

NNDSS

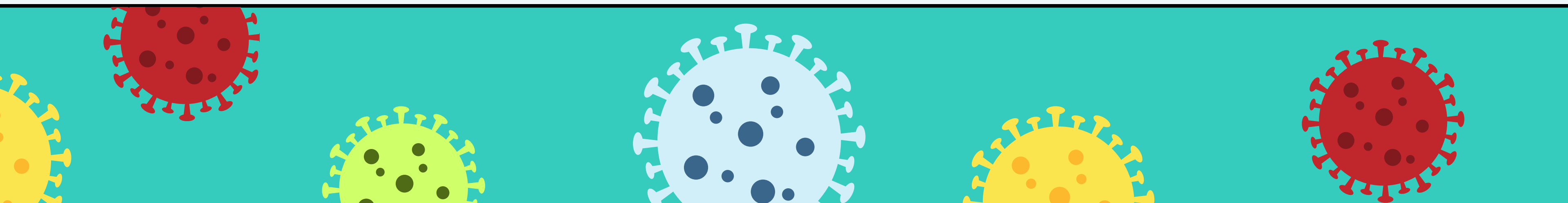
DATA

PIPELINE

PROJECT

RECAP

- Extracted the data from NNDSS website using API calls and Postman
- Data transformation using BigQuery
- EDA
- Set up the automation using Google cloud Scheduler

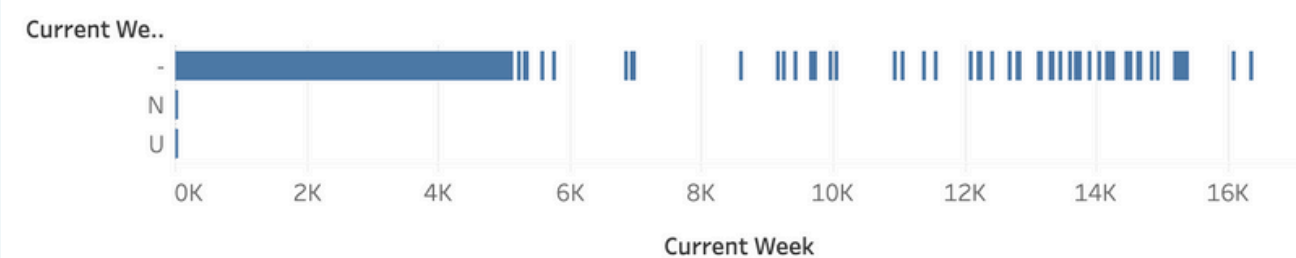


DATASET

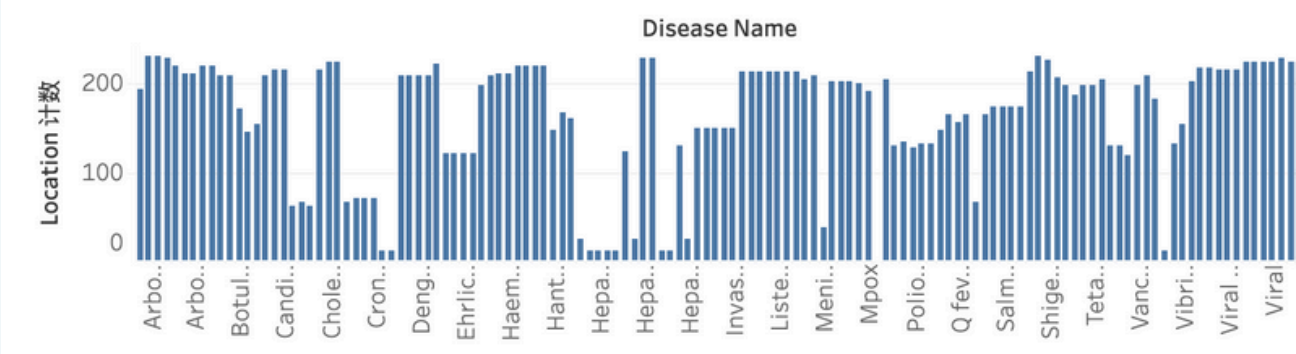
- **Dataset Source and Purpose:** National Notifiable Diseases Surveillance System (NNDSS) by the CDC, which tracks weekly reported cases of various diseases across different U.S. locations.
- Its primary purpose is to monitor disease trends, detect outbreaks, and inform public health responses.
- **Time Span and Granularity:** The dataset spans from 2022 onwards, and is organized by weekly reports (MMWR weeks). Each record provides detailed information on disease counts per location, allowing for both short-term and seasonal analysis of trends.
- **Key Variables:**
 - Disease Name: Specifies the type of disease (e.g., influenza, measles).
 - Location Information: Includes state-level data, latitude, and longitude coordinates.
 - Case Counts: Current weekly counts, along with cumulative counts and rolling averages, which are useful for trend and anomaly analysis.

