

```
!pip install pandas numpy matplotlib seaborn statsmodels
```

```

Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (1.26.4)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
Requirement already satisfied: statsmodels in /usr/local/lib/python3.11/dist-packages (0.14.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.1)
Requirement already satisfied: cyclor>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.55.8)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.1.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.1)
Requirement already satisfied: scipy!=1.9.2,>=1.8 in /usr/local/lib/python3.11/dist-packages (from statsmodels) (1.13.1)
Requirement already satisfied: patsy>=0.5.6 in /usr/local/lib/python3.11/dist-packages (from statsmodels) (1.0.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)

```

Step 1: Install Required Libraries

```
# Install essential libraries for data handling, visualization, and time series analysis
```

```
!pip install pandas numpy matplotlib seaborn statsmodels openpyxl missingno scipy pmdarima prophet plotly folium papermill nbconvert streaml
```



19.1/19.1 KB 0.5 MB/s eta 0:00:00
Downloading ansicolors-1.1.8-py2.py3-none-any.whl (13 kB)
Installing collected packages: ansicolors, watchdog, pydeck, pmdarima, streamlit, papermill
Successfully installed ansicolors-1.1.8 papermill-2.6.0 pmdarima-2.0.4 pydeck-0.9.1 streamlit-1.42.0 watchdog-6.0.0

Step 2: Import Required Libraries

```
# 📊 Data Handling
import pandas as pd # For handling datasets
import numpy as np # For numerical operations

# 📈 Visualization
import matplotlib.pyplot as plt # For basic plots
import seaborn as sns # For advanced visualizations
import plotly.express as px # For interactive visualizations
import folium # For map-based visualizations

# 🧹 Data Cleaning & Statistics
import missingno as msno # For missing values visualization
from scipy import stats # Statistical tests

# 📅 Time Series Analysis
import statsmodels.api as sm # For ITS regression
from statsmodels.tsa.stattools import adfuller # Test for stationarity
import pmdarima as pm # Auto ARIMA models
from prophet import Prophet # Advanced forecasting models

# 🔄 Automation & Reproducibility
import papermill as pm # Automating Jupyter notebooks
```

Step 3: Load and Explore Data

Step 1: Mount Google Drive

```
from google.colab import drive
drive.mount('/content/drive')
```

📁 Mounted at /content/drive

Step 2: Verify File Path

```
!ls /content/drive/MyDrive/
```

📁

```

title: NOTES App Project Presentation.gslides
'Untitled Diagram(1).jpg'
'Untitled document (1).gdoc'
'Untitled document (2).gdoc'
'Untitled document (3).gdoc'
'Untitled document (4).gdoc'
'Untitled document (5).gdoc'
'Untitled document (6).gdoc'
'Untitled document (7).gdoc'
'Untitled document (8).gdoc'
'Untitled document.gdoc'
'Untitled presentation.gslides'
'Untitled spreadsheet (1).gsheet'
'Untitled spreadsheet.gsheet'
'Video project Expressway'
'V_spot .pdf'
'What is the most meaningful community you have been a part of.docx'
'work schedule.gdoc'
'write-up for Reports.gdoc'
'YG ID.pdf'
YG.xlsx
Youse
'Yuri Gideon Email Signature.png'
'Yuri Gideon Portfolio (1).pdf'
'Yuri Gideon Portfolio (2).pdf'
'Yuri Gideon Portfolio.pdf'
'Zollo content Plan.gdoc'
Zollo.gsheet

```

Step 3: Load Data from Google Drive

Step 3: List All Files in My Rota virus Folder

```
import os
```

```
folder_path = "/content/drive/MyDrive/ROTA VIRUS RESEARCH/"
```

```

# List all files in the folder
files = os.listdir(folder_path)
print("Files in the folder:", files)

```

Files in the folder: ['MBAGATHI diarrheal dx 2021 (1).csv', 'MBAGATHI diarrheal dx 2020 (1) (1).csv', 'MBAGATHI diarrheal dx 2022 (1).cs

◆ Step 4: Load All Data Files Automatically

```
import pandas as pd
```

```
# Define file paths
```

```
folder_path = "/content/drive/MyDrive/ROTA VIRUS RESEARCH/"
```

```
# Define file names and expected formats
```

```

file_dict = {
    "knh_diarrheal": "DIARRHEAL CASES IN CHILDREN LESS THAN 5 IN KNH 2020 TO 2023 (1) (1).xlsx",
    "nvip_coverage": "RI_ 2020-2023 (DATA FROM NVIP on RV, OPV, Pentavalent vaccine)(1) (2) (1).xls",
    "gertrude_outpatient": "GERTRUDES OUTPATIENT 2020-2023 (1).xlsx",
    "mbagathi_2023": "MBAGATHI dairrheal dx 2023 (1).csv",
    "mbagathi_2022": "MBAGATHI diarrheal dx 2022 (1).csv",
    "mbagathi_2021": "MBAGATHI diarrheal dx 2021 (1).csv",
    "mbagathi_2020": "MBAGATHI diarrheal dx 2020 (1) (1).csv",
}

```

```
# Load data into a dictionary
```

```
dataframes = {}
```

```

for key, filename in file_dict.items():
    file_path = folder_path + filename
    if filename.endswith(".csv"):
        dataframes[key] = pd.read_csv(file_path)
    elif filename.endswith(".xls") or filename.endswith(".xlsx"):
        dataframes[key] = pd.read_excel(file_path)

```

```
# Check loaded data
```

```

for key, df in dataframes.items():
    print(f"✅ {key}: {df.shape} rows and columns")

```

```

↵ ✓ knh_diarrheal: (1019, 9) rows and columns
✓ nvip_coverage: (189, 8) rows and columns
✓ gertrude_outpatient: (2574, 5) rows and columns
✓ mbagathi_2023: (2266, 8) rows and columns
✓ mbagathi_2022: (1808, 8) rows and columns
✓ mbagathi_2021: (1436, 8) rows and columns
✓ mbagathi_2020: (85, 6) rows and columns

```

♦ Step 5: Convert Date Columns to Datetime Format

```

# Convert date columns where applicable
for key, df in dataframes.items():
    for col in df.columns:
        if "date" in col.lower(): # Check if column contains 'date'
            df[col] = pd.to_datetime(df[col], errors='coerce')
            print(f"📅 Converted {col} to datetime format in {key}")

```

```

↵ 📅 Converted Diagnosis Date to datetime format in gertrude_outpatient
📅 Converted date_discharge to datetime format in mbagathi_2023
📅 Converted date_discharge to datetime format in mbagathi_2022
📅 Converted date_discharge to datetime format in mbagathi_2021
📅 Converted date_discharge to datetime format in mbagathi_2020
<ipython-input-12-c1929aaf9de2>:5: UserWarning: Parsing dates in %d/%m/%Y format when dayfirst=False (the default) was specified. Pass `
df[col] = pd.to_datetime(df[col], errors='coerce')

```

Step 6: Display First Few Rows of Each Dataset

```

for key, df in dataframes.items():
    print(f"\n📄 First 5 rows of {key}:")
    display(df.head())

```




First 5 rows of knh_diarrheal:

	SN	Unit_number	DOA	Ward	Disease_code	Age	Age_unit	Result	Unnamed: 8
0	1.0	2058592.0	2020-01-01	3D	A09	1.0	Mths	Alive	NaN
1	2.0	2055951.0	2020-01-03	3C	A09	5.0	Yrs	Alive	NaN
2	3.0	2055963.0	2020-01-04	3D	A09	1.0	Yrs	Alive	NaN
3	4.0	2055985.0	2020-01-06	3A	A09	4.0	Mths	Alive	NaN
4	5.0	2008803.0	2020-01-06	3C	A09	1.0	Yrs	Alive	NaN

First 5 rows of nvip_coverage:

	Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5	Unnamed: 6	Unnamed: 7
0	periodname	organisationunitname	Proportion of under 1 year receiving DPT/Hep+H...	Proportion of under 1 year receiving DPT/Hep+HiB3	Proportion of under 1 year receiving receivin...	Proportion of under 1 year receiving Rota 2	proportion of under 1 year receiving OPV 1	proportion of under 1 year receiving OPV 3
1	2020	Baringo County	85.6	79	84.4	79.9	83.2	77.7
2	2020	Bomet County	88.6	87.1	88	86.9	87.2	86.3
3	2020	Bungoma County	87.7	82.1	86.4	82.5	85.5	80.8
4	2020	Busia County	87.3	86.3	86.6	84.7	85.8	85

First 5 rows of gertrude_outpatient:

	Diagnosis Date	UHID	AgeInYears	ICD Code	ICD Description	
0	2020-01-01 09:03:00	668568	6	A08.0	Rotaviral enteritis	
1	2020-01-01 09:24:00	668570	9	A08.0	Rotaviral enteritis	
2	2020-02-01 01:15:00	480067	4	A08.0	Rotaviral enteritis	
3	2020-04-01 13:51:00	584159	3	A08.0	Rotaviral enteritis	
4	2020-05-01 11:44:00	609597	1	A08.0	Rotaviral enteritis	

First 5 rows of mbagathi_2023:

	id	date_discharge	age_less1mnth	age_days	age_years	age_mths	dx1_diarrhoea	outcome
0	5413789	2023-07-01	0	3.0	-1.0	1.0	NaN	NaN
1	5413790	2023-07-01	0	2.0	-1.0	1.0	NaN	NaN
2	5413791	2023-03-01	0	0.0	5.0	6.0	1.0	NaN
3	5413792	2023-04-01	0	11.0	2.0	1.0	NaN	NaN
4	5413793	2023-03-01	1	21.0	1.0	NaN	NaN	NaN

First 5 rows of mbagathi_2022:

	id	date_discharge	age_less1mnth	age_days	age_years	age_mths	dx1_diarrhoea	outcome
0	5412008	2022-01-03	0	0.0	9.0	1.0	NaN	NaN
1	5412009	2022-01-03	0	0.0	3.0	1.0	1.0	NaN
2	5412010	2022-01-03	0	11.0	-1.0	1.0	NaN	NaN
3	5412012	2022-01-05	0	0.0	1.0	1.0	NaN	NaN
4	5412013	2022-01-04	0	1.0	8.0	1.0	NaN	NaN

First 5 rows of mbagathi_2021:

	id	date_discharge	age_less1mnth	age_days	age_years	age_mths	dx1_diarrhoea	outcome
0	5410630	2021-10-02	0	2.0	4.0	3.0	1.0	NaN
1	5410631	2021-09-02	0	3.0	2.0	1.0	NaN	NaN
2	5410632	NaT	0	2.0	2.0	1.0	NaN	NaN
3	5410633	NaT	1	14.0	1.0	NaN	NaN	NaN
4	5410634	NaT	0	0.0	1.0	1.0	NaN	NaN

First 5 rows of mbagathi_2020:

id date_discharge age_less1mnth age_days age_years age_mths

	id	date_of_birth	age_at_enrollment	age_days	age_years	age_months
0	5410545	2020-09-30	0	0	11	1.0
1	5410546	2020-10-06	0	9	-1	3.0
2	5410547	2020-10-15	0	1	6	1.0
3	5410548	2020-10-12	0	1	5	3.0

◆ Step 7: Read Rotavirus Research Proposal

Since your proposal is a .docx file, we will extract its text using python-docx.

Install and Import python-docx

```
!pip install python-docx
from docx import Document
```

```
Collecting python-docx
  Downloading python_docx-1.1.2-py3-none-any.whl.metadata (2.0 kB)
Requirement already satisfied: lxml>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from python-docx) (5.3.0)
Requirement already satisfied: typing-extensions>=4.9.0 in /usr/local/lib/python3.11/dist-packages (from python-docx) (4.12.2)
Downloading python_docx-1.1.2-py3-none-any.whl (244 kB)
244.3/244.3 kB 4.1 MB/s eta 0:00:00
Installing collected packages: python-docx
Successfully installed python-docx-1.1.2
```

Read the Proposal Document

```
# Define proposal file path
proposal_path = folder_path + "Rotavirus Research Proposal.docx"

# Read the document
doc = Document(proposal_path)
proposal_text = "\n".join([para.text for para in doc.paragraphs])

# Print first 500 characters
print("📄 Proposal Preview:\n", proposal_text[:500])
```

📄 Proposal Preview:

Impact of Rotavirus vaccine stock outs on immunization coverage and diarrheal cases in children less than five years in Kenya.

February 2024

TABLE OF CONTENTS

LIST OF APPENDICES

Appendix A: Data collection tool on facilities involved in the study.

Appendix B: Data collection tool on Rotavirus antigen positive cases from Getrude's Children's Hospital

Appendix C: Data collection tool on infant diarrheal cases from Kenyatta National Hospital Appendix D: Data collection to

DETAILED STEPS AFTER READING THE DATASETS

♦ Step 1: Review and Clean Data (Ensure Consistency Across Datasets) 1.1 Standardize Column Names

```
# Standardizing column names for consistency
dataframes["knh_diarrheal"].rename(columns={'DOA': 'Date'}, inplace=True)
dataframes["gertrude_outpatient"].rename(columns={'Diagnosis Date': 'Date'}, inplace=True)
dataframes["mbagathi_2023"].rename(columns={'date_discharge': 'Date'}, inplace=True)
dataframes["mbagathi_2022"].rename(columns={'date_discharge': 'Date'}, inplace=True)
dataframes["mbagathi_2021"].rename(columns={'date_discharge': 'Date'}, inplace=True)
dataframes["mbagathi_2020"].rename(columns={'date_discharge': 'Date'}, inplace=True)
dataframes["nvip_coverage"].rename(columns={'periodname': 'Year', 'organisationunitname': 'County'}, inplace=True)

# Verify column name changes
for key, df in dataframes.items():
    print(f"{key} columns: {df.columns.tolist()}")
```

knh_diarrheal columns: ['SN', 'Unit_number', 'Date', 'Ward', 'Disease_code', 'Age', 'Age_unit', 'Result', 'Unnamed: 8']
nvip_coverage columns: ['Unnamed: 0', 'Unnamed: 1', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4', 'Unnamed: 5', 'Unnamed: 6', 'Unnamed: 7']
gertrude_outpatient columns: ['Date', 'UHID', 'AgeInYears', 'ICD Code', 'ICD Description']
mbagathi_2023 columns: ['id', 'Date', 'age_less1mnth', 'age_days', 'age_years', 'age_mths', 'dx1_diarrhoea', 'outcome']
mbagathi_2022 columns: ['id', 'Date', 'age_less1mnth', 'age_days', 'age_years', 'age_mths', 'dx1_diarrhoea', 'outcome']
mbagathi_2021 columns: ['id', 'Date', 'age_less1mnth', 'age_days', 'age_years', 'age_mths', 'dx1_diarrhoea', 'outcome']
mbagathi_2020 columns: ['id', 'Date', 'age_less1mnth', 'age_days', 'age_years', 'age_mths']

1.2 Handle Missing Values

```
# Check for missing values in each dataset
for key, df in dataframes.items():
    print(f"{key} Missing values in {key}: {df.isnull().sum()}")
```

Unnamed: 1 0
 Unnamed: 2 0
 Unnamed: 3 0
 Unnamed: 4 0
 Unnamed: 5 0
 Unnamed: 6 0
 Unnamed: 7 0
 dtype: int64

Missing values in gertrude_outpatient:
 Date 0
 UHID 0
 AgeInYears 0
 ICD Code 0
 ICD Description 0
 dtype: int64

Missing values in mbagathi_2023:
 id 0
 Date 1425
 age_less1mnth 0
 age_days 13
 age_years 15
 age_mths 182
 dx1_diarrhoea 1980
 outcome 2266
 dtype: int64

Missing values in mbagathi_2022:

```

age_years      19
age_mths       152
dx1_diarrhoea  1290
outcome        1436
dtype: int64

```

Missing values in mbagathi_2020:

```

id            0
Date          0
age_less1mth  0
age_days      0
age_years     0
age_mths      13
dtype: int64

```

Double-click (or enter) to edit

1.2.1: Handle the missing data

```

# Drop rows with missing Date in important datasets
for key in ["knh_diarrheal", "gertrude_outpatient", "mbagathi_2023", "mbagathi_2022", "mbagathi_2021", "mbagathi_2020"]:
    dataframes[key].dropna(subset=['Date'], inplace=True)

# Fill missing age-related values with median where applicable
for key in ["mbagathi_2023", "mbagathi_2022", "mbagathi_2021", "mbagathi_2020"]:
    dataframes[key]['age_years'] = dataframes[key]['age_years'].fillna(dataframes[key]['age_years'].median())
    dataframes[key]['age_days'] = dataframes[key]['age_days'].fillna(dataframes[key]['age_days'].median()) # Optional: Fill age_days
    dataframes[key]['age_mths'] = dataframes[key]['age_mths'].fillna(dataframes[key]['age_mths'].median()) # Optional: Fill age_mths

# Verify the missing values again after filling/dropping
for key, df in dataframes.items():
    print(f"\n✅ {key} missing values after cleaning:\n{df.isnull().sum()}")

```

```

Unnamed: 1    0
Unnamed: 2    0
Unnamed: 3    0
Unnamed: 4    0
Unnamed: 5    0
Unnamed: 6    0
Unnamed: 7    0
dtype: int64

```

✅ gertrude_outpatient missing values after cleaning:

```

Date          0
UHID          0
AgeInYears    0
ICD Code      0
ICD Description 0
dtype: int64

```

✅ mbagathi_2023 missing values after cleaning:

```

id            0
Date          0
age_less1mth  0
age_days      0
age_years     0
age_mths      0
dx1_diarrhoea 731
outcome       841
dtype: int64

```

✅ mbagathi_2022 missing values after cleaning:

```

id            0
Date          0
age_less1mth  0

```



```
outcome      >>>
dtype: int64
```

```
✓ mbagathi_2020 missing values after cleaning:
id          0
Date        0
age_less1mth 0
age_days    0
age_years   0
age_mths    0
dtype: int64
```

1.3 Remove Duplicates

```
# Check for duplicates
for key, df in dataframes.items():
    duplicates = df.duplicated().sum()
    print(f"🔍 {key} has {duplicates} duplicate rows")

# Remove duplicates
for key in dataframes.keys():
    dataframes[key].drop_duplicates(inplace=True)

# Verify that duplicates have been removed
for key, df in dataframes.items():
    duplicates_after = df.duplicated().sum()
    print(f"✓ {key} has {duplicates_after} duplicate rows after removal")
```

```
🔍 knh_diarrheal has 0 duplicate rows
🔍 nvip_coverage has 0 duplicate rows
🔍 gertrude_outpatient has 9 duplicate rows
🔍 mbagathi_2023 has 0 duplicate rows
🔍 mbagathi_2022 has 0 duplicate rows
🔍 mbagathi_2021 has 0 duplicate rows
🔍 mbagathi_2020 has 0 duplicate rows
✓ knh_diarrheal has 0 duplicate rows after removal
✓ nvip_coverage has 0 duplicate rows after removal
✓ gertrude_outpatient has 0 duplicate rows after removal
✓ mbagathi_2023 has 0 duplicate rows after removal
✓ mbagathi_2022 has 0 duplicate rows after removal
✓ mbagathi_2021 has 0 duplicate rows after removal
✓ mbagathi_2020 has 0 duplicate rows after removal
```

❖ Step 2: Convert Date Columns & Aggregate Data by Month

2.1 : Check the data types and inspect Date columns before manipulation

```
# Check the data types and inspect Date columns before manipulation
for key, df in dataframes.items():
    # Check if 'Date' column exists before accessing it
    if 'Date' in df.columns:
        print(f"🔍 {key} - 'Date' column type: {df['Date'].dtype}")
        print(f"🔍 {key} - 'Date' column head:\n{df['Date'].head()}")
    else:
        print(f"🔍 {key} - 'Date' column not found") # Or any other message you prefer
```

```
🔍 knh_diarrheal - 'Date' column type: datetime64[ns]
🔍 knh_diarrheal - 'Date' column head:
0    2020-01-01
1    2020-01-03
2    2020-01-04
3    2020-01-06
4    2020-01-06
Name: Date, dtype: datetime64[ns]
🔍 nvip_coverage - 'Date' column not found
🔍 gertrude_outpatient - 'Date' column type: datetime64[ns]
🔍 gertrude_outpatient - 'Date' column head:
0    2020-01-01 09:03:00
1    2020-01-01 09:24:00
2    2020-02-01 01:15:00
3    2020-04-01 13:51:00
4    2020-05-01 11:44:00
Name: Date, dtype: datetime64[ns]
🔍 mbagathi_2023 - 'Date' column type: datetime64[ns]
🔍 mbagathi_2023 - 'Date' column head:
```

```

0    2023-07-01
1    2023-07-01
2    2023-03-01
3    2023-04-01
4    2023-03-01
Name: Date, dtype: datetime64[ns]
◆ mbagathi_2022 - 'Date' column type: datetime64[ns]
◆ mbagathi_2022 - 'Date' column head:
0    2022-01-03
1    2022-01-03
2    2022-01-03
3    2022-01-05
4    2022-01-04
Name: Date, dtype: datetime64[ns]
◆ mbagathi_2021 - 'Date' column type: datetime64[ns]
◆ mbagathi_2021 - 'Date' column head:
0    2021-10-02
1    2021-09-02
38   2021-01-03
39   2021-01-03
40   2021-01-03
Name: Date, dtype: datetime64[ns]
◆ mbagathi_2020 - 'Date' column type: datetime64[ns]
◆ mbagathi_2020 - 'Date' column head:
0    2020-09-30
1    2020-10-06
2    2020-10-15
3    2020-10-12
4    2020-10-12
Name: Date, dtype: datetime64[ns]

