



# EXPLORE WEATHER TRENDS

JANUARY 17, 2021

USAMA TARIQ

Udacity – Data Analyst  
Nanodegree – Project – 1





## OVERVIEW

Local temperature (°C) data of Lahore (LHR), the city of gardens -located in the north-eastern portion of Pakistan (PK), is analyzed, and compared with the global temperature (°C) data. I accessed the Client's, Udacity-Advance Your Career, database to extract the required data, manipulate, visualize, and interpret valuable insights from data for better and concrete decision making.

---

## Goals

- Extract the data from the database, transform it into comma-separated values (CSV) file.
- Load the CSV file for analysis.
- Visualize and compare the local and global temperature data.
- Identify the trends of the local and global temperature data.

## Tools

- **SQL:** for extraction of the data from the database
- **Jupyter Notebook:** as IDE for Python & Pandas
- **Pandas:** for loading the CSV file for analysis.
- **Python:** to visualize and compare the local and global temperature data.
- **Microsoft Office 365:** interpreting the insights of data and share valuable results in the form of a report.

# EXTRACT THE DATA

## Finding the Closest Big City

On the basis of my SQL expertise, I wrote the following query to find the closest city listed in the cities of the Pakistan. And I found my own city in the database.

```
1. select *
2. from city_list
3. where country = 'Pakistan';
```

Above query's results were

city	country
Faisalabad	Pakistan
Gujranwala	Pakistan
Hyderabad	Pakistan
Islamabad	Pakistan
Karachi	Pakistan
Lahore	Pakistan
Multan	Pakistan
Peshawar	Pakistan
Rawalpindi	Pakistan

## Extracting City Level data

For city-level data, the city I select is Lahore, Pakistan. The query used to extract the whole data for 'Lahore' is listed below.

```
1. select *
2. from city_data
3. where city='Lahore';
```

And the results for this query were like,

Output		198 results		<a href="#">Download CSV</a>
year	city	country	avg_temp	
1816	Lahore	Pakistan	24.69	
1817	Lahore	Pakistan	22.74	
1818	Lahore	Pakistan	23.43	
1819	Lahore	Pakistan	22.83	

Although the number of records is limited to 198, but its burdensome to show the full data.

## Extracting Global Level data

Query for retrieving the global level data,

```
1. select *  
2. from global_data;
```

This query yielded 266 number of records as a global-level data.

Output 266 results		<a href="#">Download CSV</a>
year	avg_temp	
1750	8.72	
1751	7.98	
1752	5.78	
1753	8.39	

Finally, the required data for the analysis has been extracted and transformed successfully. In the upcoming step, the data will be load for analysis.

# LOADING DATA (CSV)

For the analytical purposes of the data, csv is supported by many different tools. But I prefer the Jupyter notebooks over Spreadsheets.

In Jupyter notebooks, you can analyze data critically with the help Python and several supported libraries like [Pandas](#), [NumPy](#) etc.

So, the loading of data from csv file to a Data Frame in Python is done by,

```
1. df_CityData = pd.read_csv('City-Level-Data (CSV).csv')
2. df_GlobalData = pd.read_csv('Global-Data (CSV).csv')
```

The first line of code reads data from “City-Level-Data (CSV).csv” file and writes it to a “df\_CityData” data frame. And the second line of code reads data from “Global-Data (CSV).csv” file and writes it to a “df\_GlobalData” data frame.

To ensure whether the data has been successfully read and stored in the desired data frames, the lines of code used to examine, and the results yielded are shown side by side.

```
1. df_CityData.head()
```

	year	city	country	avg_temp
0	1816	Lahore	Pakistan	24.69
1	1817	Lahore	Pakistan	22.74
2	1818	Lahore	Pakistan	23.43
3	1819	Lahore	Pakistan	22.83
4	1820	Lahore	Pakistan	24.22

```
1. df_GlobalData.head()
```