DASHBOARD

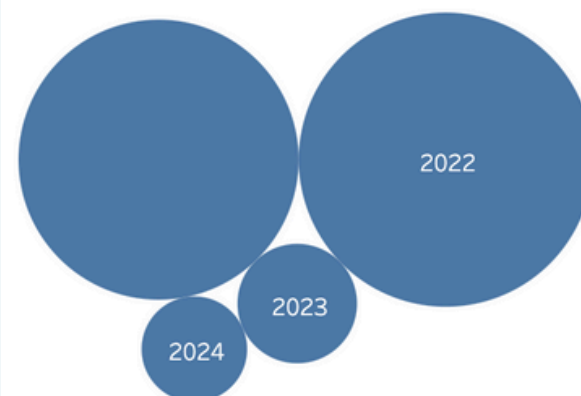
Current Week Data Distribution



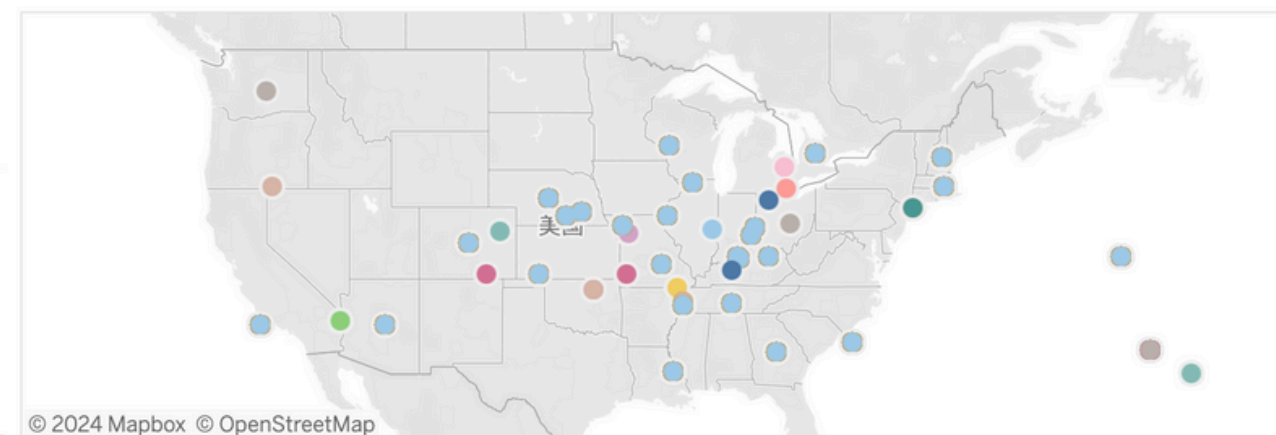
Disease Frequency by Type



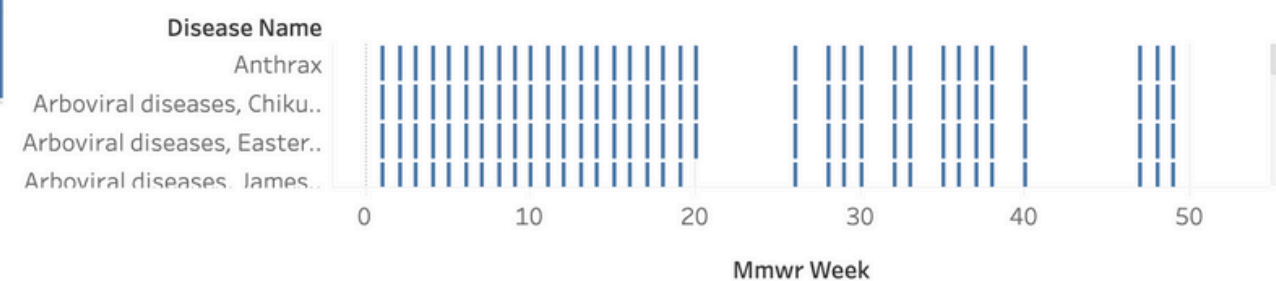
Cases Distribution by Year



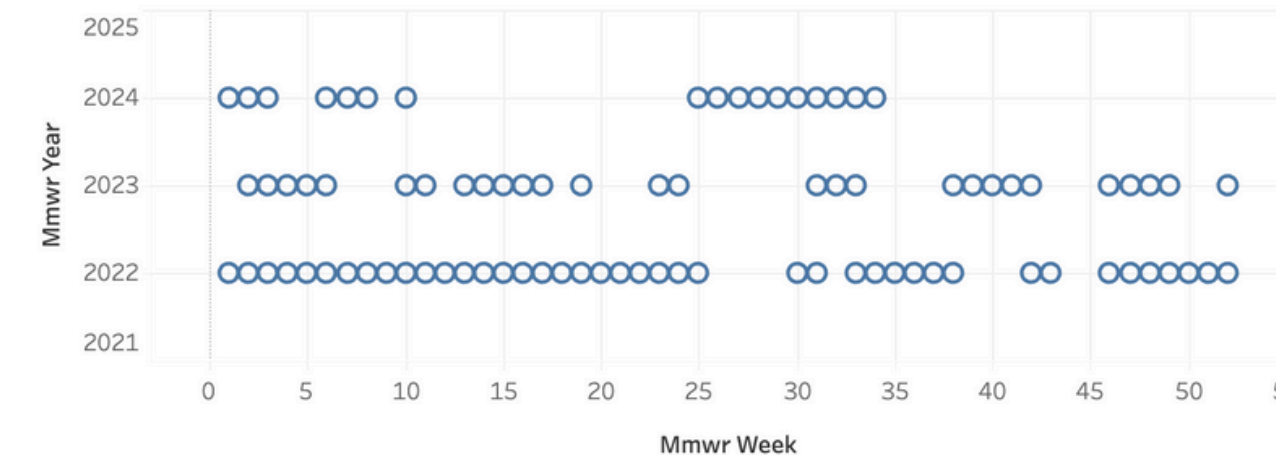
Disease Cases by Location on Map



Disease Cases Trend by MMWR Week



Yearly Disease Cases Timeline by MMWR Week



Disease Name

- ☒ (全部)
- ☒ Null
- ☒ Anthrax
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Arboviral disease...