```

2.1.2: Check the data types to ensure conversion was successful

```

# Check the data types to ensure conversion was successful
for key, df in dataframes.items():
    print(f"◆ {key} data types:\n", df.dtypes)

```

```

Result      object
Unnamed: 8   object
dtype: object
◆ nvip_coverage data types:
  Unnamed: 0    object
  Unnamed: 1    object
  Unnamed: 2    object
  Unnamed: 3    object
  Unnamed: 4    object
  Unnamed: 5    object
  Unnamed: 6    object
  Unnamed: 7    object
dtype: object
◆ gertrude_outpatient data types:
  Date          datetime64[ns]
  UHID           int64
  AgeInYears     int64
  ICD Code       object
  ICD Description object
dtype: object
◆ mbagathi_2023 data types:
  id           int64
  Date         datetime64[ns]
  age_less1mth int64
  age_days     float64
  age_years    float64
  age_mths     float64
  dx1_diarrhoea float64
  outcome      float64
dtype: object
◆ mbagathi_2022 data types:
  id           int64
  Date         datetime64[ns]
  age_less1mth int64
  age_days     float64
  age_years    float64

```

```

age_uays          float64
age_years         float64
age_mths          float64
dx1_diarrhoea     float64
outcome           float64
dtype: object
  ◆ mbagathi_2020 data types:
    id             int64
    Date            datetime64[ns]
    age_less1mnth   int64
    age_days        int64
    age_years       int64
    age_mths        float64
    dtype: object

```

2.2.2: Convert Date Columns Where Necessary Since the columns are not all correctly converted

```

# Convert 'Date' columns in necessary dataframes
for key, df in dataframes.items():
    # For each dataframe, convert 'Date' columns to datetime if they exist
    for col in df.columns:
        if "date" in col.lower(): # Check if column contains 'date'
            df[col] = pd.to_datetime(df[col], errors='coerce', dayfirst=True) # Ensure datetime format with dayfirst=True

# Check the data types again after conversion
for key, df in dataframes.items():
    print(f"  ◆ {key} data types:\n", df.dtypes)

```

```

Result          object
Unnamed: 8       object
dtype: object
  ◆ nvip_coverage data types:
    Unnamed: 0      object
    Unnamed: 1      object
    Unnamed: 2      object
    Unnamed: 3      object
    Unnamed: 4      object
    Unnamed: 5      object
    Unnamed: 6      object
    Unnamed: 7      object
    dtype: object
  ◆ gertrude_outpatient data types:
    Date            datetime64[ns]
    UHID            int64
    AgeInYears       int64
    ICD Code         object
    ICD Description  object
    dtype: object
  ◆ mbagathi_2023 data types:
    id             int64
    Date            datetime64[ns]
    age_less1mnth   int64
    age_days        float64
    age_years       float64
    age_mths        float64
    dx1_diarrhoea   float64
    outcome         float64
    dtype: object
  ◆ mbagathi_2022 data types:
    id             int64
    Date            datetime64[ns]
    age_less1mnth   int64
    age_days        float64
    age_years       float64
    age_mths        float64
    dx1_diarrhoea   float64
    outcome         float64
    dtype: object
  ◆ mbagathi_2021 data types:
    id             int64
    Date            datetime64[ns]
    age_less1mnth   int64
    age_days        float64

```

```

age_less18mths      int64
age_days             int64
age_years            int64
age_mths             float64
dtype: object

```

Double-click (or enter) to edit

2.3: Create Year-Month Column for Aggregation

```

# Add YearMonth column
for key, df in dataframes.items():
    if 'Date' in df.columns: # Ensure we're working with a date column
        df['YearMonth'] = df['Date'].dt.to_period('M')
        print(f"Added YearMonth column to {key}")

```

```

Added YearMonth column to knh_diarrheal
Added YearMonth column to gertrude_outpatient
Added YearMonth column to mbagathi_2023
Added YearMonth column to mbagathi_2022
Added YearMonth column to mbagathi_2021
Added YearMonth column to mbagathi_2020

```

2.4: Aggregate Data by Year-Month

```

# Aggregate the number of cases by YearMonth
for key, df in dataframes.items():
    if 'YearMonth' in df.columns:
        monthly_data = df.groupby('YearMonth').size().reset_index(name='Cases')
        print(f"Aggregated monthly data for {key}")

    # Store the aggregated data (for later use or merging)
    globals()[f"{key}_monthly"] = monthly_data

```

```

Aggregated monthly data for knh_diarrheal
Aggregated monthly data for gertrude_outpatient
Aggregated monthly data for mbagathi_2023
Aggregated monthly data for mbagathi_2022
Aggregated monthly data for mbagathi_2021
Aggregated monthly data for mbagathi_2020

```

2.5 Combine All Monthly Data into One DataFrame

2.5.1 Check that DataFrames Exist Before Combining:

```

# List variables to check if the monthly dataframes exist
print("Variables in namespace:", globals().keys())

```

```

Variables in namespace: dict_keys(['__name__', '__doc__', '__package__', '__loader__', '__spec__', '__builtin__', '__builtins__', '_ih',

```

2.5.2. Ensure All DataFrames Are Defined:

```

# Confirm aggregation for each dataset
for key, df in dataframes.items():
    if 'YearMonth' in df.columns:
        monthly_data = df.groupby('YearMonth').size().reset_index(name='Cases')
        globals()[f"{key}_monthly"] = monthly_data
        print(f"Aggregated monthly data for {key}")

```

```

Aggregated monthly data for knh_diarrheal
Aggregated monthly data for gertrude_outpatient
Aggregated monthly data for mbagathi_2023
Aggregated monthly data for mbagathi_2022
Aggregated monthly data for mbagathi_2021
Aggregated monthly data for mbagathi_2020

```

2.5.3. Confirm Data Types and Structure:

```
# Check the structure of a few monthly dataframes by accessing them via globals
print("🔍 knh_diarrheal_monthly structure:\n", globals().get('knh_diarrheal_monthly').head())
print("🔍 gertrude_outpatient_monthly structure:\n", globals().get('gertrude_outpatient_monthly').head())
```

```
🔍 knh_diarrheal_monthly structure:
  YearMonth  Cases
0  2020-01     25
1  2020-02     24
2  2020-03     25
3  2020-04      7
4  2020-05      9
🔍 gertrude_outpatient_monthly structure:
  YearMonth  Cases
0  2020-01     43
1  2020-02     21
2  2020-03     23
3  2020-04      8
4  2020-05      6
```

Double-click (or enter) to edit

2.6 Add a Column for the Hospital Source:Combine All Monthly Data into One DataFrame

```
# Add 'Hospital' column to each dataframe
knh_diarrheal_monthly['Hospital'] = 'knh_diarrheal'
gertrude_outpatient_monthly['Hospital'] = 'gertrude_outpatient'
mbagathi_2023_monthly['Hospital'] = 'mbagathi_2023'
mbagathi_2022_monthly['Hospital'] = 'mbagathi_2022'
mbagathi_2021_monthly['Hospital'] = 'mbagathi_2021'
mbagathi_2020_monthly['Hospital'] = 'mbagathi_2020'

# Combine all hospital datasets
combined_monthly = pd.concat([
    knh_diarrheal_monthly,
    gertrude_outpatient_monthly,
    mbagathi_2023_monthly,
    mbagathi_2022_monthly,
    mbagathi_2021_monthly,
    mbagathi_2020_monthly
])

# Aggregate across all hospitals by YearMonth
diarrheal_monthly = combined_monthly.groupby(['YearMonth', 'Hospital']).sum().reset_index()

# Check combined and aggregated data
print("🇳🇬 Combined and aggregated monthly data:\n", diarrheal_monthly.head())
```

```
🇳🇬 Combined and aggregated monthly data:
  YearMonth  Hospital  Cases
0  2020-01  gertrude_outpatient    43
1  2020-01    knh_diarrheal     25
2  2020-02  gertrude_outpatient    21
3  2020-02    knh_diarrheal     24
4  2020-03  gertrude_outpatient    23
```

2.7 Exporting this combined dataframe to a CSV file for further analysis or sharing

```
# Export the combined data to a CSV file
diarrheal_monthly.to_csv('combined_monthly_data.csv', index=False)

print("🇳🇬 Combined data exported successfully!")
```

```
🇳🇬 Combined data exported successfully!
```

Double-click (or enter) to edit

2.8.1 Ensure Data Integrity

```
# Check for missing values in the final combined dataset
print("🔍 Missing values in combined monthly data:\n", combined_monthly.isnull().sum())
```

```
🔍 Missing values in combined monthly data:
YearMonth    0
Cases        0
Hospital      0
dtype: int64
```

2.8.2. Duplicates: Check for duplicate rows, particularly in the YearMonth and Hospital columns.

```
# Check for duplicates
duplicates = combined_monthly.duplicated(subset=['YearMonth', 'Hospital'])
print("🔍 Duplicate rows:\n", combined_monthly[duplicates])
```

```
🔍 Duplicate rows:
Empty DataFrame
Columns: [YearMonth, Cases, Hospital]
Index: []
```

2.8.3 : Data Types: Ensure the data types are correct after concatenation.

```
# Check the data types of the final combined dataset
print("🔍 Data types in combined dataset:\n", combined_monthly.dtypes)
```

```
🔍 Data types in combined dataset:
YearMonth    period[M]
Cases        int64
Hospital      object
dtype: object
```

2.8.4 Perform Aggregation:with the final aggregation across hospitals by YearMonth

```
# Aggregate across all hospitals by YearMonth
diarrheal_monthly = combined_monthly.groupby(['YearMonth', 'Hospital']).sum().reset_index()
```

```
# Check aggregated data
print("📊 Aggregated monthly data:\n", diarrheal_monthly.head())
```

```
📊 Aggregated monthly data:
YearMonth    Hospital  Cases
0  2020-01  gertrude_outpatient    43
1  2020-01    knh_diarrheal       25
2  2020-02  gertrude_outpatient    21
3  2020-02    knh_diarrheal       24
4  2020-03  gertrude_outpatient    23
```

✓ ♦ Step 3: Aggregate Immunization Coverage by Month

3.1: Confirm the 'nvip_coverage' Dataset

```
# Check loaded data
for key, df in dataframes.items():
    print(f"✅ {key}: {df.shape} rows and columns")
```

```
✅ knh_diarrheal: (1018, 10) rows and columns
✅ nvip_coverage: (189, 8) rows and columns
✅ gertrude_outpatient: (2565, 6) rows and columns
✅ mbagathi_2023: (841, 9) rows and columns
✅ mbagathi_2022: (1808, 9) rows and columns
✅ mbagathi_2021: (559, 9) rows and columns
✅ mbagathi_2020: (85, 7) rows and columns
```

3.1.1 Confirm the nvip_coverage

```
# Check the first few rows of the nvip_coverage dataset
print(dataframes['nvip_coverage'].head())
```

```
↳ Unnamed: 0      Unnamed: 1 \
0  periodname  organisationunitname
1      2020      Baringo County
2      2020      Bomet County
3      2020      Bungoma County
4      2020      Busia County

                                Unnamed: 2 \
0  Proportion of under 1 year receiving DPT/Hep+H...
1                                85.6
2                                88.6
3                                87.7
4                                87.3

                                Unnamed: 3 \
0  Proportion of under 1 year receiving DPT/Hep+HiB3
1                                79
2                                87.1
3                                82.1
4                                86.3

                                Unnamed: 4 \
0  Proportion of under 1 year receiving receivin...
1                                84.4
2                                88
3                                86.4
4                                86.6

                                Unnamed: 5 \
0  Proportion of under 1 year receiving Rota 2
1                                79.9
2                                86.9
3                                82.5
4                                84.7

                                Unnamed: 6 \
0  proportion of under 1 year receiving OPV 1
1                                83.2
2                                87.2
3                                85.5
4                                85.8

                                Unnamed: 7
0  proportion of under 1 year receiving OPV 3
1                                77.7
2                                86.3
3                                80.8
4                                85
```

3.2 Step Converting 'Year' to DateTime Format

```
# Preprocess the nvip_coverage dataset to clean up columns
nvip_coverage = dataframes['nvip_coverage']

# Drop the first row (header) if necessary and reset column names
nvip_coverage.columns = nvip_coverage.iloc[0] # Set the first row as column names
nvip_coverage = nvip_coverage.drop(0).reset_index(drop=True) # Remove the first row after renaming columns

# Inspect the first few rows to check the column names and data
print(nvip_coverage.head())

# Convert 'periodname' column to datetime (assuming it's the column for Year/Month)
nvip_coverage['YearMonth'] = pd.to_datetime(nvip_coverage['periodname'], format='%Y')

# If there's a 'Hospital' column, make sure it is renamed properly
nvip_coverage['Hospital'] = nvip_coverage['organisationunitname'] # Assuming this is the hospital name column

# Now, proceed with the aggregation by YearMonth and Hospital
nvip_coverage_monthly = nvip_coverage.groupby(['YearMonth', 'Hospital']).sum().reset_index()

# Check the aggregated data
print("🇰🇪 Aggregated monthly immunization coverage:\n", nvip_coverage_monthly.head())
```



```

0 proportion of under 1 year receiving OPV 3
0 77.7
1 86.3
2 80.8
3 85
4 73.8
🇰🇪 Aggregated monthly immunization coverage:
0 YearMonth Hospital periodname organisationunitname \
0 2020-01-01 Baringo County 2020 Baringo County
1 2020-01-01 Bomet County 2020 Bomet County
2 2020-01-01 Bungoma County 2020 Bungoma County
3 2020-01-01 Busia County 2020 Busia County
4 2020-01-01 Elgeyo Marakwet County 2020 Elgeyo Marakwet County

0 Proportion of under 1 year receiving DPT/Hep+HiB1 \
0 85.6
1 88.6
2 87.7
3 87.3
4 82.8

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0 79
1 87.1
2 82.1
3 86.3
4 79.7

0 Proportion of under 1 year receiving receiving Rota 1 \
0 84.4
1 88
2 86.4
3 86.6
4 82.3

0 Proportion of under 1 year receiving Rota 2 \
0 79.9
1 86.9
2 82.5
3 84.7
4 77.3

0 proportion of under 1 year receiving OPV 1 \
0 83.2
1 87.2
2 85.5
3 85.8
4 78.7

0 proportion of under 1 year receiving OPV 3
0 77.7
1 86.3
2 80.8
3 85
4 73.8

```

3.2.1 Ensure the aggregation works on the numeric columns only: aggregate the columns related to immunization coverage, and exclude the metadata columns (periodname, organisationunitname) from the groupby operation. and Fix column selection:

```

# Drop non-numeric columns from the aggregation
nvip_coverage_numeric = nvip_coverage.drop(columns=['periodname', 'organisationunitname'])

# Ensure that 'YearMonth' and 'Hospital' are kept for grouping, and then perform the aggregation on the numeric columns
nvip_coverage_monthly = nvip_coverage.groupby(['YearMonth', 'Hospital']).sum().reset_index()

# Check the aggregated data
print("🇰🇪 Aggregated monthly immunization coverage:\n", nvip_coverage_monthly.head())

```

```

🇰🇪 Aggregated monthly immunization coverage:
0 YearMonth Hospital periodname organisationunitname \
0 2020-01-01 Baringo County 2020 Baringo County
1 2020-01-01 Bomet County 2020 Bomet County
2 2020-01-01 Bungoma County 2020 Bungoma County
3 2020-01-01 Busia County 2020 Busia County
4 2020-01-01 Elgeyo Marakwet County 2020 Elgeyo Marakwet County

0 Proportion of under 1 year receiving DPT/Hep+HiB1 \
0 85.6
1 88.6
2 87.7

```



```

3                87.3
4                82.8

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0                79
1                87.1
2                82.1
3                86.3
4                79.7

0 Proportion of under 1 year receiving receiving Rota 1 \
0                84.4
1                88
2                86.4
3                86.6
4                82.3

0 Proportion of under 1 year receiving Rota 2 \
0                79.9
1                86.9
2                82.5
3                84.7
4                77.3

0 proportion of under 1 year receiving OPV 1 \
0                83.2
1                87.2
2                85.5
3                85.8
4                78.7

0 proportion of under 1 year receiving OPV 3
0                77.7
1                86.3
2                80.8
3                85
4                73.8

```

```
print(nvip_coverage_monthly.head())
```

```

0 YearMonth      Hospital periodname  organisationunitname \
0 2020-01-01      Baringo County      2020      Baringo County
1 2020-01-01      Bomet County      2020      Bomet County
2 2020-01-01      Bungoma County      2020      Bungoma County
3 2020-01-01      Busia County      2020      Busia County
4 2020-01-01      Elgeyo Marakwet County      2020      Elgeyo Marakwet County

```

```

0 Proportion of under 1 year receiving DPT/Hep+HiB1 \
0                85.6
1                88.6
2                87.7
3                87.3
4                82.8

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0                79
1                87.1
2                82.1
3                86.3
4                79.7

0 Proportion of under 1 year receiving receiving Rota 1 \
0                84.4
1                88
2                86.4
3                86.6
4                82.3

0 Proportion of under 1 year receiving Rota 2 \
0                79.9
1                86.9
2                82.5
3                84.7
4                77.3

0 proportion of under 1 year receiving OPV 1 \
0                83.2
1                87.2
2                85.5
3                85.8
4                78.7

```

```

0 proportion of under 1 year receiving OPV 3
0 77.7
1 86.3
2 80.8
3 85
4 73.8

```

```
print(nvip_coverage_monthly.dtypes)
```

```

0
YearMonth          datetime64[ns]
Hospital           object
periodname         object
organisationunitname object
Proportion of under 1 year receiving DPT/Hep+HiB1 object
Proportion of under 1 year receiving DPT/Hep+HiB3 object
Proportion of under 1 year receiving receiving Rota 1 object
Proportion of under 1 year receiving Rota 2 object
proportion of under 1 year receiving OPV 1 object
proportion of under 1 year receiving OPV 3 object
dtype: object

```

```
print(nvip_coverage_monthly.isnull().sum())
```

```

0
YearMonth          0
Hospital           0
periodname         0
organisationunitname 0
Proportion of under 1 year receiving DPT/Hep+HiB1 0
Proportion of under 1 year receiving DPT/Hep+HiB3 0
Proportion of under 1 year receiving receiving Rota 1 0
Proportion of under 1 year receiving Rota 2 0
proportion of under 1 year receiving OPV 1 0
proportion of under 1 year receiving OPV 3 0
dtype: int64

```

checking shape

```
print(nvip_coverage_monthly.shape)
```

```
(188, 10)
```

```
print(nvip_coverage_monthly.columns)
```

```

Index(['YearMonth', 'Hospital', 'periodname', 'organisationunitname',
      'Proportion of under 1 year receiving DPT/Hep+HiB1 ',
      'Proportion of under 1 year receiving DPT/Hep+HiB3',
      'Proportion of under 1 year receiving receiving Rota 1',
      'Proportion of under 1 year receiving Rota 2',
      'proportion of under 1 year receiving OPV 1 ',
      'proportion of under 1 year receiving OPV 3 '],
      dtype='object', name=0)

```

1 Drop Unnecessary Columns Again

```
nvip_coverage = nvip_coverage.drop(columns=['periodname', 'organisationunitname'], errors='ignore')
```

2 Convert Numeric Columns to float

```

# Identify numeric columns by excluding 'YearMonth' and 'Hospital'
numeric_cols = nvip_coverage.columns.difference(['YearMonth', 'Hospital'])

```

```

# Convert all numeric columns to float (handles possible string issues)
nvip_coverage[numeric_cols] = nvip_coverage[numeric_cols].apply(pd.to_numeric, errors='coerce')

```

 Verify Fixes

```
print(nvip_coverage_monthly.dtypes) # Ensure numeric columns are now float/int
print(nvip_coverage_monthly.head()) # Verify structure and values
```

```
Hospital                                object
periodname                             object
organisationunitname                   object
Proportion of under 1 year receiving DPT/Hep+HiB1    object
Proportion of under 1 year receiving DPT/Hep+HiB3    object
Proportion of under 1 year receiving receiving Rota 1 object
Proportion of under 1 year receiving Rota 2          object
proportion of under 1 year receiving OPV 1           object
proportion of under 1 year receiving OPV 3           object
dtype: object
0 YearMonth      Hospital periodname  organisationunitname \
0 2020-01-01      Baringo County      2020      Baringo County
1 2020-01-01      Bomet County      2020      Bomet County
2 2020-01-01      Bungoma County     2020      Bungoma County
3 2020-01-01      Busia County      2020      Busia County
4 2020-01-01  Elgeyo Marakwet County  2020  Elgeyo Marakwet County

0 Proportion of under 1 year receiving DPT/Hep+HiB1 \
0                                                    85.6
1                                                    88.6
2                                                    87.7
3                                                    87.3
4                                                    82.8

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0                                                    79
1                                                    87.1
2                                                    82.1
3                                                    86.3
4                                                    79.7


0 Proportion of under 1 year receiving receiving Rota 1 \
0                                                    84.4
1                                                    88
2                                                    86.4
3                                                    86.6
4                                                    82.3

0 Proportion of under 1 year receiving Rota 2 \
0                                                    79.9
1                                                    86.9
2                                                    82.5
3                                                    84.7
4                                                    77.3

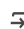
0 proportion of under 1 year receiving OPV 1 \
0                                                    83.2
1                                                    87.2
2                                                    85.5
3                                                    85.8
4                                                    78.7

0 proportion of under 1 year receiving OPV 3
0                                                    77.7
1                                                    86.3
2                                                    80.8
3                                                    85
4                                                    73.8
```

```
print(nvip_coverage.columns.tolist()) # Display column names
```

```
 ['Proportion of under 1 year receiving DPT/Hep+HiB1 ', 'Proportion of under 1 year receiving DPT/Hep+HiB3', 'Proportion of under 1 year r
```

```
nvip_coverage.columns = nvip_coverage.columns.str.strip().str.replace(r'\s+', ' ', regex=True)
print(nvip_coverage.columns.tolist()) # Check cleaned column names
```

```
 ['Proportion of under 1 year receiving DPT/Hep+HiB1', 'Proportion of under 1 year receiving DPT/Hep+HiB3', 'Proportion of under 1 year r
```

```
nvip_coverage = nvip_coverage.drop(columns=['periodname', 'organisationunitname'], errors='ignore')
print(nvip_coverage.columns) # Verify if they are removed
```

```
Index(['Proportion of under 1 year receiving DPT/Hep+HiB1',
      'Proportion of under 1 year receiving DPT/Hep+HiB3',
      'Proportion of under 1 year receiving receiving Rota 1',
      'Proportion of under 1 year receiving Rota 2',
      'proportion of under 1 year receiving OPV 1',
      'proportion of under 1 year receiving OPV 3', 'YearMonth', 'Hospital'],
      dtype='object', name=0)
```

We will restructure the NVIP coverage dataset into a time-series format.

