

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

```
In [2]: df=pd.read_csv("21_cities.csv")
df
```

Out[2]:

	id	name	state_id	state_code	state_name	country_id	country_code	country_name	
	0	52	Ashkāsham	3901	BDS	Badakhshan	1	AF	Afghanistan
	1	68	Fayzabad	3901	BDS	Badakhshan	1	AF	Afghanistan
	2	78	Jurm	3901	BDS	Badakhshan	1	AF	Afghanistan
	3	84	Khandūd	3901	BDS	Badakhshan	1	AF	Afghanistan
	4	115	Rāghistān	3901	BDS	Badakhshan	1	AF	Afghanistan

150449	131496	Redcliff	1957	MI	Midlands Province	247	ZW	Zimbabwe	
150450	131502	Shangani	1957	MI	Midlands Province	247	ZW	Zimbabwe	
150451	131503	Shurugwi	1957	MI	Midlands Province	247	ZW	Zimbabwe	
150452	131504	Shurugwi District	1957	MI	Midlands Province	247	ZW	Zimbabwe	
150453	131508	Zvishavane District	1957	MI	Midlands Province	247	ZW	Zimbabwe	

150454 rows × 11 columns



```
In [3]: df.head()
```

Out[3]:

	id	name	state_id	state_code	state_name	country_id	country_code	country_name	latitude
0	52	Ashkāsham	3901	BDS	Badakhshan	1	AF	Afghanistan	36.68333
1	68	Fayzabad	3901	BDS	Badakhshan	1	AF	Afghanistan	37.11664
2	78	Jurm	3901	BDS	Badakhshan	1	AF	Afghanistan	36.86477
3	84	Khandūd	3901	BDS	Badakhshan	1	AF	Afghanistan	36.95127
4	115	Rāghistān	3901	BDS	Badakhshan	1	AF	Afghanistan	37.66079



Data Cleaning and Data Preprocessing

In [4]:

`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150454 entries, 0 to 150453
Data columns (total 11 columns):
#   Column          Non-Null Count  Dtype
---  -
0   id               150454 non-null  int64
1   name             150454 non-null  object
2   state_id         150454 non-null  int64
3   state_code       150129 non-null  object
4   state_name       150454 non-null  object
5   country_id       150454 non-null  int64
6   country_code     150406 non-null  object
7   country_name     150454 non-null  object
8   latitude         150454 non-null  float64
9   longitude        150454 non-null  float64
10  wikiDataId       147198 non-null  object
dtypes: float64(2), int64(3), object(6)
memory usage: 12.6+ MB
```

In [5]:

`df.describe()`

Out[5]:

	id	state_id	country_id	latitude	longitude
count	150454.000000	150454.000000	150454.000000	150454.000000	150454.000000
mean	76407.091689	2678.377677	140.658460	31.556175	2.369557
std	44357.755335	1363.513591	70.666123	22.813220	68.012770
min	1.000000	1.000000	1.000000	-75.000000	-179.121980
25%	38160.250000	1451.000000	82.000000	19.000000	-58.468150
50%	75975.500000	2174.000000	142.000000	40.684720	8.669980
75%	115204.750000	3905.000000	207.000000	47.239220	27.750000
max	153528.000000	5116.000000	247.000000	73.508190	179.466000

In [6]:

`df.columns`

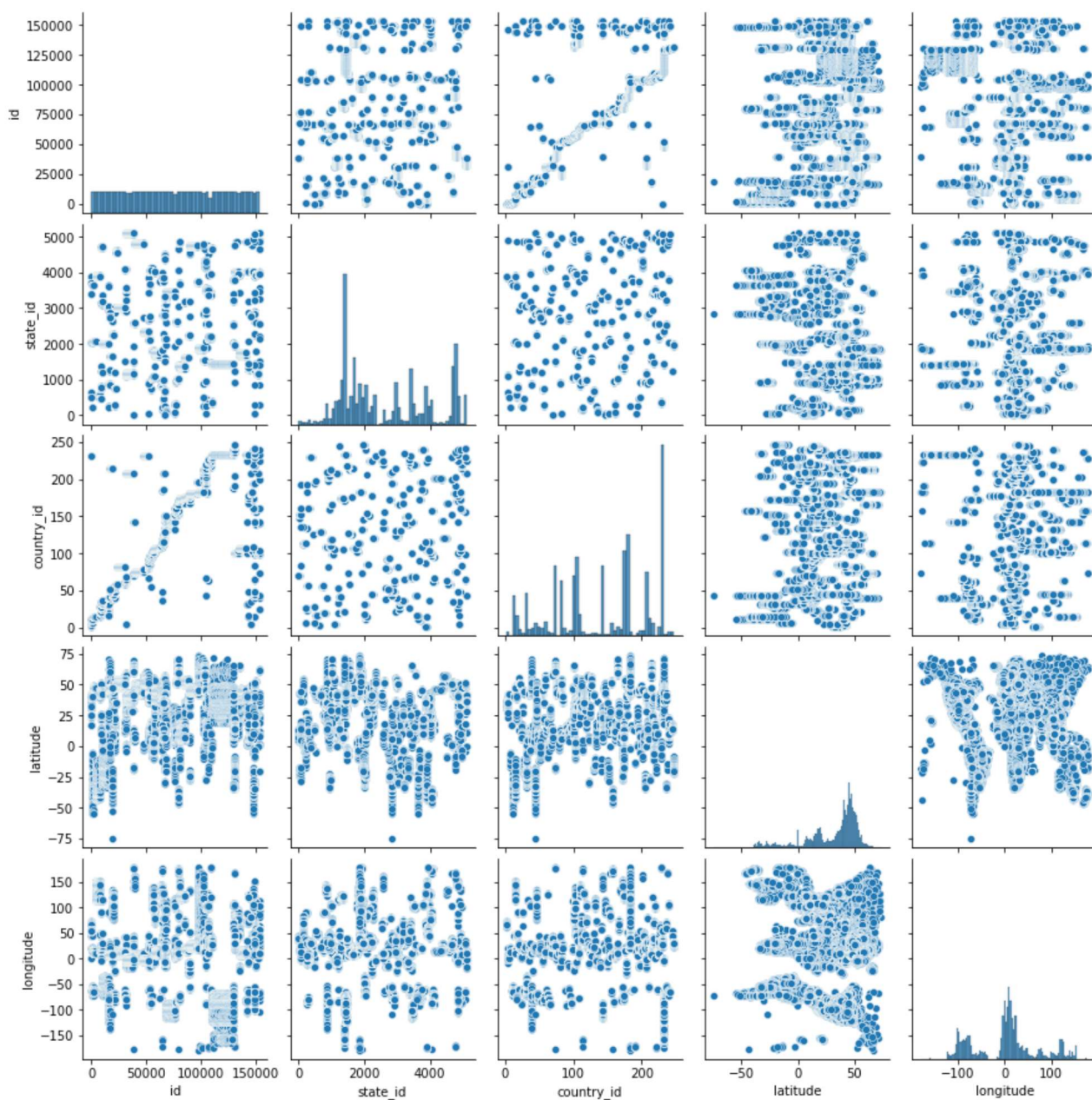
```
Out[6]: Index(['id', 'name', 'state_id', 'state_code', 'state_name', 'country_id',
              'country_code', 'country_name', 'latitude', 'longitude', 'wikiDataId'],
              dtype='object')
```

EDA and Visualization

In [7]:

`sns.pairplot(df)`

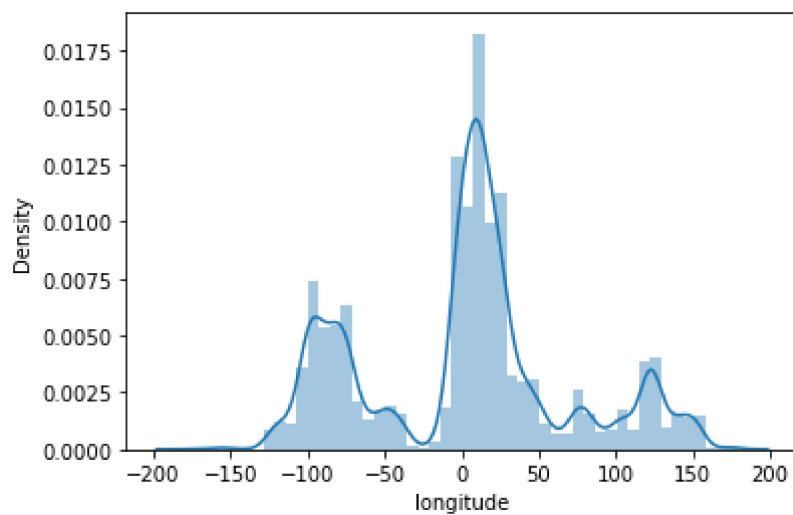
```
Out[7]: <seaborn.axisgrid.PairGrid at 0x2afbcccdb00>
```



In [8]: `sns.distplot(df['longitude'])`

C:\ProgramData\Anaconda3\lib\site-packages\seaborn\distributions.py:2557: 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)

Out[8]: <AxesSubplot:xlabel='longitude', ylabel='Density'>