- ☒ Babesiosis
- ☒ Botulism, Foodbo...
- ☒ Botulism, Infant
- ☒ Botulism, Other (...)
- ☒ Brucellosis
- ☒ Campylobacterio...

Disease Name

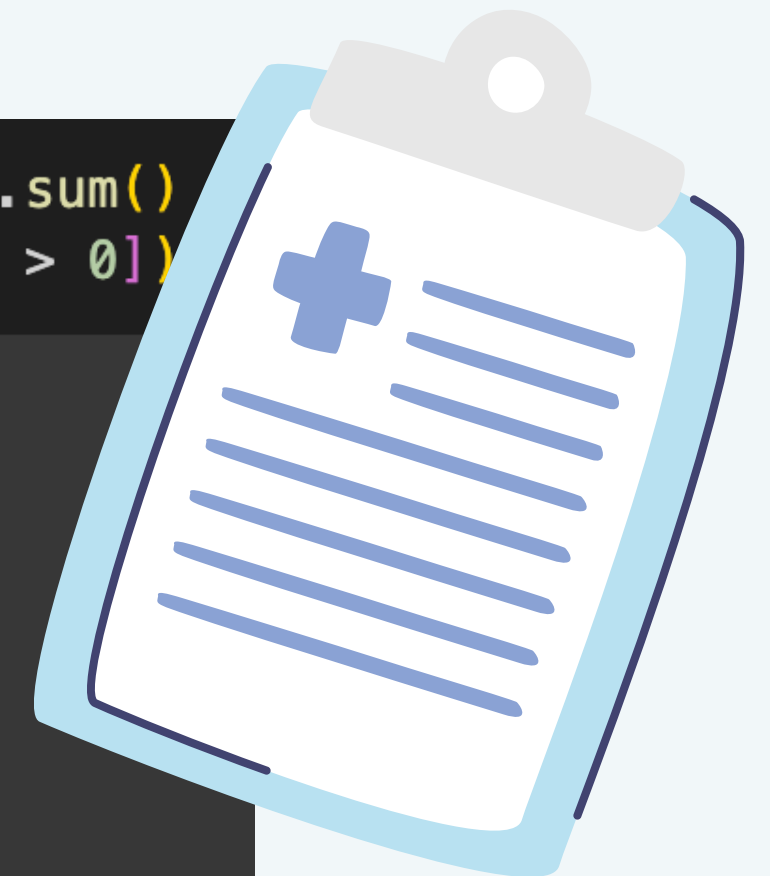
- Null
- Anthrax
- Arboviral disease...
- Arboviral disease...
- Arboviral disease...
- Arboviral disease...
- Arboviral disease...
- Arboviral disease...
- Arboviral disease...
- Babesiosis
- Botulism, Foodbo...
- Botulism, Infant
- Botulism, Other (...)
- Brucellosis
- Campylobacteriosis
- Candida auris, cli...