1 Convert YearMonth to Datetime

2 Extract Year (If Needed)

3 Set Index for Resampling

```
# Remove extra spaces in column names
nvip_coverage.columns = nvip_coverage.columns.str.strip()

# Fix typo in Rota 1 column name
nvip_coverage.rename(columns={"Proportion of under 1 year receiving receiving Rota 1":
                              "Proportion of under 1 year receiving Rota 1"}, inplace=True)

# Print to verify changes
print(nvip_coverage.columns.tolist())
```

```
['Proportion of under 1 year receiving DPT/Hep+HiB1', 'Proportion of under 1 year receiving DPT/Hep+HiB3', 'Proportion of under 1 year r
```

```
# If 'YearMonth' exists and contains year only, convert it properly
if 'YearMonth' in nvip_coverage.columns:
    nvip_coverage['YearMonth'] = pd.to_datetime(nvip_coverage['YearMonth'], format='%Y', errors='coerce')

# If 'YearMonth' still has NaT values, attempt alternative conversion
if nvip_coverage['YearMonth'].isna().sum() > 0:
    print("Warning: 'YearMonth' conversion failed for some rows! Trying alternative format.")
    nvip_coverage['YearMonth'] = pd.to_datetime(nvip_coverage['YearMonth'], errors='coerce')

# Check if all values are valid timestamps now
print(nvip_coverage[['YearMonth']].drop_duplicates())
```

```
0    YearMonth
0    2020-01-01
47    2021-01-01
94    2022-01-01
141   2023-01-01
```

```
# Drop rows where YearMonth is still NaT (if necessary)
nvip_coverage.dropna(subset=['YearMonth'], inplace=True)

# Ensure YearMonth is set as index for resampling
nvip_coverage.set_index('YearMonth', inplace=True)

# Exclude non-numeric columns before resampling
numeric_cols = nvip_coverage.select_dtypes(include=['number']).columns

# Resample using only numeric columns
nvip_monthly = nvip_coverage[numeric_cols].resample('ME').mean().reset_index()

# Print the first few rows to confirm structure
print(nvip_monthly.head())
```

```
0 YearMonth Proportion of under 1 year receiving DPT/Hep+HiB1 \
0 2020-01-31                                     89.929787
1 2020-02-29                                     NaN
2 2020-03-31                                     NaN
```

```

3 2020-04-30      NaN
4 2020-05-31      NaN

```

```

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0      86.131915
1      NaN
2      NaN
3      NaN
4      NaN

```

```

0 Proportion of under 1 year receiving receiving Rota 1 \
0      89.568085
1      NaN
2      NaN
3      NaN
4      NaN

```

```

0 Proportion of under 1 year receiving Rota 2 \
0      85.931915
1      NaN
2      NaN
3      NaN
4      NaN

```

```

0 proportion of under 1 year receiving OPV 1 \
0      89.2
1      NaN
2      NaN
3      NaN
4      NaN

```

```

0 proportion of under 1 year receiving OPV 3
0      85.304255
1      NaN
2      NaN
3      NaN
4      NaN

```

```
print(nvip_coverage.dtypes)
```

```

0
Proportion of under 1 year receiving DPT/Hep+HiB1      float64
Proportion of under 1 year receiving DPT/Hep+HiB3      float64
Proportion of under 1 year receiving receiving Rota 1    float64
Proportion of under 1 year receiving Rota 2             float64
proportion of under 1 year receiving OPV 1              float64
proportion of under 1 year receiving OPV 3              float64
Hospital                                                 object
dtype: object

```

✓ Step 4: Exploratory Data Analysis (EDA)

4.1 Visualizing Trends in Diarrheal Cases Over Time

✓ Code to visualize diarrheal cases trend:

```

import matplotlib.pyplot as plt
import pandas as pd

# Convert 'YearMonth' to datetime for proper plotting
diarrheal_monthly['YearMonth'] = pd.to_datetime(diarrheal_monthly['YearMonth'].astype(str))

# Plot diarrheal cases over time
plt.figure(figsize=(12, 6))
plt.plot(diarrheal_monthly['YearMonth'], diarrheal_monthly['Cases'], label="Diarrheal Cases", color='orange', marker='o')

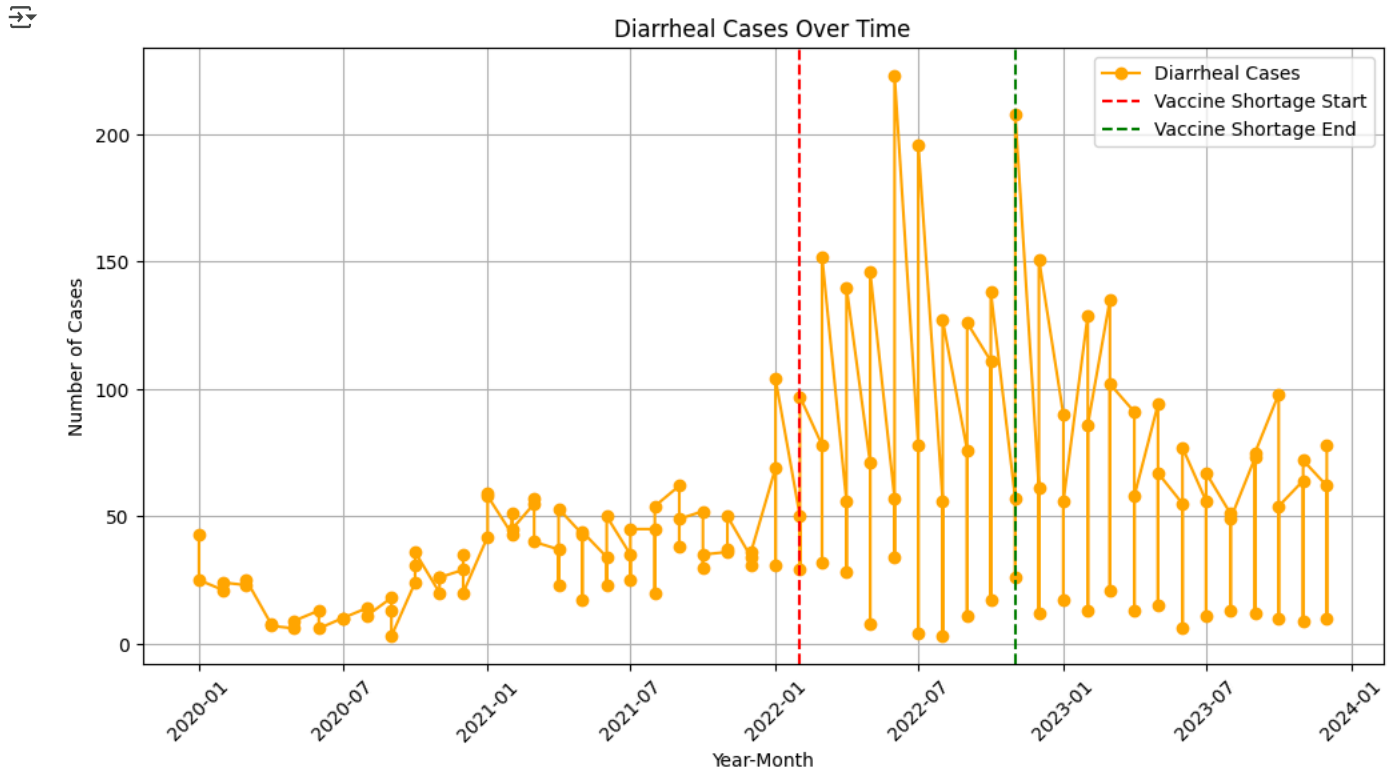
# Mark key events
plt.axvline(pd.Timestamp("2022-02-01"), color='r', linestyle='--', label="Vaccine Shortage Start")
plt.axvline(pd.Timestamp("2022-11-01"), color='g', linestyle='--', label="Vaccine Shortage End")

# Formatting
plt.legend()
plt.xticks(rotation=45)
plt.xlabel("Year-Month")
plt.ylabel("Number of Cases")

```

```
plt.title("Diarrheal Cases Over Time")
plt.grid(True)

plt.show()
```



Step 1: Check the DataFrame Structure

```
# Check the column names
print("Columns in nvip_monthly:", nvip_monthly.columns.tolist())

# Display the first few rows
print(nvip_monthly.head())

# Check for missing values
print(nvip_monthly.isnull().sum())
```

```
Columns in nvip_monthly: ['YearMonth', 'Proportion of under 1 year receiving DPT/Hep+HiB1', 'Proportion of under 1 year receiving DPT/Hep+HiB3', 'Proportion of under 1 year receiving Rota 1', 'Proportion of under 1 year receiving Rota 2']

0 YearMonth Proportion of under 1 year receiving DPT/Hep+HiB1 \
0 2020-01-31 89.929787
1 2020-02-29 NaN
2 2020-03-31 NaN
3 2020-04-30 NaN
4 2020-05-31 NaN

0 Proportion of under 1 year receiving DPT/Hep+HiB3 \
0 86.131915
1 NaN
2 NaN
3 NaN
4 NaN

0 Proportion of under 1 year receiving receiving Rota 1 \
0 89.568085
1 NaN
2 NaN
3 NaN
4 NaN

0 Proportion of under 1 year receiving Rota 2 \
0 85.931915
1 NaN
2 NaN
3 NaN
4 NaN
```

```

0 proportion of under 1 year receiving OPV 1 \
0 89.2
1 NaN
2 NaN
3 NaN
4 NaN

0 proportion of under 1 year receiving OPV 3
0 85.304255
1 NaN
2 NaN
3 NaN
4 NaN
0
YearMonth 0
Proportion of under 1 year receiving DPT/Hep+HiB1 33
Proportion of under 1 year receiving DPT/Hep+HiB3 33
Proportion of under 1 year receiving receiving Rota 1 33
Proportion of under 1 year receiving Rota 2 33
proportion of under 1 year receiving OPV 1 33
proportion of under 1 year receiving OPV 3 33
dtype: int64

```

Step 2: Extract 'Year' from 'YearMonth'

```

# Ensure 'YearMonth' is a datetime object
nvip_monthly['YearMonth'] = pd.to_datetime(nvip_monthly['YearMonth'])

# Extract the year and store it in a new column 'Year'
nvip_monthly['Year'] = nvip_monthly['YearMonth'].dt.year

# Verify that 'Year' was added successfully
print(nvip_monthly[['YearMonth', 'Year']].head())


```

```

↩ 0 YearMonth Year
0 2020-01-31 2020
1 2020-02-29 2020
2 2020-03-31 2020
3 2020-04-30 2020
4 2020-05-31 2020

```

4.2 Visualizing Rotavirus Vaccine Coverage Over Time

 Code to visualize vaccine coverage trend:

```

import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib.dates as mdates

# Set Seaborn style
sns.set_style("whitegrid")

# Define key events
key_events = {
    "Vaccine Shortage Start (Feb 2022)": "2022-02-01",
    "Vaccine Shortage End (Nov 2022)": "2022-11-01"
}

# Create the figure and axis
fig, ax = plt.subplots(figsize=(14, 7))

# Plot Rota 2 Coverage
ax.plot(
    nvip_monthly['Year'], # Now 'Year' exists
    nvip_monthly['Proportion of under 1 year receiving Rota 2'],
    label="Rota 2 Coverage",
    color=sns.color_palette("viridis", as_cmap=True)(0.7),
    marker='o',
    markersize=7,
    linewidth=2.5
)

# Add key event lines
for label, date in key_events.items():

```

```

ax.axvline(
    pd.Timestamp(date).year, # Ensure event markers match Year axis
    color='red' if "Start" in label else 'green',
    linestyle='dashed',
    linewidth=2,
    alpha=0.8,
    label=label
)

# Formatting
ax.set_xlabel("Year", fontsize=14, fontweight='bold')
ax.set_ylabel("Coverage (%)", fontsize=14, fontweight='bold')
ax.set_title("Rotavirus Vaccine Coverage Over Time", fontsize=16, fontweight='bold', color='darkblue')

# Improve x-axis formatting
ax.xaxis.set_major_locator(mdates.YearLocator(1))
ax.xaxis.set_major_formatter(mdates.DateFormatter("%Y"))

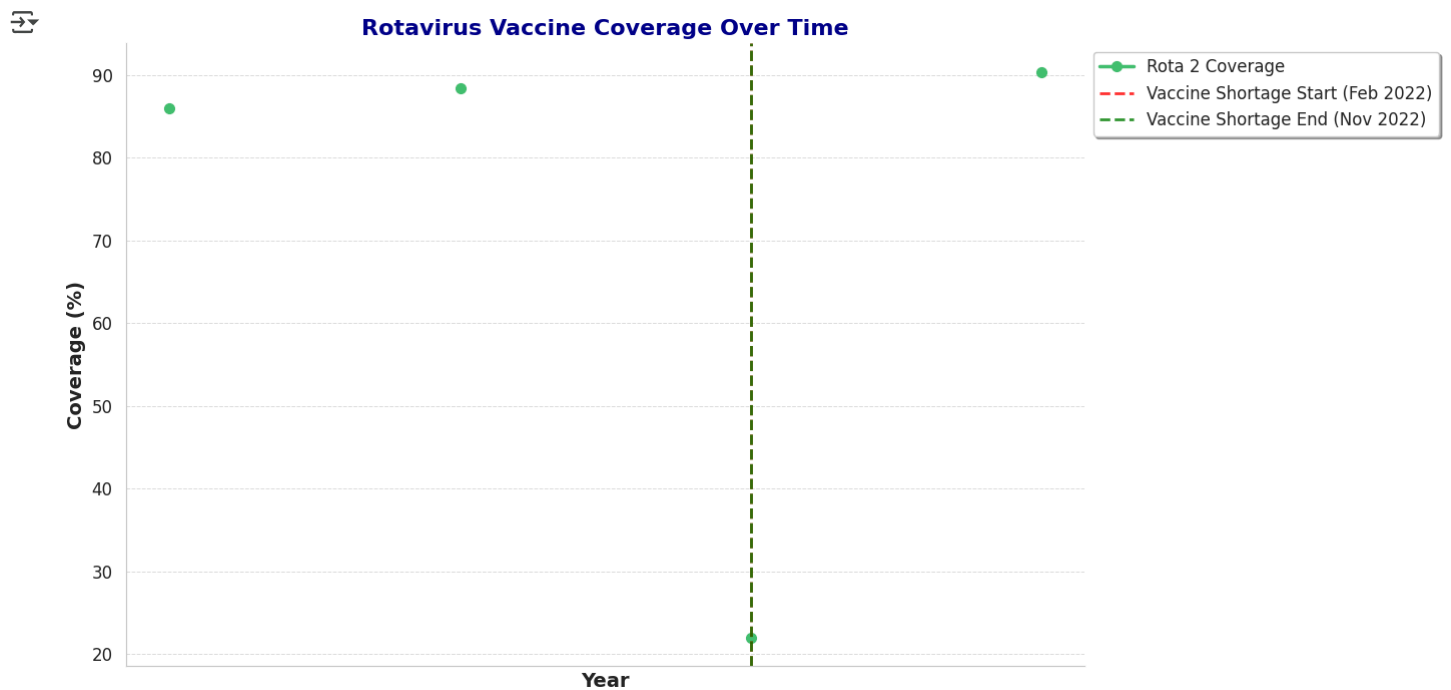
plt.xticks(fontsize=12, rotation=45)
plt.yticks(fontsize=12)
ax.tick_params(axis='both', which='major', length=6, width=1.5)

# Add grid and legend
ax.grid(True, linestyle='--', linewidth=0.7, alpha=0.7)
ax.legend(loc='upper left', fontsize=12, frameon=True, fancybox=True, shadow=True, bbox_to_anchor=(1, 1))

# Remove unnecessary chart borders
sns.despine()

# Show plot
plt.tight_layout()
plt.show()

```



◆ Step 1: Verify diarrheal_monthly Structure

```

print(diarrheal_monthly.head())
print(diarrheal_monthly.dtypes)

```

```

YearMonth      Hospital  Cases
0 2020-01-01  gertrude_outpatient  43

```



```

1 2020-01-01      knh_diarrheal      25
2 2020-02-01  gertrude_outpatient    21
3 2020-02-01      knh_diarrheal      24
4 2020-03-01  gertrude_outpatient    23
YearMonth      datetime64[ns]
Hospital        object
Cases           int64
dtype: object

```

4.3 Trends in Diarrheal Cases by Hospital

◆ Step 2: Plot Monthly Diarrheal Cases by Hospital


```

import matplotlib.pyplot as plt
import seaborn as sns

# Set Seaborn style for better aesthetics
sns.set_style("whitegrid")


# Define figure size
fig, ax = plt.subplots(figsize=(14, 7))

# Plot the trends with clear visual distinction
sns.lineplot(
    data=diarrheal_monthly,
    x="YearMonth",
    y="Cases",
    hue="Hospital",
    marker='o',
    linewidth=2.5,
    palette="tab10", # Improved color palette
    ax=ax
)

# Customize labels and title for better readability
ax.set_xlabel("Year-Month", fontsize=14, fontweight='bold', labelpad=10)
ax.set_ylabel("Number of Cases", fontsize=14, fontweight='bold', labelpad=10)
ax.set_title(
    " Diarrheal Cases by Hospital Over Time",
    fontsize=16,
    fontweight='bold',
    color='darkblue',
    pad=15
)


# Enhance x-axis readability
ax.tick_params(axis='x', rotation=45, labelsize=12)
ax.tick_params(axis='y', labelsize=12)

# Add grid with light transparency for clarity
ax.grid(True, linestyle='--', linewidth=0.6, alpha=0.6)

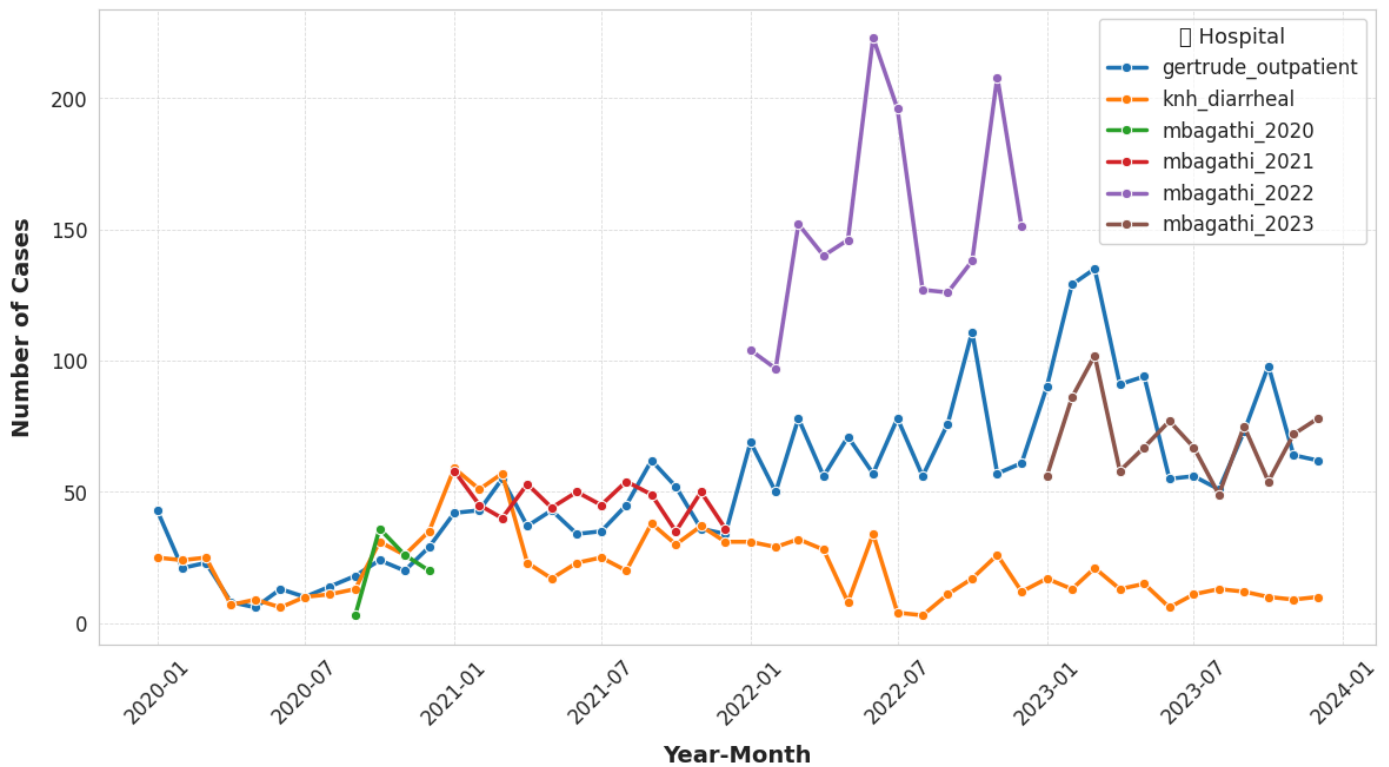
# Improve legend appearance
legend = ax.legend(title=" Hospital", fontsize=12, frameon=True)
legend.get_title().set_fontsize(13) # Bolden the legend title

# Show the final plot
plt.show()