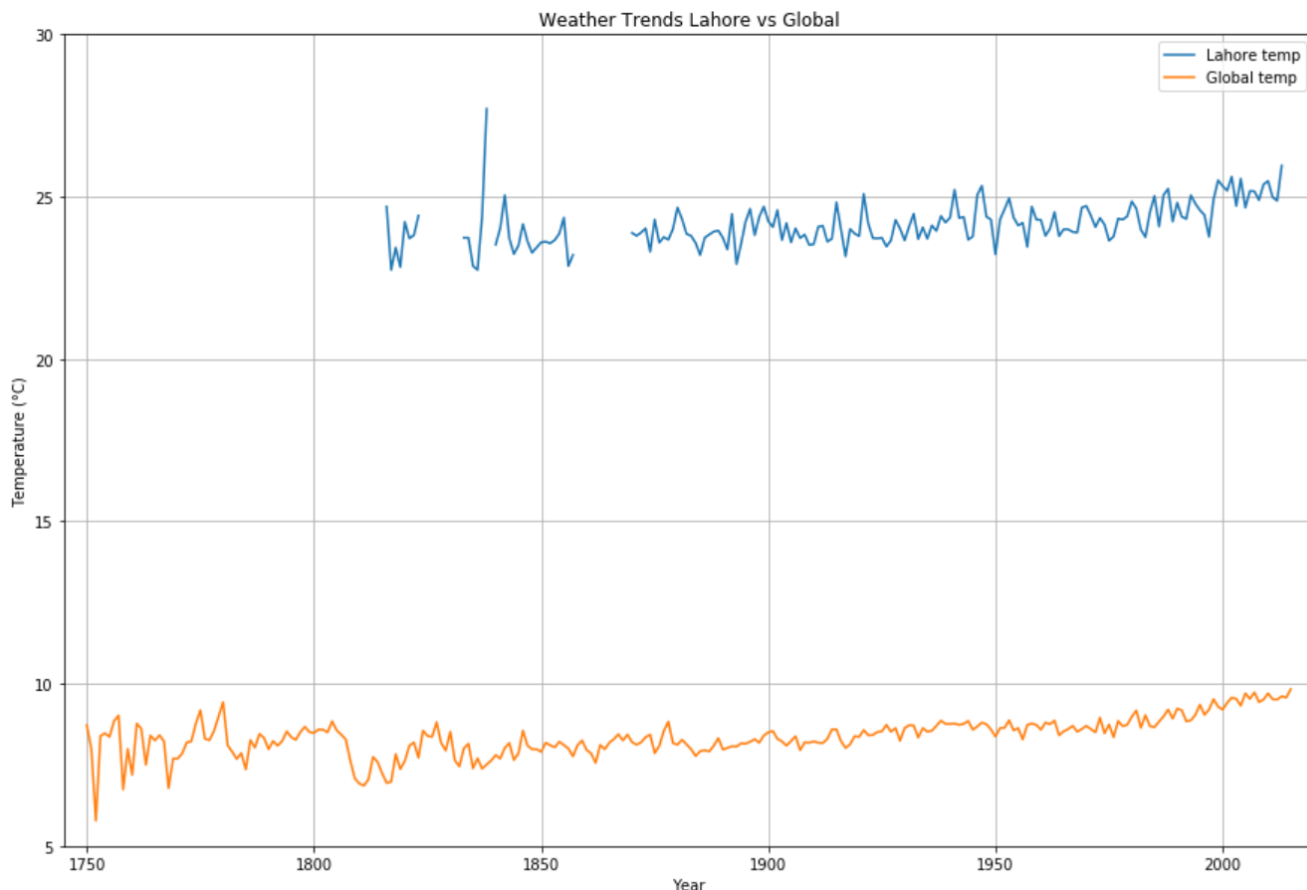
	year	avg_temp
0	1750	8.72
1	1751	7.98
2	1752	5.78
3	1753	8.39
4	1754	8.47

# VISUALIZING & CALCULATING MOVING AVERAGES

## Visualizing raw data

Firstly, to understand the data and its trends. Drawing a line chart to explore the spikes and fluctuations in the local temperatures (°C) and global temperatures (°C) for better comparisons.

```
1. plt.figure(figsize=[15,10])
2. plt.plot(df_CityData['year'],
3.          df_CityData['avg_temp'],
4.          label='Lahore temp')
5. plt.plot(df_GlobalData['year'],
6.          df_GlobalData['avg_temp'],
7.          label='Global temp')
8. plt.legend()
9. plt.grid(True)
10. plt.axis([1745, 2020, 5, 30])
11. plt.title('Weather Trends Lahore vs Global')
12. plt.xlabel('Year')
13. plt.ylabel('Temperature (°C)')
14. plt.show()
```



Lahore temperature (°C) are the blue ones, the average maximum temperature (°C) ever reached in Lahore is 27.70°C in the year 1838 and the average minimum temperature (°C) ever recorded is 22.74°C in year 1817 & 1836 as per the data. The span of data is from year 1816 to year 2013.

Global temperature (°C) are the orange ones, the average maximum temperature (°C) ever reached globally is 9.83°C in the year 2015 and the average minimum temperature (°C) ever recorded is 5.78°C in year 1752 as per the data. The span of data is from year 1750 to year 2015.