Detecting Anomalies

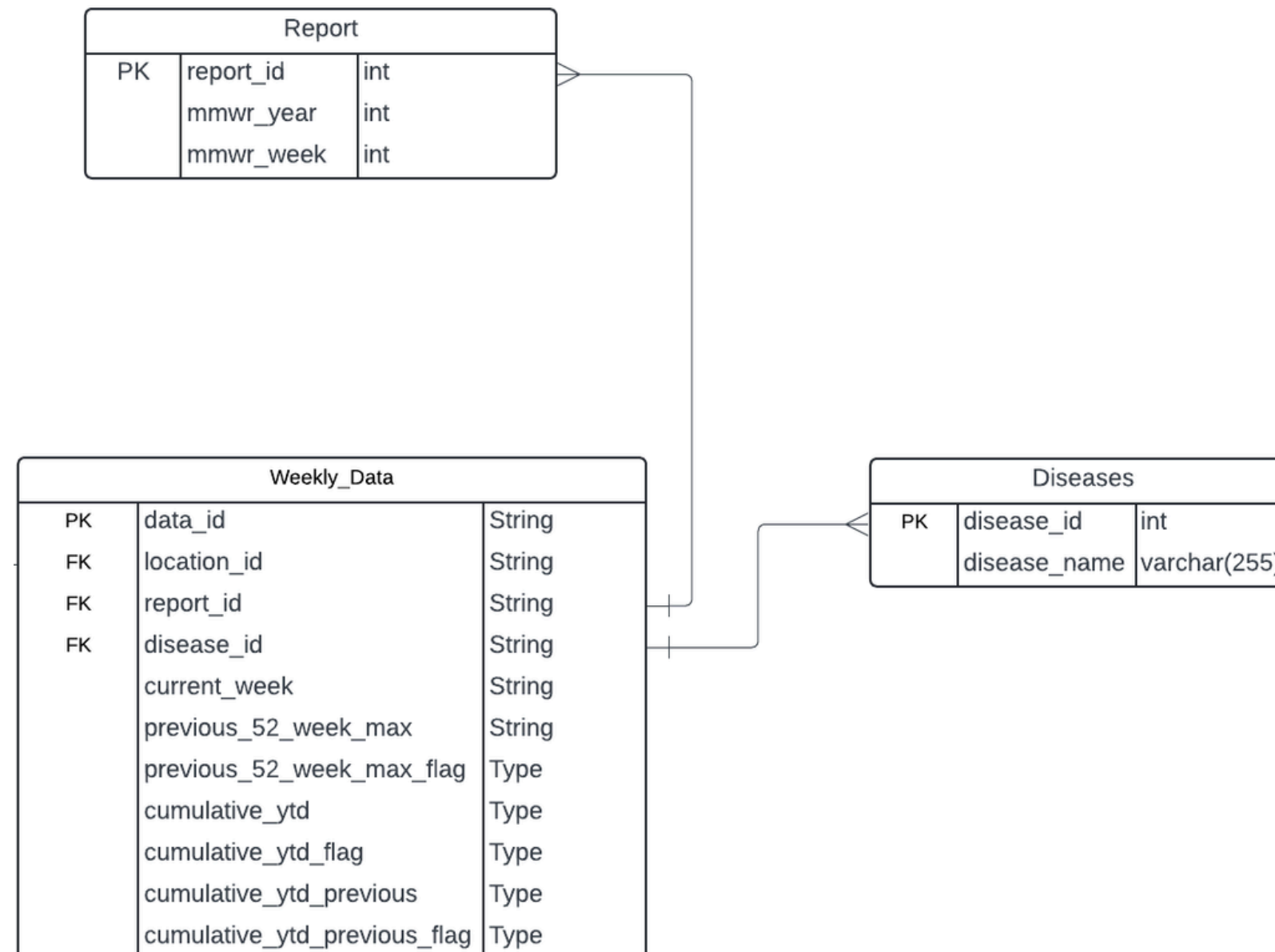
- Data Cleaning
- Trend Analysis
- Preprocessing
- Anomaly Detection

```
missing_values = merged_df.isnull().sum()  
print(missing_values[missing_values > 0])
```

