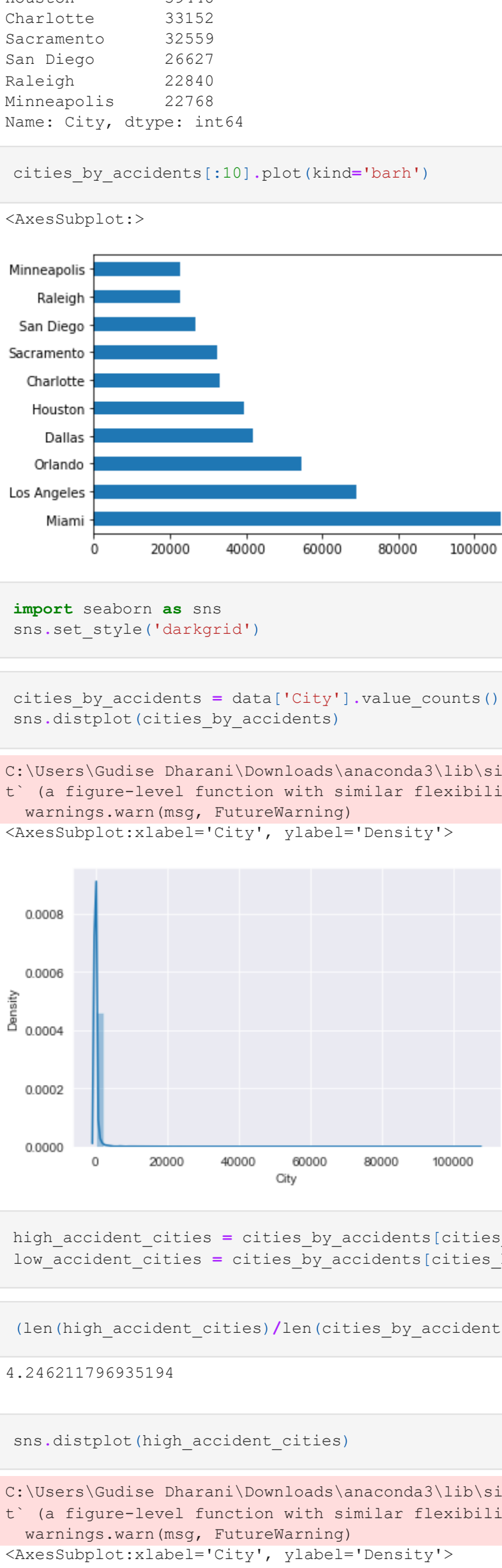


```
In [13]: import seaborn as sns
sns.set_style('darkgrid')

In [14]: cities_by_accidents = data['City'].value_counts()
sns.distplot(cities_by_accidents)

C:\Users\Goudise Bharani\Downloads\anaconda3\lib\site-packages\seaborn\distributions.py:2618: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot'
('a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:label='City', ylabel='Density'>

Out[14]:
```



```
In [15]: high_accident_cities = cities_by_accidents[cities_by_accidents >= 1000]
low_accident_cities = cities_by_accidents[cities_by_accidents < 1000]

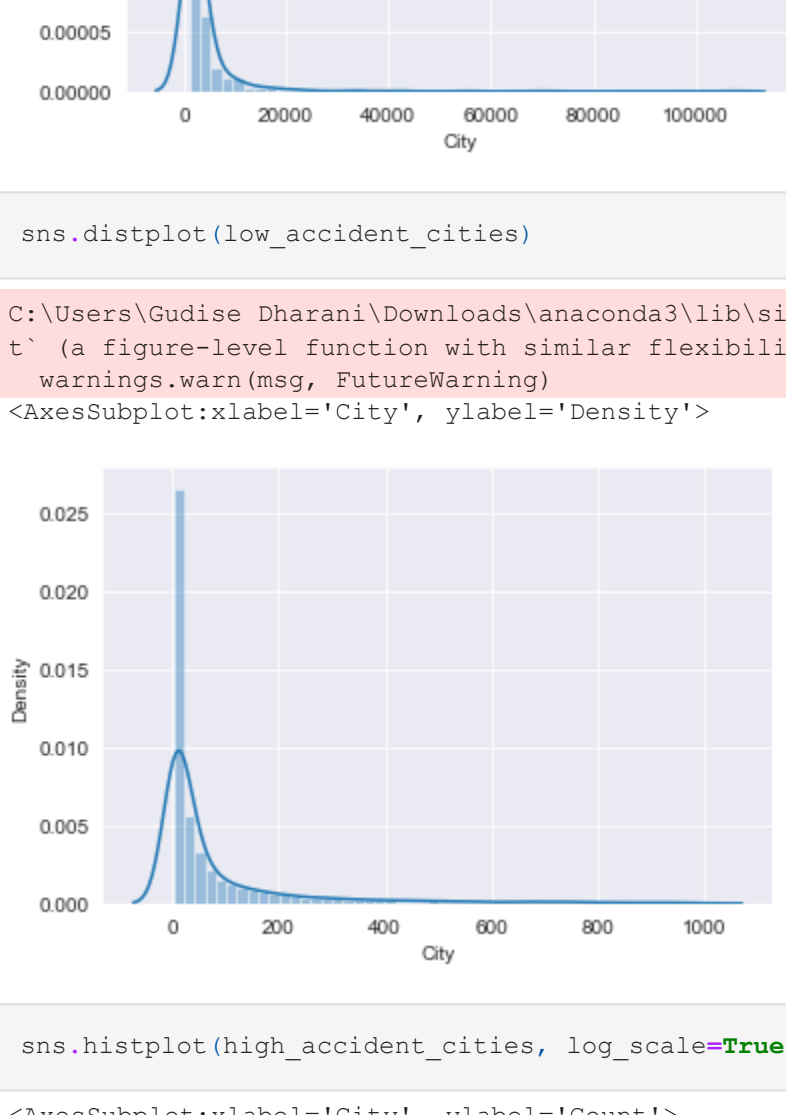
In [16]: (len(high_accident_cities)/len(cities_by_accidents))*100

Out[16]: 4.24621759335194