But the temperatures fluctuations are so rapid around the years, so we must calculate the moving average over the years temperature (°C). To predict better results.

## Calculating Moving Averages on Raw Data & Visualizing

Calculating moving averages for both local & global data temperatures (°C) for the span of a decade.

### Why the window size or moving average calculating for 10 years or decade?

As the dispersion of data will be reduced nearly by 25%, and lesser the dispersion of data the trend line will be smoother.

```
1. df_CityData['m_avg_temp'] = df_CityData['avg_temp'].rolling(10).mean()
2. df_GlobalData['m_avg_temp'] = df_GlobalData['avg_temp'].rolling(10).mean()
```

Now plotting the data of moving averages calculated for 10 years span of both local and global temperatures (°C)

```
1. plt.figure(figsize=[15,10])
2. plt.plot(df_CityData['year'],
3.          df_CityData['m_avg_temp'],
4.          label='Lahore temp')
5. plt.plot(df_GlobalData['year'],
6.          df_GlobalData['m_avg_temp'],
7.          label='Global temp')
8. plt.legend()
9. plt.grid(True)
10. plt.axis([1745, 2020, 5, 30])
11. plt.title('Weather Trends Lahore vs Global')
12. plt.xlabel('Year')
13. plt.ylabel('Temperature (°C)')
14. plt.show()
```





The moving average temperature (°C) drops for local & global data is 23.53°C & 7.2°C respectively. And the highest moving average temperature (°C) for global & local data is 9.59°C & 25.2°C respectively.

## Imputing Missing values in Raw Data & Visualizing

The inconsistencies in data leads to the misleading decisions. To minimize the risks leading to inappropriate decisions, I applied imputer to the missing data for the local temperatures (°C).

Finding the count of missing values in each column of data for local temperatures (°C) and global temperatures (°C) as well. Code and results shown side by side.

```
1. df_CityData.isnull().sum()
```

```
year      0
city      0
country   0
avg_temp  22
m_avg_temp 54
dtype: int64
```

```
1. df_GlobalData.isnull().sum()
```

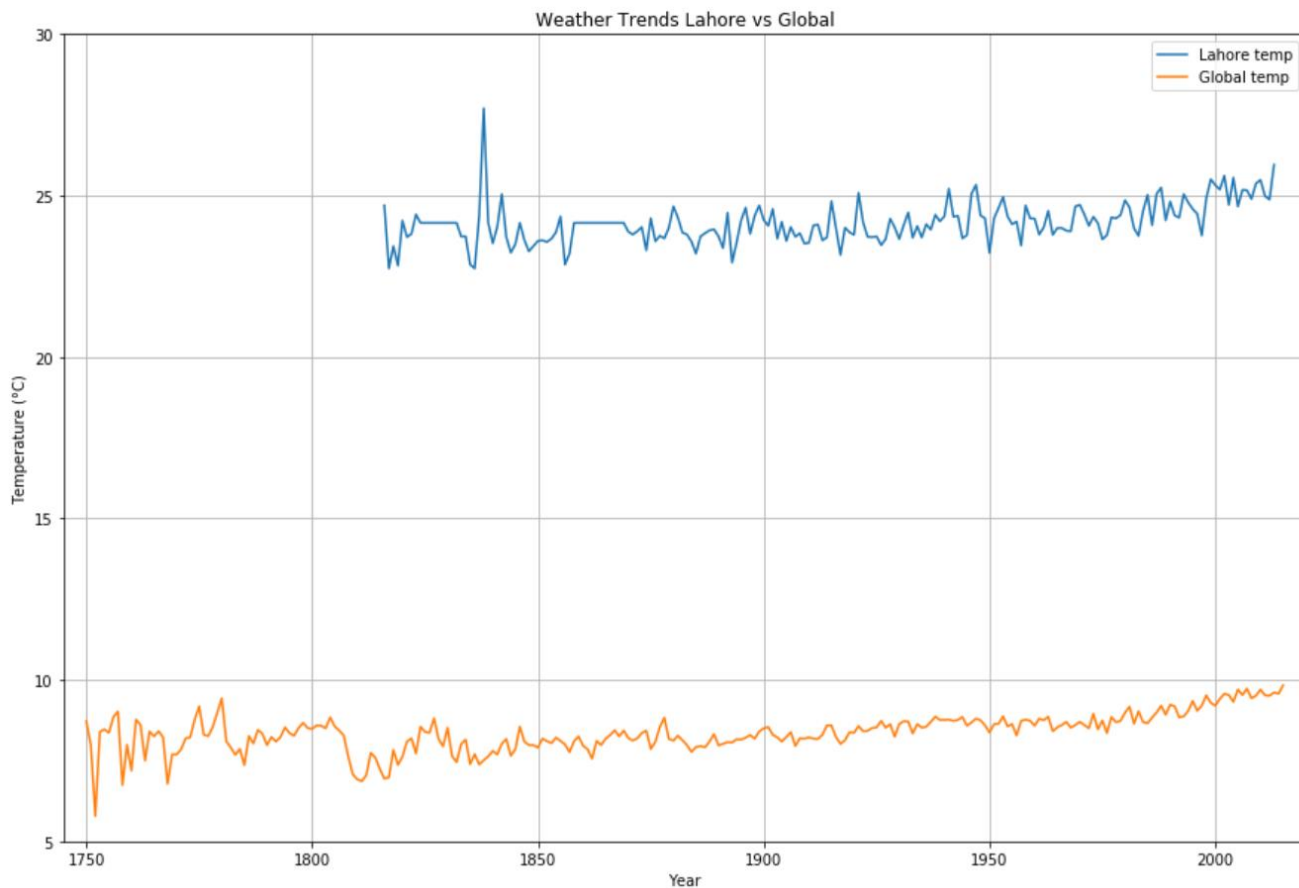
```
year      0
avg_temp  0
m_avg_temp 9
dtype: int64
```