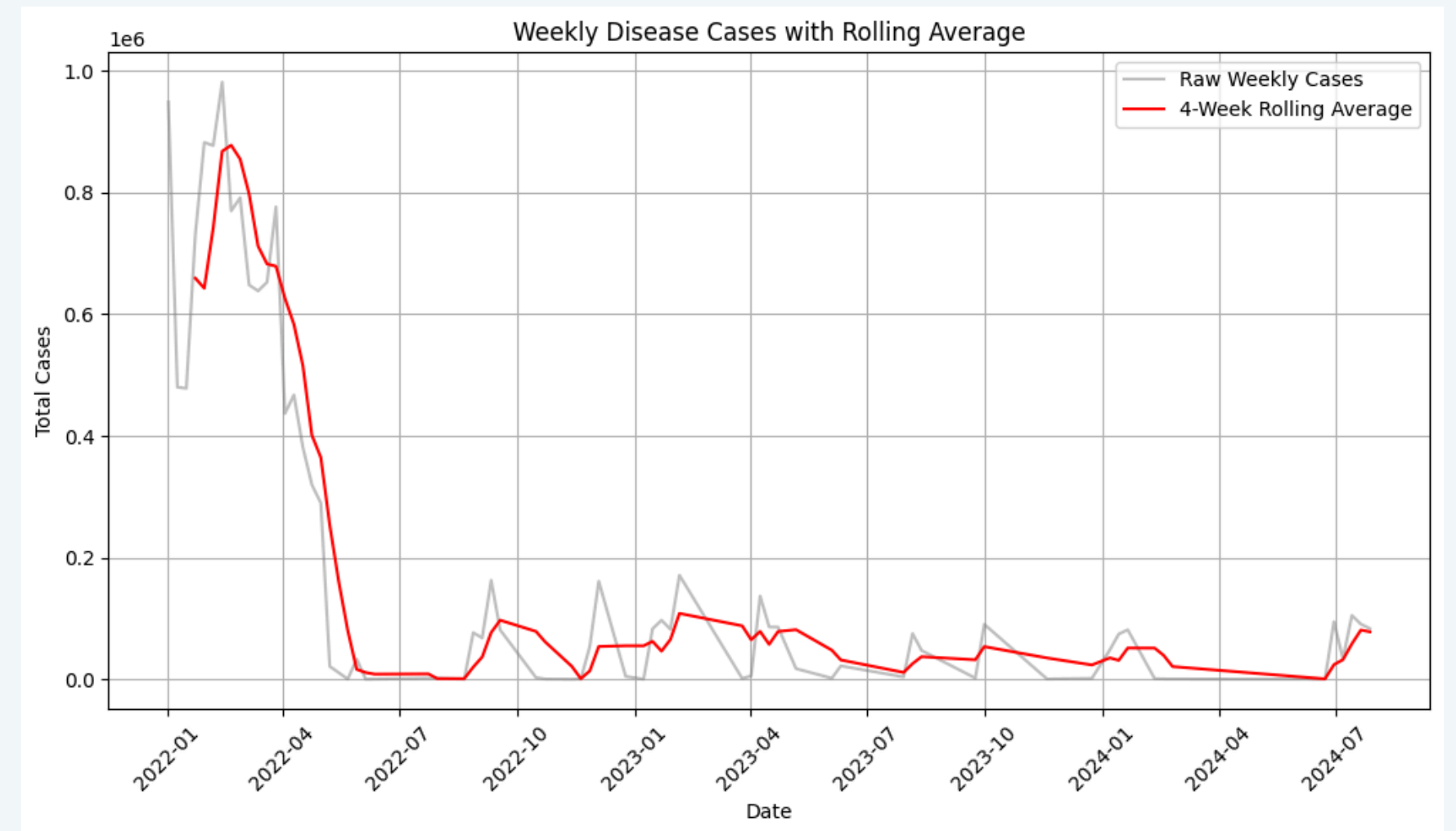
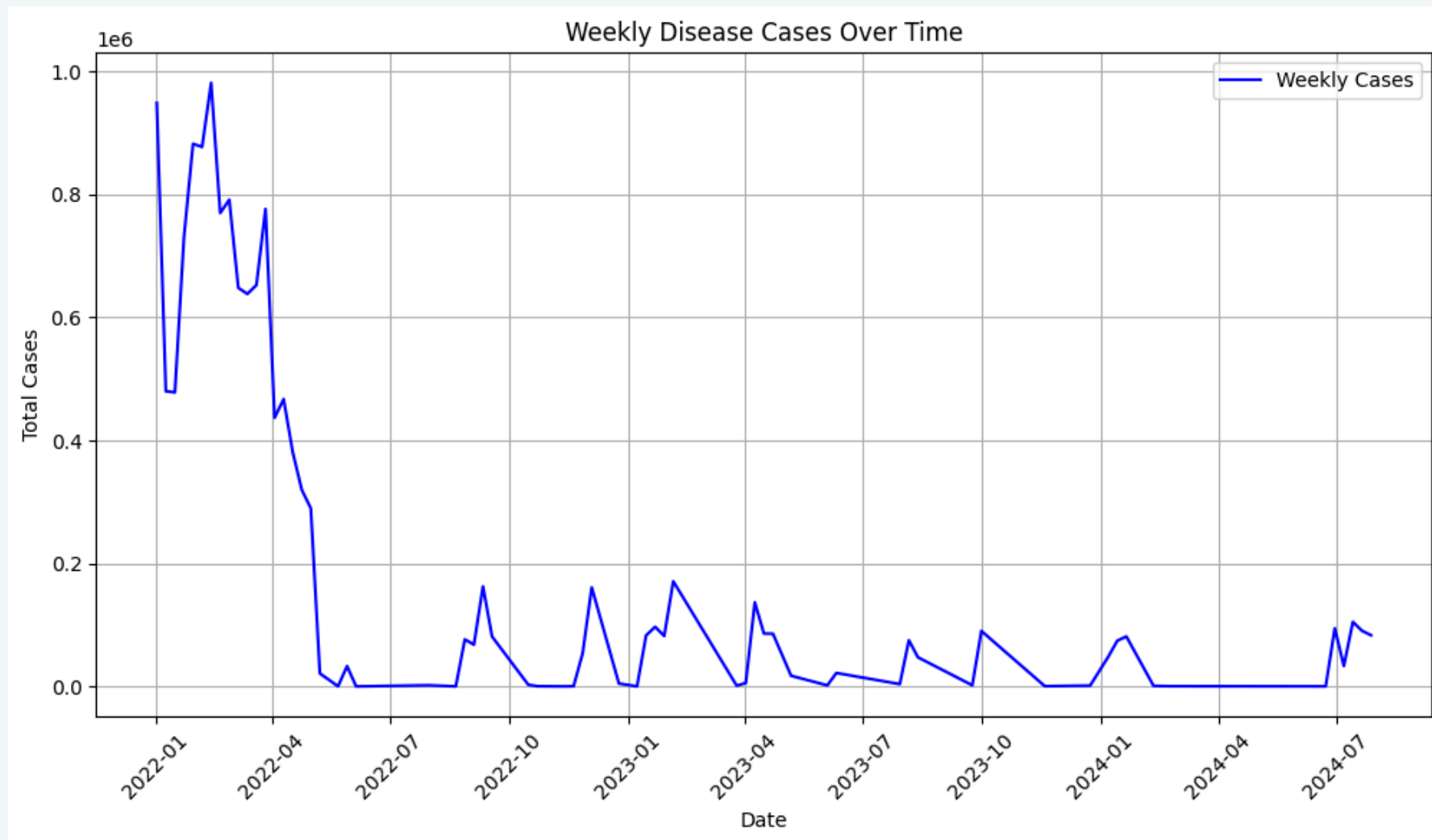
```
disease_name      66036  
location_name     2040783  
states_x          2040783  
location2         2040783  
longitude         2055217  
latitude          2055217  
mmwr_year         291068  
mmwr_week         291068  
states_y          291068  
label             291068  
dtype: int64
```



