

Understanding and Predicting Climate Change Metrics

1. Project Question

The main question guiding this project is: How can we understand and predict the impact of climate change metrics, specifically CO₂ levels and temperatures, over time?

This involves analyzing historical climate data to identify trends and using machine learning models to predict future changes in CO₂ concentrations and temperature averages.

Understanding and Predicting Climate Change Metrics

2. Data Sources

The data sources for this project were chosen based on their comprehensive coverage of climate-related metrics and their credibility. The datasets used include:

- Global Land Temperatures by City: Provides historical temperature data for various cities around the world.
- Temperatures by Country: Contains historical temperature data aggregated by country.
- Temperatures by Major City: Focuses on temperature data for major global cities.
- Temperatures by State: Includes temperature records for states within countries.
- Global Temperatures: A global dataset for temperature records.

These datasets are sourced from Kaggle's climate change datasets and contain various climate metrics such as average temperature, CO2 levels, and timestamps. Each dataset was chosen for its relevance and the granularity of the data it provides, which is essential for detailed analysis and prediction.

Licenses: Most of these datasets are available under open-data licenses, which allow for usage and distribution with proper attribution. We ensured compliance with these licenses by acknowledging the sources in our documentation and not using the data for any commercial purposes.

Understanding and Predicting Climate Change Metrics

3. Data Pipeline

Overview

The data pipeline is implemented using Python, leveraging libraries such as Pandas for data manipulation, Matplotlib and Seaborn for visualization, and Scikit-Learn for machine learning tasks.

The pipeline involves several stages:

1. Data Loading: Importing the datasets into Pandas DataFrames.
2. Data Cleaning: Dropping missing values to ensure data integrity.
3. Date Format Modification: Converting date columns to extract and use the year for aggregations.
4. Data Transformation: Aggregating temperature data by year and region for more manageable and insightful analysis.
5. Data Visualization: Creating plots to visualize trends in temperature and CO2 levels.
6. Prediction: Using linear regression models to predict future CO2 levels based on historical data.

Transformation and Cleaning Steps

- Dropping NaN Values: All datasets underwent a cleaning process where rows with NaN values were removed to prevent inaccuracies in analysis.
- Date Handling: Date columns were converted to extract the year, which simplified the aggregation of temperature data on a yearly basis.
- Aggregation: For instance, temperature data was aggregated by averaging temperatures per year and filtering specific data, like focusing on the United States.

Problems and Solutions

- Missing Values: The primary issue was handling missing values, which was resolved by dropping

Understanding and Predicting Climate Change Metrics

any rows with NaN values.

- Date Format Consistency: Different date formats were standardized to facilitate easier manipulation and analysis.
- Data Aggregation: Aggregating large datasets required efficient handling and computation, which was managed by utilizing Pandas powerful groupby and aggregation functions.

Error Handling

The pipeline includes basic error handling mechanisms to manage:

- Missing Data: Handled by dropping NaN values.
- Changing Input Data: The pipeline can be re-run with new data without requiring significant modifications, making it robust to changes in input data.

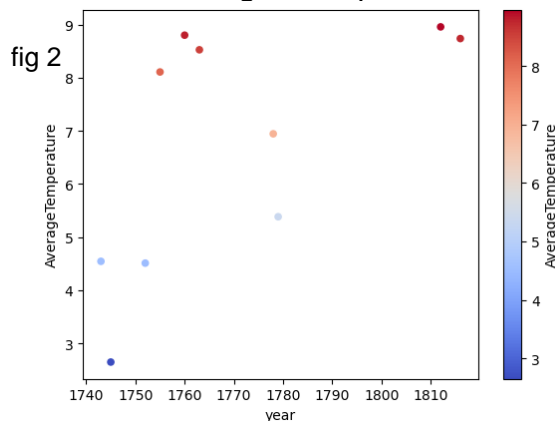


Figure 1 This scatter plot shows the average temperature over the years. There is a visible trend of increasing average temperatures over time, particularly in the most recent years.

Figure 2 This scatter plot also shows the average temperature over the years but for an earlier time period. The temperatures vary over this time period, with some years showing relatively higher temperatures and others showing lower.

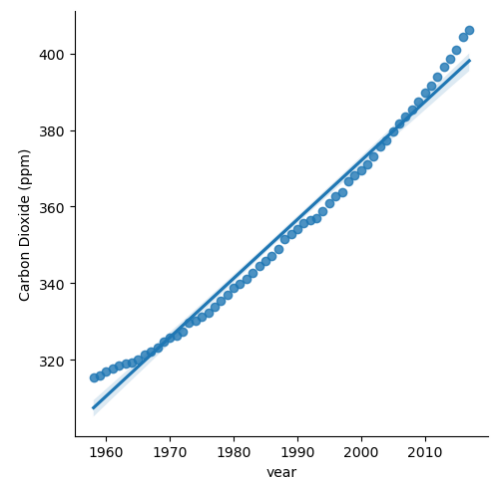
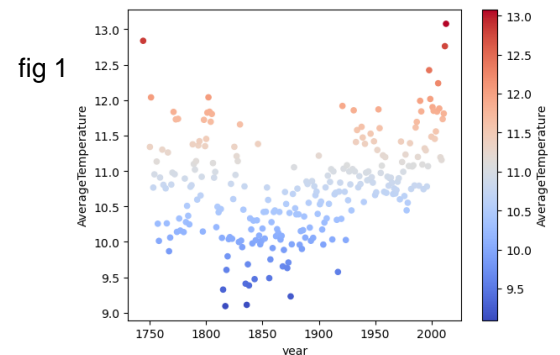


fig 3

Figure 3

This scatter plot shows the carbon dioxide (CO₂) levels over the years, with a linear regression line indicating the trend. The CO₂ levels show a consistent increase over the years, with a steep upward trend.