



Intro

Data Preprocessing

Feature Engineering

EDA

SQL

ML Model

Conclusion

BUSN 32120

COVID-19 U.S. Data Analysis

Final Project

By: Bisma Rana



Intro

Data Preprocessing

Feature Engineering

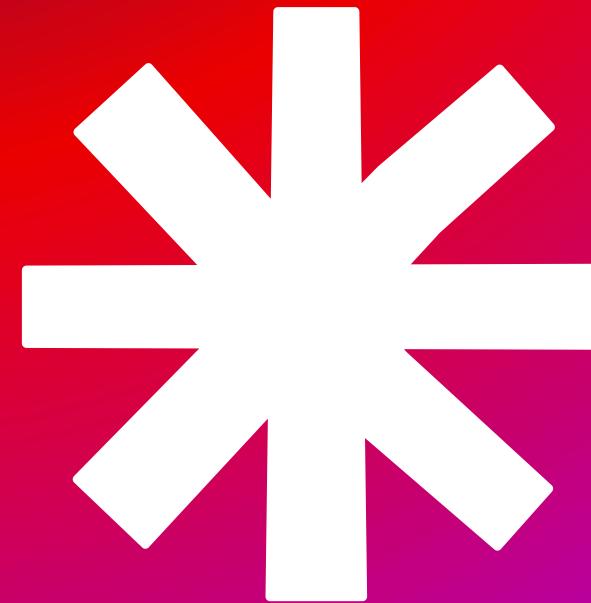
EDA

SQL

ML Model

Conclusion

Table of Contents

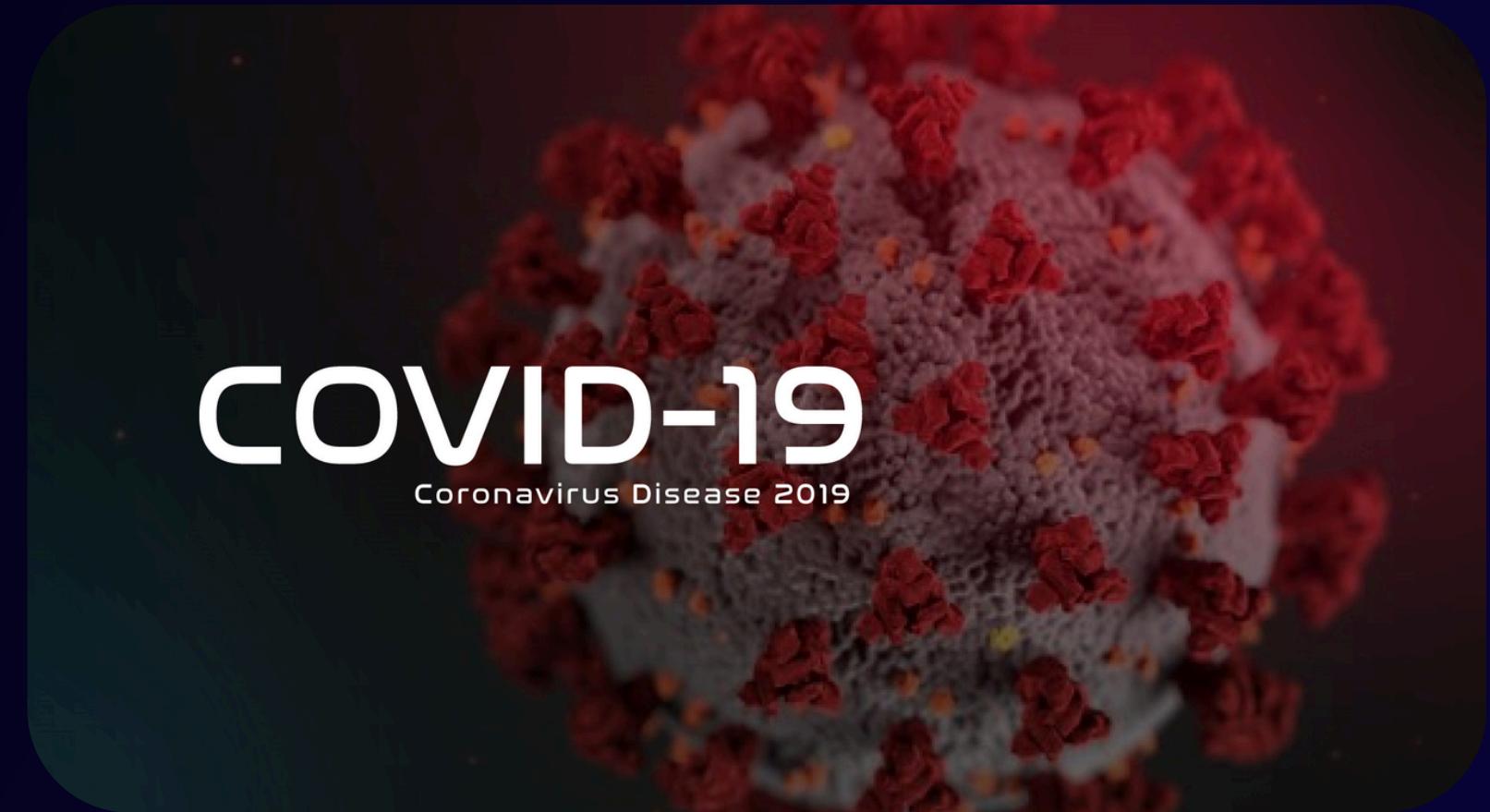


- Intro
- Data Cleaning and Preprocessing
- Feature Engineering
- EDA (Exploratory Data Analysis)
- SQL Analysis (12 Queries)
- Machine Learning Model
- Final Insights & Recommendations