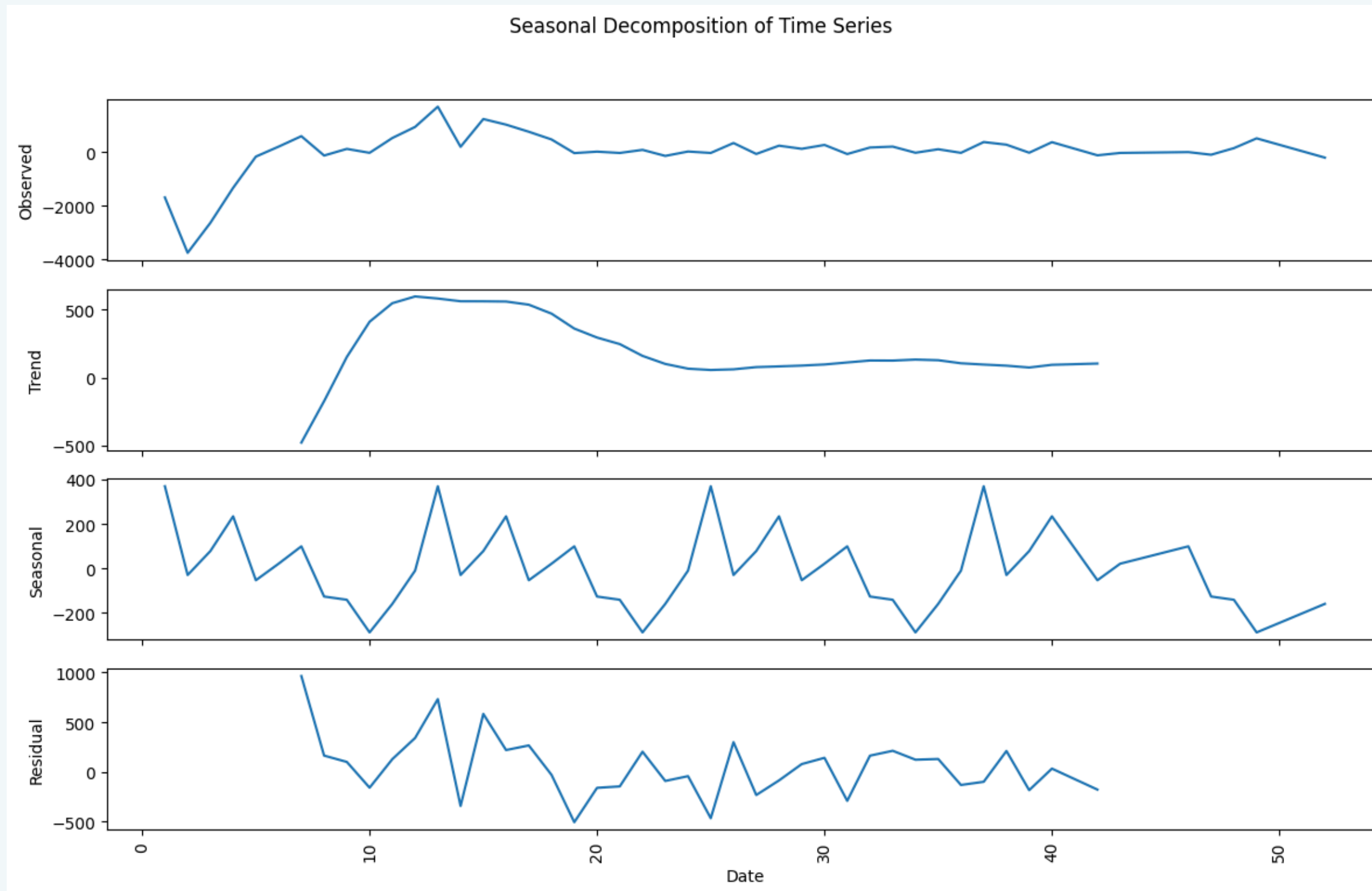
TABLES USED



Trend Analysis



Trend Analysis (Contd.)