In [17]: sns.distplot(high_accident_cities)

C:\Users\Goudise Bharani\Downloads\anaconda3\lib\site-packages\seaborn\distributions.py:2618: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot'
('a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:label='City', ylabel='Density'>

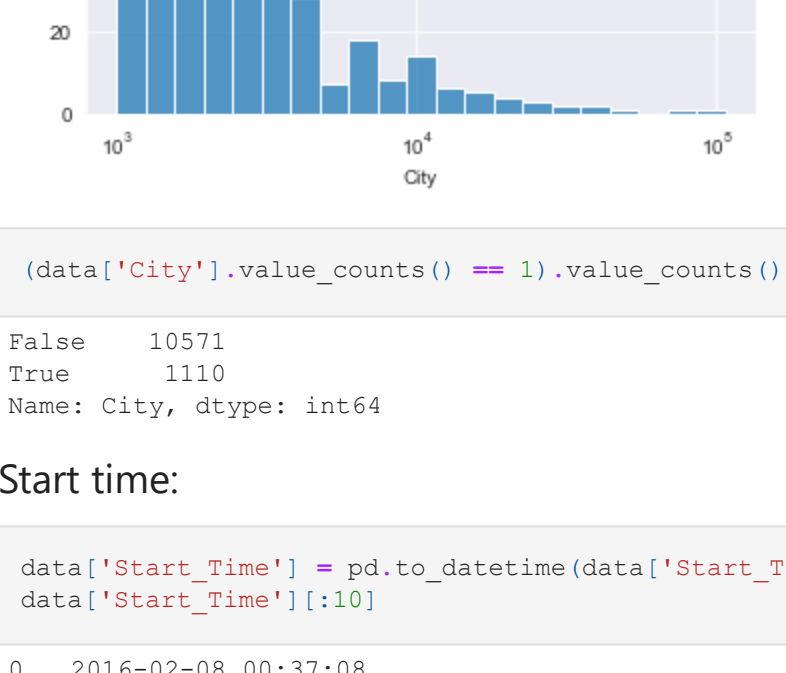
Out[17]:
```



```
In [18]: sns.distplot(low_accident_cities)

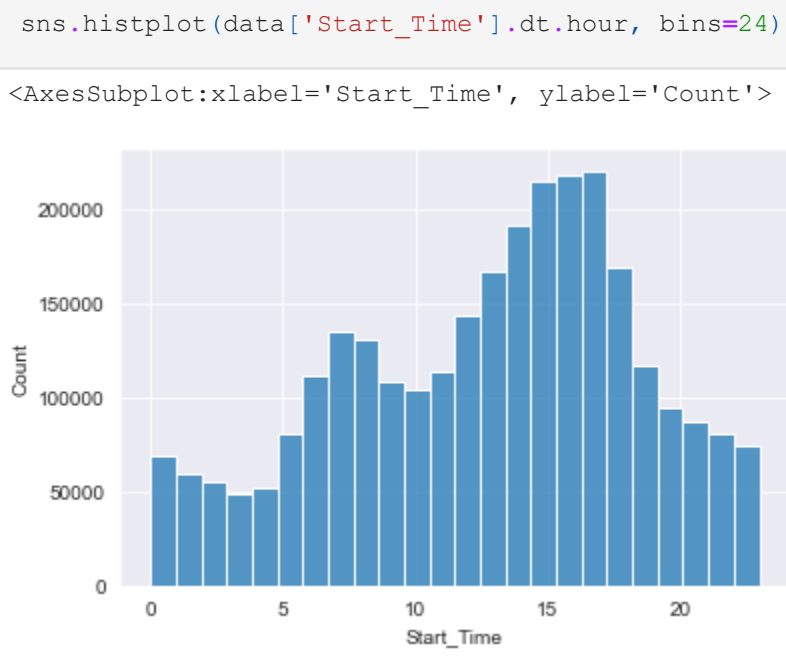
C:\Users\Goudise Bharani\Downloads\anaconda3\lib\site-packages\seaborn\distributions.py:2618: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot'
('a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:label='City', ylabel='Density'>

Out[18]:
```



```
In [19]: sns.histplot(high_accident_cities, log_scale=True)