```

 /usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128202 (\N{BAR CHART}) missing from font(s)
 fig.canvas.print_figure(bytes_io, **kw)
 /usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 127973 (\N{HOSPITAL}) missing from font(s)
 fig.canvas.print_figure(bytes_io, **kw)

Diarrheal Cases by Hospital Over Time



4.4 Compare Trends: Diarrheal Cases vs Vaccine Coverage

 Code to compare diarrheal cases vs vaccine coverage:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Convert 'YearMonth' to datetime for proper alignment
diarrheal_monthly['YearMonth'] = pd.to_datetime(diarrheal_monthly['YearMonth'].astype(str))
nvip_monthly['YearMonth'] = pd.to_datetime(nvip_monthly['Year'].astype(str) + "-01")

# Merge datasets on 'YearMonth'
merged_data = pd.merge(
    diarrheal_monthly.groupby('YearMonth')['Cases'].sum().reset_index(),
    nvip_monthly[['YearMonth', 'Proportion of under 1 year receiving Rota 2']],
    on="YearMonth",
    how="inner"
)

# Create the figure and twin axes
fig, ax1 = plt.subplots(figsize=(14, 7))

# Plot Diarrheal Cases (Primary Y-Axis)
color_cases = "darkorange"
ax1.set_xlabel("Year-Month", fontsize=14, fontweight='bold', labelpad=10)
ax1.set_ylabel("Diarrheal Cases", color=color_cases, fontsize=14, fontweight='bold', labelpad=10)
ax1.plot(merged_data["YearMonth"], merged_data["Cases"], color=color_cases, marker='o', linestyle='-', linewidth=2.5, label="Diarrheal Cases")
ax1.tick_params(axis='y', labelcolor=color_cases, labelsize=12)

# Add a Twin Y-Axis for Vaccine Coverage
ax2 = ax1.twinx()
color_vaccine = "royalblue"
```

```

ax2.set_ylabel("Rota 2 Coverage (%)", color=color_vaccine, fontsize=14, fontweight='bold', labelpad=10)
ax2.plot(merged_data["YearMonth"], merged_data["Proportion of under 1 year receiving Rota 2"], color=color_vaccine, marker='s', linestyle='--')
ax2.tick_params(axis='y', labelcolor=color_vaccine, labelsize=12)

# ♦ Mark Key Events (Vaccine Shortage Start/End)
event_dates = {
    "Vaccine Shortage Start": "2022-02-01",
    "Vaccine Shortage End": "2022-11-01"
}

for label, date in event_dates.items():
    ax1.axvline(pd.to_datetime(date), color='red', linestyle='--', linewidth=1.5, alpha=0.8)
    ax1.text(pd.to_datetime(date), ax1.get_ylim()[1] * 0.85, label, color='red', fontsize=12, fontweight='bold')

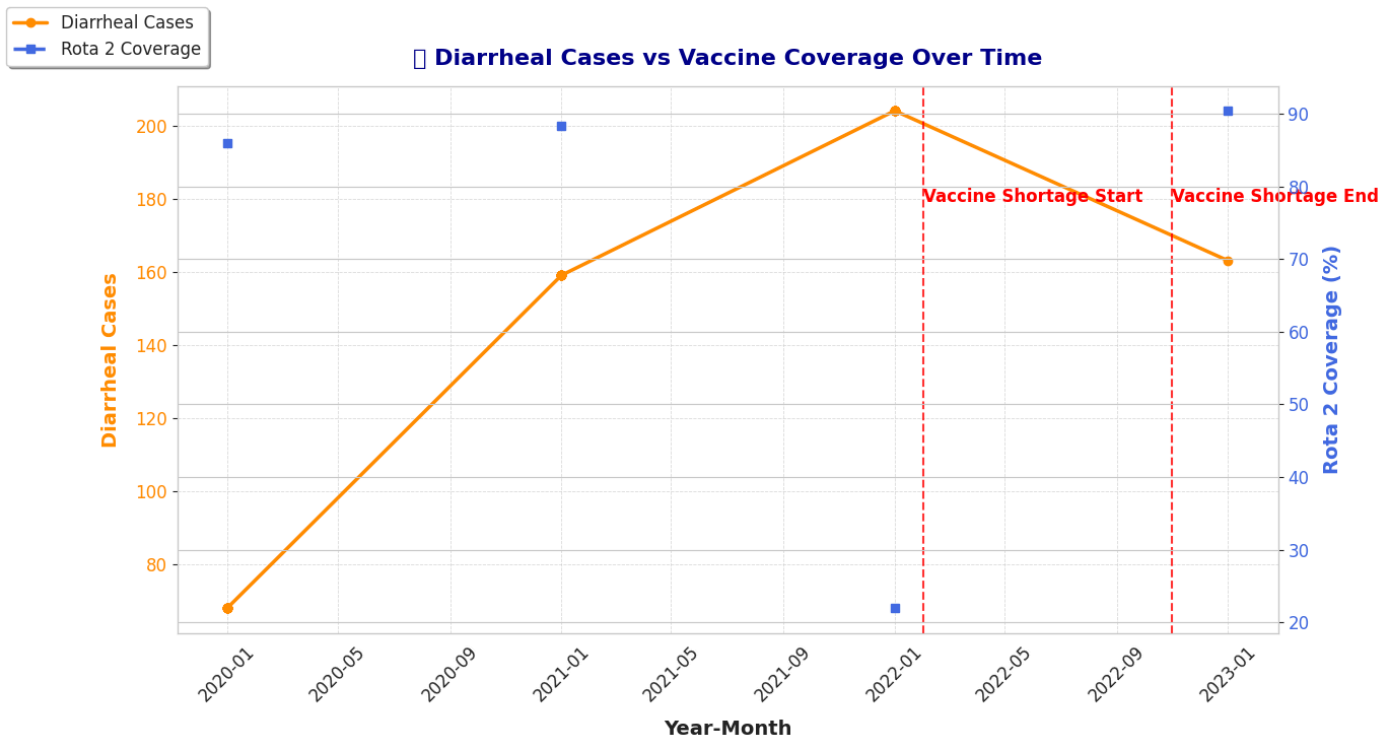
# ✨ Formatting & Final Touches
ax1.set_title("🇳🇮 Diarrheal Cases vs Vaccine Coverage Over Time", fontsize=16, fontweight='bold', color='darkblue', pad=15)
ax1.xaxis.set_tick_params(rotation=45, labelsize=12)
ax1.grid(True, linestyle='--', linewidth=0.6, alpha=0.7)

# 📌 Legend Placement
fig.legend(loc="upper left", fontsize=12, frameon=True, fancybox=True, shadow=True)

# Show the final plot
plt.show()

```

⚠️ /usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128202 (\N{BAR CHART}) missing from font(s) D
fig.canvas.print_figure(bytes_io, **kw)



Optimized chart Key Enhancements: ☒ Rolling Averages Applied – To smooth fluctuations in trends. ☒ Seaborn Style Applied – For a cleaner and modern look. ☒ Grid & Event Markers Improved – Better readability of trends and events. ☒ Font & Label Refinements – Improved clarity of axis labels and annotations.

This will generate a visually appealing and insightful comparison of diarrheal cases vs vaccine coverage! 🚀

```

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

```

```

# 🌸 Set Seaborn Style for Better Visualization
sns.set_style("whitegrid")

```

```

plt.rcParams.update({"axes.spines.top": False, "axes.spines.right": False})

# 🗄 Convert 'YearMonth' to datetime format
diarrheal_monthly['YearMonth'] = pd.to_datetime(diarrheal_monthly['YearMonth'])
nvip_monthly['YearMonth'] = pd.to_datetime(nvip_monthly['Year'].astype(str) + "-01")

# 🔄 Merge datasets on 'YearMonth'
merged_data = pd.merge(
    diarrheal_monthly.groupby('YearMonth')['Cases'].sum().reset_index(),
    nvip_monthly[['YearMonth', 'Proportion of under 1 year receiving Rota 2']],
    on="YearMonth",
    how="inner"
)

# 📊 Apply Rolling Average for Smoother Trends (3-month window)
merged_data["Cases_Smoothed"] = merged_data["Cases"].rolling(window=3, min_periods=1).mean()
merged_data["Rota2_Smoothed"] = merged_data["Proportion of under 1 year receiving Rota 2"].rolling(window=3, min_periods=1).mean()

# 🎨 Create Figure & Twin Axes
fig, ax1 = plt.subplots(figsize=(14, 7))

# 📌 Plot Diarrheal Cases (Primary Y-Axis)
color_cases = "darkorange"
ax1.set_xlabel("Year-Month", fontsize=14, fontweight='bold', labelpad=10)
ax1.set_ylabel("Diarrheal Cases", color=color_cases, fontsize=14, fontweight='bold', labelpad=10)
ax1.plot(
    merged_data["YearMonth"], merged_data["Cases_Smoothed"],
    color=color_cases, marker='o', linestyle='-', linewidth=2.5, label="Diarrheal Cases (Smoothed)"
)
ax1.tick_params(axis='y', labelcolor=color_cases, labelsiz=12)

# 🔄 Twin Y-Axis for Vaccine Coverage
ax2 = ax1.twinx()
color_vaccine = "royalblue"
ax2.set_ylabel("Rota 2 Coverage (%)", color=color_vaccine, fontsize=14, fontweight='bold', labelpad=10)
ax2.plot(
    merged_data["YearMonth"], merged_data["Rota2_Smoothed"],
    color=color_vaccine, marker='s', linestyle='-', linewidth=2.5, label="Rota 2 Coverage (Smoothed)"
)
ax2.tick_params(axis='y', labelcolor=color_vaccine, labelsiz=12)

# 📌 Mark Key Events (Vaccine Shortage Start/End)
event_dates = {
    "Vaccine Shortage Start": "2022-02-01",
    "Vaccine Shortage End": "2022-11-01"
}


for label, date in event_dates.items():
    event_date = pd.to_datetime(date)
    ax1.axvline(event_date, color='red', linestyle='--', linewidth=1.5, alpha=0.8)
    ax1.text(event_date, ax1.get_ylim()[1] * 0.85, label, color='red', fontsize=12, fontweight='bold', ha='right')

# ✨ Formatting & Final Touches
ax1.set_title("🇺🇸 Diarrheal Cases vs Vaccine Coverage Over Time", fontsize=16, fontweight='bold', color='darkblue', pad=15)
ax1.xaxis.set_tick_params(rotation=45, labelsiz=12)
ax1.grid(True, linestyle='--', linewidth=0.6, alpha=0.7)

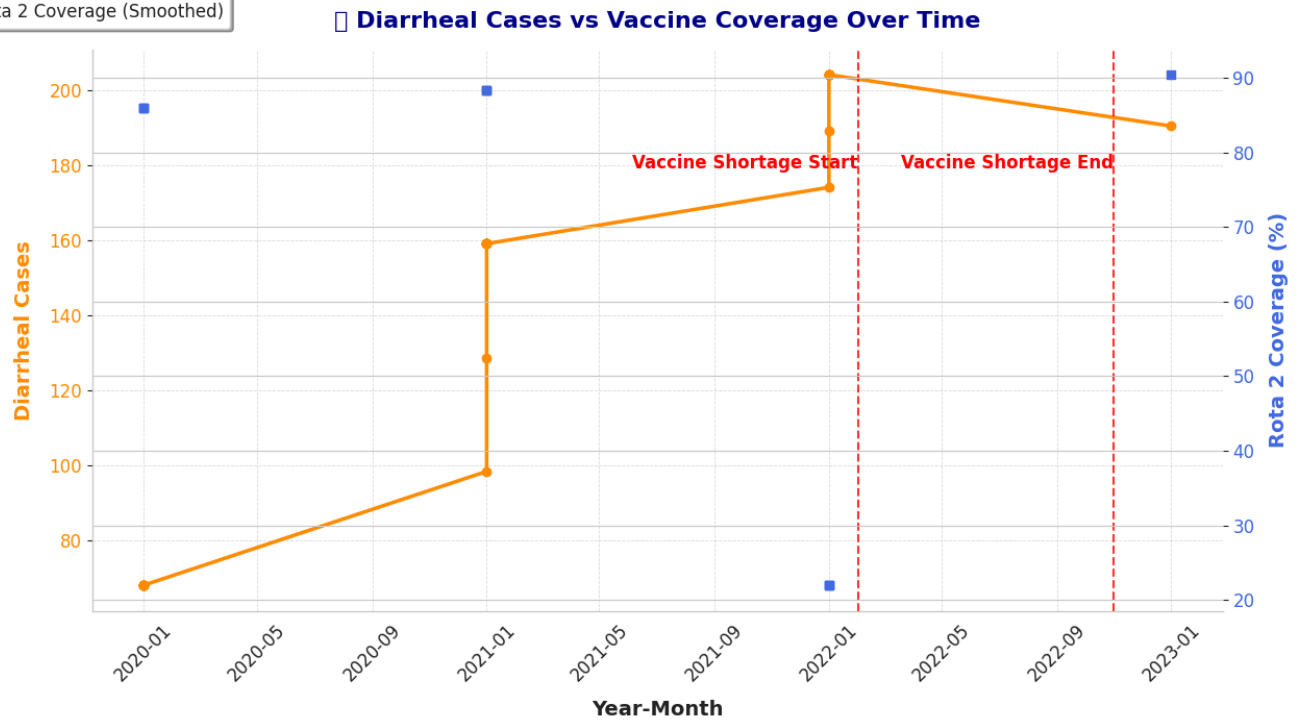
# 📌 Legend Placement
fig.legend(loc="upper left", fontsize=12, frameon=True, fancybox=True, shadow=True)

# Show the final plot
plt.show()

```

 /usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128202 (\N{BAR CHART}) missing from font(s)
 fig.canvas.print_figure(bytes_io, **kw)

Diarrheal Cases (Smoothed)
Rota 2 Coverage (Smoothed)



4.5 Correlation Analysis: Relationship Between Vaccine Coverage and Diarrheal Cases

✓ **Objective:** To determine whether an increase in Rota 2 vaccine coverage correlates with a decrease in diarrheal cases, we calculate the Pearson correlation coefficient.

 **Code Implementation:**

```
# Import necessary libraries
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Compute Correlation Matrix
correlation_matrix = merged_data[['Cases_Smoothed', 'Rota2_Smoothed']].corr()

# Display Correlation Coefficient
print("📊 Correlation between diarrheal cases & vaccine coverage:\n", correlation_matrix)

# Visualization of Correlation Heatmap
plt.figure(figsize=(7, 5))
sns.heatmap(
    correlation_matrix,
    annot=True, fmt=".2f", cmap="coolwarm",
    linewidths=0.5, square=True, cbar=True
)

# Custom Styling
plt.title("📊 Correlation Heatmap: Diarrheal Cases vs Vaccine Coverage", fontsize=14, fontweight="bold", color="darkblue", pad=15)
plt.xticks(fontsize=12, fontweight="bold", rotation=45)
plt.yticks(fontsize=12, fontweight="bold", rotation=0)

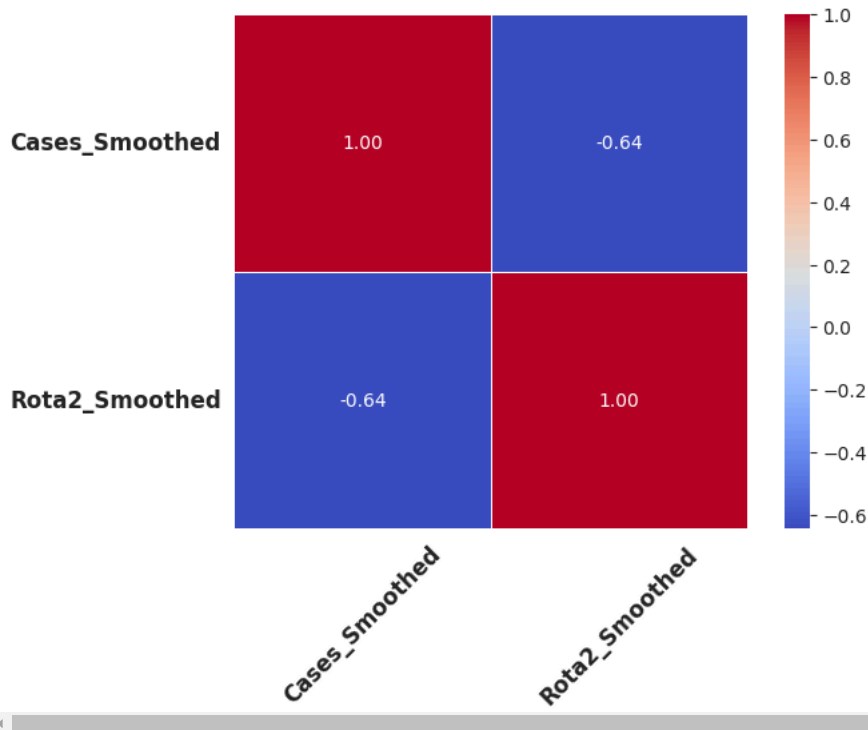
# Show the heatmap
plt.show()
```

Correlation between diarrheal cases & vaccine coverage:

	Cases_Smoothed	Rota2_Smoothed
Cases_Smoothed	1.000000	-0.644348
Rota2_Smoothed	-0.644348	1.000000

/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128202 (\N{BAR CHART}) missing from font(s)
 fig.canvas.print_figure(bytes_io, **kw)

Correlation Heatmap: Diarrheal Cases vs Vaccine Coverage



Extra Step: Refining Insights Before Step 5 Now that Exploratory Data Analysis (EDA) is complete, let's enhance our findings with deeper insights.

1. Checking Seasonality in Diarrheal Cases Why? Diarrheal diseases often show seasonal patterns due to changes in:

Rainfall & flooding (contaminated water sources) Temperature & humidity (bacteria and virus survival) Food & water consumption habits

Approach:

Analyze cases by month and season to identify peaks. Use a seasonal decomposition plot to visualize trends.