```
In [9]: df1=df[['id','state_id', 'country_id',
               'latitude','longitude']]
df1
```

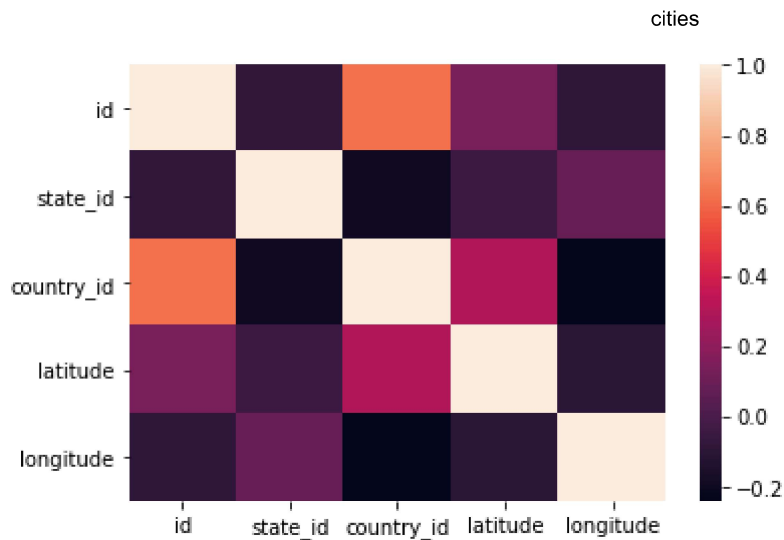
```
Out[9]:
```

	id	state_id	country_id	latitude	longitude
0	52	3901	1	36.68333	71.53333
1	68	3901	1	37.11664	70.58002
2	78	3901	1	36.86477	70.83421
3	84	3901	1	36.95127	72.31800
4	115	3901	1	37.66079	70.67346
...
150449	131496	1957	247	-19.03333	29.78333
150450	131502	1957	247	-19.78333	29.36667
150451	131503	1957	247	-19.67016	30.00589
150452	131504	1957	247	-19.75000	30.16667
150453	131508	1957	247	-20.30345	30.07514

150454 rows × 5 columns

```
In [10]: sns.heatmap(df1.corr())
```

```
Out[10]: <AxesSubplot:>
```



To Train the Model -Model Building

We are going to train Linear Regression model; We need to split out data into two variables x and y where x is independent variable (input) and y is dependent variable on x (output) we could ignore address column as it is not required for our model

```
In [11]: x=df1[['id','state_id', 'country_id',
            'latitude']]
          y=df1['longitude']
```

```
In [12]: from sklearn.model_selection import train_test_split
          x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3)
```

```
In [13]: from sklearn.linear_model import LinearRegression
          lr=LinearRegression()
          lr.fit(x_train,y_train)
```