Out[19]: <AxesSubplot:label='City', ylabel='Count'>
```



```
In [20]: (data['City'].value_counts()[0:5].value_counts())

Out[20]:
True      10571
False     1110
Name: City, dtype: int64

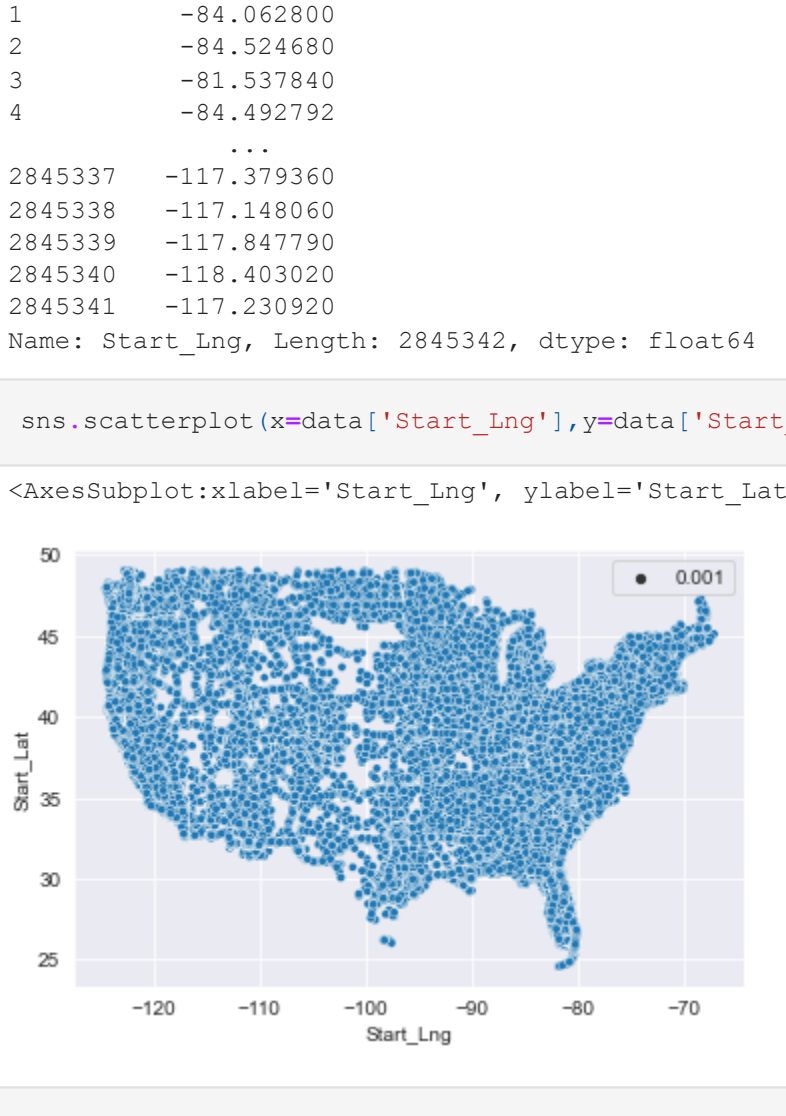
Start time:
```

```
In [21]: data['Start_Time'] = pd.to_datetime(data['Start_Time'])
data['Start_Time'].dt.hour

Out[21]:
0    2016-02-08 00:37:08
1    2016-02-08 05:58:00
2    2016-02-08 06:15:39
3    2016-02-08 06:53:45
4    2016-02-08 07:53:43
5    2016-02-08 08:16:57
6    2016-02-08 09:15:41
7    2016-02-08 11:51:46
8    2016-02-08 14:19:57
9    2016-02-08 15:16:43
Name: Start_Time, dtype: datetime64[ns]

In [22]: sns.histplot(data['Start_Time'].dt.hour, bins=24)

Out[22]: <AxesSubplot:label='Start_Time', ylabel='Count'>
```