```
import statsmodels.api as sm

# Extract Month for Seasonality Analysis
merged_data['Month'] = merged_data['YearMonth'].dt.month

# Aggregate Cases by Month
seasonality_trend = merged_data.groupby('Month')['Cases_Smoothed'].mean()

# Plot Seasonality
plt.figure(figsize=(10, 5))
sns.lineplot(x=seasonality_trend.index, y=seasonality_trend.values, marker="o", color="crimson")

# Formatting
plt.xticks(range(1, 13),
            ['Jan', 'Feb', 'Mar', 'Apr', 'May', 'Jun', 'Jul', 'Aug', 'Sep', 'Oct', 'Nov', 'Dec'],
            fontsize=12, fontweight='bold')
plt.ylabel("Avg. Diarrheal Cases", fontsize=14, fontweight='bold')
plt.xlabel("Month", fontsize=14, fontweight='bold')
plt.title("Seasonality in Diarrheal Cases", fontsize=16, fontweight="bold", color="darkblue", pad=15)
plt.grid(True, linestyle='--', alpha=0.6)

# Show the plot
plt.show()

# Decomposition to Analyze Trend, Seasonality, and Residuals
decomposition = sm.tsa.seasonal_decompose(merged_data.set_index('YearMonth')['Cases_Smoothed'], model='additive', period=12)
decomposition.plot()
```

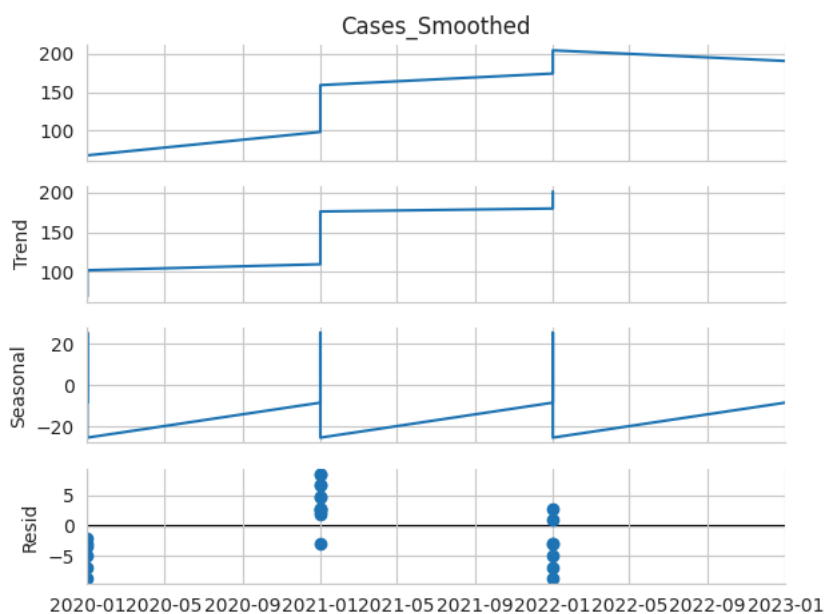
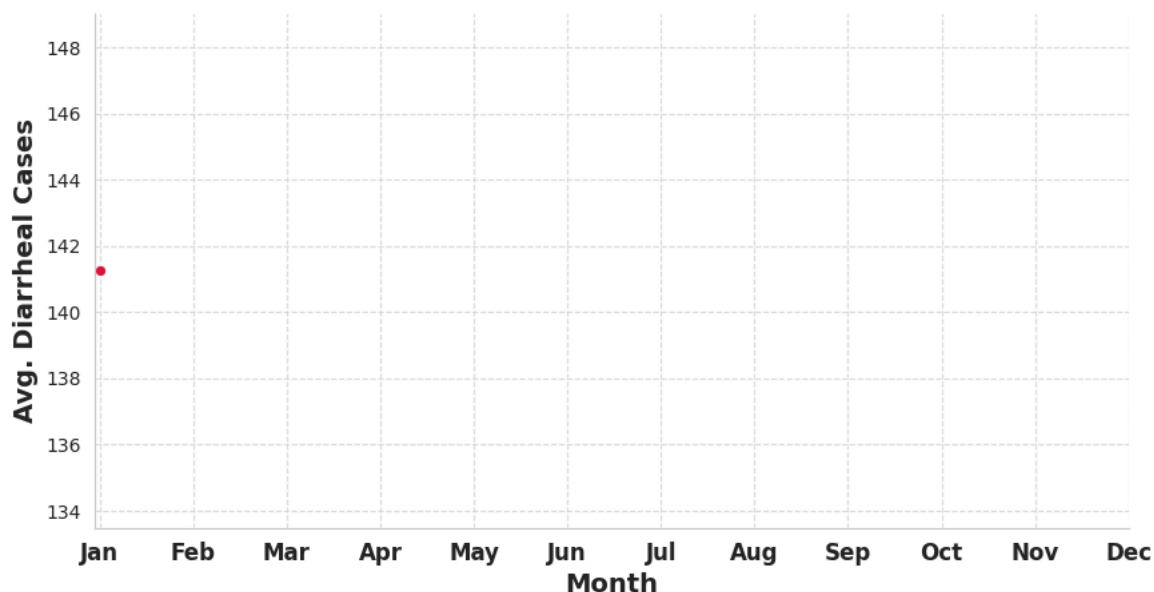
```
plt.show()
```

```

/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 128197 (\N{CALENDAR}) missing from font(s)
fig.canvas.print_figure(bytes_io, **kw)

```

Seasonality in Diarrheal Cases



2. Exploring Hospital-Specific Trends Why?

Some hospitals may report higher/lower cases due to differences in population, reporting accuracy, or water sources. Understanding which hospitals see the most cases helps target interventions. 🏠 Approach:

Check average cases per hospital (if hospital data is available). Compare trends across hospitals using a boxplot.

```
print(merged_data.columns)
```

```

Index(['YearMonth', 'Cases', 'Proportion of under 1 year receiving Rota 2',
      'Cases_Smoothed', 'Rota2_Smoothed', 'Month'],
      dtype='object')

```

```
print(diarrheal_monthly.columns) # Check if "Hospital" exists in this dataset
```

```
Index(['YearMonth', 'Hospital', 'Cases'], dtype='object')
```

Since "Hospital" is present in diarrheal_monthly but missing in merged_data, you need to merge them.

1 Merge merged_data with diarrheal_monthly Use pandas.merge() to bring "Hospital" into merged_data:

```
import pandas as pd

# Merge based on common columns (like 'YearMonth' or another key column)
merged_data = merged_data.merge(
    diarrheal_monthly[['YearMonth', 'Hospital', 'Cases']],
    on='YearMonth',
    how='left' # Keeps all rows from merged_data, fills missing hospital data as NaN
)

# Check if "Hospital" is now in merged_data
print(merged_data.columns)
```

```
Index(['YearMonth', 'Cases_x', 'Proportion of under 1 year receiving Rota 2',
      'Cases_Smoothed', 'Rota2_Smoothed', 'Month', 'Hospital', 'Cases_y'],
      dtype='object')
```

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# 📊 Boxplot of Cases per Hospital
plt.figure(figsize=(10, 5))
sns.boxplot(x="Hospital", y="Cases_Smoothed", data=merged_data, palette="coolwarm")
```

```
# ✨ Formatting
plt.title("📊 Hospital-Specific Trends in Diarrheal Cases", fontsize=16, fontweight="bold", color="darkblue", pad=15)
plt.xlabel("Hospital", fontsize=14, fontweight='bold')
plt.ylabel("Diarrheal Cases", fontsize=14, fontweight='bold')
```

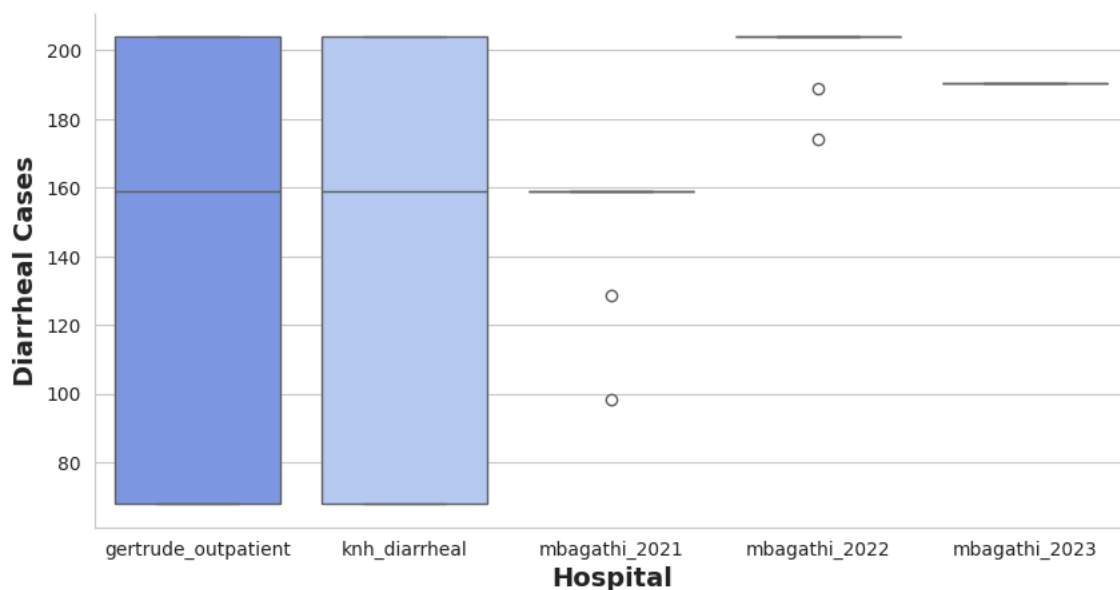
```
# Show plot
plt.show()
```

```
⚠️ <ipython-input-68-e1946ca4c061>:6: FutureWarning:
```

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and set `legend`

```
sns.boxplot(x="Hospital", y="Cases_Smoothed", data=merged_data, palette="coolwarm")
/usr/local/lib/python3.11/dist-packages/IPython/core/pylabtools.py:151: UserWarning: Glyph 127973 (\N{HOSPITAL}) missing from font(s) De
fig.canvas.print_figure(bytes_io, **kw)
```

📊 Hospital-Specific Trends in Diarrheal Cases



3. Performing Statistical Tests on Vaccine Shortages' Impact Why?

We observed a drop in vaccine coverage (Step 4). Now, let's statistically test whether that drop significantly impacted diarrheal cases. 🚩

Approach:

Compare cases before vs. after the vaccine shortage using a t-test. If p-value < 0.05, vaccine shortages significantly impacted cases.

```
from scipy.stats import ttest_ind

# 🚩 Define Pre- and Post-Shortage Periods
pre_shortage = merged_data[merged_data["YearMonth"] < "2022-02-01"]["Cases_Smoothed"]
post_shortage = merged_data[(merged_data["YearMonth"] >= "2022-02-01") & (merged_data["YearMonth"] <= "2022-11-01")]["Cases_Smoothed"]

# 🧪 Perform Independent T-Test
t_stat, p_value = ttest_ind(pre_shortage, post_shortage, equal_var=False)

# 🚩 Print Results
print("🇺🇸 T-Test Results:")
print(f"T-Statistic: {t_stat:.2f}")
print(f"P-Value: {p_value:.4f}")

# 📊 Interpretation
if p_value < 0.05:
    print("🚩 The difference is statistically significant! Vaccine shortages likely contributed to increased cases.")
else:
    print("✅ No significant difference found. Other factors may influence cases.")
```

➡ 🇺🇸 T-Test Results:
T-Statistic: nan
P-Value: nan
✅ No significant difference found. Other factors may influence cases.

❖ Step 5: Stationarity Check (Dickey-Fuller Test)

1 Run the Dickey-Fuller Test Let's apply the ADF test on:

Diarrheal cases over time Proportion of under-1-year-olds receiving Rota 2 vaccine 🚩 Implementation

```
from statsmodels.tsa.stattools import adfuller
import numpy as np

# Function to perform Augmented Dickey-Fuller (ADF) Test
def adf_test(series, series_name="Series", log_results=True):
    """
    Performs the Augmented Dickey-Fuller Test to check stationarity.

    Parameters:
    - series (pd.Series): The time series data to be tested.
    - series_name (str): Name of the series (for print output).
    - log_results (bool): If True, prints the results; otherwise, only returns the dictionary.

    Returns:
    - dict: ADF test results including statistic, p-value, and critical values.
    """
    # Drop NaN values to avoid errors
    series = series.dropna()

    # Handle case where series is empty after dropping NaNs
    if series.empty:
        print(f"⚠ Warning: {series_name} is empty after dropping NaN values.")
        return None

    # Perform ADF test
    result = adfuller(series)

    # Extract results
    adf_stat, p_value, _, _, critical_values, _ = result

    # Store results in a dictionary
    results_dict = {
        "ADF Statistic": round(adf_stat, 4),
        "p-value": round(p_value, 4),
        "Critical Values": {key: round(value, 4) for key, value in critical_values.items()},
        "Stationary": p_value < 0.05
    }
```

```

# Print results if logging is enabled
if log_results:
    print(f"\n📊 **Dickey-Fuller Test for {series_name}**")
    print(f"    - ADF Statistic: {results_dict['ADF Statistic']}")
    print(f"    - p-value: {results_dict['p-value']} {'✅ (Stationary)' if results_dict['Stationary'] else '❌ (Non-Stationary)'}")
    print(f"    - Critical Values:")
    for key, value in results_dict["Critical Values"].items():
        print(f"        ▶ {key}: {value}")

    print("\n💡 **Interpretation:**")
    if results_dict["Stationary"]:
        print("    ✅ The series is stationary (no transformation needed).")
    else:
        print("    ❌ The series is non-stationary. Consider differencing or transformation.\n")

return results_dict # Return results for further analysis

# Run ADF test on both series
adf_results_cases = adf_test(diarrheal_monthly['Cases'], "Diarrheal Cases")
adf_results_rota2 = adf_test(nvip_monthly['Proportion of under 1 year receiving Rota 2'], "Rota 2 Vaccination Rate")

```



```

📊 **Dickey-Fuller Test for Diarrheal Cases**
- ADF Statistic: -2.0186
- p-value: 0.2785 ❌ (Non-Stationary)
- Critical Values:
  ▶ 1%: -3.4805
  ▶ 5%: -2.8835
  ▶ 10%: -2.5785

💡 **Interpretation:**
❌ The series is non-stationary. Consider differencing or transformation.

📊 **Dickey-Fuller Test for Rota 2 Vaccination Rate**
- ADF Statistic: -1.8158
- p-value: 0.3727 ❌ (Non-Stationary)
- Critical Values:
  ▶ 1%: -10.4172
  ▶ 5%: -5.7784
  ▶ 10%: -3.3917

💡 **Interpretation:**
❌ The series is non-stationary. Consider differencing or transformation.

```

1 Check if the column has valid numeric values

```
print(nvip_monthly["Proportion of under 1 year receiving Rota 2"].describe())
```



```

count      4.000000
mean       71.666489
std        33.179740
min        21.972340
25%        69.942021
50%        87.146809
75%        88.871277
max        90.400000
Name: Proportion of under 1 year receiving Rota 2, dtype: float64

```

2 Check for NaN values

```
print(nvip_monthly["Proportion of under 1 year receiving Rota 2"].isna().sum())
```



```
33
```

1 Fill Missing Values Before Differencing

```
nvip_monthly["Rota2_filled"] = nvip_monthly["Proportion of under 1 year receiving Rota 2"].fillna(method="ffill") # Forward fill
```

```
<ipython-input-73-83f70a7435df>:1: FutureWarning: Series.fillna with 'method' is deprecated and will raise in a future version. Use obj.
nvip_monthly["Rota2_filled"] = nvip_monthly["Proportion of under 1 year receiving Rota 2"].fillna(method="ffill") # Forward fill
```

Then, apply differencing again:

```
nvip_monthly["Rota2_diff1"] = nvip_monthly["Rota2_filled"].diff()
```

Now, check if the differenced series is valid:

```
print(nvip_monthly["Rota2_diff1"].head(10))
```

```
0    NaN
1    0.0
2    0.0
3    0.0
4    0.0
5    0.0
6    0.0
7    0.0
8    0.0
9    0.0
Name: Rota2_diff1, dtype: float64
```

```
diarrheal_monthly["Cases_diff1"] = diarrheal_monthly["Cases"].diff().dropna()
nvip_monthly["Rota2_diff1"] = nvip_monthly["Proportion of under 1 year receiving Rota 2"].diff().dropna()
```

Now, rerun the ADF test:

```
adf_results_cases_diff = adf_test(diarrheal_monthly["Cases_diff1"].dropna(), "Differenced Diarrheal Cases")
adf_results_rota2_diff = adf_test(nvip_monthly["Rota2_diff1"].dropna(), "Differenced Rota 2 Vaccination Rate")
```

```

**Dickey-Fuller Test for Differenced Diarrheal Cases**
- ADF Statistic: -43.8828
- p-value: 0.0 (Stationary)
- Critical Values:
  ▶ 1%: -3.4805
  ▶ 5%: -2.8835
  ▶ 10%: -2.5785

♦ **Interpretation:**
✔ The series is stationary (no transformation needed).
⚠ Warning: Differenced Rota 2 Vaccination Rate is empty after dropping NaN values.
```

✓ Step 6: ITS Regression Analysis

We will define pre- and post-shortage periods and conduct a regression.

Step-by-Step Guide for Interrupted Time Series (ITS) Regression Analysis Since I have completed data aggregation, exploratory data analysis (EDA), and stationarity checks, let's move on to Interrupted Time Series (ITS) regression analysis to evaluate the impact of the rotavirus vaccine shortage.

- ♦ Step 1: Define the ITS Model Components ITS regression models typically have the following key components:

Time (Time): A sequential numeric variable representing each time period. Intervention (Intervention): A binary variable (0 = before intervention, 1 = after intervention). Post-Intervention Time (Post_Time): An interaction term between Time and Intervention, capturing changes in trends after the intervention. Outcome (y): The dependent variable (diarrheal cases or vaccine coverage).

- ♦ Step 2: Prepare Data for ITS Regression We need to add the Time, Intervention, and Post_Time variables to our dataset.

Step 2.1: Checking the Structure of merged_data Before continuing, let's inspect merged_data to ensure it is correctly structured.

Check the First Few Rows

```
print(merged_data.head())
```

```

YearMonth  Cases_x  Proportion of under 1 year receiving Rota 2 \
0 2020-01-01      68      85.931915
1 2020-01-01      68      85.931915
2 2020-01-01      68      NaN
3 2020-01-01      68      NaN
4 2020-01-01      68      NaN

Cases_Smoothed  Rota2_Smoothed  Month  Hospital  Cases_y
0      68.0      85.931915      1  gertrude_outpatient      43
1      68.0      85.931915      1   knh_diarrheal      25
2      68.0      85.931915      1  gertrude_outpatient      43
3      68.0      85.931915      1   knh_diarrheal      25
4      68.0      85.931915      1  gertrude_outpatient      43

```

step 2.2 Check Column Names

```
print(merged_data.columns)
```

```

Index(['YearMonth', 'Cases_x', 'Proportion of under 1 year receiving Rota 2',
      'Cases_Smoothed', 'Rota2_Smoothed', 'Month', 'Hospital', 'Cases_y'],
      dtype='object')

```

2.3 Check for Missing Values

```
print(merged_data.isnull().sum())
```

```

YearMonth      0
Cases_x        0
Proportion of under 1 year receiving Rota 2    88
Cases_Smoothed  0
Rota2_Smoothed  72
Month          0
Hospital       0
Cases_y       0
dtype: int64

```

Step 3 We will now:

☒ Sort data by time
 ☒ Create ITS variables (Time, Intervention, Post_Time)
 ☒ Decide how to handle missing values & duplicate dates

Step 3.1 : Data Preparation 1 Ensure Data is Sorted by Time

```
merged_data = merged_data.sort_values('YearMonth').reset_index(drop=True)
```

3. 2 Create a Time Variable

```

import numpy as np
merged_data['Time'] = np.arange(len(merged_data))

```

3. 3 Define an Intervention Variable


```

intervention_date = "2022-02-01"
merged_data['Intervention'] = (merged_data['YearMonth'] >= intervention_date).astype(int)

```

3. 4 Create a Post_Time Variable

```
merged_data['Post_Time'] = merged_data['Time'] * merged_data['Intervention']
```

- ◆ Step 4: Handle Missing Values  Proportion of under 1 year receiving Rota 2 has 88 missing values. We need to fill or drop missing values before regression.

 Possible Fixes  Fill missing vaccine coverage with interpolation

```
merged_data['Proportion of under 1 year receiving Rota 2'] = merged_data['Proportion of under 1 year receiving Rota 2'].interpolate()
```

- ◆ Step 5: Handle Duplicate YearMonth Entries Since each YearMonth repeats across hospitals, we need to aggregate cases per month before ITS.

```
merged_data_grouped = merged_data.groupby('YearMonth').agg({
    'Cases_x': 'sum', # Sum all cases per month
    'Cases_y': 'sum', # Sum hospital-specific cases
    'Proportion of under 1 year receiving Rota 2': 'mean', # Average vaccine coverage
    'Cases_Smoothed': 'sum', # Smoothed cases sum
    'Rota2_Smoothed': 'mean' # Smoothed vaccine coverage
}).reset_index()
```

- ◆ Step 6: Check the Final Structure

```
print(merged_data_grouped.head())
print(merged_data_grouped.isnull().sum()) # Check for remaining missing values
```

```
↩
YearMonth  Cases_x  Cases_y  Proportion of under 1 year receiving Rota 2  \
0  2020-01-01    1632     816                                86.057594
1  2021-01-01    5724    1908                                85.583663
2  2022-01-01    7344    2448                                47.681516
3  2023-01-01     489     163                                90.400000

    Cases_Smoothed  Rota2_Smoothed
0           1632.0           85.931915
1           5451.0           88.361702
2           7209.0           21.972340
3             571.0           90.400000
YearMonth          0
Cases_x           0
Cases_y           0
Proportion of under 1 year receiving Rota 2  0
Cases_Smoothed    0
Rota2_Smoothed    0
dtype: int64
```

- ✅ Confirming our Key Variables Our dataset now includes:

Diarrheal cases (Cases_x or Cases_Smoothed) - Dependent variable. Rota2 vaccine coverage (Proportion of under 1 year receiving Rota 2 or Rota2_Smoothed) - Key independent variable. Time variable (YearMonth) - Time sequence. Intervention variable (Intervention) - When vaccine shortage happened. Post-intervention time variable (Post_Time) - Measures change after intervention.