You can see local temperatures (°C) have '22' missing values in 'avg\_temp' which results in missing of '54' values for 'm\_avg\_temp'. These missing values breaks the continuity of data trend in chart. Now I'll impute the missing values with the **mean** of the 'avg\_temp' in the new data frame.

```
1. df_CityData_imp = df_CityData.copy()
2. df_CityData_imp['imp_avg_temp'] =
   df_CityData_imp['avg_temp'].fillna(df_CityData['avg_temp'].mean())
```

After replacing missing values with the mean temperature. Plotting the data.

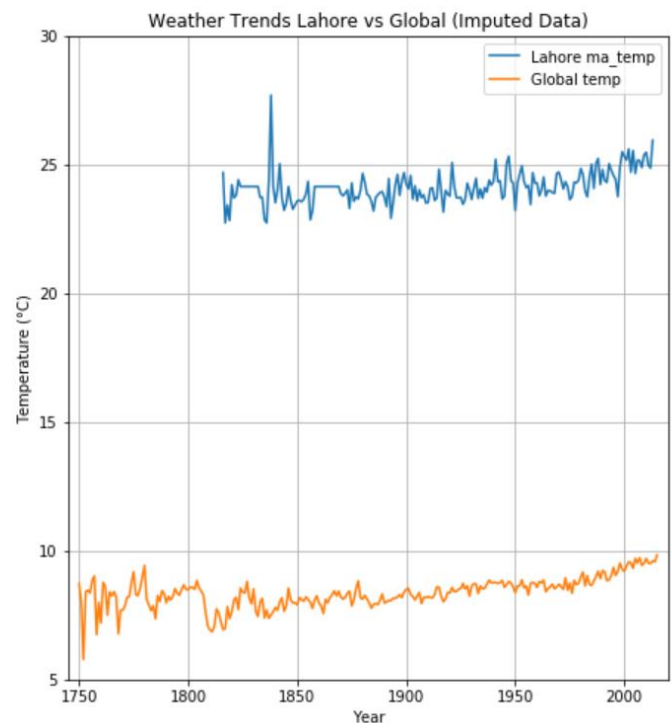
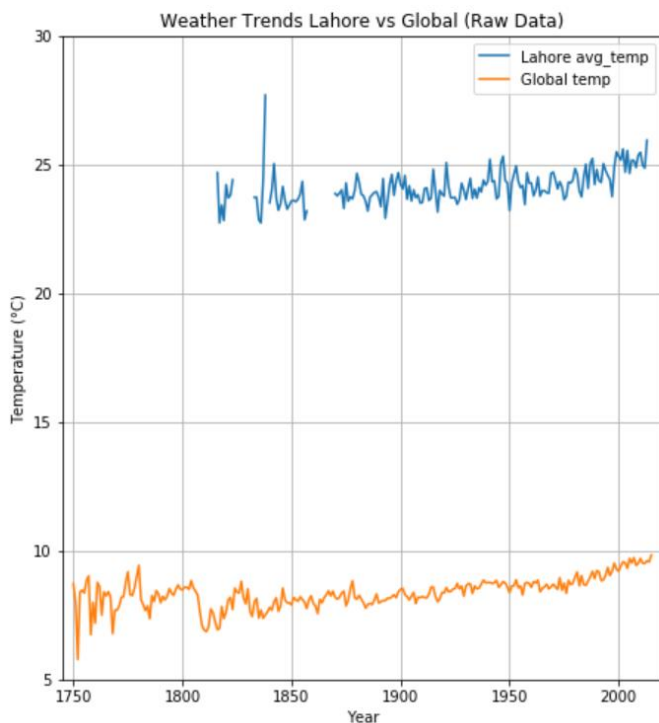
```
1. plt.figure(figsize=[15,10])
2. plt.plot(df_CityData_imp['year'],
3.         df_CityData_imp['imp_avg_temp'],
4.         label='Lahore temp')
5. plt.plot(df_GlobalData['year'],
6.         df_GlobalData['avg_temp'],
7.         label='Global temp')
8. plt.legend()
9. plt.grid(True)
10. plt.axis([1745, 2020, 5, 30])
11. plt.title('Weather Trends Lahore vs Global')
12. plt.xlabel('Year')
13. plt.ylabel('Temperature (°C)')
14. plt.show()
```



