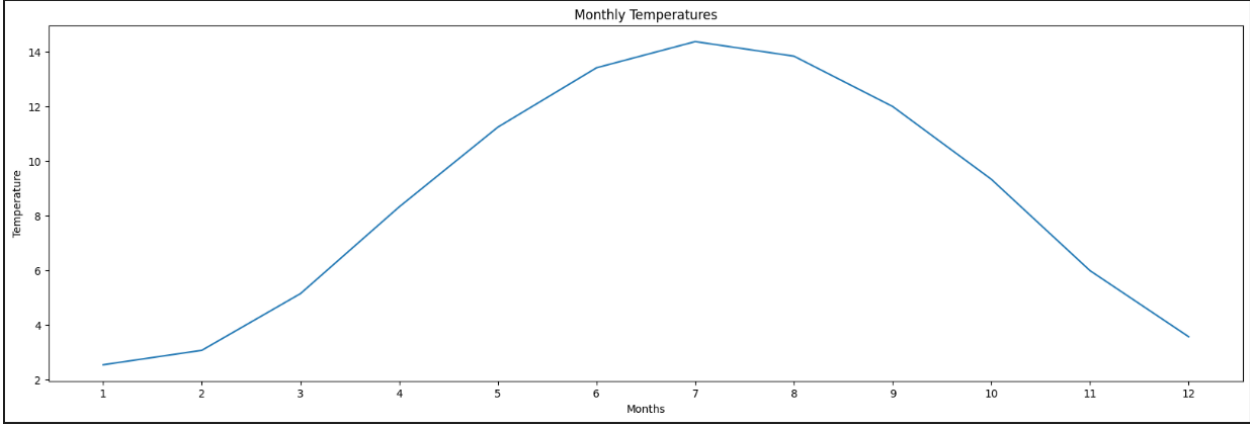


Data Collection and Preprocessing Phase

Date	9 July 2024
Team ID	SWTID1720111029
Project Title	Unveiling Climate Change Dynamics through Earth Surface Temperature Analysis
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

