

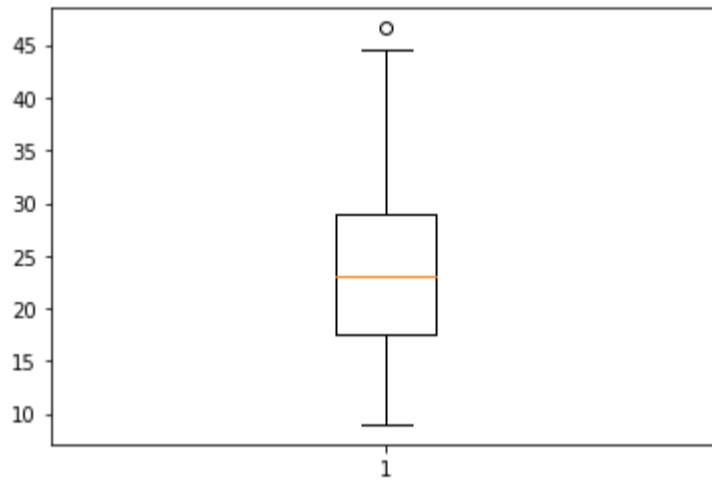
## Chap.-2 Advance Visualization

1. Box Plot (boxplot) - to identify outlier (a value that doesn't given)  
notch  
vert  
widths  
patch\_artist
2. Scatter (scatter) - relation between 2 var or 2 attribute  
c='r'  
marker='D'  
edgecolor='g'  
linewidths=2  
alpha=0.8
3. Heatmap -
4. Wafflechart - alternate of piechat
5. Regression plot.
6. Wordcloud
7. Area Plot (Area Charts)
8. Geospatial data : latitude and longitudinal

### Boxplot - to identify outlier (a value that doesn't given)

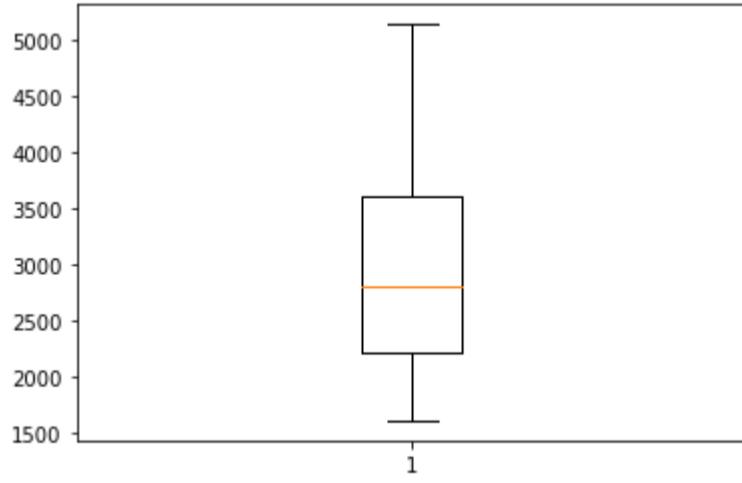
In [1]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 df = pd.read_csv('auto-mpg.csv')
4 plt.boxplot(df['mpg'])
5 plt.show()
```



In [2]:

```
1 df = pd.read_csv('auto-mpg.csv')
2 plt.boxplot(df['weight'])
3 plt.show()
```



In [3]:

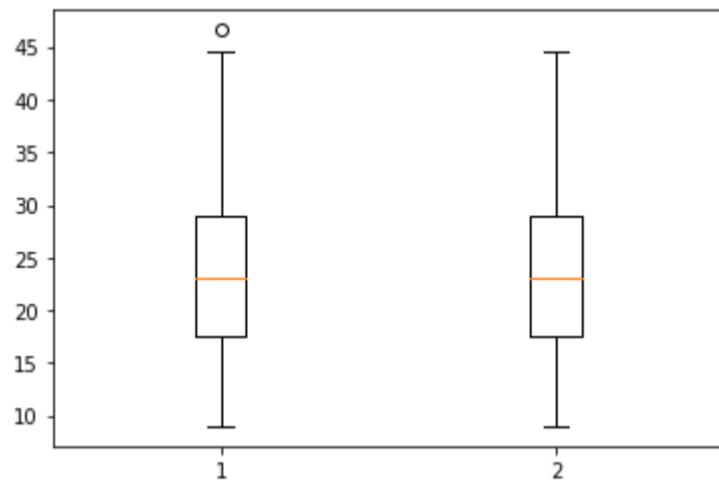
```

1 df = pd.read_csv('auto-mpg.csv')
2
3 def remove_outliers(df,col):
4     q1 = df[col].quantile(0.25)
5     q3 = df[col].quantile(0.75)
6     iqr = q3-q1
7     lower = q1-1.5*iqr
8     upper = q3+1.5*iqr
9     return (df[(df[col]>=lower)&(df[col]<=upper)])
10 ndf = remove_outliers(df, 'mpg')
11 plt.boxplot([df['mpg'],ndf['mpg']])

```

Out[3]:

```
{
'whiskers': [<matplotlib.lines.Line2D at 0x1e4cef4cc0>,
<matplotlib.lines.Line2D at 0x1e4cef59070>,
<matplotlib.lines.Line2D at 0x1e4cef654f0>,
<matplotlib.lines.Line2D at 0x1e4cef65850>],
'caps': [<matplotlib.lines.Line2D at 0x1e4cef593d0>,
<matplotlib.lines.Line2D at 0x1e4cef59730>,
<matplotlib.lines.Line2D at 0x1e4cef65bb0>,
<matplotlib.lines.Line2D at 0x1e4cef65f10>],
'boxes': [<matplotlib.lines.Line2D at 0x1e4cef4c970>,
<matplotlib.lines.Line2D at 0x1e4cef65190>],
'medians': [<matplotlib.lines.Line2D at 0x1e4cef59a90>,
<matplotlib.lines.Line2D at 0x1e4cef712b0>],
'fliers': [<matplotlib.lines.Line2D at 0x1e4cef59df0>,
<matplotlib.lines.Line2D at 0x1e4cef71610>],
'means': []}
```

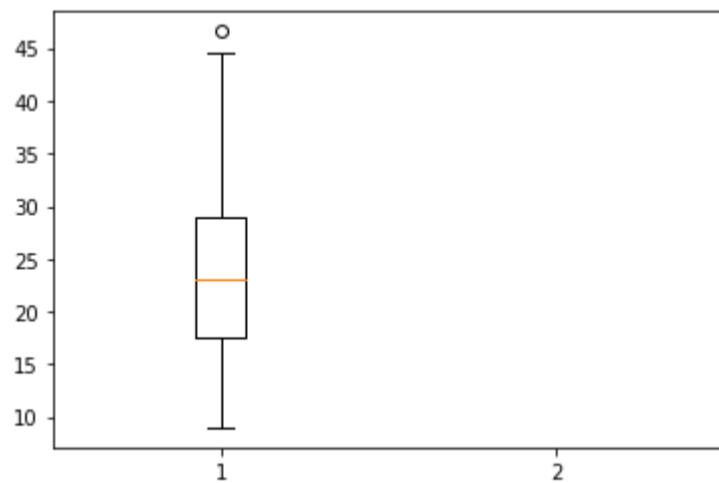


In [4]:

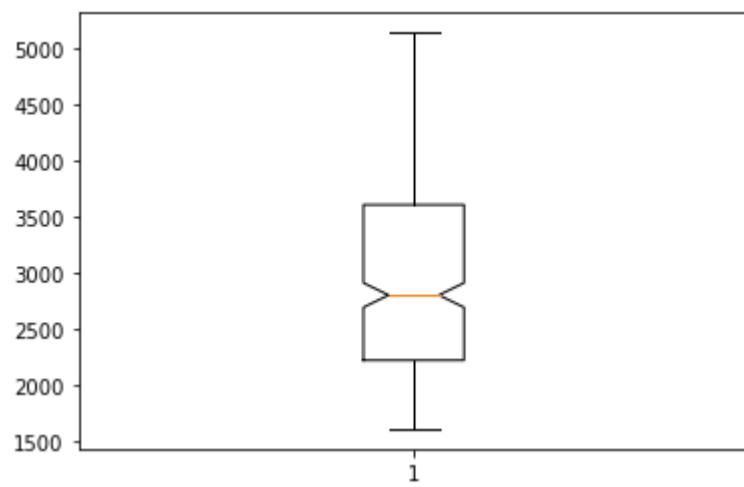
```

1 df = pd.read_csv('auto-mpg.csv')
2
3 def remove_outliers(df,col):
4     q1 = df[col].quantile(0.25)
5     q3 = df[col].quantile(0.75)
6     iqr = q3-q1
7     lower = q1-1.5*iqr
8     upper = q3+1.5*iqr
9     return (df[(df[col]<lower)&(df[col]>upper)])
10 ndf = remove_outliers(df, 'mpg')
11 plt.boxplot([df['mpg'],ndf['mpg']])
12 plt.show()

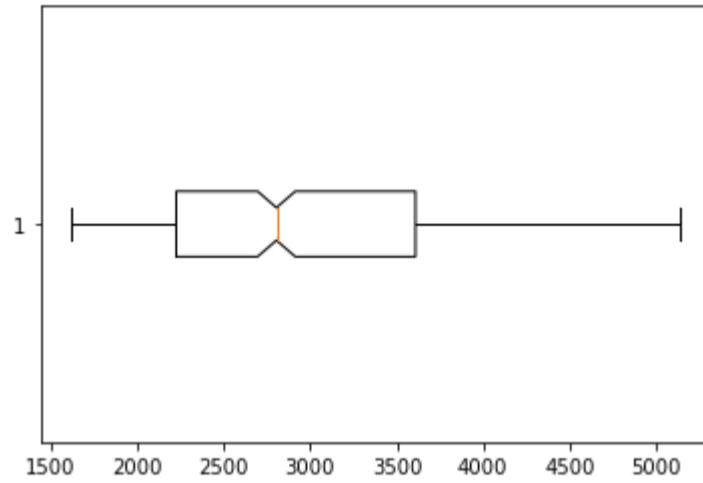
```



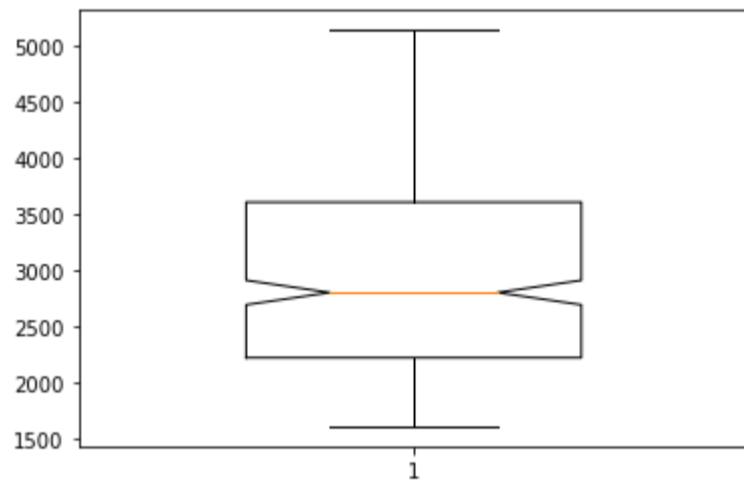
```
In [5]: 1 df = pd.read_csv('auto-mpg.csv')
2 plt.boxplot(df['weight'], notch=True)
3 plt.show()
```



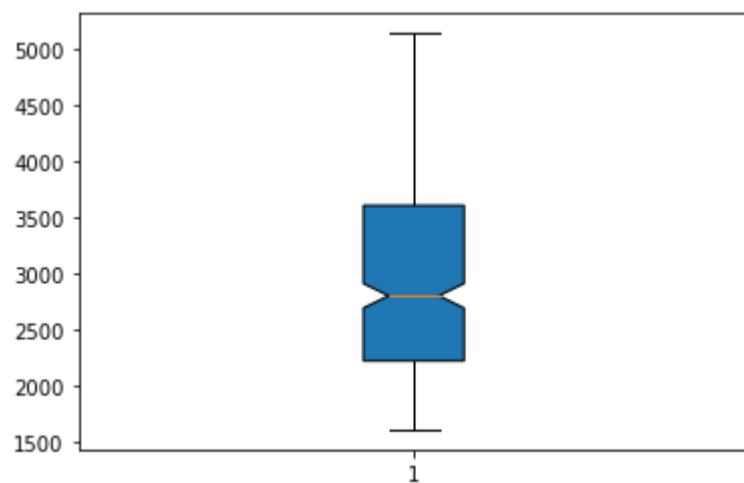
```
In [6]: 1 df = pd.read_csv('auto-mpg.csv')
2 plt.boxplot(df['weight'], notch=True, vert=False)
3 plt.show()
```