Now you can see the trend line is continuous from year 1816 to year 2013 for 'Lahore temp' despite of its spikes and fluctuations.

Now, plotting raw data along with the imputed data side by side.

```
1. plt.figure(figsize=(16,8))
2.
3. plt.subplot(121)
4. plt.plot(df_CityData['year'],
5.          df_CityData['avg_temp'],
6.          label='Lahore avg_temp')
7. plt.plot(df_GlobalData['year'],
8.          df_GlobalData['avg_temp'],
9.          label='Global temp')
10. plt.legend()
11. plt.grid(True)
12. plt.axis([1745, 2020, 5, 30])
13. plt.title('Weather Trends Lahore vs Global (Raw Data)')
14. plt.xlabel('Year')
15. plt.ylabel('Temperature (°C)')
16.
17. plt.subplot(122)
18. plt.plot(df_CityData_imp['year'],
19.          df_CityData_imp['imp_avg_temp'],
20.          label='Lahore ma_temp')
21. plt.plot(df_GlobalData['year'],
22.          df_GlobalData['avg_temp'],
23.          label='Global temp')
24. plt.legend()
25. plt.grid(True)
26. plt.axis([1745, 2020, 5, 30])
27. plt.title('Weather Trends Lahore vs Global (Imputed Data)')
28. plt.xlabel('Year')
29. plt.ylabel('Temperature (°C)')
30. plt.show()
```



## Calculating Moving Averages on Imputed Data & Visualizing

Calculating moving averages for both local & global data temperatures (°C) for the span of a decade.

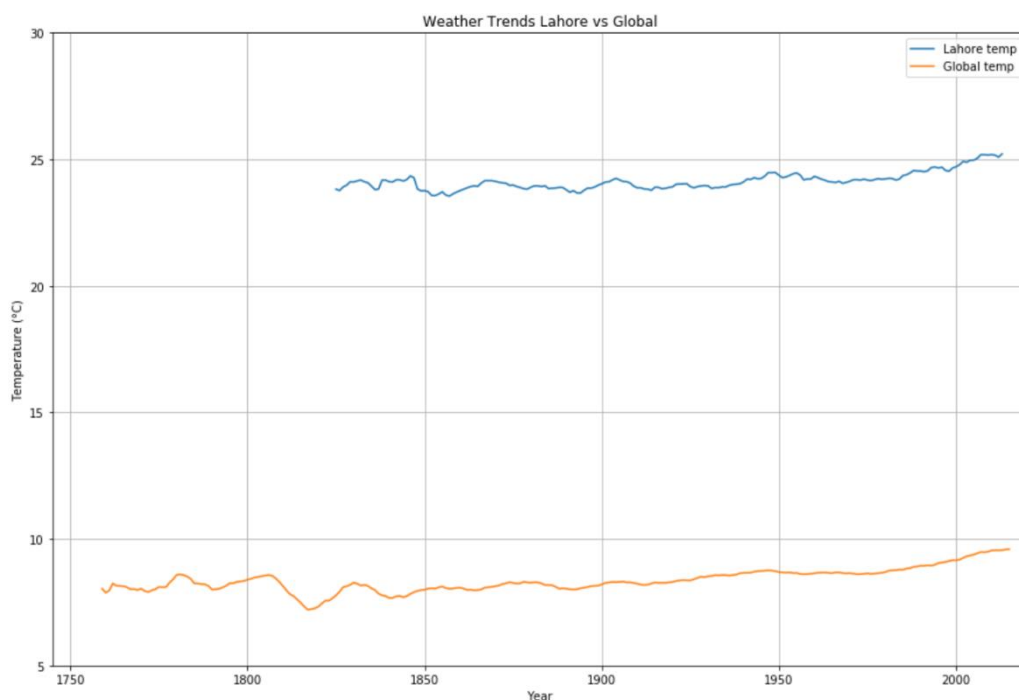
### Why the window size or moving average calculating for 10 years or decade?