Introduction

This analysis explores COVID-19 case numbers, deaths, and vaccination rates across U.S. states. The objective is to understand how vaccination coverage influenced COVID-19 fatality outcomes and to potentially uncover patterns that highlight disparities in healthcare access across different regions. To do this, I used two datasets from CDC public APIs.



COVID-19
Coronavirus Disease 2019

A detailed 3D rendering of the SARS-CoV-2 virus, showing its characteristic spike proteins on the surface of the viral envelope.

Target Audience

- Public Health Officials
- Healthcare Policymakers
- Healthcare Practitioners and Administrators

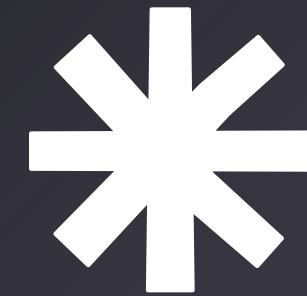


Data Cleaning & Preprocessing

DC API used for fetching two datasets of COVID-19

```
1 # Define Function
2 def pull_api_data(api_url):
3     response = requests.get(api_url)
4     if response.status_code == 200:
5         data = response.json()
6         df = pd.DataFrame(data)
7         print(f"Successfully pulled {len(df)} records from API")
8         return df
9     else:
10        print(f"Failed to pull data. Status code: {response.status_code}")
11        return None
12
13 # Define API Endpoints
14
15 # ~COVID-19 case and death data by state~
16 cases_deaths_api = 'https://data.cdc.gov/resource/pwn4-m3yp.json?$limit=50000'
17 # ~COVID-19 vaccination data by state~
18 vaccinations_api = 'https://data.cdc.gov/resource/unsk-b7fc.json?$limit=50000'
19
20 # Pull Data Into DataFrames
21 df_cases_deaths = pull_api_data(cases_deaths_api)
22 df_vaccinations = pull_api_data(vaccinations_api)
```

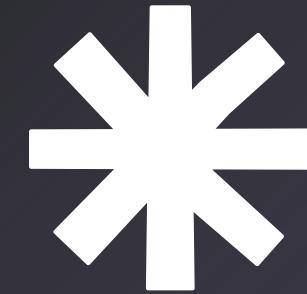
- Created a function **pull_api_data()** to fetch data from CDC APIs using **requests.get()** and convert the response to a DataFrame.



Check for Null Values

```
● ● ●  
1 # Print Header to Label Output Section  
2 print("\nMissing Values in Cases & Deaths - ", end="")  
3  
4 # Count Missing Values for Each Column in Cases & Deaths Dataset  
5 missing = df_cases_deaths.isnull().sum()  
6  
7 # Print Number of Columns with At Least One Missing Value  
8 print(len(missing[missing > 0])), "Columns with missing values: \n")  
9  
10 # Show Actual Columns and Number of Values Missing in Each  
11 print(missing[missing > 0])
```

- Uses **df.isnull().sum()** to count missing values in each column of the cases & deaths dataset.
- Prints the number of columns that have missing data.



Relevant Columns Only

```
● ● ●  
1 # Keep Relevant Columns Only  
2 df_cases_deaths = df_cases_deaths[['state', 'start_date', 'end_  
    date', 'tot_cases', 'tot_deaths', 'new_cases', 'new_deaths']]  
3  
4 # Rename Columns for Clarity  
5 df_cases_deaths = df_cases_deaths.rename(columns={  
6     'end_date': 'date',  
7     'tot_cases': 'total_cases',  
8     'tot_deaths': 'total_deaths'  
9 })  
10  
11 # Convert Date to Datetime  
12 df_cases_deaths['date'] = pd.to_datetime(df_cases_deaths['dat  
e'], errors='coerce')
```