Start Latitude & Start Longitude:

```
In [23]: data['Start_Lat']

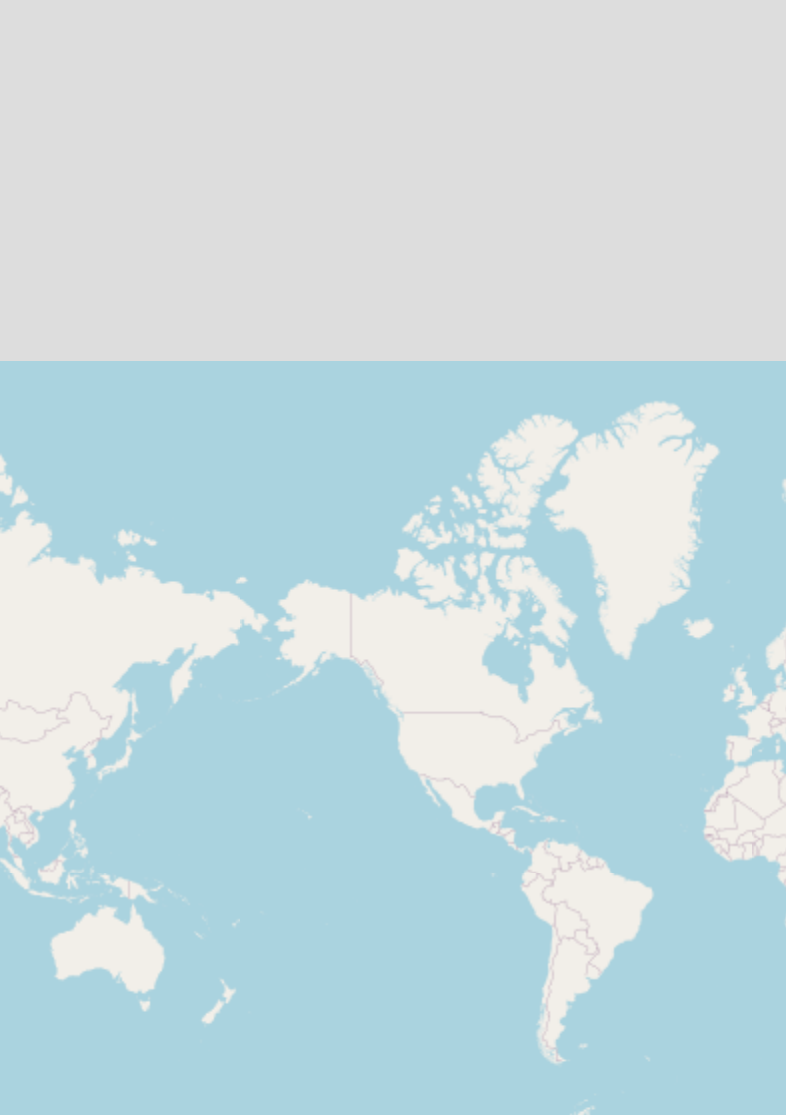
Out[23]:
0    40.108910
1    39.865420
2    39.102660
3    41.062130
4    39.172393
2845337    ...
2845338    32.766860
2845339    22.775450
2845340    33.992460
2845341    34.139350
Name: Start_Lat, Length: 2845342, dtype: float64

In [24]: data['Start_Lng']

Out[24]:
0    -81.092880
1    -81.062800
2    -80.524880
3    -81.537840
4    -84.492792
2845337    -117.378360
2845338    -117.140560
2845339    -117.847790
2845340    -118.480020
2845341    -117.203920
Name: Start_Lng, Length: 2845342, dtype: float64

In [25]: sns.scatterplot(x=data['Start_Time'],y=data['Start_Lat'], size=0.001)

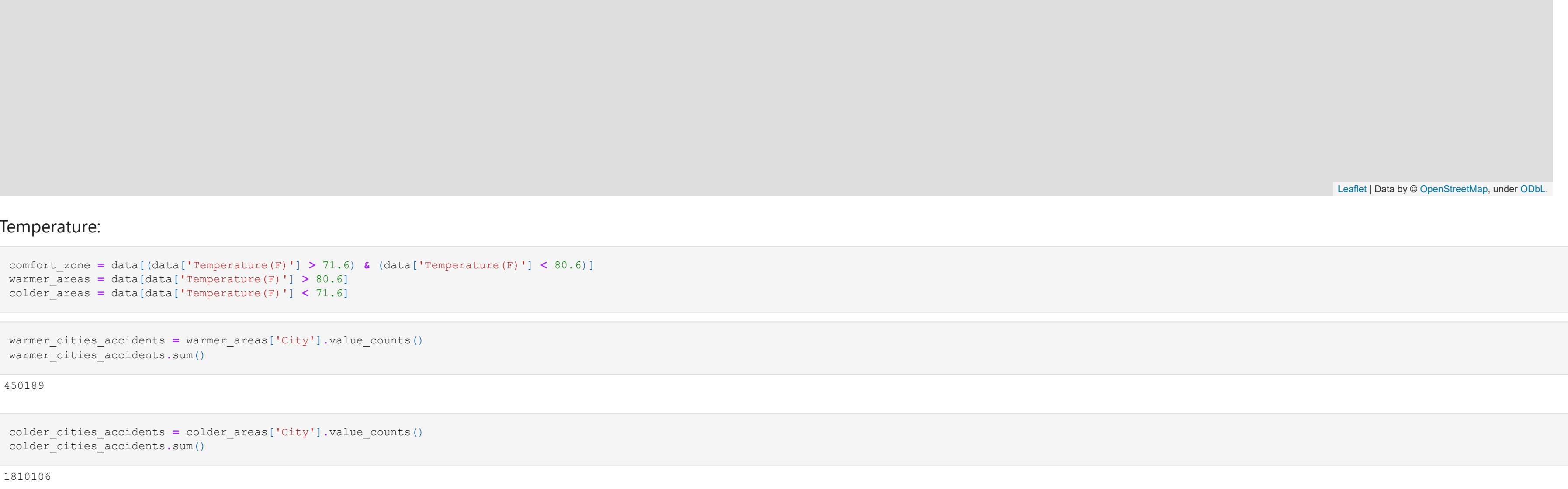
Out[25]: <AxesSubplot:label='Start_Lng', ylabel='Start_Lat'>
```



```
In [26]: import folium
from folium.plugins import HeatMap

In [27]: map = folium.Map()
HeatMap(list(zip(list(data['Start_Lat']), list(data['Start_Lng']))))).add_to(map)
map

Out[27]:
```



```
Temperature:

In [28]: warmer_zone = data[(data['Temperature(F)'] > 71.6) & (data['Temperature(F)'] < 80.6)]
warmer_area = data[(data['Temperature(F)'] > 80.6)]
colider_area = data[(data['Temperature(F)'] < 71.6)]

In [29]: warmer_cities_accidents = warmer_area['City'].value_counts()
warmer_cities_accidents.sum()

Out[29]: 450189

In [30]: colider_cities_accidents = colider_area['City'].value_counts()
colider_cities_accidents.sum()

Out[30]: 1810106

Ask & Answer questions:
```