```
In [7]: 1 df = pd.read_csv('auto-mpg.csv')
2 plt.boxplot(df['weight'], notch=True, widths=0.5)
3 plt.show()
```



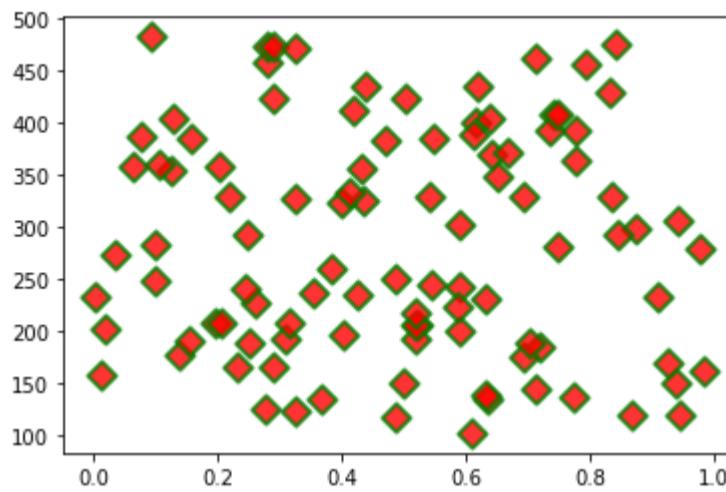
```
In [8]: 1 df = pd.read_csv('auto-mpg.csv')
2 plt.boxplot(df['weight'], notch=True, patch_artist=True)
3 plt.show()
```



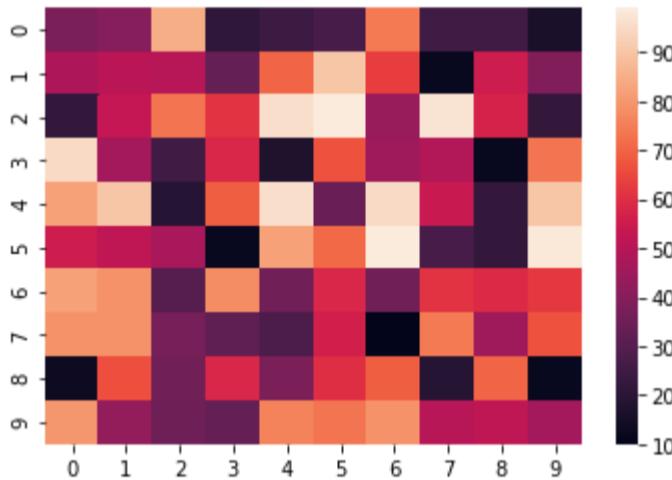
## Scatter - relation between 2 var or 2 attribute

```
In [9]: 1 import numpy as np
2 x = np.random.rand(100)
3 y = np.random.randint(low=100,high=500,size=100)
4 plt.scatter(x,y,s=100,c='r',marker='D',edgecolor='g',linewidths=2,alpha=0.8)
```

Out[9]: <matplotlib.collections.PathCollection at 0x1e4ceb14be0>



```
In [10]: 1 import seaborn as sn
2 import numpy as np
3 import matplotlib.pyplot as plt
4 data = np.random.randint(low=10,high=100,size=(10,10))
5 sn.heatmap(data)
6 plt.show()
```



```
In [11]: 1 import seaborn as sn
2 import numpy as np
3 import matplotlib.pyplot as plt
4 data=np.random.randint(low=10,high=100,size=(10,10))
5 sn.heatmap(data,annot=True)#annot=Display number
6 plt.show()
```



In [12]:

```

1 import seaborn as sn
2 import numpy as np
3 import matplotlib.pyplot as plt
4 data=np.random.randint(low=10,high=100,size=(10,10))
5 sn.heatmap(data,annot=True,vmin=10,vmax=50,cbar=False)#cbar =False remove cbar
6 plt.show()

```

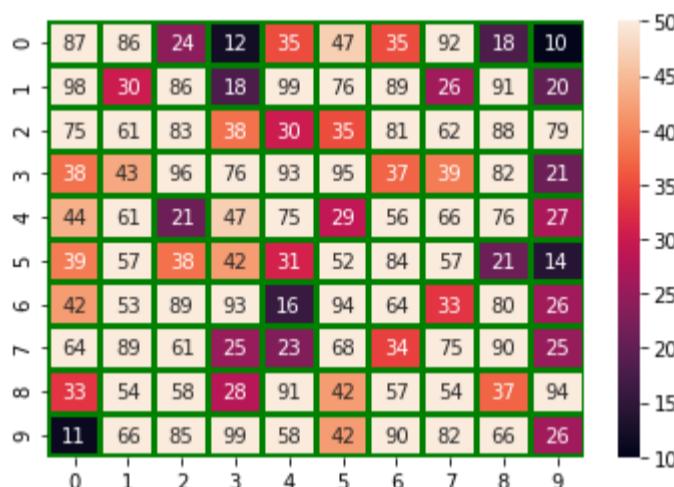


In [13]:

```

1 import seaborn as sn
2 import numpy as np
3 import matplotlib.pyplot as plt
4 data=np.random.randint(low=10,high=100,size=(10,10))
5 sn.heatmap(data,annot=True,vmin=10,vmax=50,linewidths=2)
6 plt.show()

```



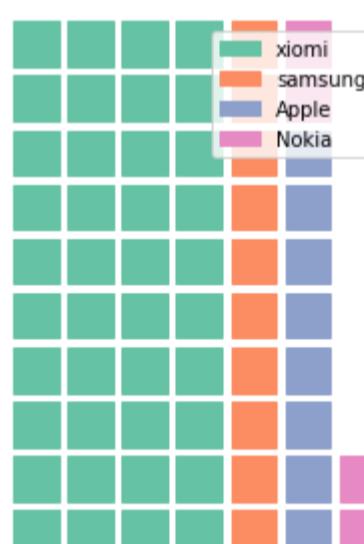
## Wafflechart - alternate of piechat

In [18]:

```

1 import pandas as pd
2 from pywaffle import Waffle
3 import matplotlib.pyplot as plt
4
5 df=pd.read_csv('auto-mpg.csv')
6 df=pd.DataFrame(data)
7 data={'phone': ['xiomi', 'samsung', 'Apple', 'Nokia'], 'stock':[40,10,8,4]}
8 new_fig=plt.figure(FigureClass=Waffle,rows=10,values=df.stock,labels=list(df.phone))

```



In [15]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from pywaffle import Waffle
4 data={'phone':['samsung','iphone','nokia','redmi','vivo','oppo','google','microsoft'],
5 'stock':[44,35,67,89,10,69,45,34]}
6 df=pd.DataFrame(data)
7 df.head() #It prints first five rows
8 fig=plt.figure(FigureClass=Waffle,rows=10,values=df.stock,labels=list(df.phone))

```

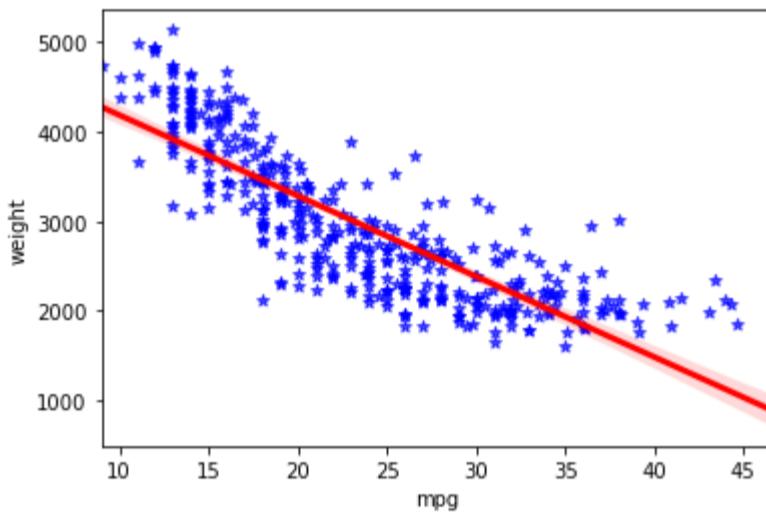


In [19]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 df=pd.read_csv('auto-mpg.csv')
5 sns.regplot(y=df['weight'],x=df['mpg'],marker='*',
6             scatter_kws={'color':'blue','alpha':0.7},
7             line_kws={'color':'red','linewidth':3})
8 plt.show()

```

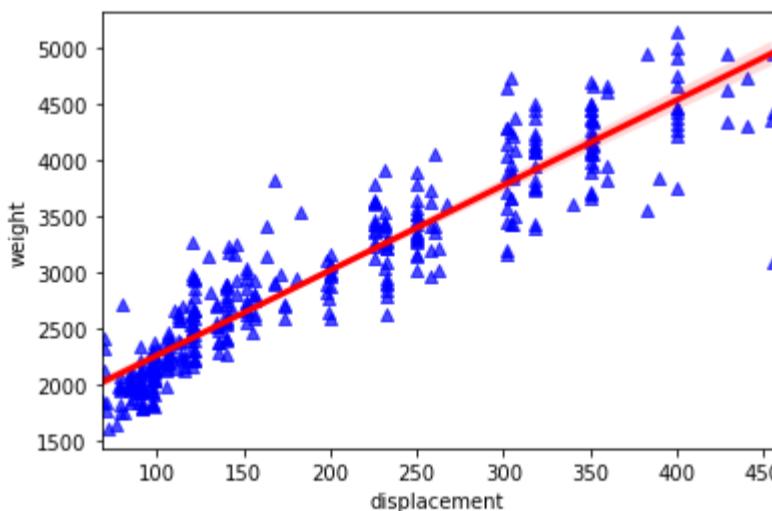


In [20]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 df=pd.read_csv('auto-mpg.csv')
5 sns.regplot(y=df['weight'],x=df['displacement'],marker='^',
6             scatter_kws={'color':'blue','alpha':0.7},
7             line_kws={'color':'red','linewidth':3})
8 plt.show()

```



In [ ]:

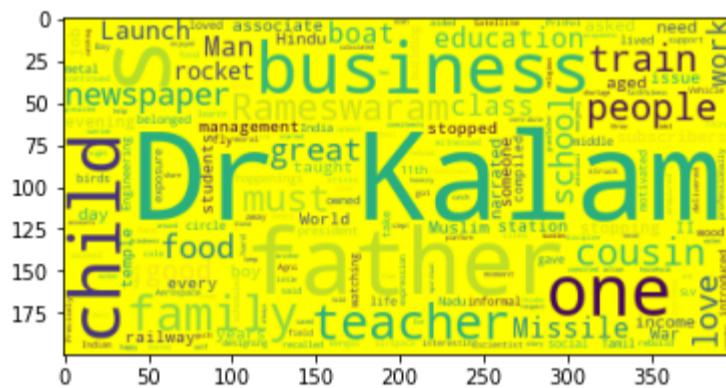
```
1 pip install wordcloud
```

```
In [21]: 1 from wordcloud import WordCloud, STOPWORDS  
2 stopwords = set(STOPWORDS)  
3 print(stopwords)  
4 # stopwords will not give any meaning it is only used for joining as suppose "the" we are using much so it  
5 # it should be stopped and removed from respective textfile
```

{"how's", 'a', 'are', 'because', "wouldn't", 'over', 'yourselves', 'otherwise', 'himself', 'an', "they're", 'where', "mustn't", 'am', 'is', 'themselves', "when's", 'between', 'each', "that's", 'further', "you'll", 'also', "we'll", "they'll", 'who', 'some', "he'd", 'what', 'to', 'nor', 'all', "he'll", 'do', 'would', 'the', 'whom', 'than', "she'll", 'into', 'hers', 'had', 'ours', "hadn't", 'these', 'its', 'few', "we'd", 'again', 'but', 'theirs', 'this', 'herself', 'why', 'my', "haven't", 'on', "what's", "we've", 'only', 'out', "we're", "here's", 'has', "you're", 'there', "can't", 'doing', 'can', 'other', 'she', 'before', 'being', 'it', 'such', 'most', "she's", 'you', 'during', 'i', 'below', 'for', "weren't", 'their', 'own', "there's", 'them', "i'm", "they'd", 'which', 'same', "i'll", 'above', 'however', 'that', 'http', "shouldn't", 'with', 'him', 'those', 'yourself', 'hence', 'until', 'as', 'did', 'ever', 'from', 'while', "shan't", "aren't", 'any', 'so', "doesn't", 'under', "don't", 'were', 'up', 'should', 'cannot', 'if', 'once', 'com', 'have', "where's", 'he', 'www', "you'd", 'therefore', "wasn't", 'after', 'too', 'ought', "isn't", 'k', 'how', 'against', 'here', 'get', "hasn't", 'not', 'else', "it's", 'at', 'of', 'itself', 'me', 'was', 'having', 'your', "why's", 'they', 'r', 'both', 'be', 'her', "couldn't", 'our', 'by', 'or', 'very', 'his', 'shall', 'just', "who's", 'and', 'off', 'yours', 'been', 'does', 'could', 'myself', 'then', 'when', 'we', "she'd", 'in', 'about', 'ourselves', 'since', "he's", "let's", "you've", 'more', 'no', "won't", 'down', "didn't", 'like', "i've", 'through', 'i'd", "they've"}

