

Deforestation Correlation with Air Pollution in Amazon

Introduction

Human activity has led to environmental challenges like deforestation and air pollution. The Amazon plays a crucial role in climate regulation, but deforestation threatens its functions. This led to the following question:

How does deforestation in the Amazon correlate with pollution outcomes in Brazil, specifically in São Paulo?

Data Sources

Global Forest Watch - Amazon Deforestation (SAD Alerts)

- **Source:** [Global Forest Watch](#)
- **Data Description:** The dataset contains geo-referenced deforestation and degradation alerts for the Brazilian Amazon from 2008 to 2018, with attributes like date, area affected, and spatial data.
- **Reason for Choosing:** relevance to understanding deforestation's impact on global environmental changes, including pollution in São Paulo.
- **Data Structure and Quality:** It includes time-stamped records with geographic coordinates and may have some inaccuracies, addressed during preprocessing. The dataset is updated hourly and available as a CSV.
- **License:** Open Data Commons Attribution License (ODC-By), which requires proper attribution in any published work.

Kaggle- Air Pollution at São Paulo, Brazil, Since 2013

- **Source:** [Kaggle](#)
- **Data Description:** This dataset provides hourly air pollution data for São Paulo, Brazil, covering pollutants like MP10, O3, NO2, CO, MP2.5, SO2, benzene, and toluene, spanning from 2013 onward.
- **Reason for Choosing:** It's relevant for understanding the impact of Amazon deforestation to pollution levels.
- **Data Structure and Quality:** The data includes time-stamped readings with some gaps, stored in a CSV format for easy integration with the deforestation data.
- **License:** Open Data License (Kaggle dataset terms apply) which requires proper attribution in any published material.

Data Pipeline Overview

The pipeline automates the download, cleaning, transformation, and storage of two datasets: deforestation from Global Forest Watch and air pollution from Kaggle. It is implemented in Python, using libraries such as pandas, requests, matplotlib, seaborn, and sqlite3.

Pipeline Steps:

1. Data Acquisition:

Deforestation data is fetched via the Global Forest Watch API, and air pollution data is retrieved using the Kaggle API.

2. Data Cleaning & Transformation:

- **Deforestation Data:** Filtered for major deforestation events, columns renamed, time format standardized, and data aggregated by month. Missing values were interpolated.
- **Pollution Data:** Columns renamed, missing rows dropped, and pollutant levels averaged monthly. Time format standardized, and missing pollutant values filled with zero.

3. Data Storage:

The cleaned data is saved into SQLite databases: deforestation.db for aggregated deforestation data and air_pollution.db for averaged pollutant levels.

4. Visualizations:

Several plots were created to analyze trends and correlations:

- **Deforestation Trend Plot:** A line plot of the total affected area over time.
- **Pollutant Trend Plots:** Line plots for each pollutant (e.g., PM10, CO, NO2).
- **Correlation Heatmap:** The correlation between deforestation and pollutants.
- **Scatter Plots:** Deforestation vs. pollutants (e.g., PM10, PM2.5) to explore correlations.

5. Meta-Quality Measures & Error Handling:

- **Error Handling:** The pipeline retries failed downloads to ensure reliability.
- **Changing Input Data:** Pipeline fetches the latest data, avoiding issues.
- **Quality Assurance:** Data is validated after transformations to handle missing values and anomalies.

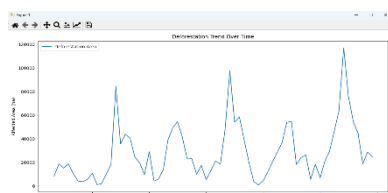


Fig. 1. Deforestation trend (monthly)

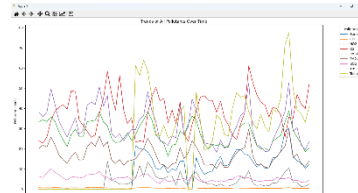


Fig. 2. Trend of all pollutants



Fig. 3. Correlation heatmap

Results, Limitations, and future prospect

Output Data Structure and Quality:

- **Structure:** The final output consists of two SQLite databases: deforestation.db and air_pollution.db. Each stores cleaned and aggregated data in structured tables.
 - **Deforestation Data:** Includes columns for date (“Date”) and total affected area (“AffectedArea”) from 2013 to 2018.
 - **Pollution Data:** Contains pollutants (e.g., PM10, PM2.5, CO, NO2) with their monthly averages from 2013 to 2018.
- **Quality:** The data is clean, with missing values interpolated and aggregated monthly. Data quality depends on the integrity of the source, especially with missing values.

Output Data Format and Reasoning:

- **Format:** The output is stored in SQLite format, chosen for its simplicity and efficiency in handling structured data. It enables easy querying and integration with other tools, making it ideal for local data storage.
- **Reasoning:** SQLite is used because it doesn't require a server, is easily queried, and ensures data consistency and reliability.