Are there more accidents in warmer or colder areas?

Ans) Colder areas have more accidents compared to warmer areas CODE Above mentioned

Which 5 states have highest no of accidents?

Ans) Miami, Los Angeles, Orlando, Dallas, Houston CODE Above mentioned

Does New York show up in the data? If yes, why is the count lower if it is one of the most populated city?

Ans) No, the dataset consists of only 49 states. There is no data about New York.

```
In [31]: 'NY' in data

Out[31]: False

What time of the day are accidents most frequent in?

Ans) From 2 pm to 5 pm accidents are highest and from 6 am to 10 am they are second highest CODE Above mentioned



Which days of the week have most accidents?

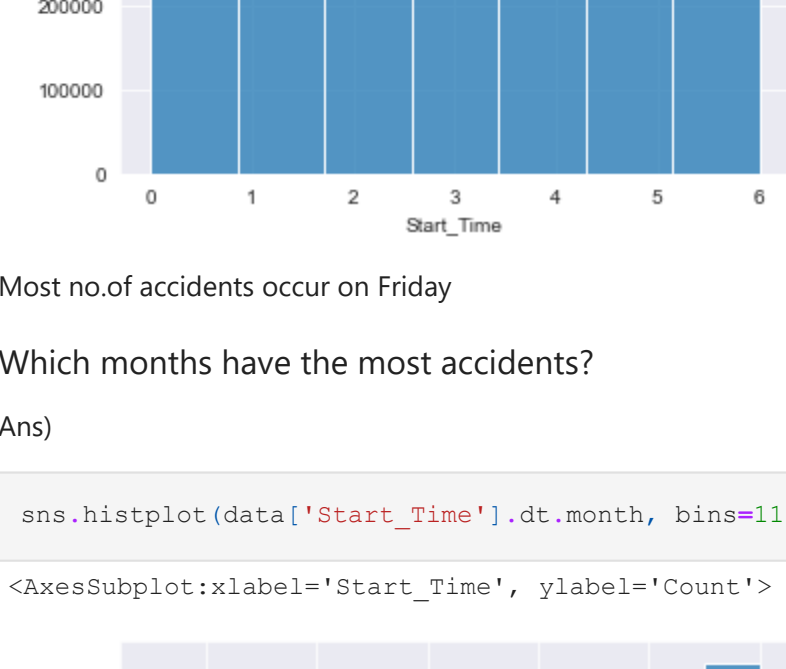


Ans)


```

```
In [32]: sns.histplot(data['Start_Time'].dt.dayofweek, bins=7)

Out[32]: <AxesSubplot:label='Start_Time', ylabel='Count'>
```



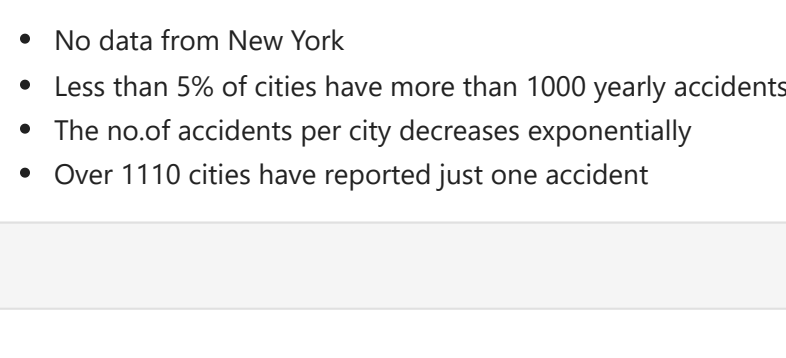
Most no of accidents occur on Friday

Which months have the most accidents?

Ans)

```
In [33]: sns.histplot(data['Start_Time'].dt.month, bins=11)

Out[33]: <AxesSubplot:label='Start_Time', ylabel='Count'>
```



Most of the accidents occur in December, whereas November is having second highest no of accidents. But there is no data regarding January in the dataset.

Summary and Conclusion:

Insights:

- No data from New York
- Less than 5% of cities have more than 1000 yearly accidents
- The no of accidents per city decreases exponentially
- Over 1110 cities have reported just one accident