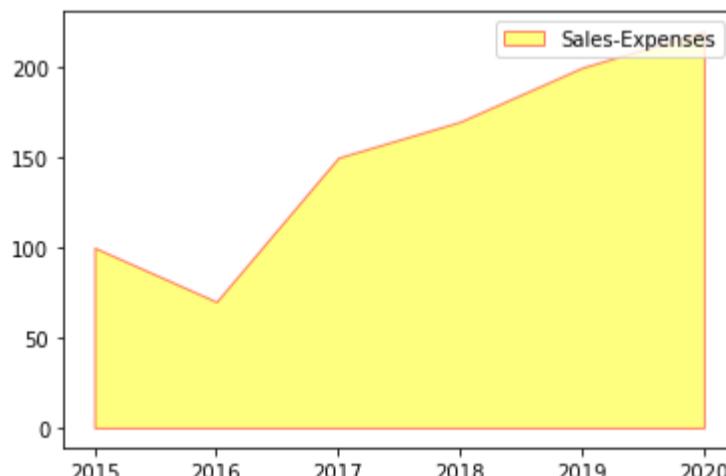
```
In [22]: 1 f = open('demo.txt', 'r', encoding='utf-8')
2 apj = f.read()
3 wc = WordCloud(background_color='yellow', max_words=200, stopwords=set(STOPWORDS))
4 #max_words to limit the words
5 wc.generate(apj)
6 plt.imshow(wc)
```

**Out[22]:** <matplotlib.image.AxesImage at 0x1e4d15cd340>



## Area Plot (Area Charts)

```
In [23]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 data = {'Year':[2015,2016,2017,2018,2019,2020],
5         'Sales':[100,70,150,170,200,220],
6         'Expenses':[70,80,90,100,110,120]}
7 df = pd.DataFrame(data)
8 plt.fill_between(df['Year'],df['Sales'],
9                  facecolor='yellow',
10                 edgecolor='red',
11                 alpha=0.5,
12                 label='Sales-Expenses')
13
14 plt.legend()
15 plt.show()
```

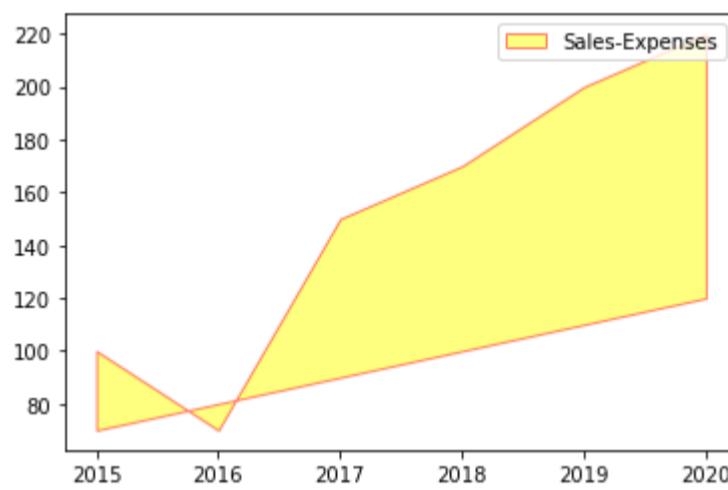


In [25]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 data = {'Year':[2015,2016,2017,2018,2019,2020],
5         'Sales':[100,70,150,170,200,220],
6         'Expenses':[70,80,90,100,110,120]}
7 df = pd.DataFrame(data)
8 plt.fill_between(df['Year'],df['Sales'],df['Expenses'],
9                   facecolor='yellow',
10                  edgecolor='red',
11                  alpha=0.5,
12                  label='Sales-Expenses')
13
14 plt.legend()
15 plt.show()

```



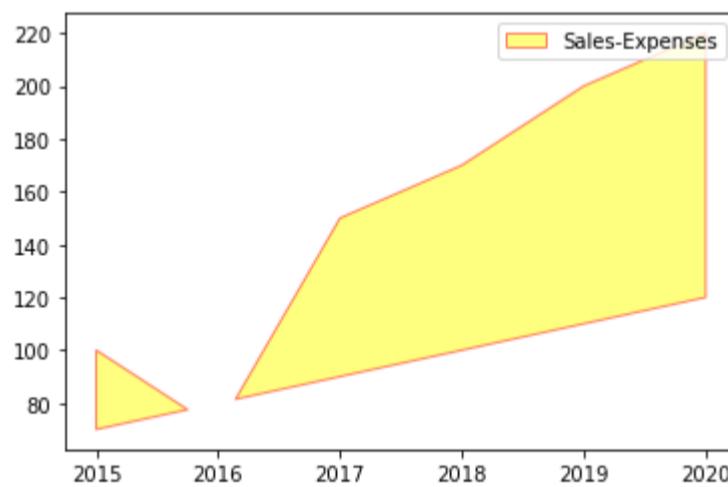
- **interpolate - to indentify intersection**

In [26]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 data = {'Year':[2015,2016,2017,2018,2019,2020],
5         'Sales':[100,70,150,170,200,220],
6         'Expenses':[70,80,90,100,110,120]}
7 df = pd.DataFrame(data)
8 plt.fill_between(df['Year'],df['Sales'],df['Expenses'],
9                   where=df['Sales']>df['Expenses'],interpolate=True,
10                  facecolor='yellow',
11                  edgecolor='red',
12                  alpha=0.5,
13                  label='Sales-Expenses')
14
15 plt.legend()
16 plt.show()

```

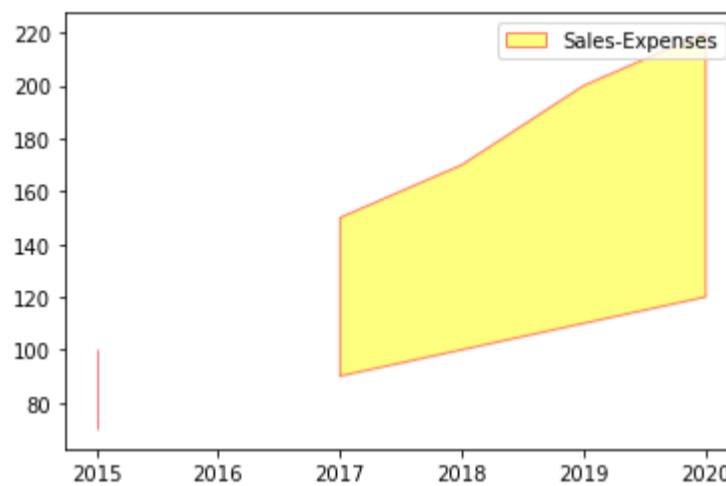


In [27]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 data = {'Year':[2015,2016,2017,2018,2019,2020],
5         'Sales':[100,70,150,170,200,220],
6         'Expenses':[70,80,90,100,110,120]}
7 df = pd.DataFrame(data)
8 plt.fill_between(df['Year'],df['Sales'],df['Expenses'],
9                   where=df['Sales']>df['Expenses'],
10                  facecolor='yellow',
11                  edgecolor='red',
12                  alpha=0.5,
13                  label='Sales-Expenses')
14
15 plt.legend()
16 plt.show()

```



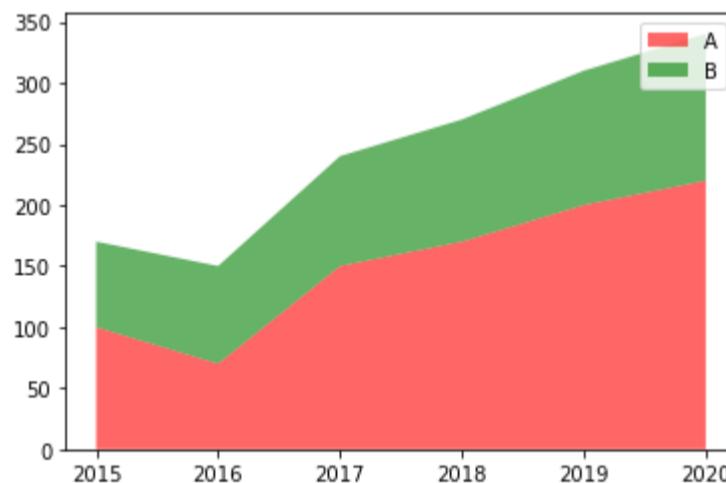
## Stackplot - A variation of Area Plot

In [28]:

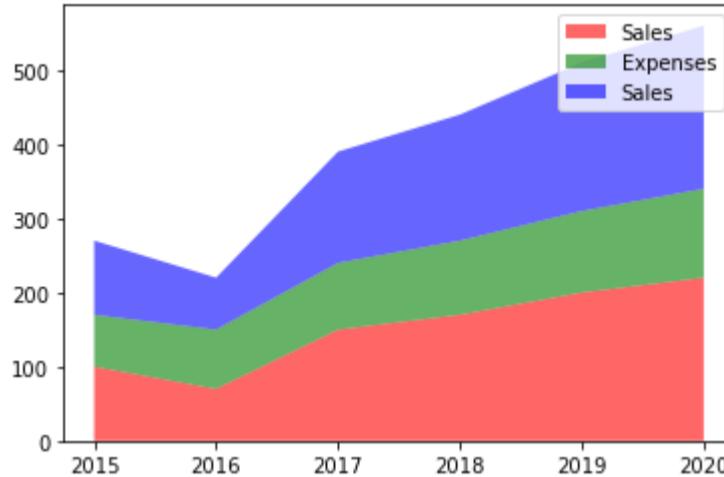
```

1 data = {'Year':[2015,2016,2017,2018,2019,2020],
2         'Sales':[100,70,150,170,200,220],
3         'Expenses':[70,80,90,100,110,120]}
4 df = pd.DataFrame(data)
5 plt.stackplot(df['Year'],df['Sales'],df['Expenses'],
6                 labels=['A','B'],
7                 colors=['r','g'],
8                 alpha=0.6)
9
10 plt.legend()
11 plt.show()

```



```
In [29]: 1 data = {'Year':[2015,2016,2017,2018,2019,2020],
2           'Sales':[100,70,150,170,200,220],
3           'Expenses':[70,80,90,100,110,120]}
4 df = pd.DataFrame(data)
5 plt.stackplot(df['Year'],df['Sales'],df['Expenses'],df['Sales'],
6                 labels=['Sales','Expenses','Sales'],
7                 colors=['r','g','b'],
8                 alpha=0.6)
9
10 plt.legend()
11 plt.show()
```

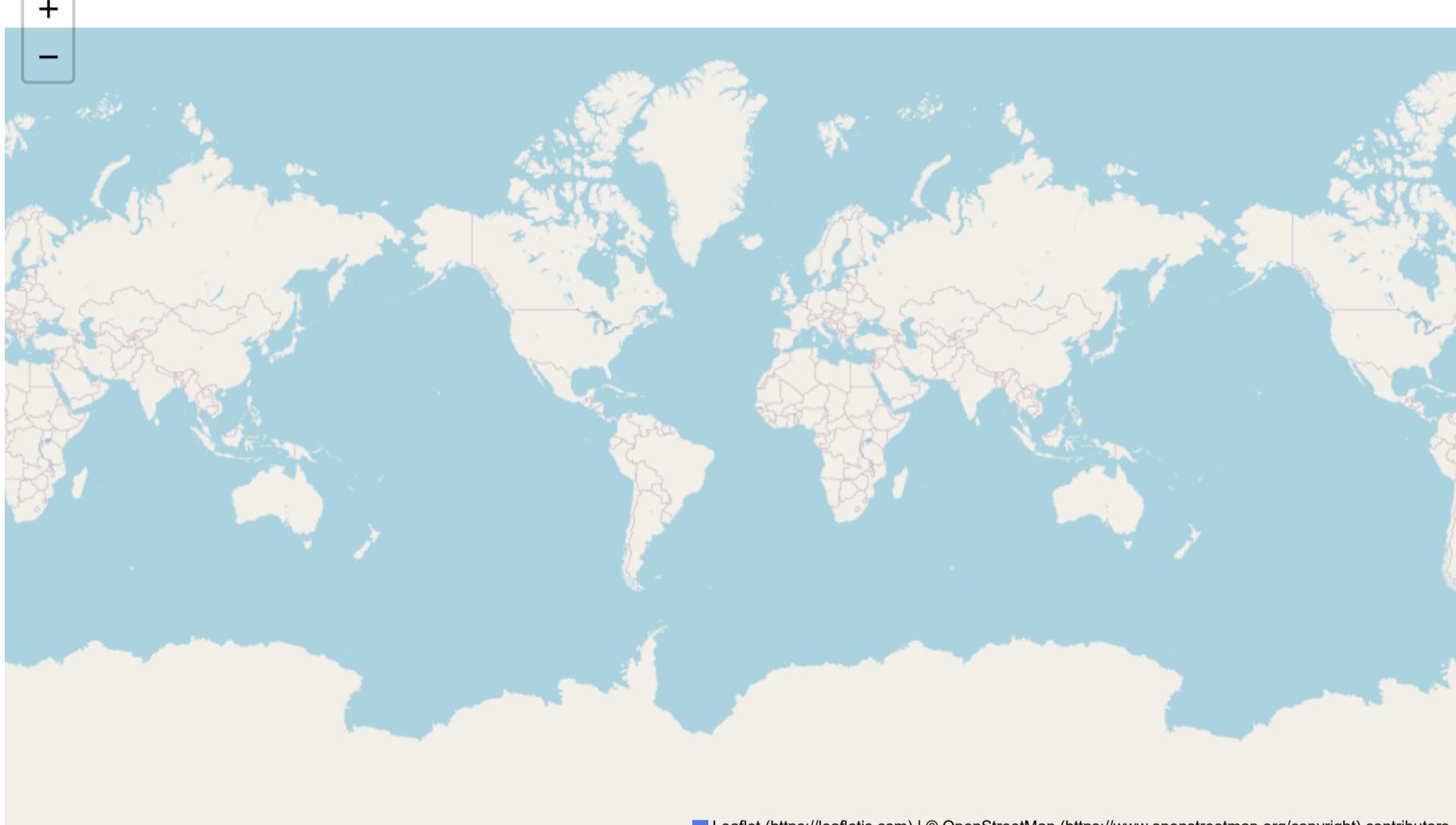