- Selected relevant and rename using **df.rename()** for clearer names.
- Convert dates using **pd.to_datetime()** for proper date handling (errors='coerce').



```
● Converts vaccination columns to numeric  
with pd.to_numeric() and creates a  
Wednesday-aligned weekly date using  
pd.to_timedelta() and pd.Timedelta().  
  
● Groups by state and week with  
.groupby() and .agg('mean'), then  
renames the grouped week column to  
'date' using .rename()
```

```
1 # Convert Vaccination Columns to Numeric, Force Errors to NaN
2 cols_to_convert = [
3     'fully_vaccinated_pct',
4     'one_dose_pct',
5     'booster_pct',
6     'fully_vaccinated_65plus_pct',
7     'booster_65plus_pct'
8 ]
9 for col in cols_to_convert:
10     df_vaccinations[col] = pd.to_numeric(df_vaccinations[col], errors='coerce')
11
12
13 # Create 'Week' Column
14 df_vaccinations['week'] = df_vaccinations['date'] - pd.to_timedelta(df_vaccinations['date'].dt.weekday, unit='d') + pd.Timedelta(days=2) #Realign to Wednesday
15
16 # Group by 'State' and 'Week' and Take Mean
17 df_vaccinations_weekly = df_vaccinations.groupby(['state', 'week']).agg({
18     'fully_vaccinated_pct': 'mean',
19     'one_dose_pct': 'mean',
20     'booster_pct': 'mean',
21     'fully_vaccinated_65plus_pct': 'mean',
22     'booster_65plus_pct': 'mean'
23 }).reset_index()
24
25 # Rename 'Week' Column Back to Date
26 df_vaccinations_weekly = df_vaccinations_weekly.rename(columns={'week': 'date'})
```

Aggregate Vaccination Data to Weekly Level



Check for Outliers



```
1 numeric_cols = ['total_cases', 'total_deaths', 'new_cases', 'new_deaths']
2
3 # Convert Columns to Numeric
4 for col in numeric_cols:
5     df_cases_deaths[col] = pd.to_numeric(df_cases_deaths[col], errors='coerce')
6
7 # Create Boxplot
8 for col in numeric_cols:
9     plt.figure(figsize=(6, 4))
10    sns.boxplot(x=df_cases_deaths[col], color='orange')
11    plt.title(f'Boxplot of {col}')
12    plt.xlabel(col)
13    plt.show()
```

- Converts specified columns to numeric using **pd.to_numeric()** with error coercion to handle invalid data
- Generates boxplots for each numeric column with **sns.boxplot()** to visually identify outliers and distribution.



Handle Outliers



```
1 # Filter to Check Rows with 'Negative' New Deaths
2 df_cases_deaths[df_cases_deaths['new_deaths'] < 0]
3
4 # Manually Inspect Random Rows
5 df_cases_deaths.sample(10)
6
7 # Convert Negative Values in New Deaths to Positive (Absolute Value)
8 df_cases_deaths['new_deaths'] = abs(df_cases_deaths['new_deaths'])
9
10 # Recheck to Confirm
11 df_cases_deaths[df_cases_deaths['new_deaths'] < 0]
12
13 # View Sample
14 df_cases_deaths.sample(5)
```

- Uses filtering and sampling **df_cases_deaths.sample()** to identify and inspect negative new_deaths values as outliers.
- Corrects outliers by applying the absolute value function **abs()** to new_deaths, then verifies the fix by re-filtering and sampling the cleaned data.



- Connects to an SQLite database (or creates it if not exists) using **sqlite3.connect()**.
- Saves dataframes **df_cases_deaths** and **df_vaccinations_weekly** into tables (**cases_deaths** and **vaccinations_weekly**) with **.to_sql()**, replacing existing tables if needed.
- Closes the connection with **conn.close()** and confirms successful storage with a print statement.

```
● ● ●  
1 # Store Data into SQLite Database  
2 import sqlite3  
3  
4 # Define Database File Name  
5 db_filename = "covid_analysis.db"  
6  
7 # Connect to SQLite Database  
8 conn = sqlite3.connect(db_filename)  
9  
10 # Store Cases & Deaths Dataframe as Table  
11 df_cases_deaths.to_sql("cases_deaths", conn, if_exists="replace", index=False)  
12  
13 # Store Weekly Aggregated Vaccinations Dataframe as Table  
14 df_vaccinations_weekly.to_sql("vaccinations_weekly", conn, if_exists="replace", index=False)  
15  
16 # Close Database Connection  
17 conn.close()  
18  
19 # Confirm Storage Completion  
20 print(f"Data stored in {db_filename} (tables: cases_deaths, vaccinations_weekly)")
```

Store Cleaned Datasets into SQLite Database



Merging Datasets on State and Date



```
1 # Merging on State and Date Using an Inner Join
2 df_merged = pd.merge(df_cases_deaths, df_vaccinations_weekly, on=
   ['state', 'date'], how='inner')
3
4 # View Shape and Sample of Merged Dataset
5 print("Merged dataset shape:", df_merged.shape)
6 print("\nSample of merged data:")
7 display(df_merged.head())
8
9 # Check for Missing Values After Merging
10 print("\nMissing values after merging:")
11 print(df_merged.isnull().sum())
```

- Supports deeper insights into how vaccination coverage relates to COVID-19 case and death trends over time and across different states.



Key Findings

- Merged cases, deaths, and vaccination data on state and date – resulting in a clean, unified dataset with **7,127** records across key public health metrics.
- Vaccination metrics are consistently populated; missing values in **booster_pct** appear only in early periods, reflecting the later rollout of booster campaigns.
- Final dataset is complete and analysis-ready – ideal for exploring relationships between vaccination coverage and COVID-19 outcomes over time.