Critical Reflection on Data and Limitations:

- **Data Limitations:**
 - Deforestation data has gaps or errors due to missing alerts or delayed reporting, affecting trend accuracy.
 - The overlapping period between deforestation and pollution data may not fully capture the long-term effects of deforestation on pollution.
 - Pollution data may be incomplete or inconsistent, addressed by filling gaps with zeros or interpolation.
- **Anticipated Issues:**
 - **Correlation Validity:** External variables like transportation or seasonal weather effects may distort the correlation between deforestation and pollution. Possible solutions include studying longer periods, filtering industry-related pollutants, and using better statistical methods or Machine Learning approaches.
 - **Data Resolution:** Monthly aggregation may obscure short-term variations. Solutions include studying different temporal resolutions, using advanced interpolation methods, and investigating specific events like forest fires.

Output Visualization:

The final output is stored in deforestation.db and air_pollution.db, and a sample visualization will be provided.

PM2.5											
Year	Month	PM2.5	PM10	CO	NO2	SO2	O3	Temp	Humidity	Wind	Pressure
2013	1	12.5	25.0	1.2	0.05	0.01	0.02	15.0	65.0	10.0	1013.25
2013	2	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0	1012.50
2013	3	8.0	18.0	0.9	0.03	0.01	0.02	22.0	75.0	15.0	1011.75
2013	4	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0	1011.00
2013	5	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0	1010.25
2013	6	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0	1009.50
2013	7	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0	1008.75
2013	8	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0	1008.00
2013	9	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0	1007.25
2013	10	30.0	55.0	2.1	0.12	0.01	0.02	40.0	110.0	30.0	1006.50
2013	11	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0	1007.25
2013	12	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0	1008.00
2014	1	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0	1008.75
2014	2	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0	1009.50
2014	3	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0	1010.25
2014	4	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0	1011.00
2014	5	12.5	25.0	1.2	0.05	0.01	0.02	22.0	75.0	15.0	1011.75
2014	6	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0	1012.50
2014	7	8.0	18.0	0.9	0.03	0.01	0.02	15.0	65.0	10.0	1013.25
2014	8	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0	1012.50
2014	9	12.5	25.0	1.2	0.05	0.01	0.02	22.0	75.0	15.0	1011.75
2014	10	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0	1011.00
2014	11	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0	1010.25
2014	12	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0	1009.50
2015	1	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0	1008.75
2015	2	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0	1008.00
2015	3	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0	1007.25
2015	4	30.0	55.0	2.1	0.12	0.01	0.02	40.0	110.0	30.0	1006.50
2015	5	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0	1007.25
2015	6	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0	1008.00
2015	7	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0	1008.75
2015	8	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0	1009.50
2015	9	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0	1010.25
2015	10	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0	1011.00
2015	11	12.5	25.0	1.2	0.05	0.01	0.02	22.0	75.0	15.0	1011.75
2015	12	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0	1012.50

Fig. 5. Pollution.db

Deforestation											
Year	Month	Deforestation	PM2.5	PM10	CO	NO2	SO2	O3	Temp	Humidity	Wind
2013	1	10.0	12.5	25.0	1.2	0.05	0.01	0.02	15.0	65.0	10.0
2013	2	8.0	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0
2013	3	5.0	8.0	18.0	0.9	0.03	0.01	0.02	22.0	75.0	15.0
2013	4	12.0	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0
2013	5	15.0	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0
2013	6	18.0	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0
2013	7	20.0	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0
2013	8	22.0	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0
2013	9	25.0	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0
2013	10	28.0	30.0	55.0	2.1	0.12	0.01	0.02	40.0	110.0	30.0
2013	11	25.0	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0
2013	12	22.0	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0
2014	1	20.0	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0
2014	2	18.0	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0
2014	3	15.0	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0
2014	4	12.0	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0
2014	5	10.0	12.5	25.0	1.2	0.05	0.01	0.02	22.0	75.0	15.0
2014	6	8.0	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0
2014	7	5.0	8.0	18.0	0.9	0.03	0.01	0.02	15.0	65.0	10.0
2014	8	12.0	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0
2014	9	15.0	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0
2014	10	18.0	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0
2014	11	20.0	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0
2014	12	22.0	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0
2015	1	25.0	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0
2015	2	28.0	30.0	55.0	2.1	0.12	0.01	0.02	40.0	110.0	30.0
2015	3	25.0	28.0	50.0	2.0	0.11	0.01	0.02	38.0	105.0	28.0
2015	4	22.0	25.0	45.0	1.8	0.10	0.01	0.02	35.0	100.0	26.0
2015	5	20.0	22.0	40.0	1.7	0.09	0.01	0.02	32.0	95.0	24.0
2015	6	18.0	20.0	38.0	1.6	0.08	0.01	0.02	30.0	90.0	22.0
2015	7	15.0	18.0	35.0	1.5	0.07	0.01	0.02	28.0	85.0	20.0
2015	8	12.0	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0
2015	9	10.0	12.5	25.0	1.2	0.05	0.01	0.02	22.0	75.0	15.0
2015	10	8.0	10.0	20.0	1.0	0.04	0.01	0.02	18.0	70.0	12.0
2015	11	5.0	8.0	18.0	0.9	0.03	0.01	0.02	15.0	65.0	10.0
2015	12	12.0	15.0	30.0	1.3	0.06	0.01	0.02	25.0	80.0	18.0

Fig. 4. Deforestation.db