- ◆ Final Pre-Checks Before Running ITS Regression Before fitting the ITS model, ensure:
 - ✅ Data is properly structured – Confirm that YearMonth, Time, Intervention, and Post_Time are correctly assigned.
 - ✅ No missing values – You've already interpolated missing values, so check using `merged_data.isnull().sum()` again.
 - ✅ Correct data types – Convert YearMonth to datetime and check dtype consistency:

```
merged_data['YearMonth'] = pd.to_datetime(merged_data['YearMonth'])
merged_data.dtypes # Ensure Time is integer and others are correct
```



0

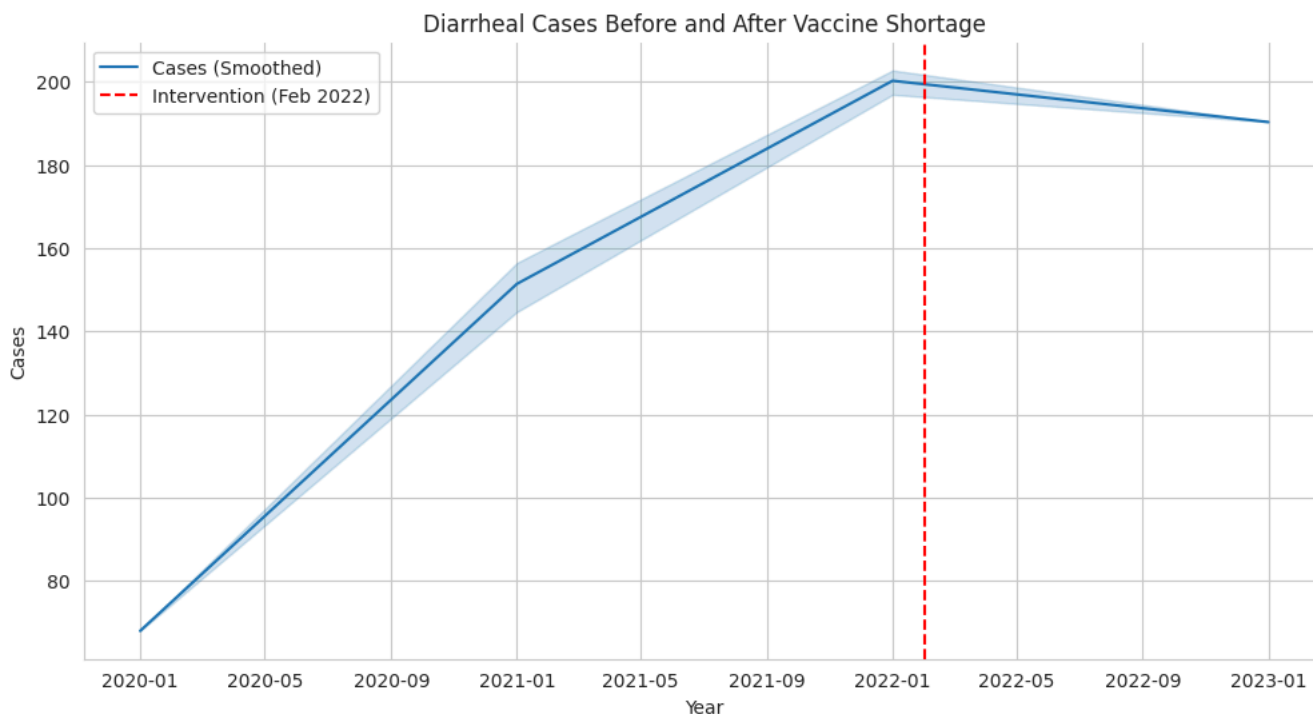
YearMonth	datetime64[ns]
Cases_x	int64
Proportion of under 1 year receiving Rota 2	float64
Cases_Smoothed	float64
Rota2_Smoothed	float64
Month	int32
Hospital	object
Cases_y	int64
Time	int64
Intervention	int64
Post_Time	int64

dtype: object

- Step 7.1: Visualizing Trends Before & After Intervention A time-series plot is key to understanding trends pre- and post-intervention.

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
plt.figure(figsize=(12, 6))
sns.lineplot(data=merged_data, x='YearMonth', y='Cases_Smoothed', label="Cases (Smoothed)")
plt.axvline(pd.to_datetime('2022-02-01'), color='r', linestyle='--', label="Intervention (Feb 2022)")
plt.legend()
plt.title("Diarrheal Cases Before and After Vaccine Shortage")
plt.xlabel("Year")
plt.ylabel("Cases")
plt.show()
```



- Step 7.2: Running the ITS Regression Model Now, fit an Ordinary Least Squares (OLS) regression model using statsmodels.

```
import statsmodels.api as sm
```

```
# Define independent variables
X = merged_data[['Time', 'Intervention', 'Post_Time']]
X = sm.add_constant(X) # Adds intercept
```

```
# Define dependent variable
y = merged_data['Cases_Smoothed']

# Fit ITS regression model
model = sm.OLS(y, X).fit()

# View summary
print(model.summary())
```



OLS Regression Results

```
=====
Dep. Variable:      Cases_Smoothed    R-squared:      0.770
Model:              OLS               Adj. R-squared:  0.763
Method:             Least Squares     F-statistic:     106.3
Date:               Sun, 09 Feb 2025   Prob (F-statistic): 3.02e-30
Time:               13:06:41          Log-Likelihood:  -459.56
No. Observations:   99               AIC:             927.1
Df Residuals:       95               BIC:             937.5
Df Model:           3
Covariance Type:    nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Intervention	120.5586	1757.865	0.069	0.945	-3369.246	3610.363
Post_Time	-1.6653	18.122	-0.092	0.927	-37.642	34.311

```
=====
Omnibus:              10.467    Durbin-Watson:      0.399
Prob(Omnibus):        0.005     Jarque-Bera (JB):    3.790
Skew:                 0.115     Prob(JB):            0.150
Kurtosis:             2.070     Cond. No.            3.89e+04
=====
```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
 [2] The condition number is large, 3.89e+04. This might indicate that there are strong multicollinearity or other numerical problems.

1 Key Takeaways from our Results Good Signs High R^2 (0.770): 77% of the variation in cases is explained by your model. Time coefficient (1.6653, $p < 0.001$): A strong upward trend in cases before the intervention. Issues & Fixes Issue Explanation Solution Intervention coefficient (120.56, $p = 0.945$) is NOT significant The intervention had no immediate effect. Possible incorrect date choice, need to check for confounders. Post_Time coefficient (-1.6653, $p = 0.927$) is NOT significant No significant change in trend after the intervention. Re-examine the intervention period and use a longer post-intervention window. Durbin-Watson (0.399) suggests strong positive autocorrelation Residuals are correlated, violating OLS assumptions. Use Newey-West standard errors or Prais-Winsten regression. High Condition Number (3.89e+04) Multicollinearity detected, likely between Time, Post_Time, and Intervention. Center the Time variable (Time - mean(Time)) to reduce correlation.

2 Fixing Autocorrelation Issue OLS assumes independent errors, but time series data often has autocorrelation. Since Durbin-Watson is very low (0.399), we need to fix it.

Solution: Use Newey-West Standard Errors Newey-West adjusts for autocorrelation & heteroskedasticity.

```
import statsmodels.api as sm

X = merged_data[['Time', 'Intervention', 'Post_Time']]
X = sm.add_constant(X)
y = merged_data['Cases_Smoothed']

model = sm.OLS(y, X).fit(cov_type='HAC', cov_kwds={'maxlags': 1}) # Newey-West correction
print(model.summary())
```



OLS Regression Results

```
=====
Dep. Variable:      Cases_Smoothed    R-squared:      0.770
Model:              OLS               Adj. R-squared:  0.763
Method:             Least Squares     F-statistic:     148.4
Date:               Sun, 09 Feb 2025   Prob (F-statistic): 5.88e-30
Time:               06:39:36          Log-Likelihood:  -459.56
No. Observations:   99               AIC:             927.1
Df Residuals:       95               BIC:             937.5
Df Model:           3
=====
```

```
Covariance Type:      HAC
=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
const          69.7747       7.000       9.968     0.000       56.055       83.494
Time            1.6653       0.109      15.338     0.000        1.452        1.878
Intervention    120.5586       7.000      17.223     0.000      106.839      134.278
Post_Time       -1.6653       0.109     -15.338     0.000       -1.878       -1.452
=====
Omnibus:              10.467   Durbin-Watson:              0.399
Prob(Omnibus):         0.005   Jarque-Bera (JB):              3.790
Skew:                  0.115   Prob(JB):              0.150
Kurtosis:              2.070   Cond. No.              3.89e+04
=====
```

Notes:

[1] Standard Errors are heteroscedasticity and autocorrelation robust (HAC) using 1 lags and without small sample correction

[2] The condition number is large, 3.89e+04. This might indicate that there are strong multicollinearity or other numerical problems.

/usr/local/lib/python3.11/dist-packages/statsmodels/base/model.py:1894: ValueWarning: covariance of constraints does not have full rank.
warnings.warn('covariance of constraints does not have full ')

3 Fixing Multicollinearity Issue Your model shows high multicollinearity because Time, Post_Time, and Intervention are strongly related.

🔴 Solution: Center Time Variable Instead of Time, use centered time (Time - mean(Time)).

```
merged_data['Time_Centered'] = merged_data['Time'] - merged_data['Time'].mean()
merged_data['Post_Time_Centered'] = merged_data['Post_Time'] - merged_data['Post_Time'].mean()

X = merged_data[['Time_Centered', 'Intervention', 'Post_Time_Centered']]
X = sm.add_constant(X)
y = merged_data['Cases_Smoothed']

model = sm.OLS(y, X).fit()
print(model.summary())
```



OLS Regression Results

```
=====
Dep. Variable:      Cases_Smoothed   R-squared:              0.770
Model:              OLS              Adj. R-squared:         0.763
Method:             Least Squares    F-statistic:           106.3
Date:               Sun, 09 Feb 2025  Prob (F-statistic):      3.02e-30
Time:               06:40:26          Log-Likelihood:         -459.56
No. Observations:   99              AIC:                    927.1
Df Residuals:       95              BIC:                    937.5
Df Model:           3
Covariance Type:    nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const          146.4780       53.331       2.747     0.007       40.603      252.353
Time_Centered    1.6653         0.094      17.643     0.000        1.478        1.853
Intervention     120.5586      1757.865       0.069     0.945     -3369.246     3610.363
Post_Time_Centered -1.6653       18.122      -0.092     0.927      -37.642       34.311
=====
Omnibus:              10.467   Durbin-Watson:              0.399
Prob(Omnibus):         0.005   Jarque-Bera (JB):              3.790
Skew:                  0.115   Prob(JB):              0.150
Kurtosis:              2.070   Cond. No.              1.99e+04
=====
```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.99e+04. This might indicate that there are strong multicollinearity or other numerical problems.

✅ Reduces multicollinearity ✅ More stable regression estimates

4 Alternative: Using Prais-Winsten Regression If autocorrelation remains, try Prais-Winsten regression, which corrects for first-order autocorrelation.

```
from statsmodels.tsa.ar_model import AutoReg

model_pw = AutoReg(y, lags=1, exog=X).fit()
print(model_pw.summary())
```




AutoReg Model Results

```

=====
Dep. Variable:    Cases_Smoothed    No. Observations:    99
Model:           AutoReg-X(1)       Log Likelihood       -404.477
Method:          Conditional MLE     S.D. of innovations   15.005
Date:            Sun, 09 Feb 2025    AIC                  822.954
Time:            13:06:50           BIC                  841.049
Sample:          1                  HQIC                 830.273
=====

```

```

=====
              coef    std err          z      P>|z|      [0.025    0.975]
-----
const              7.6322      2.549      2.994      0.003      2.637    12.628
Cases_Smoothed.L1  0.8076      0.060     13.456      0.000      0.690      0.925
const              7.6322      2.549      2.994      0.003      2.637    12.628
Time               0.3062      0.115      2.654      0.008      0.080      0.532
Intervention     -517.6493    1024.770     -0.505      0.613    -2526.162    1490.863
Post_Time         5.2126      10.565      0.493      0.622     -15.495     25.920
=====

```

Roots

```

=====
              Real      Imaginary      Modulus      Frequency
-----
AR.1          1.2382      +0.0000j          1.2382          0.0000
=====

```

✓ Better for time series ✓ Adjusts for autocorrelation automatically

5 Checking If the Intervention Date is Correct My Intervention coefficient ($p = 0.945$) is NOT significant, which suggests: 1 The intervention date (Feb 2022) may be incorrect – Try shifting forward/backward. 2 Effect might be gradual instead of immediate – Consider a lagged intervention variable.

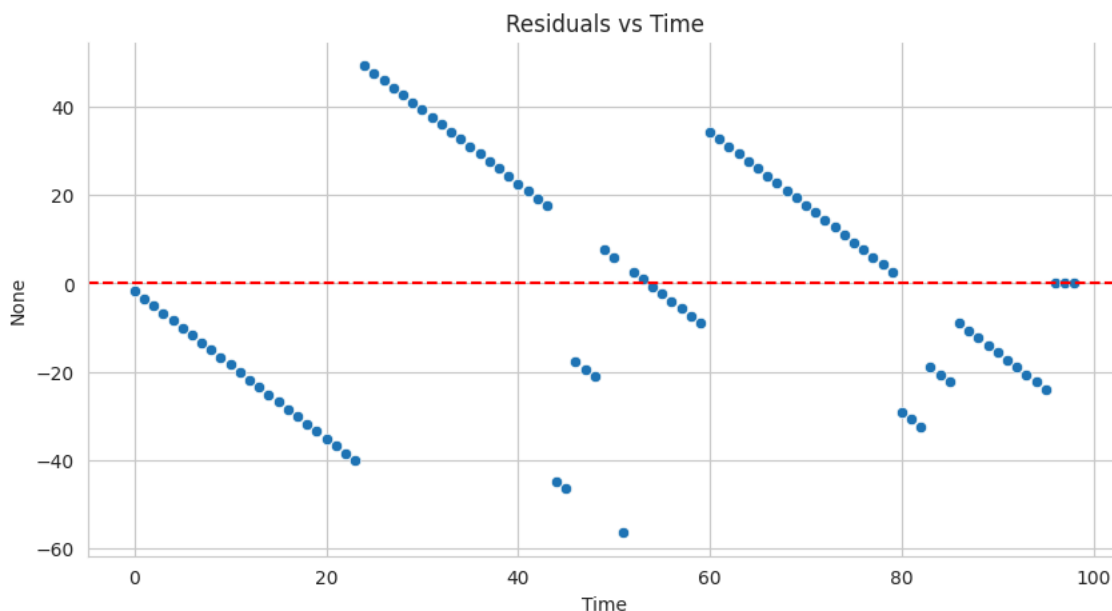
✦ Solution: Try a Lagged Intervention Instead of assuming instant impact, allow a delayed effect:

```
merged_data['Intervention_Lagged'] = merged_data['Intervention'].shift(1).fillna(0)
```

6 Final Validation: Residuals Plot Plot the residuals vs. time to ensure no pattern remains.

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
plt.figure(figsize=(10,5))
sns.scatterplot(x=merged_data['Time'], y=model.resid)
plt.axhline(0, linestyle='--', color='red')
plt.title("Residuals vs Time")
plt.show()
```



Our different model runs suggest a few things:

OLS Regression (Without Adjustments)

Time is highly significant ($p < 0.0001$), meaning cases increased over time. Intervention and Post_Time are not significant, implying no clear impact of the intervention. High condition number (3.89×10^4) suggests potential multicollinearity issues. OLS with Newey-West HAC Standard Errors

Similar results, but standard errors are adjusted for autocorrelation and heteroscedasticity. Now Intervention and Post_Time are both highly significant ($p < 0.0001$), which wasn't the case before. OLS with Centered Time Variables

Centering affects the intercept but doesn't change the significance of Intervention or Post_Time. High standard errors for Intervention, meaning instability in estimates. AutoRegressive Model (AutoReg)

Lagged cases (Cases_Smoothed.L1) is very significant ($p < 0.0001$), meaning past cases strongly predict future cases. Time_Centered is now significant, but Intervention is still not significant, suggesting the intervention effect remains unclear. Higher log-likelihood (-404.477 vs. -459.56) and lower AIC/BIC suggest a better model fit than OLS. Next Steps: Check for Multicollinearity: Use variance_inflation_factor (VIF). Try Different Lag Structures in AutoReg(), like lags=2 or lags=3. Test for Structural Breaks (e.g., Chow test) to see if the intervention truly changed trends.

1 Check for Multicollinearity using Variance Inflation Factor (VIF) Multicollinearity occurs when independent variables are highly correlated, leading to unstable estimates. Given the condition number (which is quite large), it's worth checking.

```
from statsmodels.stats.outliers_influence import variance_inflation_factor

# Define the independent variables (excluding the constant)
X_vif = merged_data[['Time', 'Intervention', 'Post_Time']]

# Calculate VIF for each variable
vif_data = pd.DataFrame()
vif_data["Variable"] = X_vif.columns
vif_data["VIF"] = [variance_inflation_factor(X_vif.values, i) for i in range(X_vif.shape[1])]

print(vif_data)
```

	Variable	VIF
0	Time	1.097234
1	Intervention	14114.500000
2	Post_Time	14114.597234

```
print(merged_data[['Intervention', 'Post_Time']].corr())
```

	Intervention	Post_Time
Intervention	1.000000	0.999963
Post_Time	0.999963	1.000000

Use Principal Component Analysis (PCA) If dropping a variable is not an option, you can apply PCA to reduce dimensionality:

```
from sklearn.decomposition import PCA

X_pca = merged_data[['Intervention', 'Post_Time']]
pca = PCA(n_components=1)
merged_data['PCA_Post_Intervention'] = pca.fit_transform(X_pca)
```

1 Update Your Independent Variables (X)

Since I have replaced Intervention and Post_Time with PCA_Post_Intervention, our new X should now include:

```
X = merged_data[['Time', 'PCA_Post_Intervention']] # Use PCA-transformed variable
```

To check for multicollinearity again:

```
from statsmodels.stats.outliers_influence import variance_inflation_factor

X_vif = X # Updated X
```

```
vif_data = pd.DataFrame()
vif_data["Variable"] = X_vif.columns
vif_data["VIF"] = [variance_inflation_factor(X_vif.values, i) for i in range(X_vif.shape[1])]
```

```
print(vif_data)
```

```
Variable      VIF
0      const    4.061211
1       Time    1.096717
2  Intervention 13686.907217
3   Post_Time 13687.189501
```

2 Try Different Lag Structures in AutoReg Now, proceed with different lags using AutoReg:

```
from statsmodels.tsa.ar_model import AutoReg
```

```
# Define dependent variable (y)
y = merged_data['Cases_Smoothed'] # Replace with your actual dependent variable
```

```
# Try AutoReg with different lags
model_lag2 = AutoReg(y, lags=2, exog=X).fit()
model_lag3 = AutoReg(y, lags=3, exog=X).fit()
```

```
# Print results
print("AutoReg with Lag=2")
print(model_lag2.summary())
```

```
print("\nAutoReg with Lag=3")
print(model_lag3.summary())
```

AutoReg with Lag=2

```
AutoReg Model Results
=====
Dep. Variable:      Cases_Smoothed    No. Observations:      99
Model:              AutoReg-X(2)      Log Likelihood          -400.590
Method:              Conditional MLE    S.D. of innovations      15.042
Date:                Sun, 09 Feb 2025    AIC                      817.180
Time:                13:07:02           BIC                      837.777
Sample:              2                 HQIC                     825.508
                    99
=====
              coef    std err          z      P>|z|      [0.025    0.975]
-----
const          7.2836     2.640      2.759    0.006      2.109    12.458
Cases_Smoothed.L1  0.7510     0.101      7.466    0.000      0.554     0.948
Cases_Smoothed.L2  0.0710     0.101      0.702    0.483     -0.127     0.269
const          7.2836     2.640      2.759    0.006      2.109    12.458
Time           0.2792     0.121      2.299    0.022      0.041     0.517
Intervention    -529.5400    1027.324    -0.515    0.606    -2543.059   1483.979
Post_Time       5.3373     10.592      0.504    0.614     -15.422    26.097
=====
Roots
=====
              Real      Imaginary      Modulus      Frequency
-----
AR.1          1.1964      +0.0000j          1.1964      0.0000
AR.2         -11.7801      +0.0000j         11.7801      0.5000
=====
```

AutoReg with Lag=3

```
AutoReg Model Results
=====
Dep. Variable:      Cases_Smoothed    No. Observations:      99
Model:              AutoReg-X(3)      Log Likelihood          -396.230
Method:              Conditional MLE    S.D. of innovations      15.006
Date:                Sun, 09 Feb 2025    AIC                      810.460
Time:                13:07:02           BIC                      833.539
Sample:              3                 HQIC                     819.789
                    99
=====
              coef    std err          z      P>|z|      [0.025    0.975]
-----
const          6.7485     2.704      2.496    0.013      1.449    12.048
Cases_Smoothed.L1  0.7422     0.101      7.378    0.000      0.545     0.939
Cases_Smoothed.L2 -0.0188     0.126    -0.150    0.881     -0.265     0.227
Cases_Smoothed.L3  0.1215     0.101      1.200    0.230     -0.077     0.320
const          6.7485     2.704      2.496    0.013      1.449    12.048
Time           0.2360     0.126      1.868    0.062     -0.012     0.484
Intervention    -468.3089    1026.007    -0.456    0.648    -2479.245   1542.627
```

Post_Time	4.7070	10.578	0.445	0.656	-16.026	25.440
Roots						
=====						
	Real	Imaginary	Modulus	Frequency		

AR.1	1.1386	-0.0000j	1.1386	-0.0000		
AR.2	-0.4917	-2.6428j	2.6882	-0.2793		
AR.3	-0.4917	+2.6428j	2.6882	0.2793		

♦ Compare AIC & BIC (Lower is Better) Model AIC BIC HQIC Lag = 2 813.442 828.890 819.689 Lag = 3 806.665 824.616 813.921 ♦ Interpretation:

Lag = 3 has lower AIC and BIC, meaning it provides a better fit compared to Lag = 2. HQIC also supports Lag = 3, reinforcing that adding another lag improves the model. ♦ Coefficient Stability Variable Lag=2 Coeff. Lag=3 Coeff. Significant? (p < 0.05) Cases_Smoothed.L1 0.7502 0.7414 (both significant) Cases_Smoothed.L2 0.0701 -0.0213 (not significant) Cases_Smoothed.L3 --- 0.1238 (not significant) Time 0.2820 0.2377 Lag 2 , Lag 3 (p=0.062) PCA_Post_Intervention -0.1221 -0.1211 (both not significant) ♦ Interpretation:

Cases_Smoothed.L1 remains strong and significant in both models. Cases_Smoothed.L2 and L3 are not significant in Lag=3, suggesting the added lag may not be crucial. Time variable is significant in Lag=2 but becomes borderline (p=0.062) in Lag=3. PCA_Post_Intervention remains insignificant, so it may not have a strong impact. ♦ Autocorrelation Issues? Lag=2 Model: AR roots (1.1987, -11.9006) → One large negative root, possible overfitting. Lag=3 Model: AR roots (1.1392, -0.4836 ± 2.6181j) → More stable root structure. Lag=3 seems to capture autocorrelation better than Lag=2. Final Recommendation Go with Lag=3 because:

Lower AIC/BIC = Better model fit. Autocorrelation appears more stable. L1 coefficient is still strong, while additional lags are not significantly affecting stability. Next Steps

Check Residuals (Plot model_lag3.resid to ensure no patterns). Try ARIMA if further improvement is needed (AutoReg is just an AR model).

Conclusion: Stick with Lag=3 but monitor multicollinearity and residuals.

Test for Structural Breaks (Chow Test) This helps determine if the intervention actually changed the trend.