Feature Engineering *

Case Fatality Rate

Calculated as
 $(\text{total_deaths} / \text{total_cases}) * 100$,
this metric indicates the severity of
COVID-19 by showing the proportion
of confirmed cases that resulted in
death.



```
1 # Convert 'total_cases' and 'total_deaths' to Numeric
2 df_merged['total_cases'] = pd.to_numeric(df_merged['total_cases'], errors='coerce')
3 df_merged['total_deaths'] = pd.to_numeric(df_merged['total_deaths'], errors='coerce')
4
5 df_merged['case_fatality_rate'] = (df_merged['total_deaths'] / df_merged['total_cases']) * 100
```



Feature Engineering *

Vaccination Coverage Category

astype(float) converts values to float type, and **apply(categorize_vaccination)** runs the function on each value to categorize vaccination coverage into 'High', 'Medium', or 'Low'.

```
● ● ●  
1 def categorize_vaccination(coverage_pct):  
2     if coverage_pct >= 70:  
3         return 'High'  
4     elif coverage_pct >= 50:  
5         return 'Medium'  
6     else:  
7         return 'Low'  
8  
9 df_merged['vaccination_category'] = df_merged['fully_vaccinated_pct'].astype(float).apply(categorize_vaccination)
```



Intro

Data Preprocessing

Feature Engineering

EDA

SQL

ML Model

Conclusion

Feature Engineering *

Elderly Protection Index

Calculates the elderly protection index as the average of **fully_vaccinated_65plus_pct** and **booster_65plus_pct** (converted to **float**), indicating overall vaccine coverage among seniors.



```
1 df_merged['elderly_protection_index'] = (
2     df_merged['fully_vaccinated_65plus_pct'].astype(
3         float) + df_merged['booster_65plus_pct'].astype(float)
4 ) / 2
```

[Intro](#)[Data Preprocessing](#)[**Feature Engineering**](#)[EDA](#)[SQL](#)[ML Model](#)[Conclusion](#)

Feature Engineering *

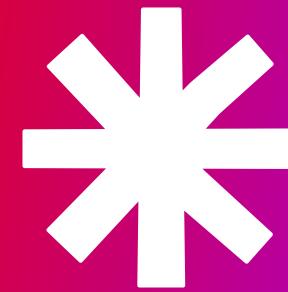
Fatality Risk Classification

This code creates a new column **fatality_risk** by applying the **classify_fatality()** function to each row's case_fatality_rate, classifying it as 'Low', 'Moderate', or 'High' based on thresholds.

```
● ● ●  
1 def classify_fatality(r):  
2     if r < 0.9:  
3         return 'Low'  
4     elif r < 1.52:  
5         return 'Moderate'  
6     else:  
7         return 'High'  
8  
9 df_merged["fatality_risk"] = df_merged["case_fatality_rat  
e"].apply(classify_fatality)
```

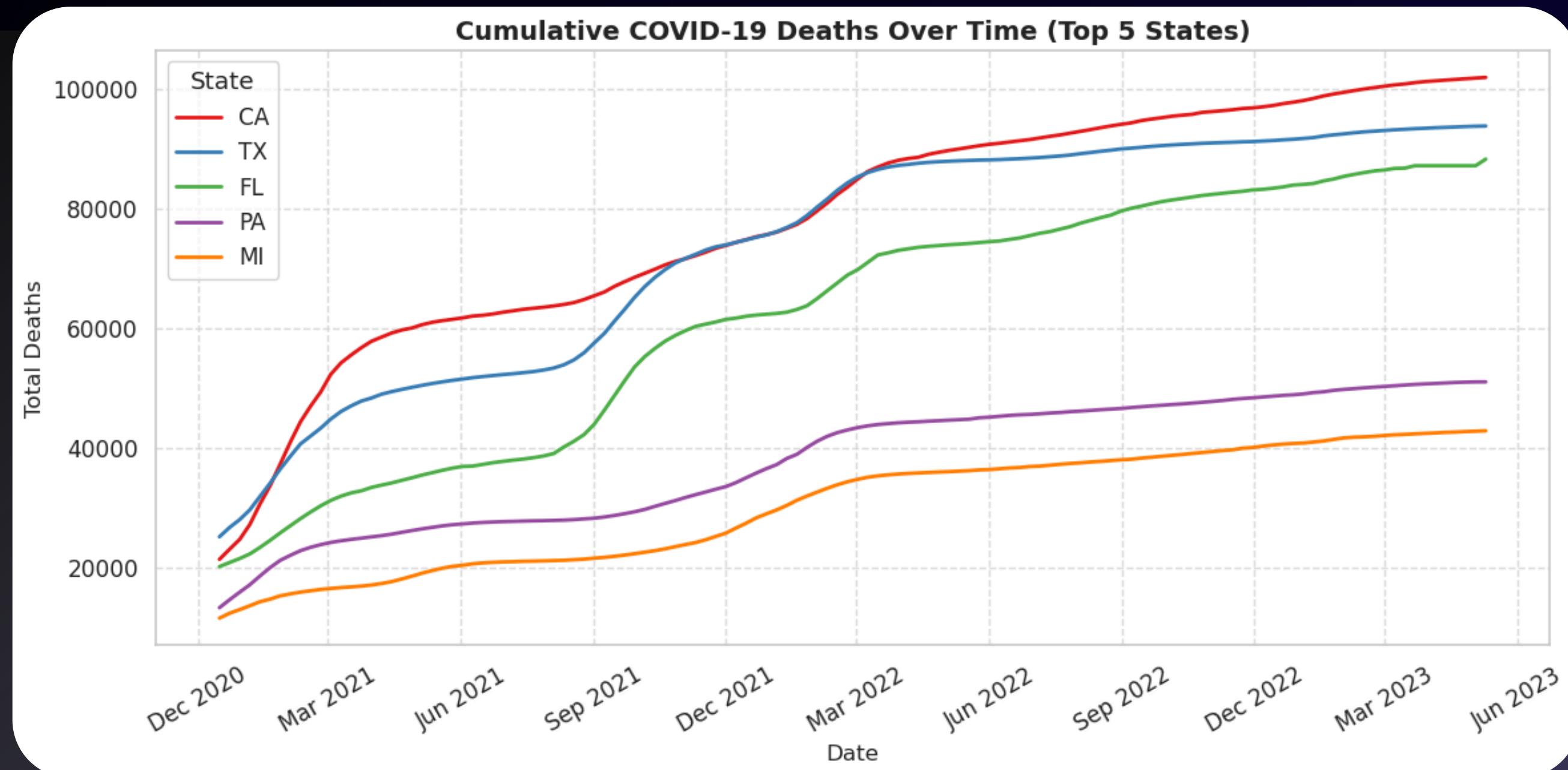
[Intro](#)[Data Preprocessing](#)[Feature Engineering](#)[EDA](#)[SQL](#)[ML Model](#)[Conclusion](#)