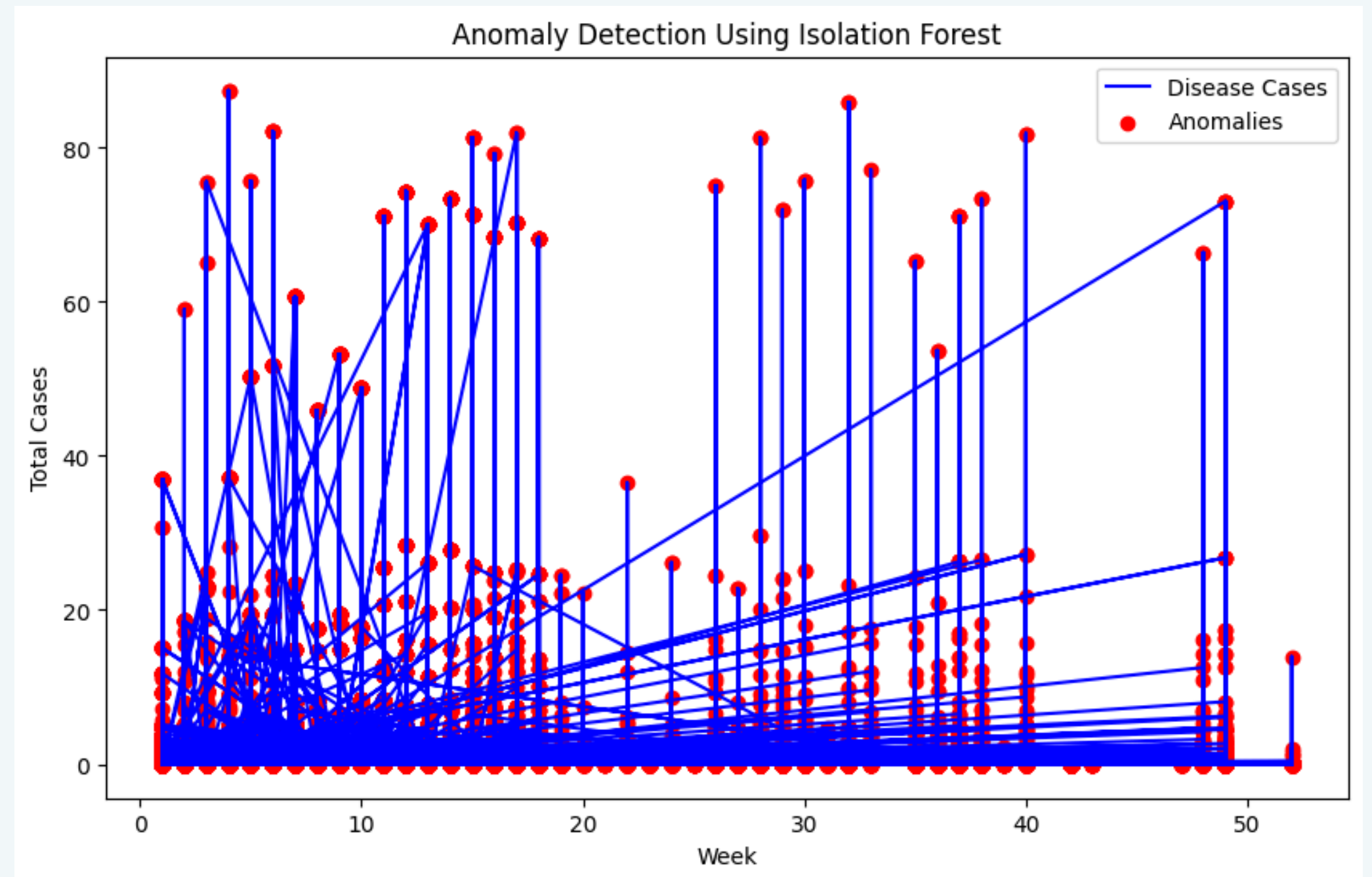
- **Observed:** Shows raw disease case data over time, highlighting overall fluctuations.
- **Trend:** Illustrates the long-term direction of cases, helping to identify sustained increases or decreases.
- **Seasonal:** Captures regular patterns, revealing potential cycles like weekly or yearly spikes in cases.
- **Residual:** Displays remaining variation, potentially indicating rare events, anomalies, or data irregularities.

Anomaly Detection

Methods used:

1. Z-Score
2. Isolation Forest
3. DBSCAN (Did not work)

- The blue line represents the overall weekly trend
- Red dots indicate anomalies—data points flagged as outliers by the model.
- With a contamination setting of 5%, the algorithm expects about 5% of the points to be anomalies, identifying weeks with case counts that significantly differ from the overall pattern.



Model Evaluation

```
Confusion Matrix:  
[[1821555  14124]  
 [         0   4264]]
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	0.99	1.00	1835679
1	0.23	1.00	0.38	4264
accuracy			0.99	1839943
macro avg	0.62	1.00	0.69	1839943
weighted avg	1.00	0.99	0.99	1839943

- The model performs well in identifying normal weeks, with most non-anomalous weeks correctly classified (True Negatives: 1,821,555)
- A small number of normal weeks misclassified as anomalies (False Positives: 14,124).
- Successfully captures all true anomalies (True Positives: 4,264), achieving perfect recall, meaning no actual anomalies were missed.
- The precision for anomalies (0.23) indicates that, while some flagged cases are normal, most are meaningful deviations. This balance shows the model's ability to capture unusual patterns while keeping false positives low, making it effective for monitoring potential disease outbreaks.

CHALLENGES

- Using Apache Superset for Visualisations
- Dealing with Date column during preprocessing
- DBSCAN - cannot work with large dataset
- Error during cloud function setup for ML deployment

NEXT STEPS

- Correctly deploying our ML using cloud functions
- Attempt neural network for more advanced modeling and prediction



THANK YOU