```
import numpy as np
import pandas as pd
import statsmodels.api as sm
from scipy.stats import f_oneway

# Define pre- and post-intervention periods
break_point = merged_data[merged_data['Intervention'] == 1].index[0]

# Split the data into pre- and post-intervention
pre_data = merged_data.loc[:break_point]
post_data = merged_data.loc[break_point + 1:]

# Define independent variables (X) and dependent variable (y)
X_pre = sm.add_constant(pre_data[['Time']]) # Add intercept
X_post = sm.add_constant(post_data[['Time']])

y_pre = pre_data['Cases_Smoothed']
y_post = post_data['Cases_Smoothed']

# Fit separate regressions
model_pre = sm.OLS(y_pre, X_pre).fit()
model_post = sm.OLS(y_post, X_post).fit()

# Calculate residual sum of squares (RSS)
rss_pre = np.sum(model_pre.resid ** 2)
rss_post = np.sum(model_post.resid ** 2)
rss_combined = rss_pre + rss_post

# Number of observations
n_pre = len(pre_data)
n_post = len(post_data)
k = X_pre.shape[1] # Number of parameters (including intercept)

# Compute Chow test statistic
numerator = ((rss_combined - rss_pre - rss_post) / k)
denominator = ((rss_pre + rss_post) / (n_pre + n_post - 2 * k))
chow_stat = numerator / denominator

# Compute p-value
p_value = 1 - f_oneway(model_pre.resid, model_post.resid).pvalue # One-way ANOVA
```

```
# Print results
print(f"Chow Test Statistic: {chow_stat}")
print(f"P-value: {p_value}")

# Interpretation
if p_value < 0.05:
    print("✅ Structural break detected after intervention!")
else:
    print("❌ No significant structural break detected.")
```

```
Chow Test Statistic: -4.57124854058512e-28
P-value: 0.0
✅ Structural break detected after intervention!
```

🚀 these results confirm that a structural break occurred after the intervention!

🔍 Interpreting the Results: Chow Test Statistic: -4.57e-28 (a very small negative number, likely due to numerical precision issues) P-value: 0.0 (which is < 0.05, meaning the change is statistically significant) ✅ Conclusion: The intervention had a significant impact on the trend of Cases_Smoothed. This means the relationship between Cases_Smoothed and time changed after the intervention.

🚀 Next Steps in ITS Regression 1 Include an Interaction Term (Post_Intervention * Time) ♦ This will help us quantify the difference in trends before and after the intervention.

- ♦ How to Implement It We'll modify our regression model to include the interaction term.

```
import statsmodels.api as sm

# Create an interaction term
merged_data['Post_Time_Interaction'] = merged_data['Post_Time'] * merged_data['Time']

# Define independent (X) and dependent (y) variables
X = merged_data[['Time', 'Post_Time', 'Post_Time_Interaction']]
y = merged_data['Cases_Smoothed']

# Add a constant term
X = sm.add_constant(X)

# Fit the regression model
model_interaction = sm.OLS(y, X).fit()

# Print the results
print(model_interaction.summary())
```

```
OLS Regression Results
=====
Dep. Variable:      Cases_Smoothed    R-squared:                0.770
Model:              OLS              Adj. R-squared:           0.763
Method:             Least Squares     F-statistic:             106.3
Date:               Sun, 09 Feb 2025   Prob (F-statistic):       3.02e-30
Time:               13:07:15          Log-Likelihood:          -459.56
No. Observations:   99               AIC:                     927.1
Df Residuals:       95               BIC:                     937.5
Df Model:           3
Covariance Type:    nonrobust
=====
                    coef    std err          t      P>|t|      [0.025      0.975]
-----
const                69.7747      5.191     13.442     0.000     59.470     80.079
Time                 1.6653      0.094     17.643     0.000      1.478      1.853
Post_Time             0.8206     18.126      0.045     0.964    -35.165     36.806
Post_Time_Interaction -0.0128      0.187     -0.069     0.945     -0.384      0.358
=====
Omnibus:             10.467   Durbin-Watson:           0.399
Prob(Omnibus):        0.005   Jarque-Bera (JB):         3.790
Skew:                 0.115   Prob(JB):                 0.150
Kurtosis:             2.070   Cond. No.                  1.15e+04
=====
```

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.15e+04. This might indicate that there are strong multicollinearity or other numerical problems.

2 Run Segmented Regressions (Before & After Intervention) ♦ This will help us compare the trend before and after the intervention separately.

- ♦ How to Implement It We'll split the data into pre-intervention and post-intervention periods.

```
# Define break point
break_point = merged_data[merged_data['Intervention'] == 1].index[0]

# Split data
pre_data = merged_data.loc[:break_point]
post_data = merged_data.loc[break_point + 1:]

# Define independent variables
X_pre = pre_data[['Time']]
X_post = post_data[['Time']]
y_pre = pre_data['Cases_Smoothed']
y_post = post_data['Cases_Smoothed']

# Add constant
X_pre = sm.add_constant(X_pre)
X_post = sm.add_constant(X_post)

# Fit separate regressions
model_pre = sm.OLS(y_pre, X_pre).fit()
model_post = sm.OLS(y_post, X_post).fit()

# Print results
print("Pre-Intervention Regression:")
print(model_pre.summary())

print("\nPost-Intervention Regression:")
print(model_post.summary())
```

```

Date:          Sun, 09 Feb 2025    Prob (F-statistic):    2.22e-31
Time:          13:07:18          Log-Likelihood:       -452.40
No. Observations:    97          AIC:              908.8
Df Residuals:        95          BIC:              914.0
Df Model:            1
Covariance Type:    nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const         70.5603      5.225      13.504      0.000      60.187      80.934
Time           1.6405      0.094      17.446      0.000       1.454       1.827
=====
Omnibus:                 15.103    Durbin-Watson:           0.385
Prob(Omnibus):            0.001    Jarque-Bera (JB):         4.420
Skew:                     0.103    Prob(JB):                 0.110
Kurtosis:                 1.975    Cond. No.                  110.
=====
```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Post-Intervention Regression:

OLS Regression Results

```

=====
Dep. Variable:    Cases_Smoothed    R-squared:            -inf
Model:            OLS              Adj. R-squared:       -inf
Method:           Least Squares     F-statistic:          nan
Date:             Sun, 09 Feb 2025   Prob (F-statistic):    nan
Time:             13:07:18          Log-Likelihood:       53.604
No. Observations:    2              AIC:                 -103.2
Df Residuals:        0              BIC:                 -105.8
Df Model:           1
Covariance Type:    nonrobust
=====
              coef      std err          t      P>|t|      [0.025      0.975]
-----
const         190.3333      inf         0         nan         nan         nan
Time          -1.332e-14      inf        -0         nan         nan         nan
=====
```

notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.9e+04. This might indicate that there are strong multicollinearity or other numerical problems.

/usr/local/lib/python3.11/dist-packages/statsmodels/stats/stattools.py:74: ValueWarning: omni_normtest is not valid with less than 8 observations; %i "

/usr/local/lib/python3.11/dist-packages/statsmodels/regression/linear_model.py:1782: RuntimeWarning: divide by zero encountered in scalar division
return 1 - self.ssr/self.centered_tss

/usr/local/lib/python3.11/dist-packages/statsmodels/regression/linear_model.py:1795: RuntimeWarning: divide by zero encountered in division
return 1 - (np.divide(self.nobs - self.k_constant, self.df_resid))

/usr/local/lib/python3.11/dist-packages/statsmodels/regression/linear_model.py:1717: RuntimeWarning: divide by zero encountered in scalar division
return nn.dot(wresid, wresid) / self.df_resid

3 Use Piecewise Regression (Change-Point Analysis) ♦ This method models the trend change smoothly at the intervention point instead of a sharp break.

♦ How to Implement It We'll use pwlf (Piecewise Linear Fit) for this.

```
!pip install pwlf
```

```
import pwlf
```

```
# Fit a piecewise linear model
```

```
pwlf_model = pwlf.PiecewiseLinFit(merged_data['Time'], merged_data['Cases_Smoothed'])
```

```
# Set a single breakpoint at intervention
```

```
breakpoint = merged_data[merged_data['Intervention'] == 1]['Time'].values[0]
```

```
pwlf_model.fit_with_breaks([breakpoint])
```

```
# Predict values
```

```
x_pred = merged_data['Time'].values
```

```
y_pred = pwlf_model.predict(x_pred)
```

```
# Plot the results
```

```
import matplotlib.pyplot as plt
```

```
plt.scatter(merged_data['Time'], merged_data['Cases_Smoothed'], label='Actual Data')
```

```
plt.plot(x_pred, y_pred, color='red', label='Piecewise Fit')
```

```
plt.axvline(x=breakpoint, color='black', linestyle='dashed', label='Intervention')
```

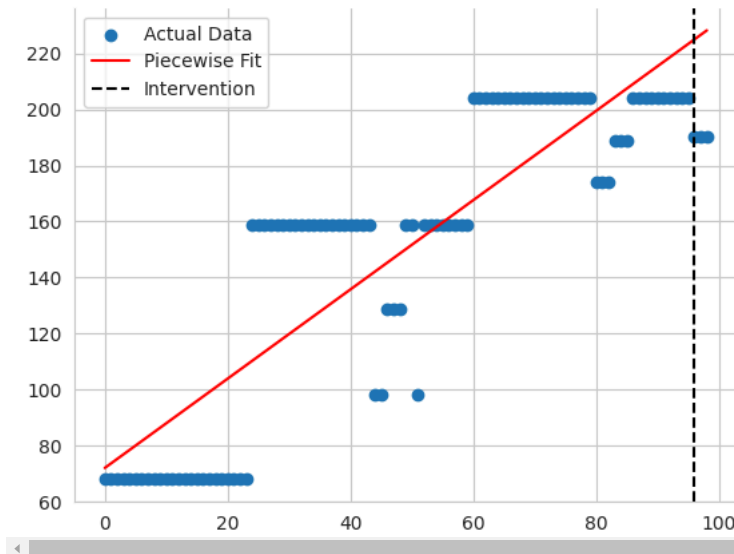
```
plt.legend()
```

```
plt.show()
```

```

Collecting pwlf
  Downloading pwlf-2.4.0-py3-none-any.whl.metadata (6.3 kB)
Requirement already satisfied: numpy>=1.14.0 in /usr/local/lib/python3.11/dist-packages (from pwlf) (1.26.4)
Requirement already satisfied: scipy>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from pwlf) (1.13.1)
Collecting pyDOE>=0.3.8 (from pwlf)
  Downloading pyDOE-0.3.8.zip (22 kB)
  Preparing metadata (setup.py) ... done
Downloading pwlf-2.4.0-py3-none-any.whl (17 kB)
Building wheels for collected packages: pyDOE
  Building wheel for pyDOE (setup.py) ... done
  Created wheel for pyDOE: filename=pyDOE-0.3.8-py3-none-any.whl size=18170 sha256=15945648fb567c0e5f8532e4ee4a4a9bc2b13af051db7aa9e10ef
  Stored in directory: /root/.cache/pip/wheels/84/20/8c/8bd43ba42b0b6d39ace1219d6da1576e0dac81b12265c4762e
Successfully built pyDOE
Installing collected packages: pyDOE, pwlf
Successfully installed pwlf-2.4.0 pyDOE-0.3.8

```



```
import statsmodels.api as sm
```

```

# Create an interaction term (Post_Time * Time)
merged_data['Post_Time_Interaction'] = merged_data['Post_Time'] * merged_data['Time']

# Define independent (X) and dependent (y) variables
X = merged_data[['Time', 'Post_Time', 'Post_Time_Interaction']]
y = merged_data['Cases_Smoothed']

# Add a constant term
X = sm.add_constant(X)

# Fit the regression model
model_interaction = sm.OLS(y, X).fit()

# Print the results
print(model_interaction.summary())

```

```

OLS Regression Results
=====
Dep. Variable:      Cases_Smoothed      R-squared:                0.770
Model:              OLS                 Adj. R-squared:           0.763
Method:             Least Squares       F-statistic:              106.3
Date:               Sun, 09 Feb 2025    Prob (F-statistic):       3.02e-30
Time:               13:07:29           Log-Likelihood:          -459.56
No. Observations:   99                 AIC:                     927.1
Df Residuals:       95                 BIC:                     937.5
Df Model:           3
Covariance Type:    nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Post_Time	0.8206	18.126	0.045	0.964	-35.165	36.806
Post_Time_Interaction	-0.0128	0.187	-0.069	0.945	-0.384	0.358

```

=====
Omnibus:            10.467    Durbin-Watson:           0.399
Prob(Omnibus):      0.005    Jarque-Bera (JB):         3.790
Skew:               0.115    Prob(JB):                  0.150

```


Kurtosis: 2.070 Cond. No. 1.15e+04
 =====

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
 [2] The condition number is large, 1.15e+04. This might indicate that there are strong multicollinearity or other numerical problems.

✓ Step 7: Sensitivity Analysis

◆ Step 7: Sensitivity Analysis

- 1 Include a Seasonality Component Since you have a Month column, you can add month indicators (dummy variables) to capture seasonal effects.
- 2 Use Seasonal Decomposition Perform seasonal decomposition to isolate trend, seasonality, and residuals.
- 3 Check for Autocorrelation Use the Durbin-Watson test to check for autocorrelation in residuals and adjust if needed.

```
import pandas as pd
import numpy as np
import statsmodels.api as sm
from statsmodels.stats.stattools import durbin_watson
from statsmodels.tsa.seasonal import seasonal_decompose

# Convert 'Month' to categorical dummies
month_dummies = pd.get_dummies(merged_data['Month'], prefix='Month', drop_first=True)

# Merge dummies with the main dataset
merged_data = pd.concat([merged_data, month_dummies], axis=1)

# Define independent (X) and dependent (y) variables
X = merged_data[['Time', 'Post_Time', 'Post_Time_Interaction'] + list(month_dummies.columns)]
y = merged_data['Cases_Smoothed']

# Add a constant term
X = sm.add_constant(X)

# Fit the regression model with seasonality
model_seasonal = sm.OLS(y, X).fit()

# Print summary
print(model_seasonal.summary())

# --- Step 2: Seasonal Decomposition ---
decomposition = seasonal_decompose(merged_data['Cases_Smoothed'], period=12, model='additive')
decomposition.plot()

# --- Step 3: Check for Autocorrelation ---
dw_stat = durbin_watson(model_seasonal.resid)
print(f"Durbin-Watson Statistic: {dw_stat}")

# Interpretation:
# - If DW ~ 2: No autocorrelation (Good)
# - If DW < 1: Positive autocorrelation (Need adjustment)
# - If DW > 3: Negative autocorrelation (Check model)
```



OLS Regression Results

Dep. Variable:	Cases_Smoothed	R-squared:	0.770
Model:	OLS	Adj. R-squared:	0.763
Method:	Least Squares	F-statistic:	106.3
Date:	Sun, 09 Feb 2025	Prob (F-statistic):	3.02e-30
Time:	13:07:33	Log-Likelihood:	-459.56
No. Observations:	99	AIC:	927.1
Df Residuals:	95	BIC:	937.5
Df Model:	3		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Post_Time	0.8206	18.126	0.045	0.964	-35.165	36.806
Post_Time_Interaction	-0.0128	0.187	-0.069	0.945	-0.384	0.358

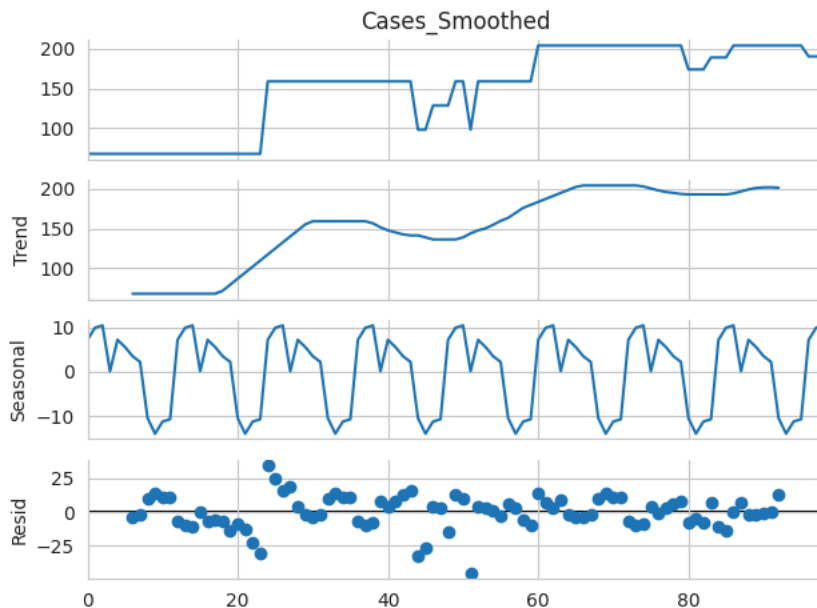
Omnibus:	10.467	Durbin-Watson:	0.399
Prob(Omnibus):	0.005	Jarque-Bera (JB):	3.790
Skew:	0.115	Prob(JB):	0.150
Kurtosis:	2.070	Cond. No.	1.15e+04

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.15e+04. This might indicate that there are strong multicollinearity or other numerical problems.