As the dispersion of data will be reduced nearly by 25%, and lesser the dispersion of data the trend line will be smoother.

```
1. df_CityData_imp['m_imp_avg_temp'] =  
   df_CityData_imp['imp_avg_temp'].rolling(10).mean()
```

Now plotting the data with no missing values and with a moving average span of a decade.

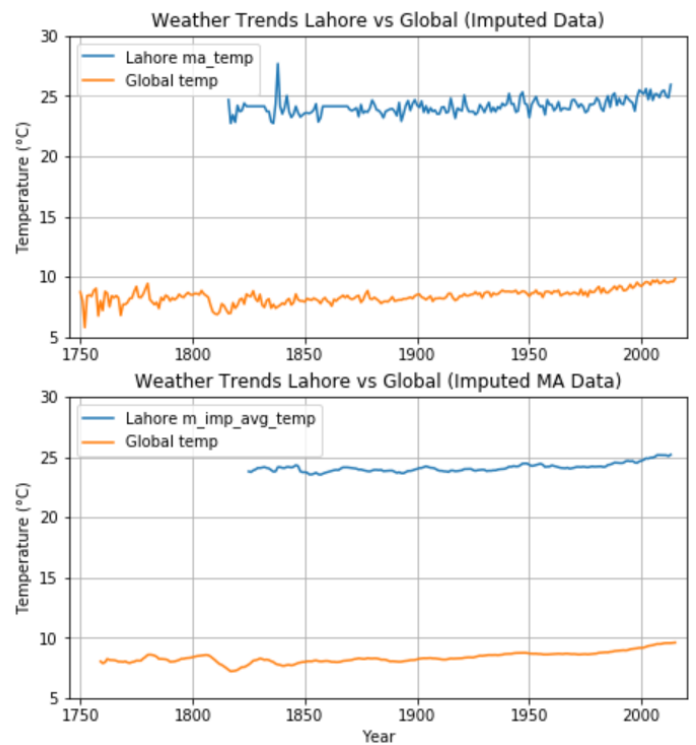
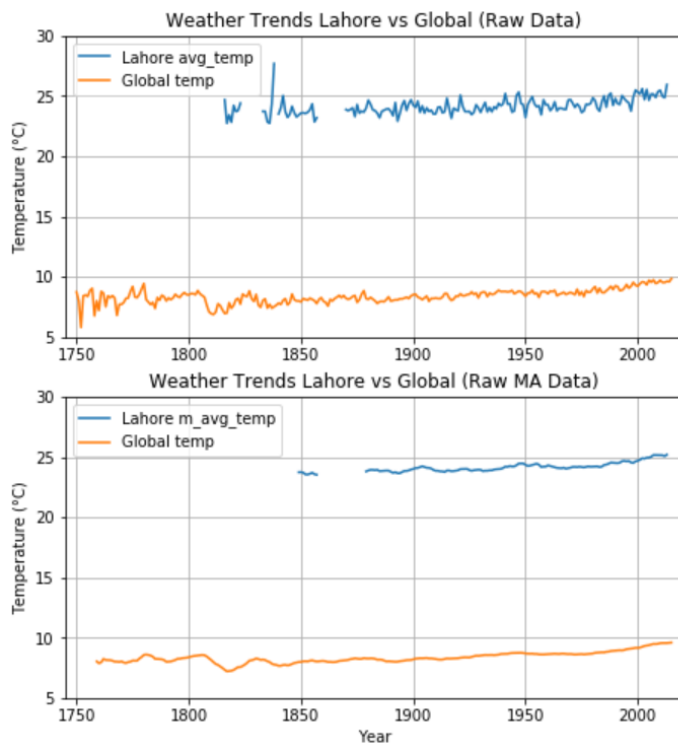
```
1. plt.figure(figsize=[15,10])  
2. plt.plot(df_CityData_imp['year'],  
3.         df_CityData_imp['m_imp_avg_temp'],  
4.         label='Lahore temp')  
5. plt.plot(df_GlobalData['year'],  
6.         df_GlobalData['m_avg_temp'],  
7.         label='Global temp')  
8. plt.legend()  
9. plt.grid(True)  
10. plt.axis([1745, 2020, 5, 30])  
11. plt.title('Weather Trends Lahore vs Global')  
12. plt.xlabel('Year')  
13. plt.ylabel('Temperature (°C)')  
14. plt.show()
```