Exploratory Data Analysis



Most Important Insights



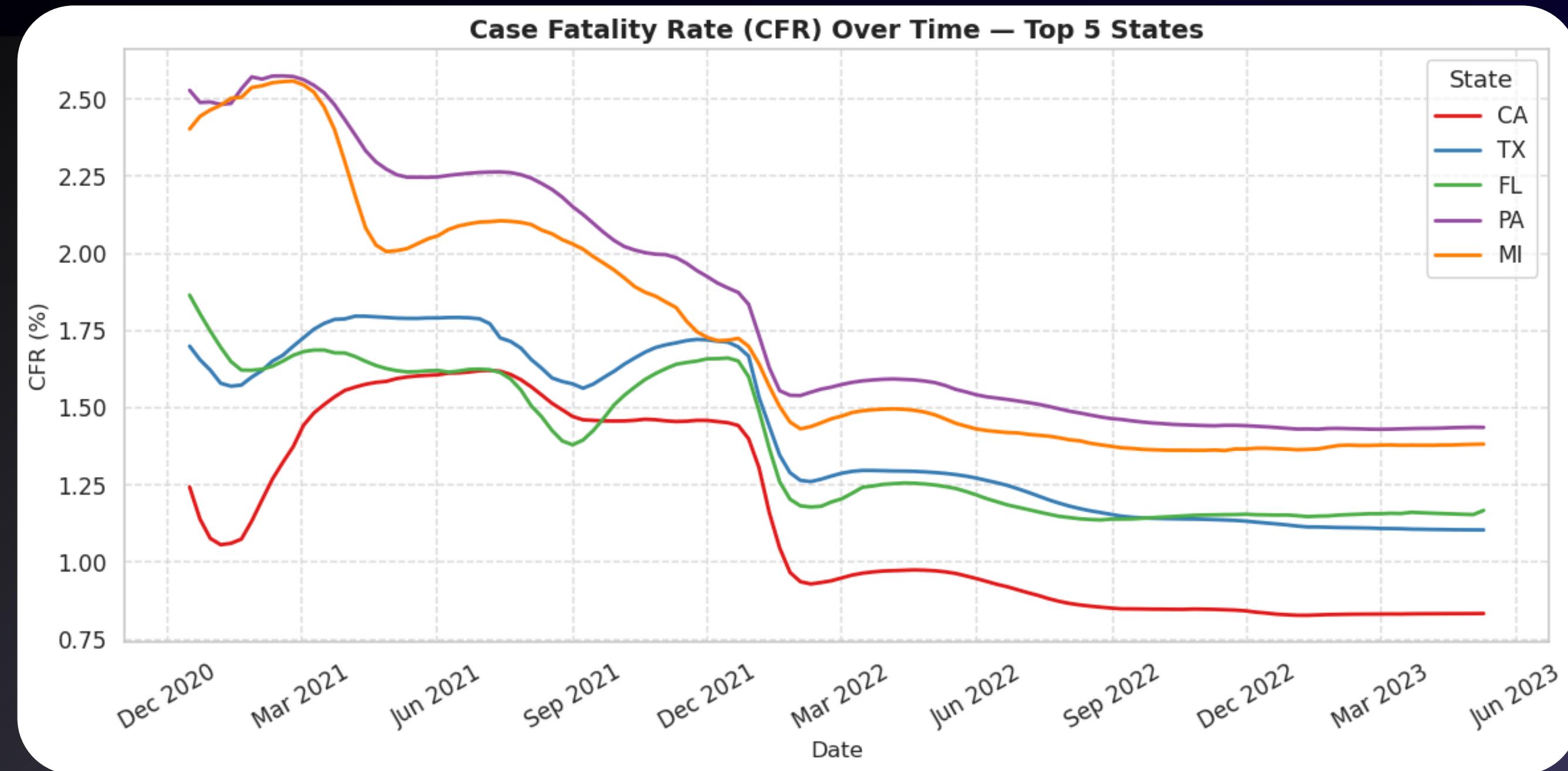


- California and Texas had the highest total COVID-19 deaths, reflecting large population sizes and overall burden.
- Cumulative death totals don't fully capture outbreak severity relative to case counts.
- All states showed steady increases in deaths over time, highlighting persistent vulnerability.

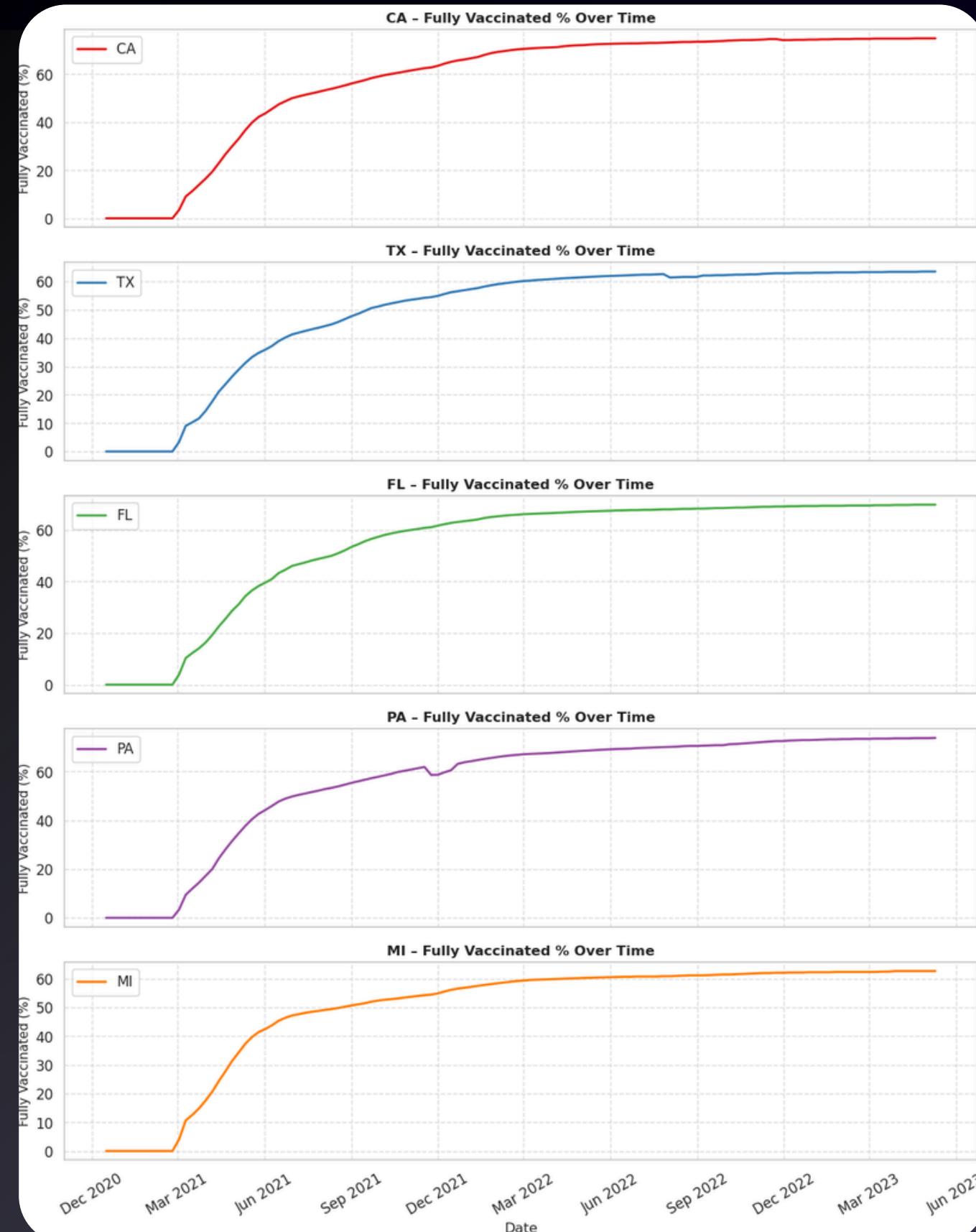


```
● ● ●  
1 import matplotlib.dates as mdates  
2 plt.figure(figsize=(12, 6))  
3  
4 # Define Top 5 States  
5 top5_states = ['CA', 'TX', 'FL', 'PA', 'MI']  
6  
7 # Plot Total Deaths Over Time  
8 for state in top5_states:  
9     subset = df_merged[df_merged['state'] == state]  
10    plt.plot(  
11        subset['date'],  
12        subset['total_deaths'],  
13        linewidth=2,  
14        label=state  
15    )  
16  
17 # Add Title and Labels  
18 plt.title('Cumulative COVID-19 Deaths Over Time (Top 5 States)', fontsize=14, fontweight='bold')  
19 plt.xlabel('Date', fontsize=12)  
20 plt.ylabel('Total Deaths', fontsize=12)  
21  
22 # Add Legend, Grid, and Layout  
23 plt.legend(title='State')  
24 plt.grid(True, linestyle='--', alpha=0.6)  
25 plt.gca().xaxis.set_major_locator(mdates.MonthLocator(interval=3))  
26 plt.gca().xaxis.set_major_formatter(mdates.DateFormatter('%b %Y'))  
27 plt.xticks(rotation=30)  
28 plt.tight_layout()  
29  
30 # Show Plot  
31 plt.show()
```