```
Out[13]: LinearRegression()
```

```
In [14]: print(lr.intercept_)
```

```
27.702993611184347
```

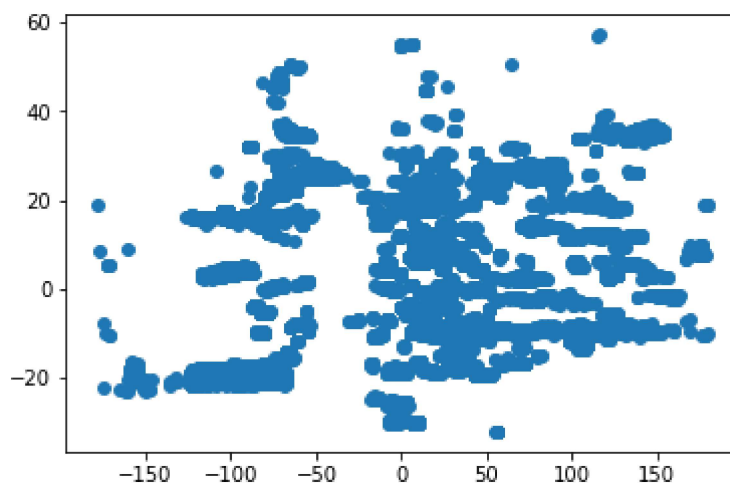
```
In [15]: coeff=pd.DataFrame(lr.coef_,x.columns,columns=['Co-efficient'])
          coeff
```

```
Out[15]:
```

	Co-efficient
id	0.000151
state_id	0.001927
country_id	-0.276142
latitude	-0.098076

```
In [16]: prediction =lr.predict(x_test)
plt.scatter(y_test,prediction)
```

Out[16]: <matplotlib.collections.PathCollection at 0x2afbf2d2d60>



```
In [17]: lr.score(x_test,y_test)
```

Out[17]: 0.06664691164703307

```
In [18]: lr.score(x_train,y_train)
```

Out[18]: 0.06704026928081253

```
In [19]: from sklearn.linear_model import Ridge,Lasso
```

```
In [20]: rr=Ridge(alpha=10)
rr.fit(x_train,y_train)
```

Out[20]: Ridge(alpha=10)

```
In [21]: rr.score(x_test,y_test)
```

Out[21]: 0.06664691169704884

```
In [22]: rr.score(x_train,y_train)
```

Out[22]: 0.06704026928081264

```
In [23]: la=Lasso(alpha=10)
la.fit(x_train,y_train)
```

Out[23]: Lasso(alpha=10)

In [24]: `la.score(x_test,y_test)`

Out[24]: 0.06664246983810551

In [25]: `la.score(x_train,y_train)`

Out[25]: 0.06699683818998758

In [26]: `from sklearn.linear_model import ElasticNet
en=ElasticNet()
en.fit(x_train,y_train)`

Out[26]: ElasticNet()

In [27]: `en.coef_`

Out[27]: array([1.51194520e-04, 1.92799793e-03, -2.76037633e-01, -9.70639850e-02])

In [28]: `en.intercept_`

Out[28]: 27.667315907034485

In [29]: `prediction=en.predict(x_test)`

In [30]: `en.score(x_test,y_test)`

Out[30]: 0.0666490216933937

Evaluation Metrics

In [31]: `from sklearn import metrics`

In [32]: `print("Mean Absolute Error:",metrics.mean_absolute_error(y_test,prediction))`

Mean Absolute Error: 51.4130859753299

In [33]: `print("Mean Squared Error:",metrics.mean_squared_error(y_test,prediction))`

Mean Squared Error: 4299.514633968312

In [34]: `print("Root Mean Squared Error:",np.sqrt(metrics.mean_squared_error(y_test,prediction)))`

Root Mean Squared Error: 65.57068425728309

Model Saving

```
In [35]: import pickle
```

```
In [36]: filename="prediction"
pickle.dump(lr,open(filename,'wb'))
```

```
In [37]: import pandas as pd
import pickle
```

```
In [38]: filename="prediction"
model=pickle.load(open(filename,'rb'))
```

```
In [39]: real=[[10,20,30,40],[11,45,10,25]]
result=model.predict(real)
```

```
In [40]: result
```

```
Out[40]: array([15.53575002, 22.57805637])
```