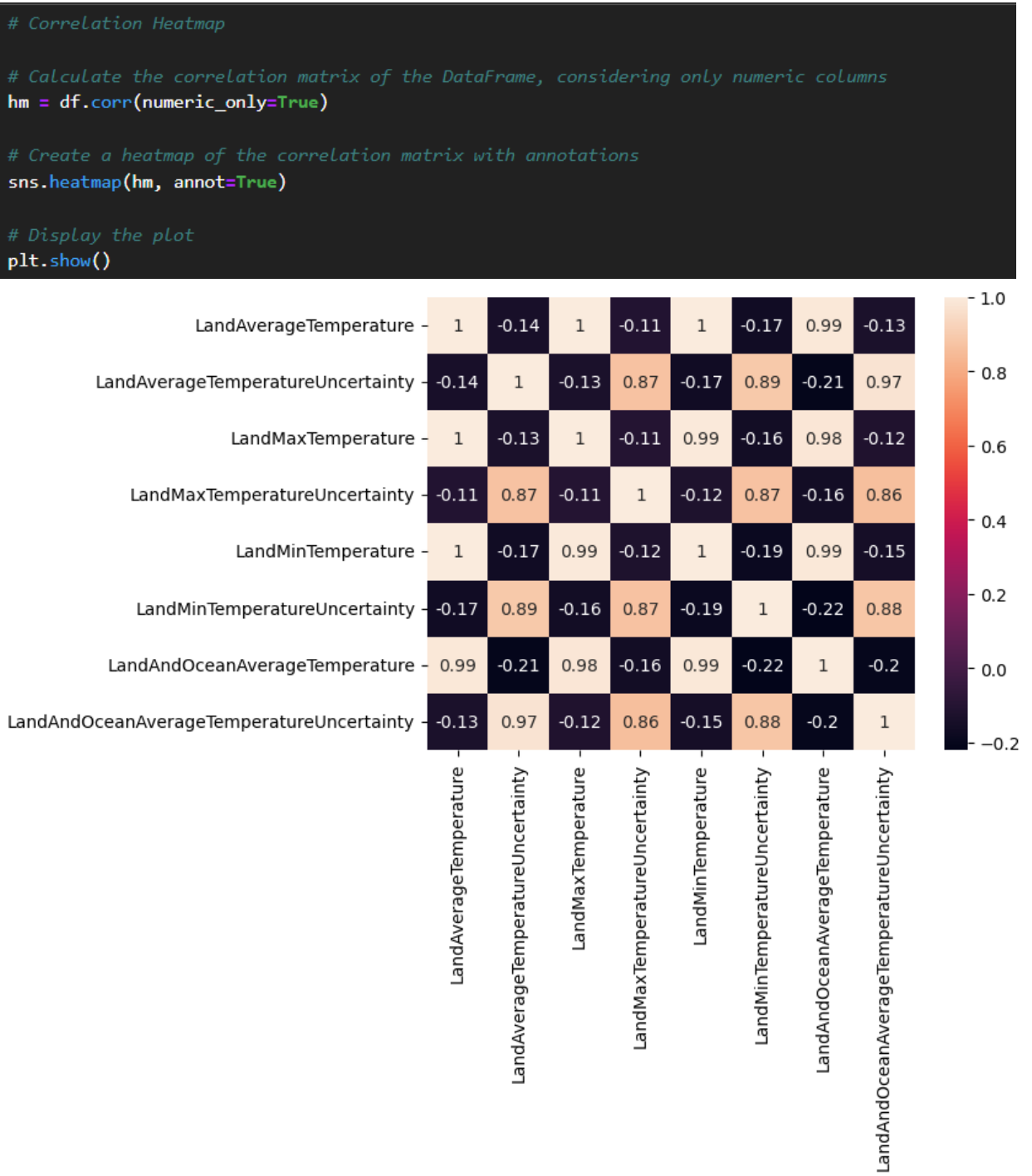
Section	Description																																																							
Data Overview	<div>RangeIndex: 3192 entries, 0 to 3191</div> <div>Data columns (total 9 columns):</div> <table><tr><th>#</th><th>Column</th><th>Non-Null</th><th>Count</th><th>Dtype</th></tr><tr><td>---</td><td>-----</td><td>-----</td><td>-----</td><td></td></tr><tr><td>0</td><td>dt</td><td>3192 non-null</td><td>object</td><td></td></tr><tr><td>1</td><td>LandAverageTemperature</td><td>3180 non-null</td><td>float64</td><td></td></tr><tr><td>2</td><td>LandAverageTemperatureUncertainty</td><td>3180 non-null</td><td>float64</td><td></td></tr><tr><td>3</td><td>LandMaxTemperature</td><td>1992 non-null</td><td>float64</td><td></td></tr><tr><td>4</td><td>LandMaxTemperatureUncertainty</td><td>1992 non-null</td><td>float64</td><td></td></tr><tr><td>5</td><td>LandMinTemperature</td><td>1992 non-null</td><td>float64</td><td></td></tr><tr><td>6</td><td>LandMinTemperatureUncertainty</td><td>1992 non-null</td><td>float64</td><td></td></tr><tr><td>7</td><td>LandAndOceanAverageTemperature</td><td>1992 non-null</td><td>float64</td><td></td></tr><tr><td>8</td><td>LandAndOceanAverageTemperatureUncertainty</td><td>1992 non-null</td><td>float64</td><td></td></tr></table> <div>dtypes: float64(8), object(1)</div> <div>memory usage: 224.6+ KB</div>	#	Column	Non-Null	Count	Dtype	---	-----	-----	-----		0	dt	3192 non-null	object		1	LandAverageTemperature	3180 non-null	float64		2	LandAverageTemperatureUncertainty	3180 non-null	float64		3	LandMaxTemperature	1992 non-null	float64		4	LandMaxTemperatureUncertainty	1992 non-null	float64		5	LandMinTemperature	1992 non-null	float64		6	LandMinTemperatureUncertainty	1992 non-null	float64		7	LandAndOceanAverageTemperature	1992 non-null	float64		8	LandAndOceanAverageTemperatureUncertainty	1992 non-null	float64	
#	Column	Non-Null	Count	Dtype																																																				
---	-----	-----	-----																																																					
0	dt	3192 non-null	object																																																					
1	LandAverageTemperature	3180 non-null	float64																																																					
2	LandAverageTemperatureUncertainty	3180 non-null	float64																																																					
3	LandMaxTemperature	1992 non-null	float64																																																					
4	LandMaxTemperatureUncertainty	1992 non-null	float64																																																					
5	LandMinTemperature	1992 non-null	float64																																																					
6	LandMinTemperatureUncertainty	1992 non-null	float64																																																					
7	LandAndOceanAverageTemperature	1992 non-null	float64																																																					
8	LandAndOceanAverageTemperatureUncertainty	1992 non-null	float64																																																					
Univariate Analysis	<div><pre>pivot = data.pivot_table(values='LandAverageTemperature', index=data.index.year, columns=data.index.month)  # Plot the monthly seasonality monthly_seasonality = pivot.mean(axis=0) monthly_seasonality.plot(figsize=(20, 6)) plt.title('Monthly Temperatures') plt.xlabel('Months') plt.ylabel('Temperature') plt.xticks(range(1, 13)) plt.show()</pre></div> <p>Seasonal fluctuations in temperature on a monthly basis. Monthly temperature variations reflecting seasonal patterns.</p>																																																							