```
In [ ]: 1 pip install folium
```

## Geospatial data : latitude and longitudinal

```
In [30]: 1 import folium
2 world_map = folium.Map()
3 world_map
```

Out[30]: Make this Notebook Trusted to load map: File -> Trust Notebook



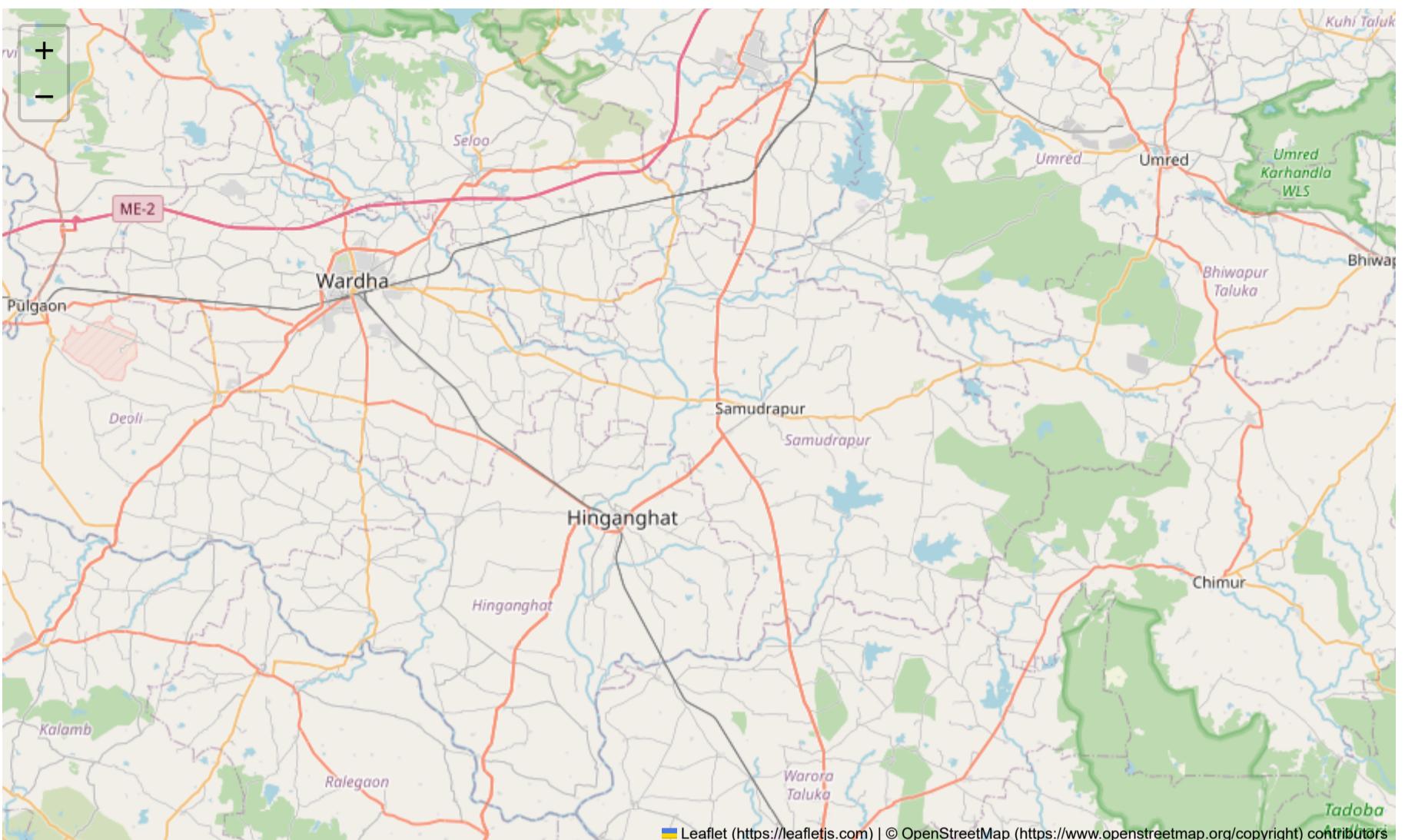
In [31]:

```

1 import folium
2 world_map=folium.Map(location=[20.5937,78.9629],zoom_start=10)
3 #North positive value, South negative value, east positive value, west negative value
4 world_map

```

Out[31]:



In [36]:

```

1 import folium
2 mapObj=folium.Map(location=[23,72],zoom_start=10)
3 folium.Circle(location=[23,72],radius=50000).add_to(mapObj)
4 mapObj.save('Output.html')

```

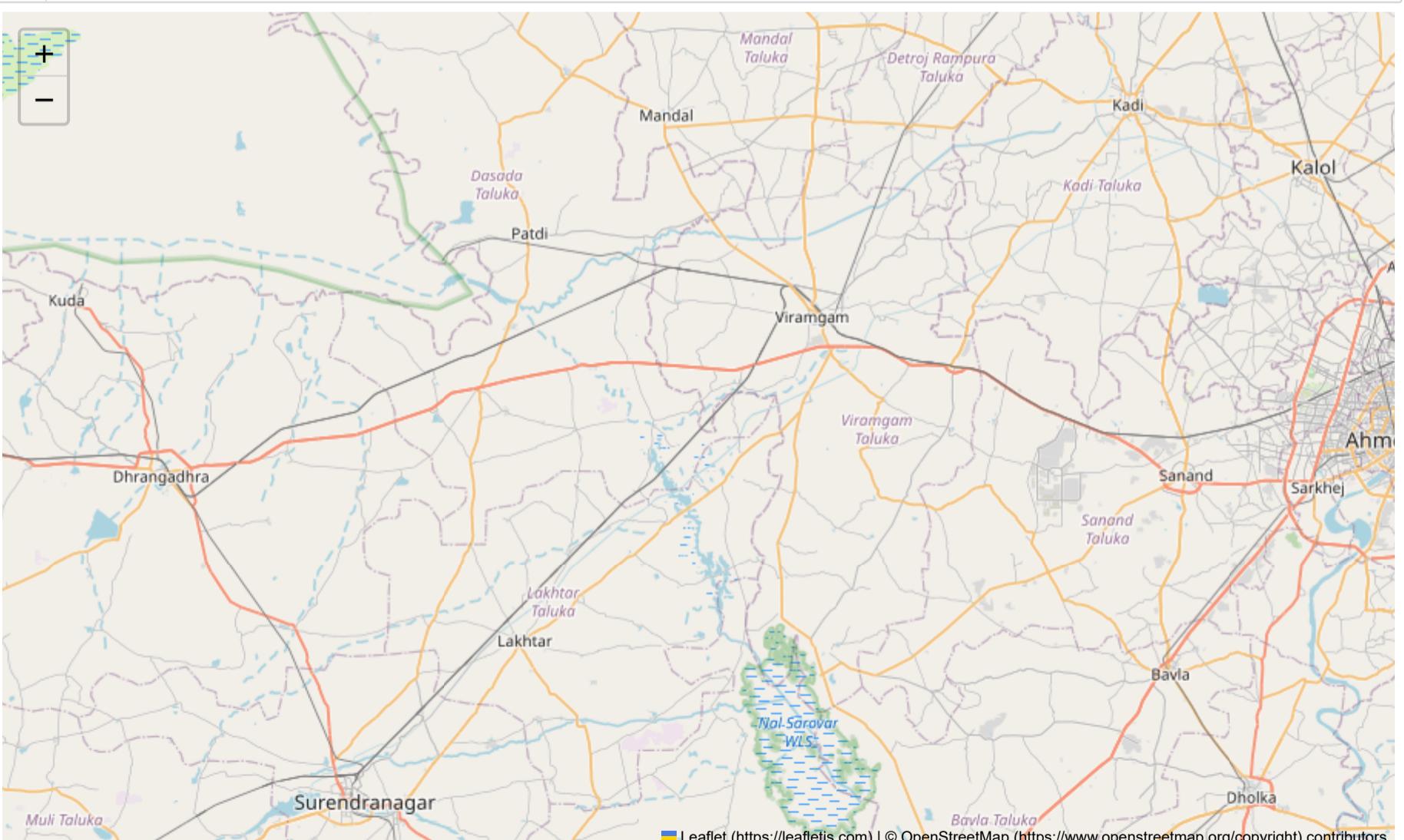
In [40]:

```

1 import folium
2 mapObj=folium.Map(location=[23,72],zoom_start=10)
3 folium.Circle(location=[23,72],radius=50000,color='green',weight=6,fill_color='red',
               opacity=0.1)
4 mapObj

```

Out[40]:



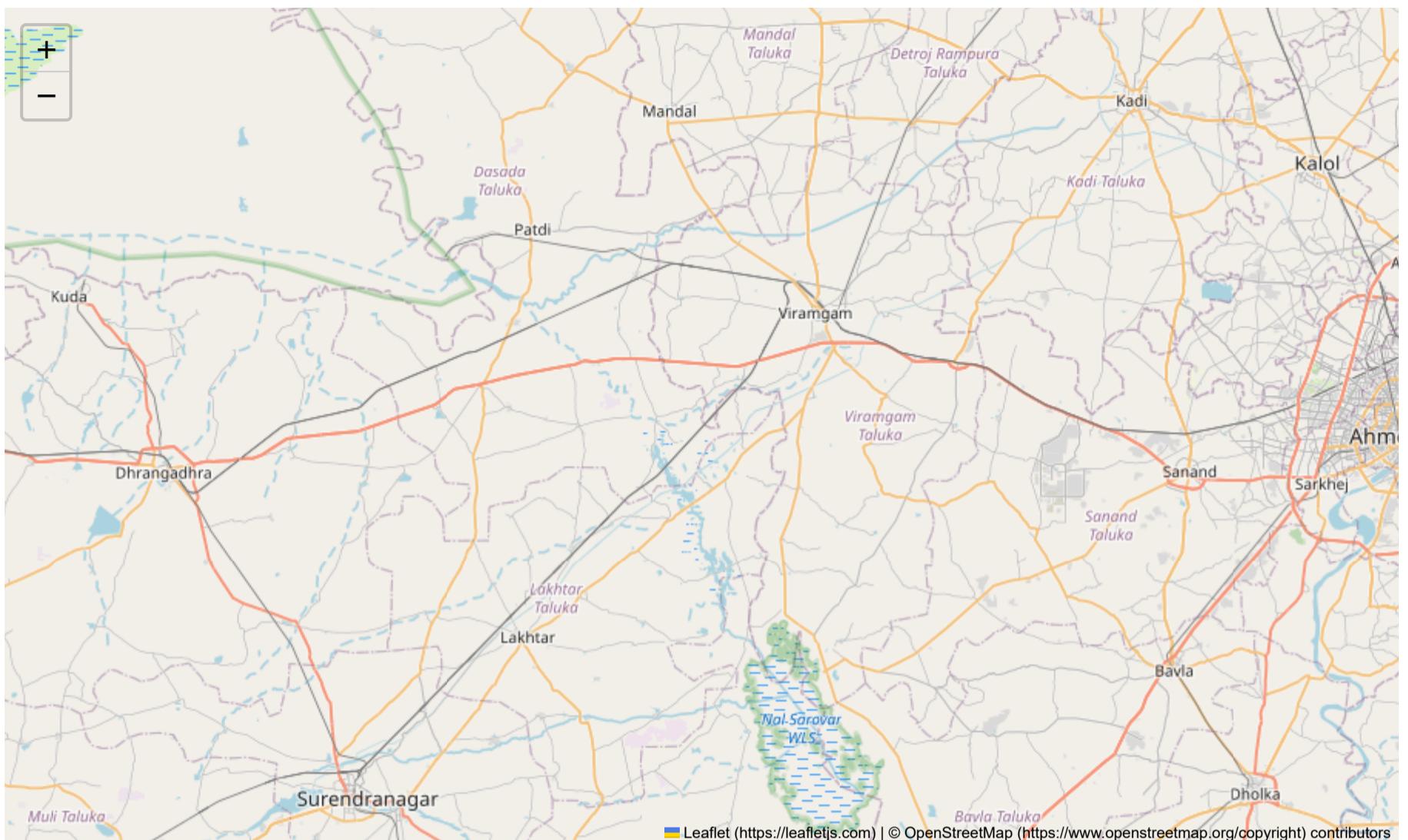
In [41]:

```

1 import folium
2 mapObj=folium.Map(location=[23,72],zoom_start=10)
3 folium.Circle(location=[23,72],radius=50000,tooltip="This is tooltip text",
4                 popup=folium.Popup("""<h1>Hiie</h1>"""))
5 mapObj

```

Out[41]:



In [43]:

```

1 import folium
2 mapObj=folium.Map(location=[23,72],zoom_start=10)
3 folium.Circle(location=[23,72],radius=50000,tooltip="This is tooltip text",
4                 popup=folium.Popup("""<h1>Hiie</h1>"""))
5 folium.LayerControl().add_to(mapObj)
6 mapObj.save('output.html')

```

In [44]:

```

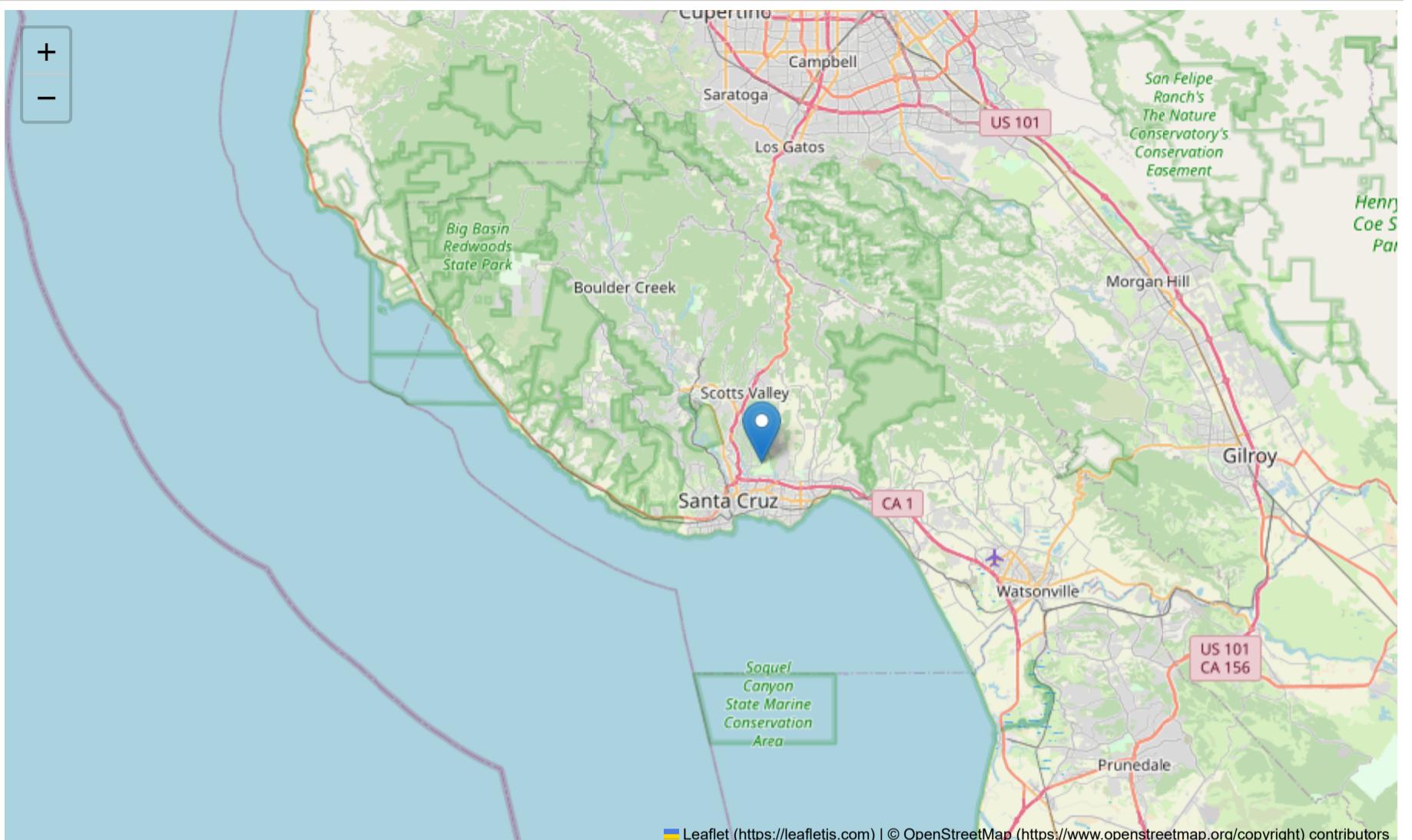
1 import folium
2 mapObj=folium.Map(location=[23,72],zoom_start=10)
3 folium.CircleMarker(location=[23,72],radius=50).add_to(mapObj)
4 mapObj.save('Output4.html')

```

In [45]:

```
1 import folium
2 mapObj = folium.Map(location=[37, -122])
3 icon=folium.features.Icon(icon='cloud')
4 marker=folium.Marker(location=[37, -122])
5 marker.add_to(mapObj)
6 mapObj
```

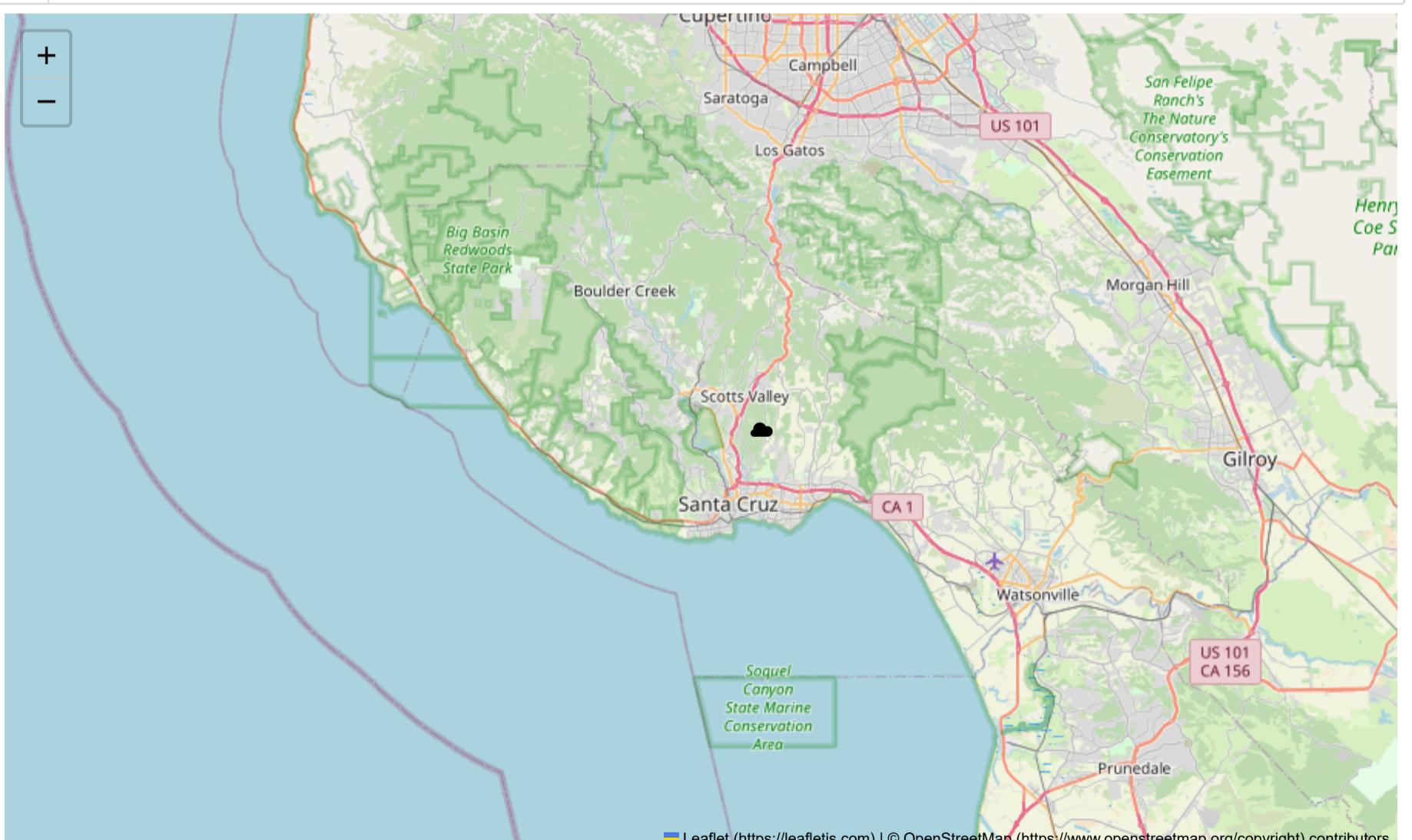
Out[45]:



In [46]:

```
1 import folium
2 mapObj = folium.Map(location=[37, -122])
3 icon=folium.features.Icon(icon='cloud')
4 marker=folium.Marker(location=[37, -122],icon=icon)
5 marker.add_to(mapObj)
6 mapObj
```

Out[46]:



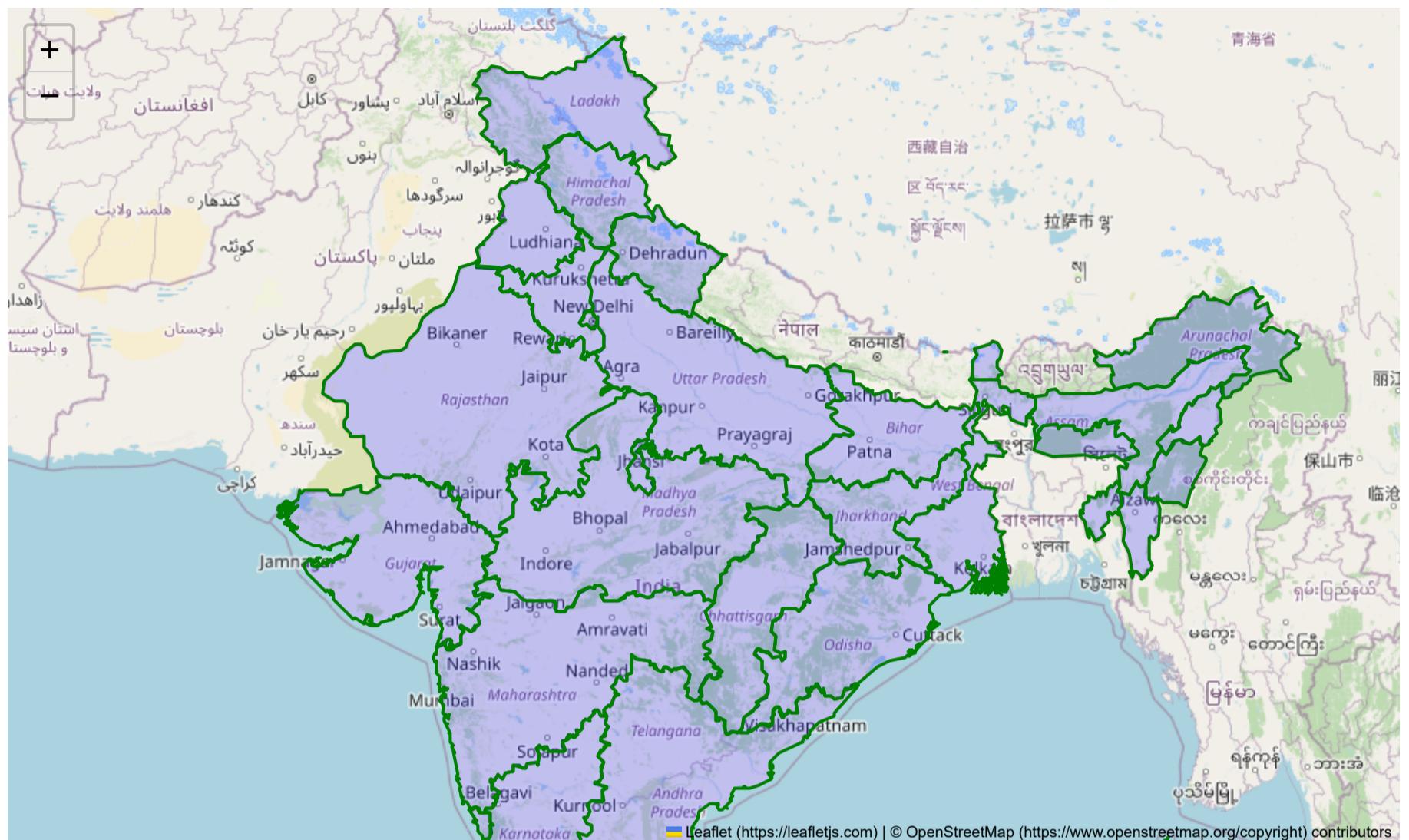
In [50]:

```

1 import folium
2 mob = folium.Map(location=[23,72])
3 folium.GeoJson(data=open("Indian_States", 'r').read(),
4                 name="India", color='green', weight=2,
5                 fillColor='blue', fillOpacity=0.2).add_to(mob)
6 mob

```

Out[50]:



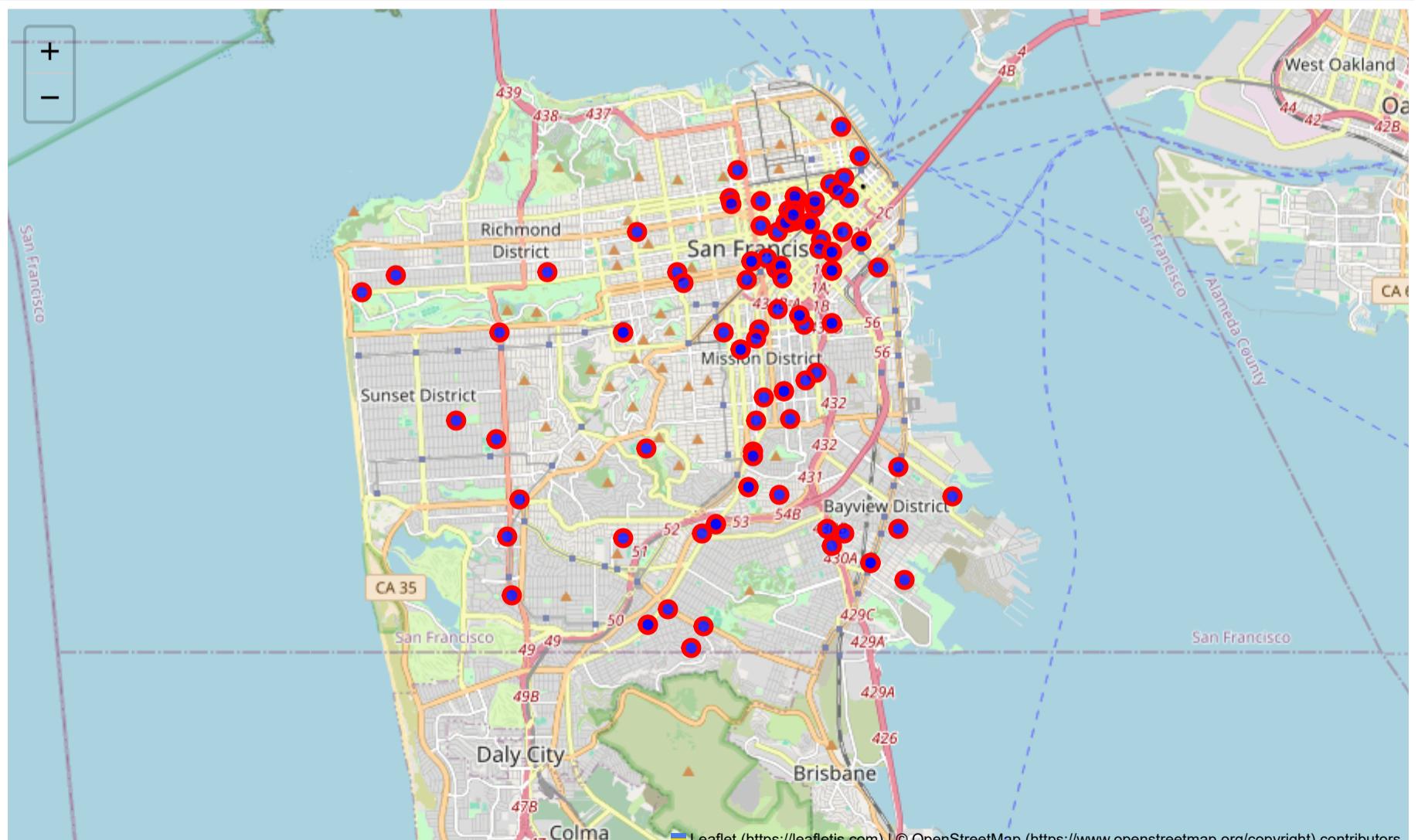
In [53]:

```

1 df = pd.read_csv('Datasets/Police_Department_Incidents_-_Previous_Year__2016_.csv')
2 df = df.iloc[0:100,:]
3 mob = folium.Map(location=[37,-122])
4 incidents = folium.map.FeatureGroup()
5 for lat,lon in zip(df.Y,df.X):
6     incidents.add_child(folium.features.CircleMarker([lat,lon],
7             radius=5,color='red',
8             fill=True,fill_color='blue',
9             fill_opacity=0.8,))
10 mob.add_child(incidents)
11
12 #zip can be done using df.X and df.Y if both have some number of datas to create a child node

```

Out[53]:



In [57]:

```

1 #police file for san fransico
2 import folium
3 import pandas as pd
4 df_incidents=pd.read_csv('Datasets/Police_Department_Incidents_-_Previous_Year__2016_.csv')
5 print(df_incidents.head())
6 df_incidents.shape
7 df=df_incidents.iloc[0:100,:]
8 df.shape
9 sanfran_map=folium.Map(location=[37.77,-122.44],zoom_start=12)
10 incidents=folium.map.FeatureGroup()
11 for lat,lon in zip(df.Y,df.X):
12     incidents.add_child(folium.features.CircleMarker([lat,lon],
13             radius=5,color='red',
14             fill=True,fill_color='blue',
15             fill_opacity=0.8,))
```

#to show circle with filled color

```

19 #iloc in folium for limiting data
20 sanfran_map.add_child(incidents)
```

	IncidntNum	Category	Descript
0	120058272	WEAPON LAWS	POSS OF PROHIBITED WEAPON
1	120058272	WEAPON LAWS	FIREARM, LOADED, IN VEHICLE, POSSESSION OR USE
2	141059263	WARRANTS	WARRANT ARREST
3	160013662	NON-CRIMINAL	LOST PROPERTY
4	160002740	NON-CRIMINAL	LOST PROPERTY

	DayOfWeek	Date	Time	PdDistrict	Resolution
0	Friday	01/29/2016	12:00:00 AM	11:00	SOUTHERN ARREST, BOOKED
1	Friday	01/29/2016	12:00:00 AM	11:00	SOUTHERN ARREST, BOOKED
2	Monday	04/25/2016	12:00:00 AM	14:59	BAYVIEW ARREST, BOOKED
3	Tuesday	01/05/2016	12:00:00 AM	23:50	TENDERLOIN NONE
4	Friday	01/01/2016	12:00:00 AM	00:30	MISSION NONE

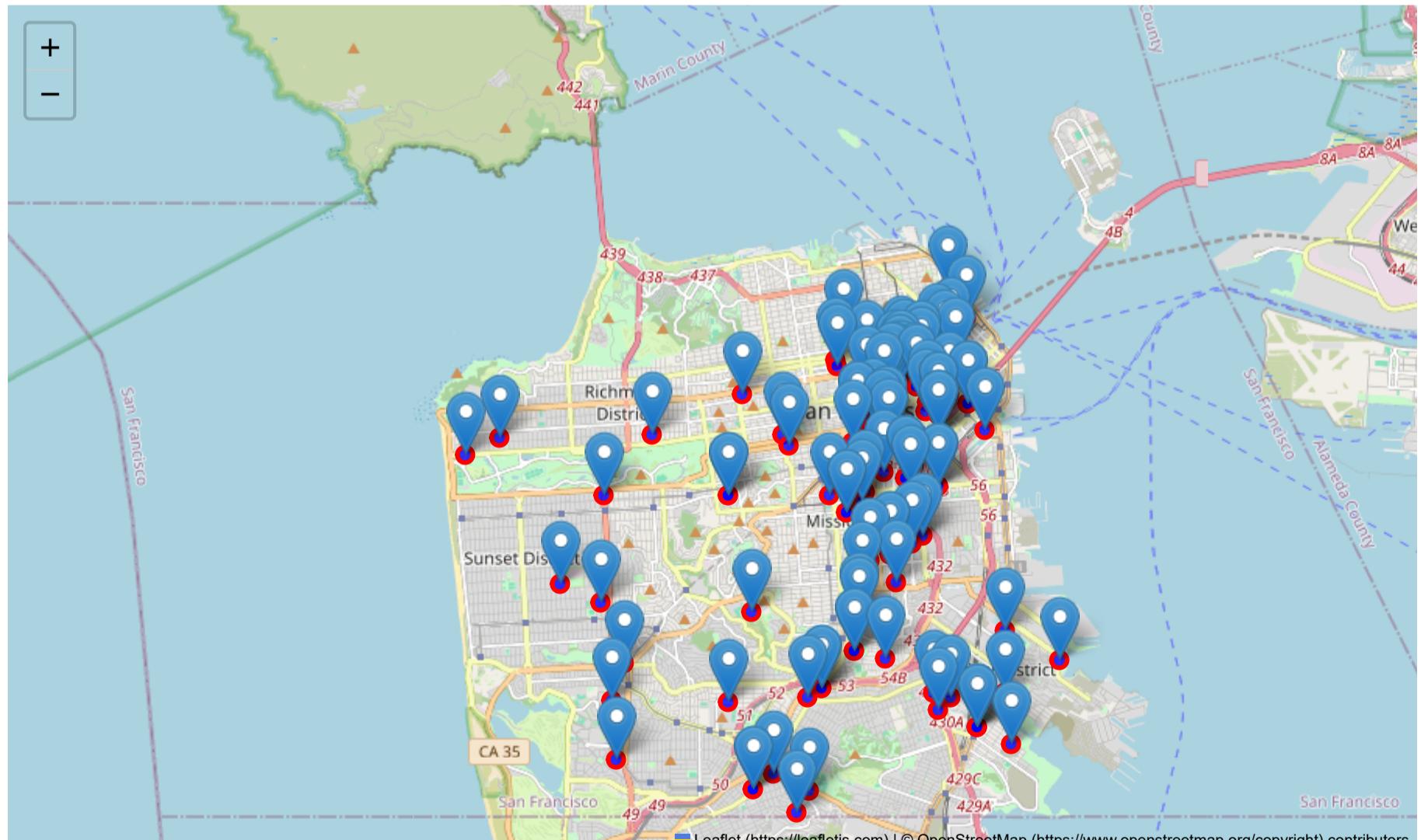
  

	Address	X	Y
0	800 Block of BRYANT ST	-122.403405	37.775421
1	800 Block of BRYANT ST	-122.403405	37.775421
2	KEITH ST / SHAFTER AV	-122.388856	37.729981
3	JONES ST / OFARRELL ST	-122.412971	37.785788
4	16TH ST / MISSION ST	-122.419672	37.765050

	Location	PdId
0	(37.775420706711, -122.403404791479)	12005827212120
1	(37.775420706711, -122.403404791479)	12005827212168
2	(37.7299809672996, -122.388856204292)	14105926363010
3	(37.7857883766888, -122.412970537591)	16001366271000
4	(37.7650501214668, -122.419671780296)	16000274071000

Out[57]:



**Choropleth Maps : to make such type of heat maps in folium (world map)**

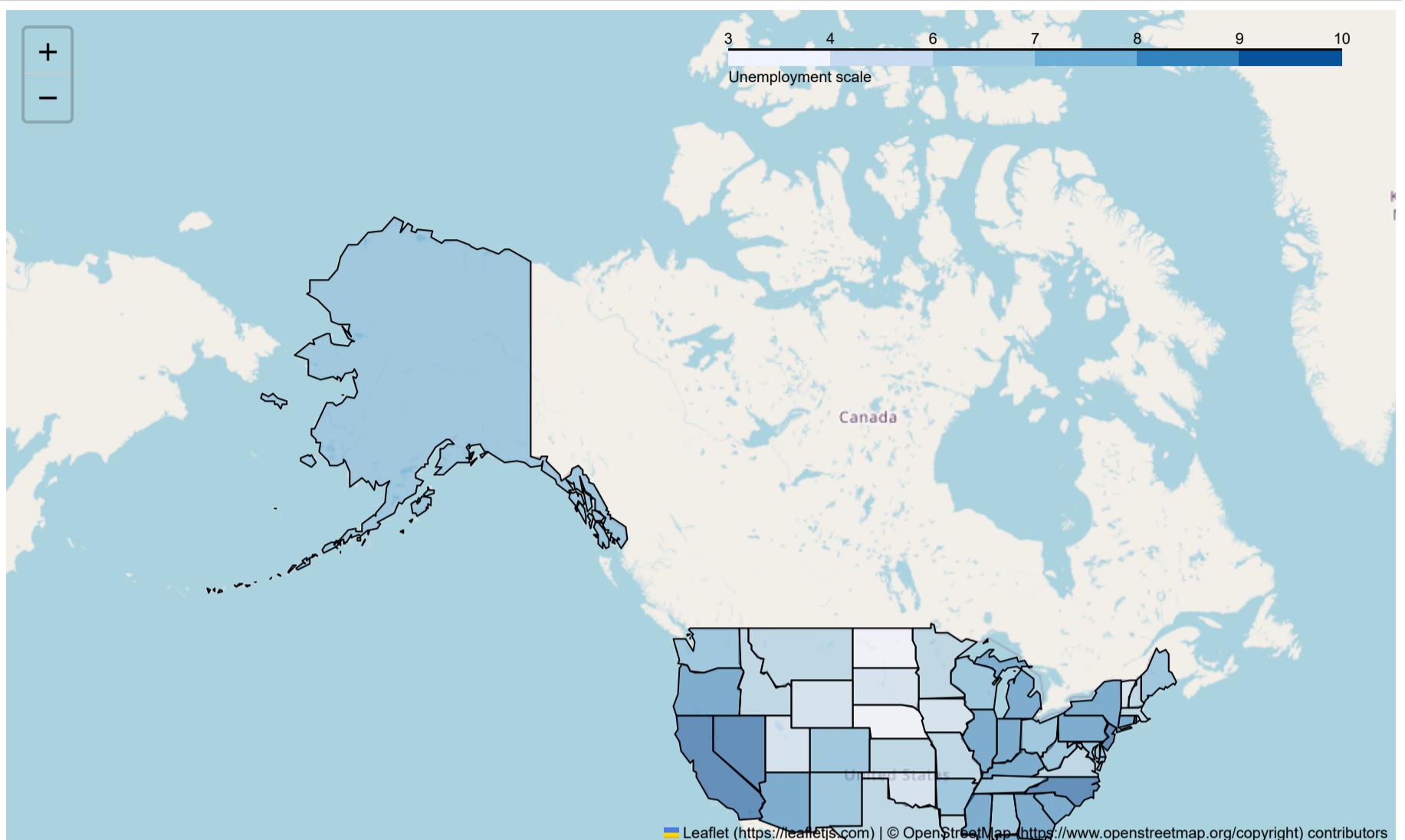
In [65]:

```

1 import folium
2 import pandas as pd
3 state_unemp=pd.read_csv('Datasets/US_Unemployment_Oct2012.csv')
4 state_geo='us-states.json'
5 m=folium.Map(location=[48,-102])
6 folium.Choropleth(geo_data=state_geo,name='choropleth',data=state_unemp,
7 columns=['State','Unemployment'],key_on='feature.id',
8 legend_name='Unemployment scale').add_to(m)
9 m

```

Out[65]:



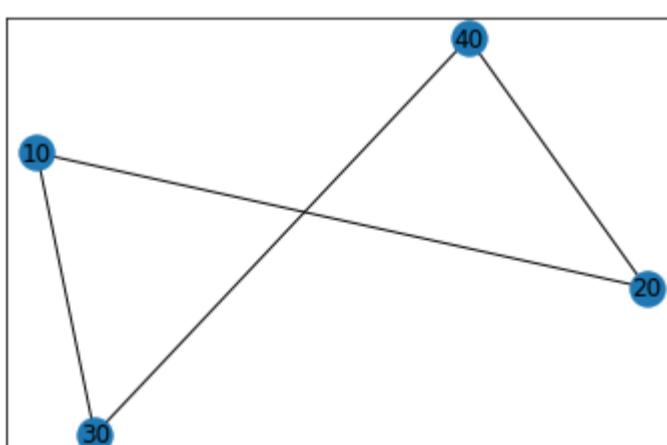
## NetworkX

In [5]:

```

1 import networkx as nx
2 G=nx.Graph()
3 G.add_node(10)
4 G.add_nodes_from([20,30,40])
5 G.add_edge(10,20)
6 G.add_edges_from([(10,30),(20,40),(30,40)])
7 nx.draw_networkx(G)

```

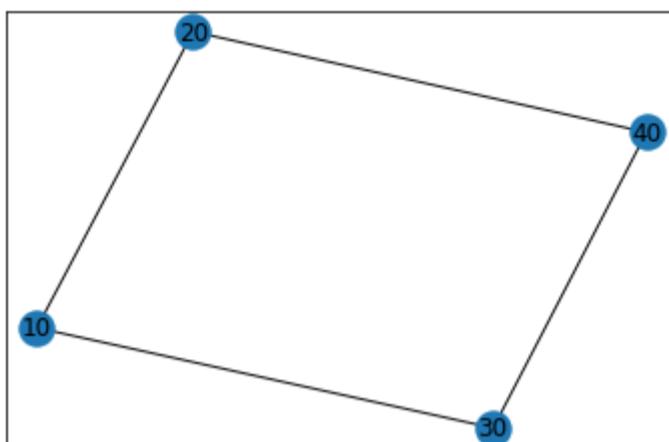


In [6]:

```

1 import networkx as nx
2 G=nx.Graph()
3 G.add_node(10)
4 G.add_nodes_from([20,30,40])
5 G.add_edge(10,20)
6 G.add_edges_from([(10,30),(20,40),(30,40)])
7 nx.draw_networkx(G)

```

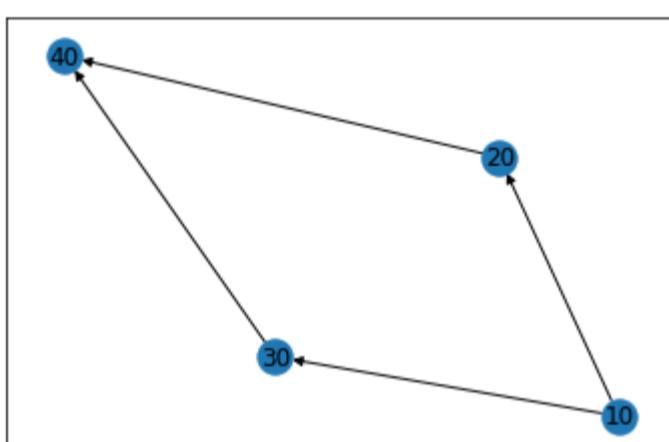


In [7]:

```

1 import networkx as nx
2 G=nx.DiGraph()
3 G.add_node(10)
4 G.add_nodes_from([20,30,40])
5 G.add_edge(10,20)
6 G.add_edges_from([(10,30),(20,40),(30,40)])
7 nx.draw_networkx(G)

```

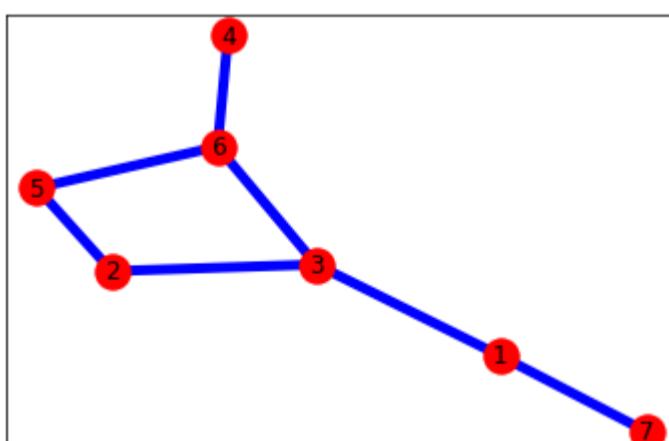


In [16]:

```

1 #undirected graphs
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 G=nx.Graph()
5 G.add_node(1)
6 G.add_nodes_from([2,3]) #to take nodes from list
7 G.add_nodes_from(range(4,7)) #to take nodes in range 7 not included
8 G.add_edge(1,3)
9 G.add_edge(2,5)
10 G.add_edge(1,7)
11 G.add_edge(3,3)
12 G.add_edges_from([(2,3),(3,6),(4,6),(5,6)]) #to draw edges in list
13 nx.draw_networkx(G,node_size=300,node_color='red',edge_color='blue',width=5)
14 plt(figsize=[150,150]
15 plt.show()

```

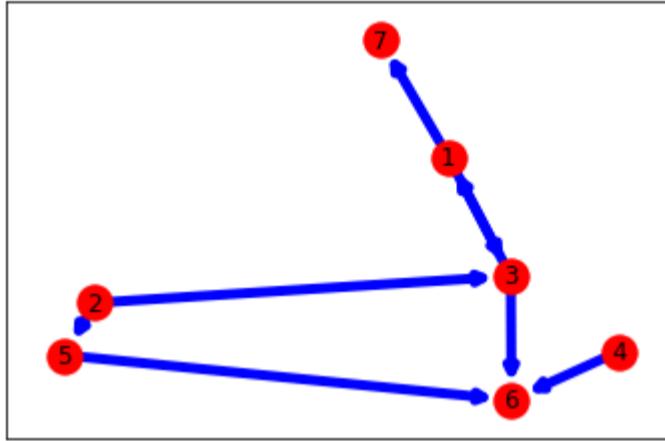


In [20]:

```

1 #Directed edges graph
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 G=nx.DiGraph()
5 G.add_node(1)
6 G.add_nodes_from([2,3]) #to take nodes from list
7 G.add_nodes_from(range(4,7)) #to take nodes in range 7 not included
8 G.add_edge(1,3)
9 G.add_edge(3,1)
10 G.add_edge(2,5)
11 G.add_edge(1,7)
12 G.add_edge(3,3)
13 G.add_edges_from([(2,3),(3,6),(4,6),(5,6)]) #to draw edges in list
14 nx.draw_networkx(G,node_size=300,node_color='red',edge_color='blue',width=5)
15 plt(figsize=[150,150]
16 plt.show()
17

```

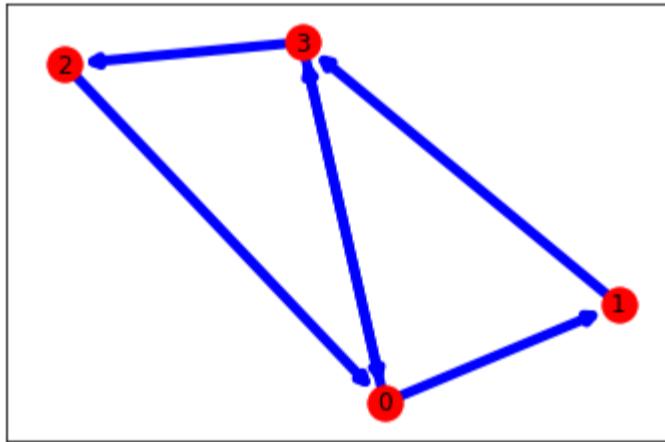


In [22]:

```

1 #Creating graph from adjacency matrix
2 #Directed edges graph
3 import networkx as nx
4 import matplotlib.pyplot as plt
5 G=nx.DiGraph()
6 G.add_nodes_from(range(4))
7 L=[[0,1,0,1],[0,0,0,1],[1,0,0,0],[1,0,1,0]]
8 for i in range(4):
9     for j in range(4):
10         if L[i][j]==1:
11             G.add_edge(i,j)
12 nx.draw_networkx(G,node_size=300,node_color='red',edge_color='blue',width=5)
13 plt(figsize=[150,150]
14 plt.show()
15

```

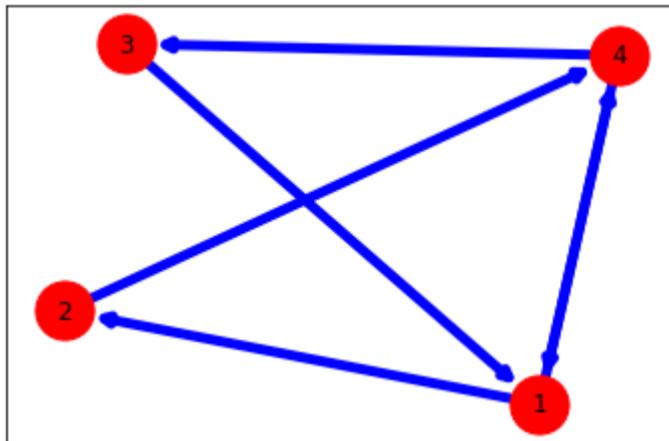


In [23]:

```

1 #Creating graph from adjacency matrix
2 #Directed edges graph
3 import networkx as nx
4 import matplotlib.pyplot as plt
5 G=nx.DiGraph()
6 G.add_nodes_from(range(1,5))
7 L=[[0,1,0,1],[0,0,0,1],[1,0,0,0],[1,0,1,0]]
8 for i in range(1,5):
9     for j in range(1,5):
10        if L[i-1][j-1]==1:
11            G.add_edge(i,j)
12 nx.draw_networkx(G,node_size=850,node_color='red',edge_color='blue',width=5)
13 plt.show()

```



## P.b. 88

In [ ]:

1

In [13]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from pandas.plotting import parallel_coordinates
5
6 # Load dataset
7 df = pd.read_csv("Datasets\supermarket_sales.csv") # Replace with correct path if needed
8
9 # 1. First 8 rows
10 print("1. First 8 rows:")
11 print(df.head(8))

```

1. First 8 rows:

	Invoice ID	Branch	City	Customer type	Gender	\
0	750-67-8428	A	Yangon	Member	Female	
1	226-31-3081	C	Naypyitaw	Normal	Female	
2	631-41-3108	A	Yangon	Normal	Male	
3	123-19-1176	A	Yangon	Member	Male	
4	373-73-7910	A	Yangon	Normal	Male	
5	699-14-3026	C	Naypyitaw	Normal	Male	
6	355-53-5943	A	Yangon	Member	Female	
7	315-22-5665	C	Naypyitaw	Normal	Female	

	Product line	Unit price	Quantity	Tax %	Total	\
0	Health and beauty	74.69	7	26.1415	548.9715	
1	Electronic accessories	15.28	5	3.8200	80.2200	
2	Home and lifestyle	46.33	7	16.2155	340.5255	
3	Health and beauty	58.22	8	23.2880	489.0480	
4	Sports and travel	86.31	7	30.2085	634.3785	
5	Electronic accessories	85.39	7	29.8865	627.6165	
6	Electronic accessories	68.84	6	20.6520	433.6920	
7	Home and lifestyle	73.56	10	36.7800	772.3800	

	Date	Time	Payment	cogs	gross margin	percentage	\
0	01-05-2019	13:08	Ewallet	522.83		4.761905	
1	03-08-2019	10:29	Cash	76.40		4.761905	
2	03-03-2019	13:23	Credit card	324.31		4.761905	
3	1/27/2019	20:33	Ewallet	465.76		4.761905	
4	02-08-2019	10:37	Ewallet	604.17		4.761905	
5	3/25/2019	18:30	Ewallet	597.73		4.761905	
6	2/25/2019	14:36	Ewallet	413.04		4.761905	
7	2/24/2019	11:38	Ewallet	735.60		4.761905	

	gross income	Rating
0	26.1415	9.1
1	3.8200	9.6
2	16.2155	7.4
3	23.2880	8.4
4	30.2085	5.3
5	29.8865	4.1
6	20.6520	5.8
7	36.7800	8.0

```
In [14]: 1 # 2. Check and fill missing values with mean
2 print("\n2. Missing values before filling:")
3 print(df.isnull().sum())
4 df.fillna(df.mean(numeric_only=True), inplace=True)
5 print("Missing values after filling:\n", df.isnull().sum())
```