- Uses a for loop to go through each state in top5_states and plots a line for its total_deaths over time with plt.plot().
- Uses **mdates.MonthLocator()** and **mdates.DateFormatter('%b %Y')** to format the x-axis dates in clean 3-month intervals (like "Jan 2021").
- Improves readability with **plt.title()**, **plt.legend()**, **plt.grid()**, and **plt.tight_layout()** so the chart looks clean and easy to follow.

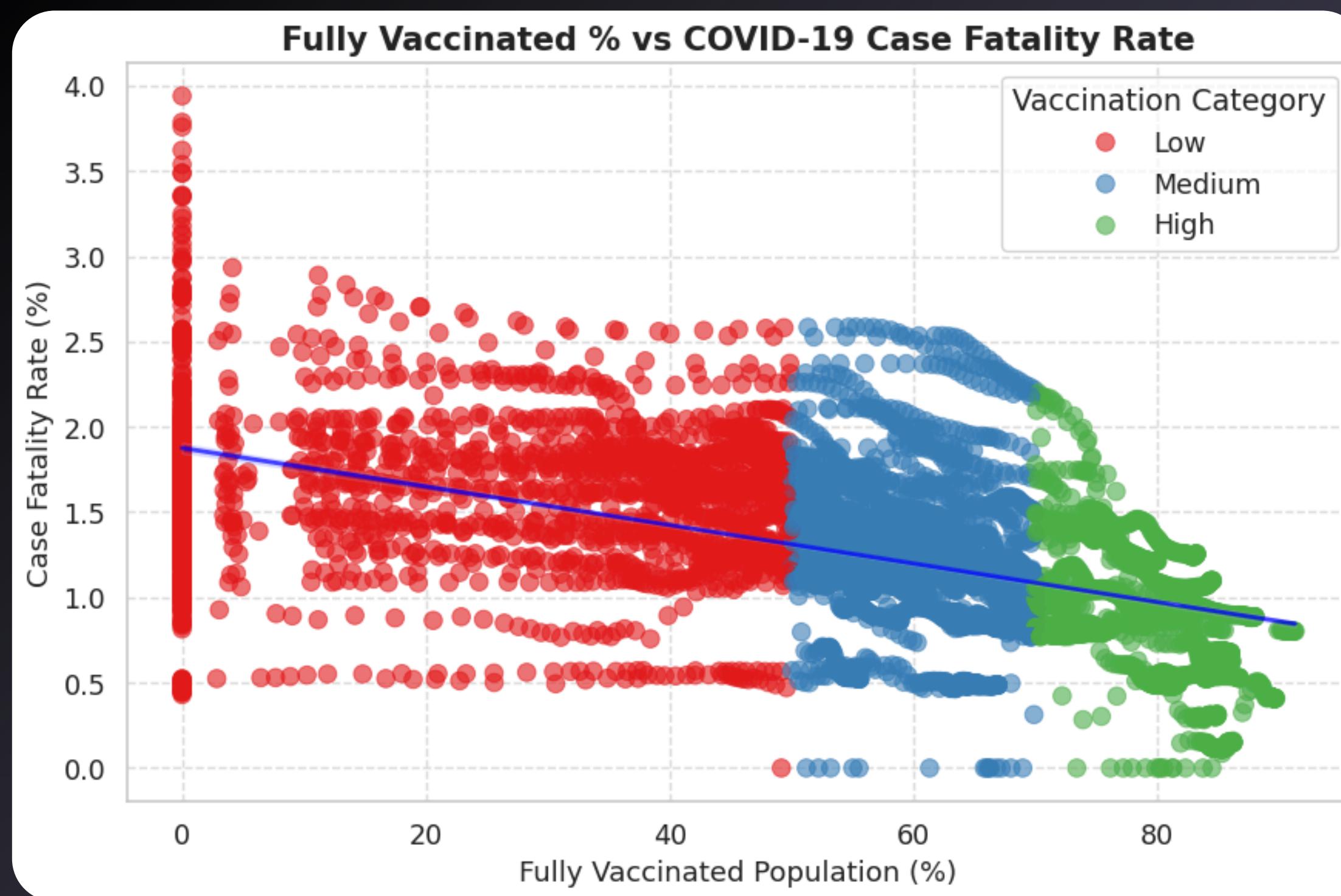


- CFR is a better measure than death counts alone, showing how deadly the virus was relative to infections.
- CFR decreased over time in all states, likely due to better treatments and healthcare response.
- Pennsylvania and Michigan started with higher CFRs ($>2.5\%$), while California maintained the lowest fatality rate.