All the plots in one chart row data along with its moving averages and imputed data along with its moving averages.

```
1. plt.figure(figsize=(16,8))
2.
3. plt.subplot(221)
4. plt.plot(df_CityData['year'],
5.          df_CityData['avg_temp'],
6.          label='Lahore avg_temp')
7. plt.plot(df_GlobalData['year'],
8.          df_GlobalData['avg_temp'],
9.          label='Global temp')
10. plt.legend()
11. plt.grid(True)
12. plt.axis([1745, 2020, 5, 30])
13. plt.title('Weather Trends Lahore vs Global (Raw Data)')
14. plt.ylabel('Temperature (°C)')
15.
16. plt.subplot(222)
17. plt.plot(df_CityData_imp['year'],
18.          df_CityData_imp['imp_avg_temp'],
19.          label='Lahore ma_temp')
20. plt.plot(df_GlobalData['year'],
21.          df_GlobalData['avg_temp'],
22.          label='Global temp')
23. plt.legend()
24. plt.grid(True)
25. plt.axis([1745, 2020, 5, 30])
26. plt.title('Weather Trends Lahore vs Global (Imputed Data)')
27. plt.ylabel('Temperature (°C)')
28.
29. plt.subplot(223)
30. plt.plot(df_CityData['year'],
31.          df_CityData['m_avg_temp'],
32.          label='Lahore m_avg_temp')
33. plt.plot(df_GlobalData['year'],
34.          df_GlobalData['m_avg_temp'],
35.          label='Global temp')
36. plt.legend()
37. plt.grid(True)
38. plt.axis([1745, 2020, 5, 30])
39. plt.title('Weather Trends Lahore vs Global (Raw MA Data)')
40. plt.xlabel('Year')
41. plt.ylabel('Temperature (°C)')
42.
43. plt.subplot(224)
44. plt.plot(df_CityData_imp['year'],
45.          df_CityData_imp['m_imp_avg_temp'],
46.          label='Lahore m_imp_avg_temp')
47. plt.plot(df_GlobalData['year'],
48.          df_GlobalData['m_avg_temp'],
49.          label='Global temp')
50. plt.legend()
51. plt.grid(True)
52. plt.axis([1745, 2020, 5, 30])
53. plt.title('Weather Trends Lahore vs Global (Imputed MA Data)')
54. plt.xlabel('Year')
55. plt.ylabel('Temperature (°C)')
56. plt.show()
```