Bivariate Analysis



The mean temperatures across the land for each of the 12 months over multiple years.

Multivariate Analysis



Data Preprocessing Code Screenshots	
Loading Data	<pre>df = pd.read_csv('GlobalTemperatures.csv')</pre>
Handling Missing Data	<p>Missing values in LandAverageTemperature and LandAverageTemperatureUncertainty columns were replaced with the mean and missing values in other columns were dropped.</p> <pre>df.isnull().sum()</pre> <pre>dt LandAverageTemperature LandAverageTemperatureUncertainty LandMaxTemperature LandMaxTemperatureUncertainty LandMinTemperature LandMinTemperatureUncertainty LandAndOceanAverageTemperature LandAndOceanAverageTemperatureUncertainty dtype: int64</pre> <pre># Impute LandAverageTemperature and LandAverageTemperatureUncertainty with mean # Fill missing values in the 'LandAverageTemperature' column with the mean of that column df['LandAverageTemperature'].fillna(df['LandAverageTemperature'].mean(), inplace=True) # Fill missing values in the 'LandAverageTemperatureUncertainty' column with the mean of that column df['LandAverageTemperatureUncertainty'].fillna(df['LandAverageTemperatureUncertainty'].mean(), inplace=True)  # For columns with 1200 missing values, drop those rows # List of columns to check for missing values cols_to_dropna = ['LandMaxTemperature', 'LandMaxTemperatureUncertainty', 'LandMinTemperature', 'LandMinTemperatureUncertainty', 'LandAndOceanAverageTemperature', 'LandAndOceanAverageTemperatureUncertainty'] # Loop through each column in the list for col in cols_to_dropna:     # Drop rows where the specified column has missing values     df.dropna(subset=[col], inplace=True)  # Verify if there are any remaining missing values print(df.isnull().sum())</pre> <pre>dt LandAverageTemperature LandAverageTemperatureUncertainty LandMaxTemperature LandMaxTemperatureUncertainty LandMinTemperature LandMinTemperatureUncertainty LandAndOceanAverageTemperature LandAndOceanAverageTemperatureUncertainty dtype: int64</pre>
Data Transformation	<p>Scaling the Xtrain, Xtest , Ytrain &amp; Ytest values to be suitable for LSTM model.</p> <pre># Scale X and y using MinMaxScaler scaler_x = MinMaxScaler() scaler_y = MinMaxScaler()  # Fit and transform the training data X_train_scaled = scaler_x.fit_transform(X_train) y_train_scaled = scaler_y.fit_transform(y_train.values.reshape(-1, 1)) # Reshape y_train to a 2D array  # Only transform the testing data X_test_scaled = scaler_x.transform(X_test) y_test_scaled = scaler_y.transform(y_test.values.reshape(-1, 1)) # Reshape y_test to a 2D array</pre>
Feature Engineering	<p>Additionally adding year and month columns on the basis of the date column.</p> <pre># Add Year and Month columns based on 'dt' column df['Year'] = pd.to_datetime(df['dt']).dt.year df['Month'] = pd.to_datetime(df['dt']).dt.month</pre>

Save Processed Data

```
# Save the scaled data
joblib.dump({
    'X_train_scaled': X_train_scaled,
    'y_train_scaled': y_train_scaled,
    'X_test_scaled': X_test_scaled,
    'y_test_scaled': y_test_scaled
}, 'scaled_data.pkl')
```