2. Missing values before filling:

Invoice ID	0
Branch	0
City	0
Customer type	0
Gender	0
Product line	0
Unit price	0
Quantity	0
Tax 5%	0
Total	0
Date	0
Time	0
Payment	0
cogs	4
gross margin percentage	0
gross income	0
Rating	5
dtype: int64	

Missing values after filling:

Invoice ID	0
Branch	0
City	0
Customer type	0
Gender	0
Product line	0
Unit price	0
Quantity	0
Tax 5%	0
Total	0
Date	0
Time	0
Payment	0
cogs	0
gross margin percentage	0
gross income	0
Rating	0
dtype: int64	

```
In [15]: 1 # 3. Orders with Quantity < 3 and (Rating > 8.5 or Total > 600)
2 condition = (df['Quantity'] < 3) & ((df['Rating'] > 8.5) | (df['Total'] > 600))
3 num_orders = df[condition].shape[0]
4 print(f"\n3. Number of matching orders: {num_orders}")
```

3. Number of matching orders: 45

```
In [16]: 1 # 4. Total purchase sum by Customer type
2 total_by_customer_type = df.groupby('Customer type')['Total'].sum()
3 print("\n4. Total purchase by Customer Type:")
4 print(total_by_customer_type)
```

4. Total purchase by Customer Type:

Customer type	
Member	164223.444
Normal	158743.305
Name: Total, dtype: float64	

```
In [17]: 1 # 5. % of gross income by Payment method
2 gross_income_by_payment = df.groupby('Payment')['gross income'].sum()
3 gross_income_percent = (gross_income_by_payment / gross_income_by_payment.sum()) * 100
4 print("\n5. Gross income percentage by Payment method:")
5 print(gross_income_percent)
```

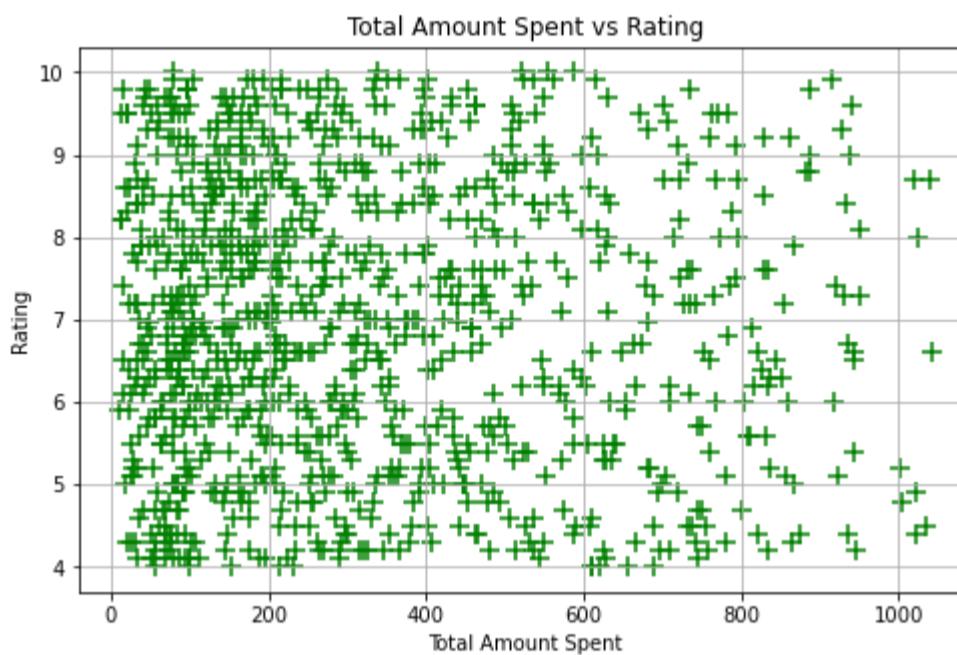
5. Gross income percentage by Payment method:

Payment	
Cash	34.742453
Credit card	31.200448
Ewallet	34.057099
Name: gross income, dtype: float64	

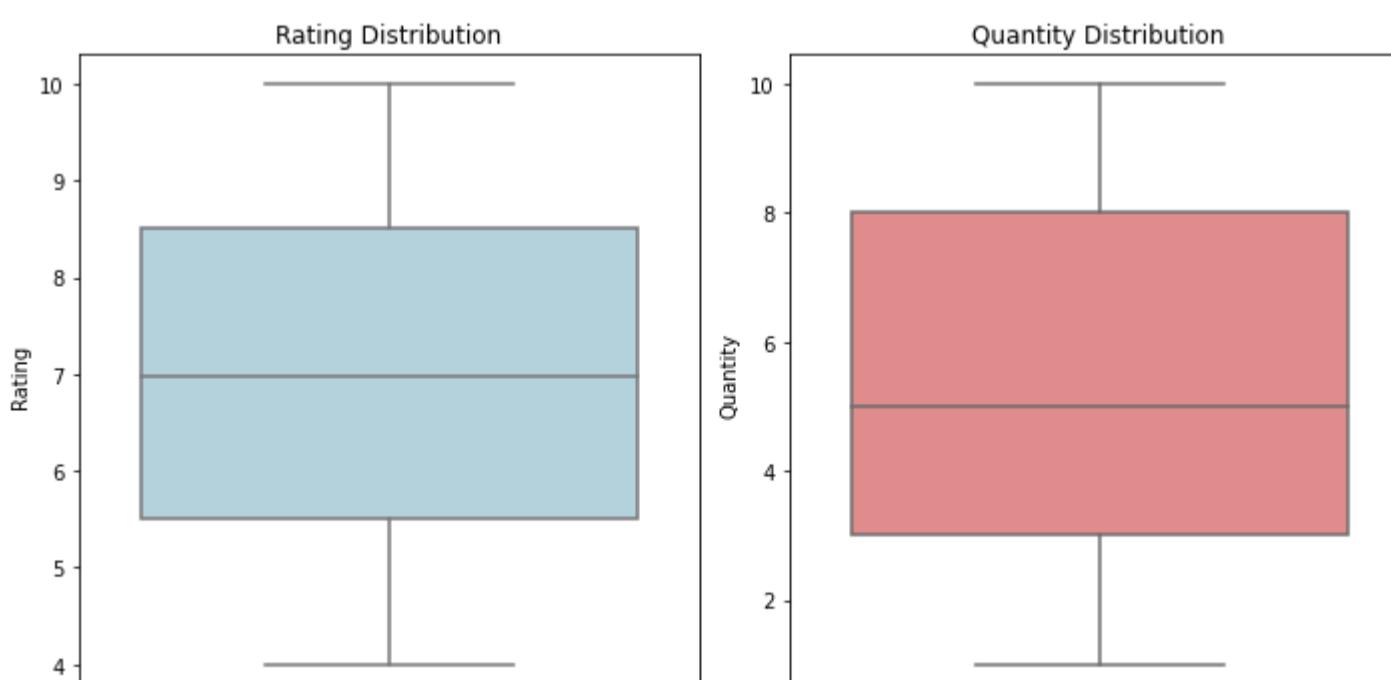
```
In [18]: 1 # 6. Average purchase price by Gender
2 avg_total_by_gender = df.groupby('Gender')['Total'].mean()
3 print("\n6. Average purchase price by Gender:")
4 print(avg_total_by_gender)
```

6. Average purchase price by Gender:  
 Gender  
 Female 335.095659  
 Male 310.789226  
 Name: Total, dtype: float64

```
In [20]: 1 # 7. Scatter plot: Total vs Rating
2 plt.figure(figsize=(8, 5))
3 plt.scatter(df['Total'], df['Rating'], marker='+', color='green', s=100)
4 plt.title('Total Amount Spent vs Rating')
5 plt.xlabel('Total Amount Spent')
6 plt.ylabel('Rating')
7 plt.grid(True)
8 plt.show()
```



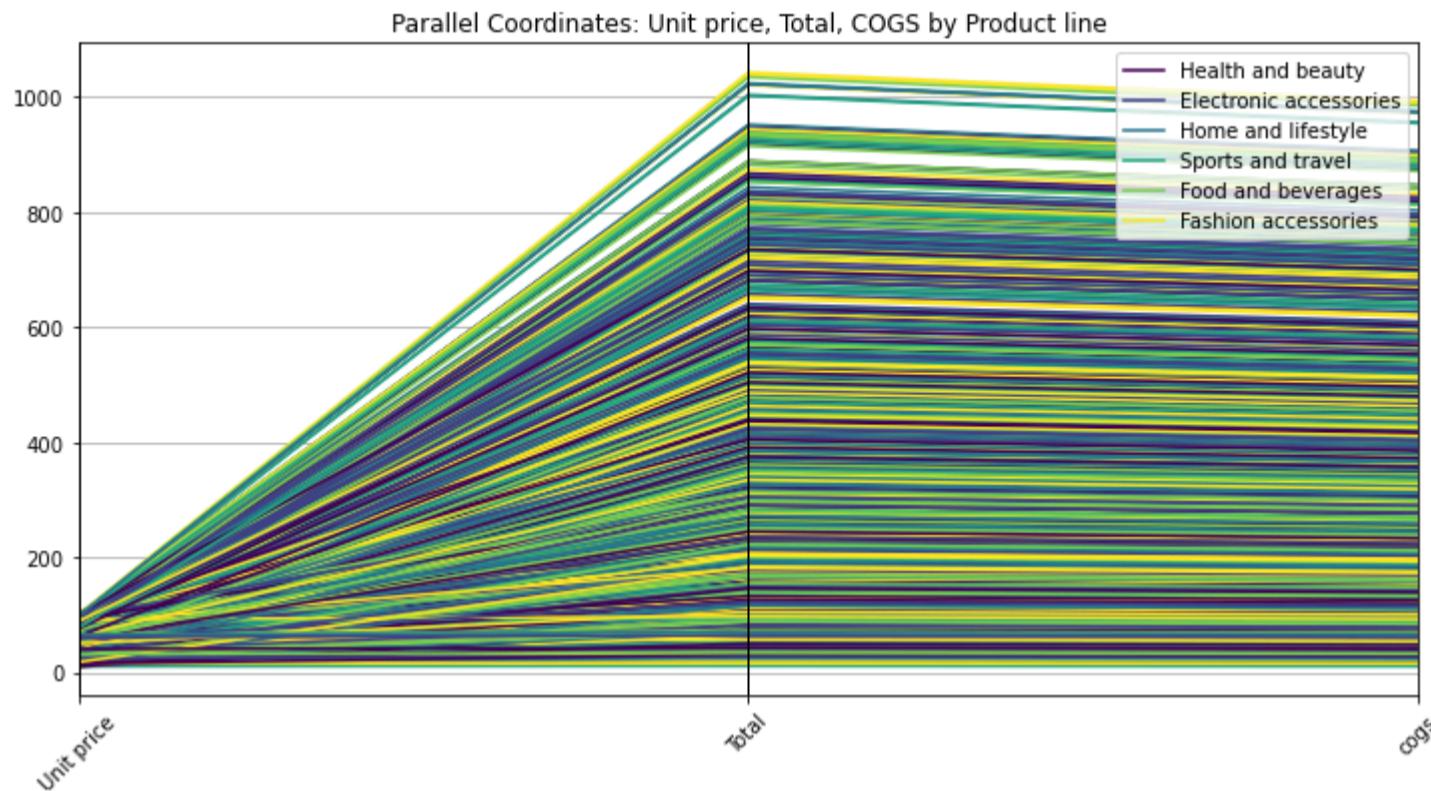
```
In [21]: 1 # 8. Box plots for Rating and Quantity
2 plt.figure(figsize=(10, 5))
3
4 plt.subplot(1, 2, 1)
5 sns.boxplot(y=df['Rating'], color='lightblue')
6 plt.title("Rating Distribution")
7
8 plt.subplot(1, 2, 2)
9 sns.boxplot(y=df['Quantity'], color='lightcoral')
10 plt.title("Quantity Distribution")
11
12 plt.tight_layout()
13 plt.show()
14
15 print("\n8. Outlier analysis:")
16 print("- Look for dots beyond the whiskers in the box plots.")
17 print("- 'Rating' is between 4 and 10 mostly, any lower may be outlier.")
18 print("- 'Quantity' outliers might appear at very high or very low ends.")
```



8. Outlier analysis:  
 - Look for dots beyond the whiskers in the box plots.  
 - 'Rating' is between 4 and 10 mostly, any lower may be outlier.  
 - 'Quantity' outliers might appear at very high or very low ends.

In [22]:

```
1 # 9. Parallel Coordinates Plot
2 subset = df[['Product line', 'Unit price', 'Total', 'cogs']]
3 plt.figure(figsize=(12, 6))
4 parallel_coordinates(subset, class_column='Product line', colormap='viridis')
5 plt.title("Parallel Coordinates: Unit price, Total, COGS by Product line")
6 plt.xticks(rotation=45)
7 plt.grid(True)
8 plt.show()
```



In [ ]:

1