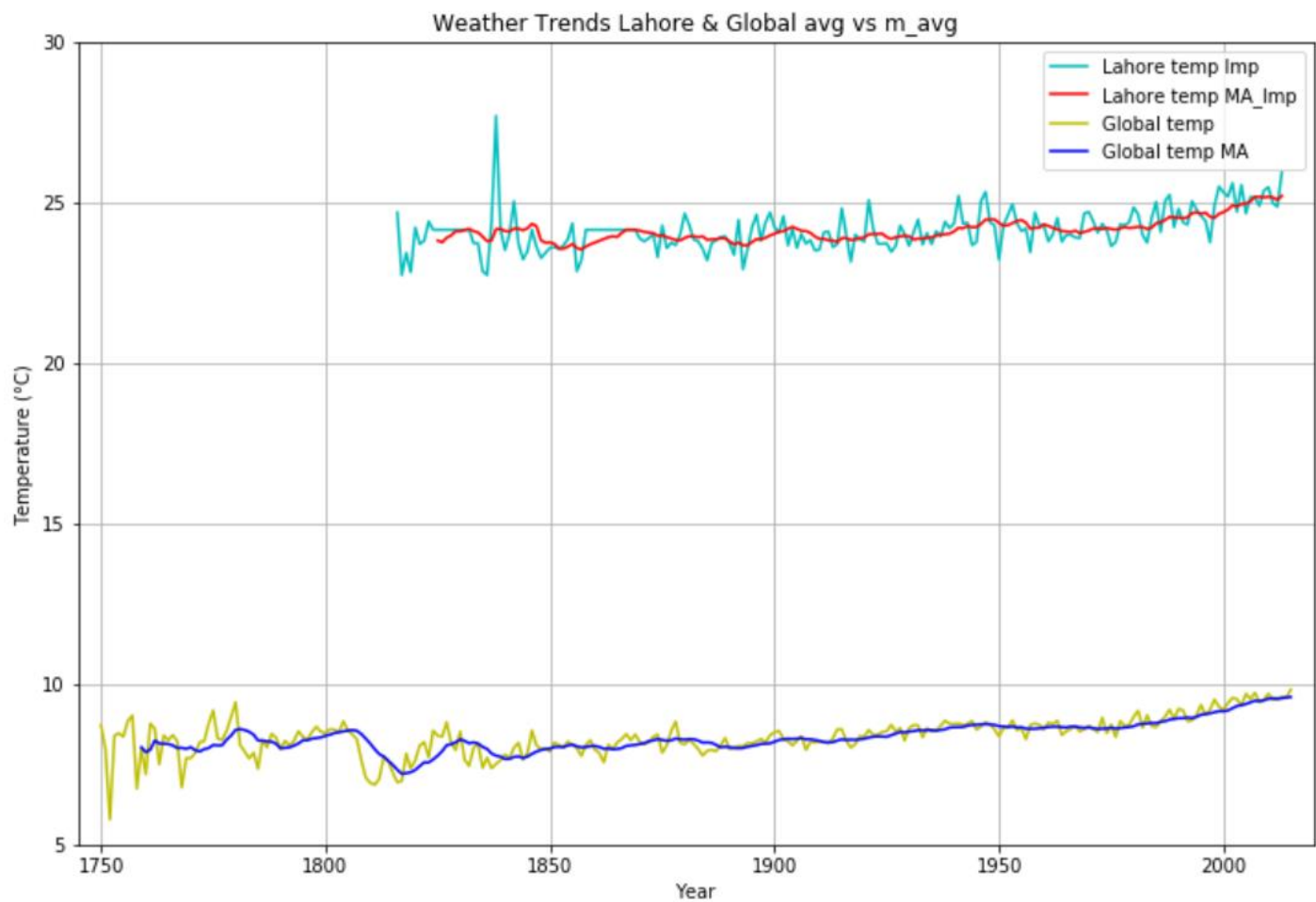


Last, but not least. Plotting the data with no missing values and average temperatures along with calculated moving averages.

```

1. plt.figure(figsize=[12,8])
2. plt.plot(df_CityData_imp['year'],
3.          df_CityData_imp['imp_avg_temp'],
4.          label='Lahore temp Imp', color= 'c')
5. plt.plot(df_CityData_imp['year'],
6.          df_CityData_imp['m_imp_avg_temp'],
7.          label='Lahore temp MA_Imp', color= 'r')
8. plt.plot(df_GlobalData['year'],
9.          df_GlobalData['avg_temp'],
10.         label='Global temp', color= 'y')
11. plt.plot(df_GlobalData['year'],
12.          df_GlobalData['m_avg_temp'],
13.          label='Global temp MA', color= 'b')
14. plt.legend()
15. plt.grid(True)
16. plt.axis([1745, 2020, 5, 30])
17. plt.title('Weather Trends Lahore & Global avg vs m_avg')
18. plt.xlabel('Year')
19. plt.ylabel('Temperature (°C)')
20. plt.show()

```



---

# CONCLUSION

From the extraction of the data from the database, transforming it into a csv file and loading it into a data-frame. Then performing the analysis with several perspectives. I come with some discrete finding from the data which will govern your decisions.

- Today, we are standing at the highest level of temperatures (°C) globally from last two centuries.
- Lahore's temperatures (°C) are relatively twice than the global average temperatures (°C).
- Over few decades, local temperatures (°C) are increasing consistently.
- Over time, the city temperatures have changed similarly with the global average.
- The world is getting hotter day by day, as the global and local temperatures (°C) exhibit an uphill trend.
- Over few decades, global temperatures (°C) are increasing consistently.



---

# REFERENCES

1. <https://www.datacamp.com/community/tutorials/moving-averages-in-pandas> (moving averages)
2. <https://towardsdatascience.com/stop-using-mean-to-fill-missing-data-678c0d396e22> (imputing data)
3. <https://matplotlib.org/3.1.1/tutorials/introductory/pyplot.html#sphx-glr-tutorials-introductory-pyplot-py> (plotting the data)
4. <https://knowledge.udacity.com/questions/372272> (Plotting data)