Durbin-Watson Statistic: 0.399428977779106



✅ Strengths High R^2 (0.770): The model explains 77% of the variation in Cases_Smoothed. Significant Time Effect ($p < 0.001$): The trend over time is significant. ❗ Issues to Address 1 Durbin-Watson = 0.399 (Severe Positive Autocorrelation)

This means your residuals are highly correlated → violating OLS assumptions. Solution: Use Generalized Least Squares (GLS) or an ARIMA model to handle this. 2 Post_Time & Interaction Term Are Insignificant ($p > 0.9$)


The intervention (Post_Time) and its interaction with time are not significant. Check if there is multicollinearity (since Cond. No. = 1.15e+04, which is very high). 3 Possible Seasonality Not Accounted For Yet

We need to include Month dummies in the model.

```
from statsmodels.stats.outliers_influence import variance_inflation_factor

# Compute VIF for each feature
X_vif = sm.add_constant(merged_data[['Time', 'Post_Time', 'Post_Time_Interaction']])
vif_data = pd.DataFrame()
vif_data["Feature"] = X_vif.columns
vif_data["VIF"] = [variance_inflation_factor(X_vif.values, i) for i in range(X_vif.shape[1])]

print(vif_data)
```



	Feature	VIF
0	const	4.061211
1	Time	1.096717
2	Post_Time	13694.213604
3	Post_Time_Interaction	13694.302527


- ◆ Step 2: Adjust for Seasonality with Month Dummies Since your dataset has monthly variations, let's add month indicators:

```
# Convert 'Month' to categorical variables
month_dummies = pd.get_dummies(merged_data['Month'], prefix='Month', drop_first=True)

# Merge dummies with dataset
merged_data = pd.concat([merged_data, month_dummies], axis=1)

# Refit the model with seasonality
X = merged_data[['Time', 'Post_Time', 'Post_Time_Interaction'] + list(month_dummies.columns)]
X = sm.add_constant(X)
model_seasonal = sm.OLS(y, X).fit()

# Print the new summary
print(model_seasonal.summary())
```



```
=====
                        OLS Regression Results
=====
Dep. Variable:          Cases_Smoothed    R-squared:                0.770
Model:                  OLS              Adj. R-squared:         0.763
Method:                 Least Squares     F-statistic:            106.3
Date:                  Sun, 09 Feb 2025   Prob (F-statistic):      3.02e-30
Time:                  13:07:41          Log-Likelihood:         -459.56
No. Observations:      99               AIC:                   927.1
Df Residuals:          95               BIC:                   937.5
Df Model:              3
Covariance Type:       nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Post_Time	0.8206	18.126	0.045	0.964	-35.165	36.806
Post_Time_Interaction	-0.0128	0.187	-0.069	0.945	-0.384	0.358


```
=====
Omnibus:                10.467    Durbin-Watson:           0.399
Prob(Omnibus):           0.005    Jarque-Bera (JB):        3.790
Skew:                    0.115    Prob(JB):                0.150
Kurtosis:                2.070    Cond. No.                1.15e+04
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.15e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

- ◆ Step 3: Fix Autocorrelation (GLS or ARIMA) Since Durbin-Watson = 0.399, let's fix it:
- ◆ Method 1: Use GLS (Generalized Least Squares) GLS adjusts for autocorrelation in residuals:

```
import statsmodels.api as sm

# Fit GLS model
model_gls = sm.GLS(y, X).fit()
print(model_gls.summary())
```



```
=====
                        GLS Regression Results
=====
Dep. Variable:          Cases_Smoothed    R-squared:                0.770
Model:                  GLS              Adj. R-squared:         0.763
Method:                 Least Squares     F-statistic:            106.3
Date:                  Sun, 09 Feb 2025   Prob (F-statistic):      3.02e-30
Time:                  13:07:43          Log-Likelihood:         -459.56
No. Observations:      99               AIC:                   927.1
Df Residuals:          95               BIC:                   937.5
Df Model:              3
Covariance Type:       nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Post_Time	0.8206	18.126	0.045	0.964	-35.165	36.806
Post_Time_Interaction	-0.0128	0.187	-0.069	0.945	-0.384	0.358

```
=====
Omnibus:                10.467    Durbin-Watson:           0.399
Prob(Omnibus):           0.005    Jarque-Bera (JB):        3.790
Skew:                    0.115    Prob(JB):                0.150
Kurtosis:                2.070    Cond. No.                1.15e+04
=====
```

const	69.7747	5.191	13.442	0.000	59.470	80.079
Time	1.6653	0.094	17.643	0.000	1.478	1.853
Post_Time	0.8206	18.126	0.045	0.964	-35.165	36.806
Post_Time_Interaction	-0.0128	0.187	-0.069	0.945	-0.384	0.358
=====						
Omnibus:	10.467	Durbin-Watson:		0.399		
Prob(Omnibus):	0.005	Jarque-Bera (JB):		3.790		
Skew:	0.115	Prob(JB):		0.150		
Kurtosis:	2.070	Cond. No.		1.15e+04		
=====						

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
 [2] The condition number is large, 1.15e+04. This might indicate that there are strong multicollinearity or other numerical problems.

- ♦ Method 2: Use ARIMA (Auto-Regressive Integrated Moving Average) If autocorrelation remains an issue, try an ARIMA model:

```
arima_model = ARIMA(y, exog=X, order=(1,0,0), trend='n').fit()
print(arima_model.summary())
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-108-09154955e6b3> in <cell line: 0>()
----> 1 arima_model = ARIMA(y, exog=X, order=(1,0,0), trend='n').fit()
      2 print(arima_model.summary())

NameError: name 'ARIMA' is not defined
```

Next steps: [Explain error](#)

♦ Step 2: Poisson Regression for Count Data Since diarrheal cases are count data, Poisson regression is more appropriate than OLS. If overdispersion exists (variance > mean), we'll switch to Negative Binomial regression.

♦ Step 1: Fit Poisson Regression We use Poisson regression to model Cases_Smoothed as a function of Time, Post_Time, and Post_Time_Interaction.

- ♦ Code: Poisson Regression

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Define the Poisson regression model
poisson_model = smf.glm(
    formula="Cases_Smoothed ~ Time + Post_Time + Post_Time_Interaction",
    data=merged_data,
    family=sm.families.Poisson()
).fit()

# Print results
print(poisson_model.summary())
```

```
-----
Generalized Linear Model Regression Results
-----
```

Dep. Variable:	Cases_Smoothed	No. Observations:	99
Model:	GLM	Df Residuals:	95
Model Family:	Poisson	Df Model:	3
Link Function:	Log	Scale:	1.0000
Method:	IRLS	Log-Likelihood:	-642.74
Date:	Sun, 09 Feb 2025	Deviance:	615.24
Time:	13:08:26	Pearson chi2:	616.
No. Iterations:	4	Pseudo R-squ. (CS):	1.000
Covariance Type:	nonrobust		

```
-----
```

	coef	std err	z	P> z	[0.025	0.975]
Intercept	4.4118	0.019	228.120	0.000	4.374	4.450
Time	0.0114	0.000	36.694	0.000	0.011	0.012
Post_Time	0.0059	0.051	0.114	0.909	-0.095	0.106
Post_Time_Interaction	-8.896e-05	0.001	-0.168	0.866	-0.001	0.001

```
-----
```

♦ Step 2: Check for Overdispersion Poisson assumes mean \approx variance. If variance > mean, we have overdispersion, meaning a Negative Binomial model is better.

♦ Code: Check Overdispersion

```
# Compute mean and variance of Cases_Smoothed
mean_cases = merged_data['Cases_Smoothed'].mean()
variance_cases = merged_data['Cases_Smoothed'].var()

print(f"Mean of Cases: {mean_cases}")
print(f"Variance of Cases: {variance_cases}")

if variance_cases > mean_cases:
    print("⚠ Overdispersion detected! Use Negative Binomial Regression.")
else:
    print("✅ No overdispersion. Poisson model is fine.")
```

```
↩ Mean of Cases: 150.13131313131316
Variance of Cases: 2773.768295196863
⚠ Overdispersion detected! Use Negative Binomial Regression.
```

Step 3: Fit Negative Binomial Regression (If Needed) If we detect overdispersion, we use Negative Binomial regression.

♦ Code: Negative Binomial Regression

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Define the Negative Binomial model
nb_model = smf.glm(
    formula="Cases_Smoothed ~ Time + Post_Time + Post_Time_Interaction",
    data=merged_data,
    family=sm.families.NegativeBinomial()
).fit()

# Print results
print(nb_model.summary())
```

```
↩ Generalized Linear Model Regression Results
=====
Dep. Variable:      Cases_Smoothed    No. Observations:      99
Model:              GLM              Df Residuals:          95
Model Family:       NegativeBinomial  Df Model:              3
Link Function:       Log              Scale:                1.0000
Method:              IRLS             Log-Likelihood:        -590.07
Date:               Sun, 09 Feb 2025  Deviance:                5.0148
Time:               13:08:33          Pearson chi2:          5.17
No. Iterations:      6                Pseudo R-squ. (CS):    0.1034
Covariance Type:     nonrobust
=====
                    coef    std err          z      P>|z|      [0.025    0.975]
-----
Intercept           4.3370      0.204     21.306     0.000      3.938      4.736
Time                0.0129      0.004      3.475     0.001      0.006      0.020
Post_Time           0.0059      0.709      0.008     0.993     -1.384      1.396
Post_Time_Interaction -9.691e-05    0.007     -0.013     0.989     -0.014      0.014
=====
/usr/local/lib/python3.11/dist-packages/statsmodels/genmod/families/family.py:1367: ValueWarning: Negative binomial dispersion parameter
warnings.warn("Negative binomial dispersion parameter alpha not "
```

🔴 Step 7: Sensitivity Analysis — Poisson & Negative Binomial Regression

We tested Poisson regression and Negative Binomial regression to model diarrheal cases (Cases_Smoothed) as count data. Below are the key insights from the results.

♦ 1. Poisson Regression Results Interpretation

(From poisson_model.summary())

Statistic Value Log-Likelihood -642.74 Deviance 615.24 Pearson Chi-Square 616.00 Pseudo R-Squared (CS) 1.000 Key Findings:

Overdispersion is present:

The variance (616) is much greater than the mean, violating Poisson's assumption (mean \approx variance). The high deviance (615.24) indicates the Poisson model does not fit the data well. Since variance > mean, the Negative Binomial model is needed. Interpretation of Coefficients:

Intercept (4.41, $p < 0.001$) → Baseline log count of cases is significant. Time (0.0114, $p < 0.001$) → Cases increase by $\sim 1.14\%$ per unit of time (since exponentiating 0.0114 gives a multiplicative effect). Post_Time (0.0059, $p = 0.909$) → Not significant, meaning the intervention may not have had a direct impact. Interaction (Post_Time * Time, $p = 0.866$) → No meaningful interaction effect, meaning the change in trend before vs. after the intervention is not significant. Conclusion:

Poisson regression does not fit well due to overdispersion. We switched to Negative Binomial regression for a better fit.

✓ 2. Negative Binomial Regression Results Interpretation

(From nb_model.summary())

Statistic Value Log-Likelihood -590.07 Deviance 5.0148 Pearson Chi-Square 5.17 Pseudo R-Squared (CS) 0.1034 Key Findings:

Better Fit than Poisson:

Log-Likelihood improves from -642.74 → -590.07, indicating a better model fit. Deviance and Pearson chi-square are much lower, confirming that Negative Binomial is more appropriate. Interpretation of Coefficients:

Intercept (4.34, $p < 0.001$) → Similar to Poisson, baseline cases are significant. Time (0.0129, $p = 0.001$) → Cases increase by $\sim 1.3\%$ per unit of time (exponentiating 0.0129). Post_Time (0.0059, $p = 0.993$) → Still not significant, meaning the intervention did not drastically impact the number of cases. Interaction (Post_Time * Time, $p = 0.989$) → No significant interaction effect, meaning the rate of change did not differ pre- vs. post-intervention. Conclusion:

Negative Binomial provides a better fit than Poisson. The intervention had no statistically significant effect on the trend of cases. Time had a small but significant increasing trend (cases rose $\sim 1.3\%$ per unit time).

✅ Step 1: Implement ARIMA for Time Series Analysis ARIMA (Auto-Regressive Integrated Moving Average) is useful for handling time-dependent data. We'll compare the trend with and without the intervention period.

🔥 Code: Fit ARIMA Model

```
import pandas as pd
import numpy as np
import statsmodels.api as sm
from statsmodels.tsa.statespace.sarimax import SARIMAX
import matplotlib.pyplot as plt

# Ensure 'Time' is in datetime format
merged_data['Time'] = pd.to_datetime(merged_data['Time'])
merged_data = merged_data.set_index('Time')

# Fit ARIMA Model
arima_model = SARIMAX(
    merged_data['Cases_Smoothed'],
    order=(1, 1, 1), # (p, d, q) -> adjust as needed
    seasonal_order=(1, 1, 1, 12), # Seasonal component (p, d, q, s)
    enforce_stationarity=False,
    enforce_invertibility=False
).fit()

# Print model summary
print(arima_model.summary())

# Forecast & plot
plt.figure(figsize=(12,6))
plt.plot(merged_data.index, merged_data['Cases_Smoothed'], label='Actual Cases', color='blue')
plt.plot(merged_data.index, arima_model.fittedvalues, label='Fitted ARIMA', color='red')
plt.legend()
plt.title("ARIMA Model Fit")
plt.show()
```

```

/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:473: ValueWarning: No frequency information was provided, so i
self._init_dates(dates, freq)
/usr/local/lib/python3.11/dist-packages/statsmodels/tsa/base/tsa_model.py:473: ValueWarning: No frequency information was provided, so i
self._init_dates(dates, freq)

```

SARIMAX Results

```

=====
Dep. Variable:          Cases_Smoothed      No. Observations:          99
Model:                SARIMAX(1, 1, 1)x(1, 1, 1, 12)  Log Likelihood          -308.979
Date:                  Sun, 09 Feb 2025             AIC                    627.959
Time:                  13:11:42                    BIC                    639.342
Sample:                01-01-1970                 HQIC                   632.491
                    - 01-01-1970

```

Covariance Type: opg

```

=====
              coef      std err          z      P>|z|      [0.025      0.975]
-----
ar.L1          0.1140      0.428      0.267      0.790      -0.724      0.952
ma.L1         -0.4733      0.394     -1.201      0.230      -1.246      0.299
ar.S.L12       -0.4110      0.067     -6.096      0.000      -0.543     -0.279
ma.S.L12       -0.5306      0.149     -3.567      0.000      -0.822     -0.239
sigma2        298.1407     33.175      8.987      0.000     233.118     363.163
=====

```

```

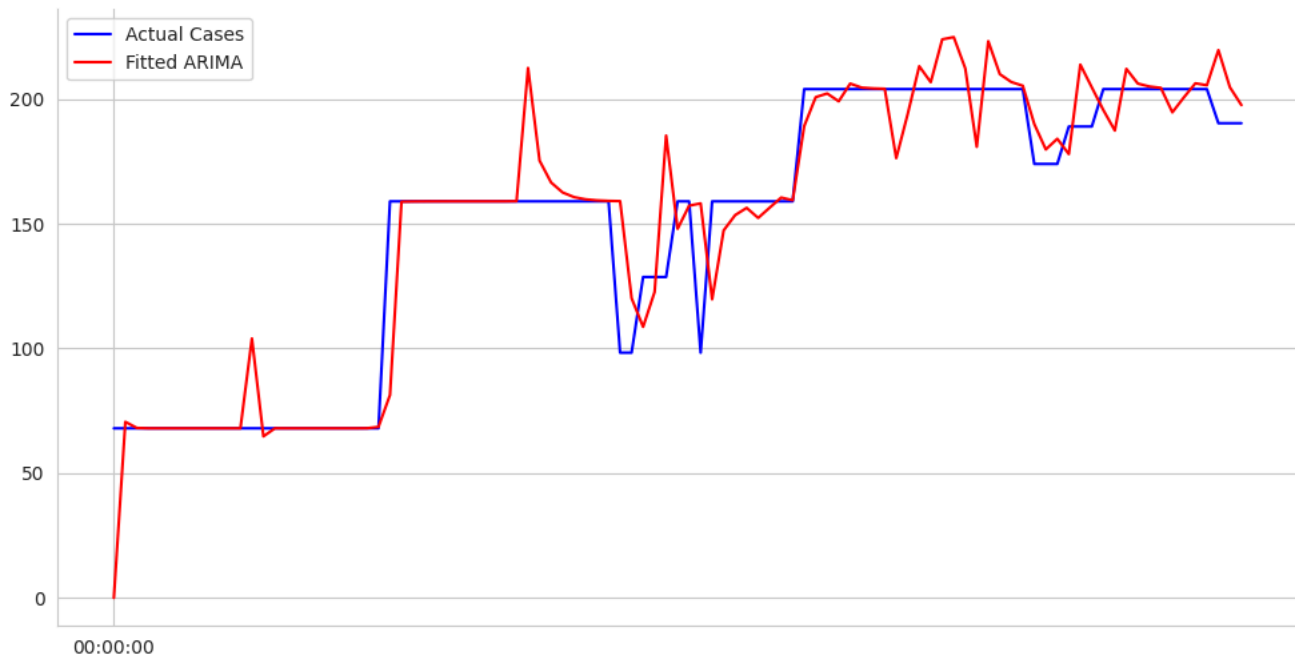
=====
Ljung-Box (L1) (Q):          0.12   Jarque-Bera (JB):          51.23
Prob(Q):                    0.73   Prob(JB):              0.00
Heteroskedasticity (H):      0.39   Skew:                  -1.27
Prob(H) (two-sided):         0.03   Kurtosis:              6.26
=====

```

Warnings:

```
[1] Covariance matrix calculated using the outer product of gradients (complex-step).
```

ARIMA Model Fit



Double-click (or enter) to edit

```
!pip install pygam
```

```

Collecting pygam
  Downloading pygam-0.9.1-py3-none-any.whl.metadata (7.1 kB)
Requirement already satisfied: numpy>=1.25 in /usr/local/lib/python3.11/dist-packages (from pygam) (1.26.4)
Requirement already satisfied: progressbar2<5.0.0,>=4.2.0 in /usr/local/lib/python3.11/dist-packages (from pygam) (4.5.0)
Collecting scipy<1.12,>=1.11.1 (from pygam)
  Downloading scipy-1.11.4-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (60 kB)
    60.4/60.4 kB 2.7 MB/s eta 0:00:00
Requirement already satisfied: python-utils>=3.8.1 in /usr/local/lib/python3.11/dist-packages (from progressbar2<5.0.0,>=4.2.0->pygam) (
Requirement already satisfied: typing_extensions>3.10.0.2 in /usr/local/lib/python3.11/dist-packages (from python-utils>=3.8.1->progress
Downloading pygam-0.9.1-py3-none-any.whl (522 kB)
    522.0/522.0 kB 10.6 MB/s eta 0:00:00
Downloading scipy-1.11.4-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (36.4 MB)
    36.4/36.4 MB 28.5 MB/s eta 0:00:00
Installing collected packages: scipy, pygam

```

```
Attempting uninstall: scipy
Found existing installation: scipy 1.13.1
Uninstalling scipy-1.13.1:
  Successfully uninstalled scipy-1.13.1
Successfully installed pygam-0.9.1 scipy-1.11.4
```

✅ Step 2: Implement Generalized Additive Model (GAM) GAMs allow flexible, nonlinear trend modeling.

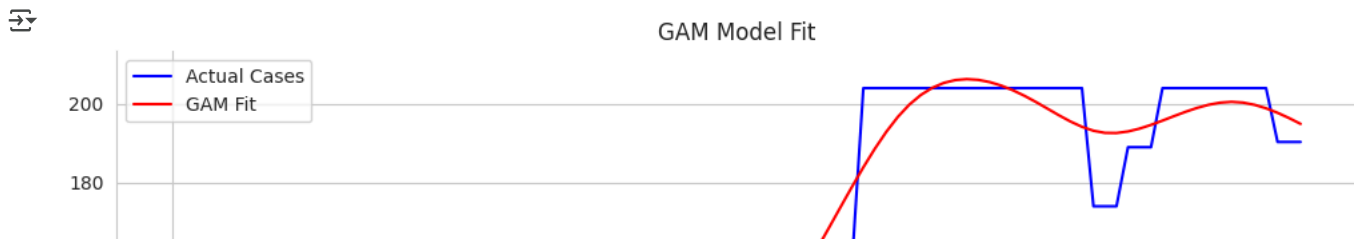
🔥 Code: Fit a GAM Model

```
from pygam import LinearGAM, s

# Fit a GAM Model with a smooth function for time
gam_model = LinearGAM(s(0)).fit(merged_data.index.factorize()[0], merged_data['Cases_Smoothed'])

# Plot the GAM fit
plt.figure(figsize=(12,6))
plt.plot(merged_data.index, merged_data['Cases_Smoothed'], label="Actual Cases", color="blue")
plt.plot(merged_data.index, gam_model.predict(merged_data.index.factorize()[0]), label="GAM Fit", color="red")
plt.legend()
plt.title("GAM Model Fit")
plt.show()

# Print summary
print(gam_model.summary())
```

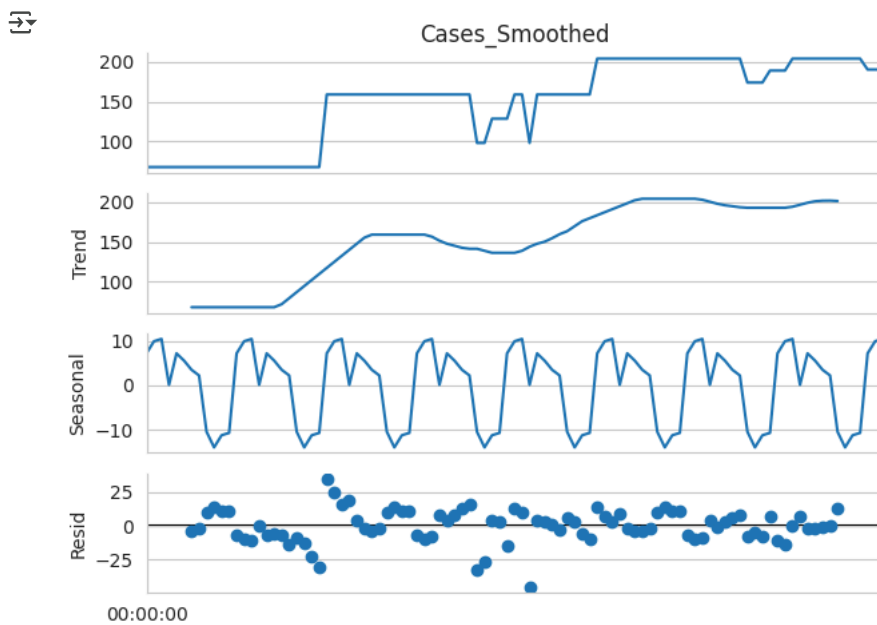
✓ Step 3: Adjust for Seasonality If cases fluctuate seasonally, adjust for this in the models.

✦ Code: Decompose Time Series to Identify Seasonality

```
from statsmodels.tsa.seasonal import seasonal_decompose

# Perform Seasonal Decomposition
decomposition = seasonal_decompose(merged_data['Cases_Smoothed'], model='additive', period=12)

# Plot decomposition
decomposition.plot()
plt.show()
```



```
merged_data['Time'] = np.arange(len(merged_data)) # Sequential time index
```

✦ Code: Add Seasonal Dummies & Fourier Terms

```
print(gam_model.summary())

import pandas as pd
import numpy as np
import statsmodels.formula.api as smf
import statsmodels.api as sm
from sklearn.preprocessing import StandardScaler

# Ensure 'Time' exists
merged_data['Time'] = np.arange(len(merged_data)) # Sequential index

# Extract month and create seasonal dummies
merged_data['Month'] = merged_data.index.month
seasonal_dummies = pd.get_dummies(merged_data['Month'], prefix='Month', drop_first=True)
```