Fully Vaccinated % Over Time

- All five states show a clear increase in vaccination rates starting early 2021, marking the mass vaccine rollout.
- **CA** and **PA** achieved the highest vaccination rates (~75%), while **TX**, **FL**, and **MI** plateaued around **60–65%**.
- Higher vaccination rates, especially in **CA** and **PA**, correlate with lower case fatality rates (CFR), indicating vaccines helped reduce fatality risk.

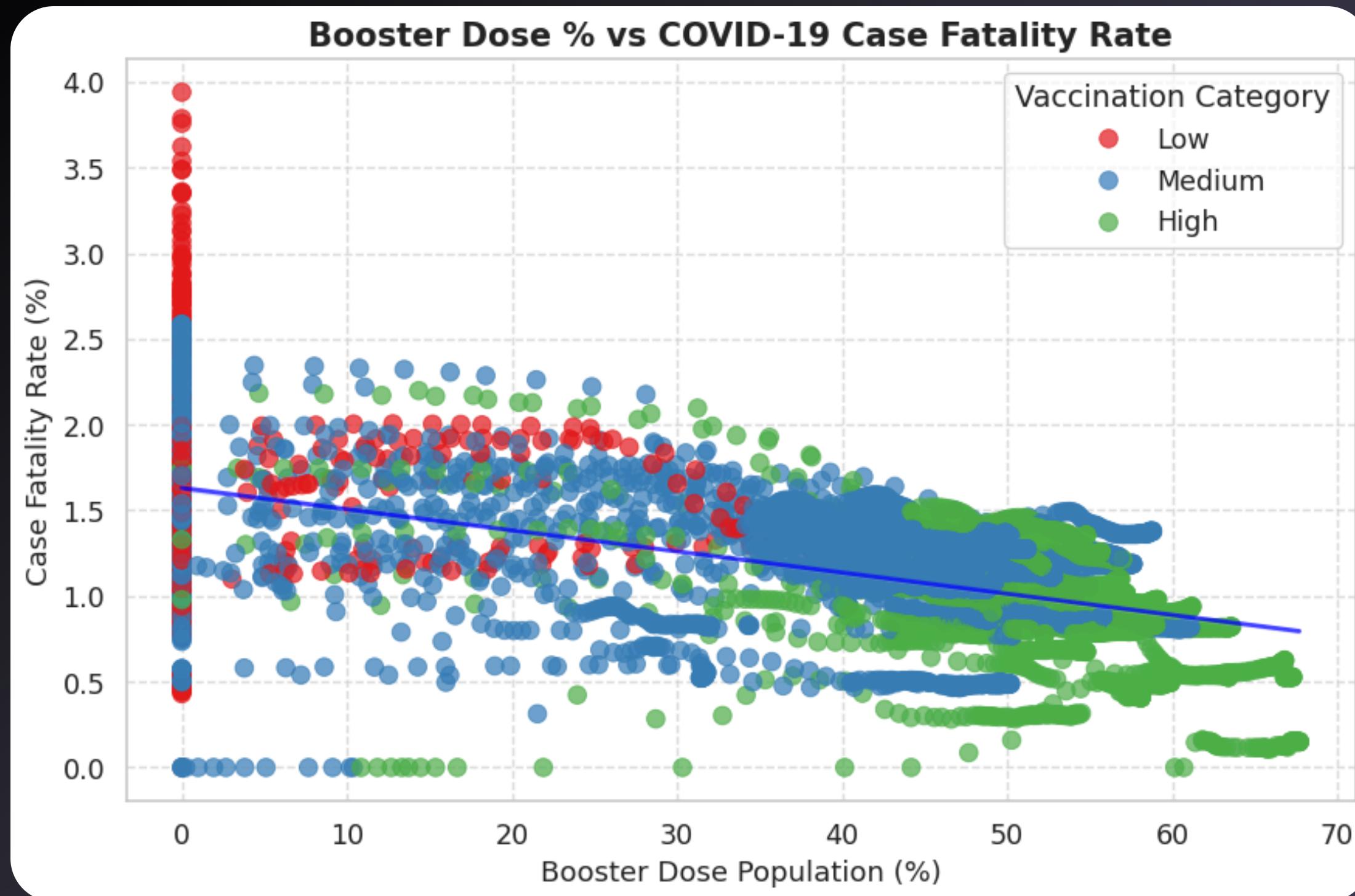


- Increasing vaccination rates are linked to a decrease in case fatality rates (CFR), showing a modest negative correlation.
- Once vaccination surpasses 60–70%, CFR mostly stays below 1%, indicating stronger protection against severe outcomes.

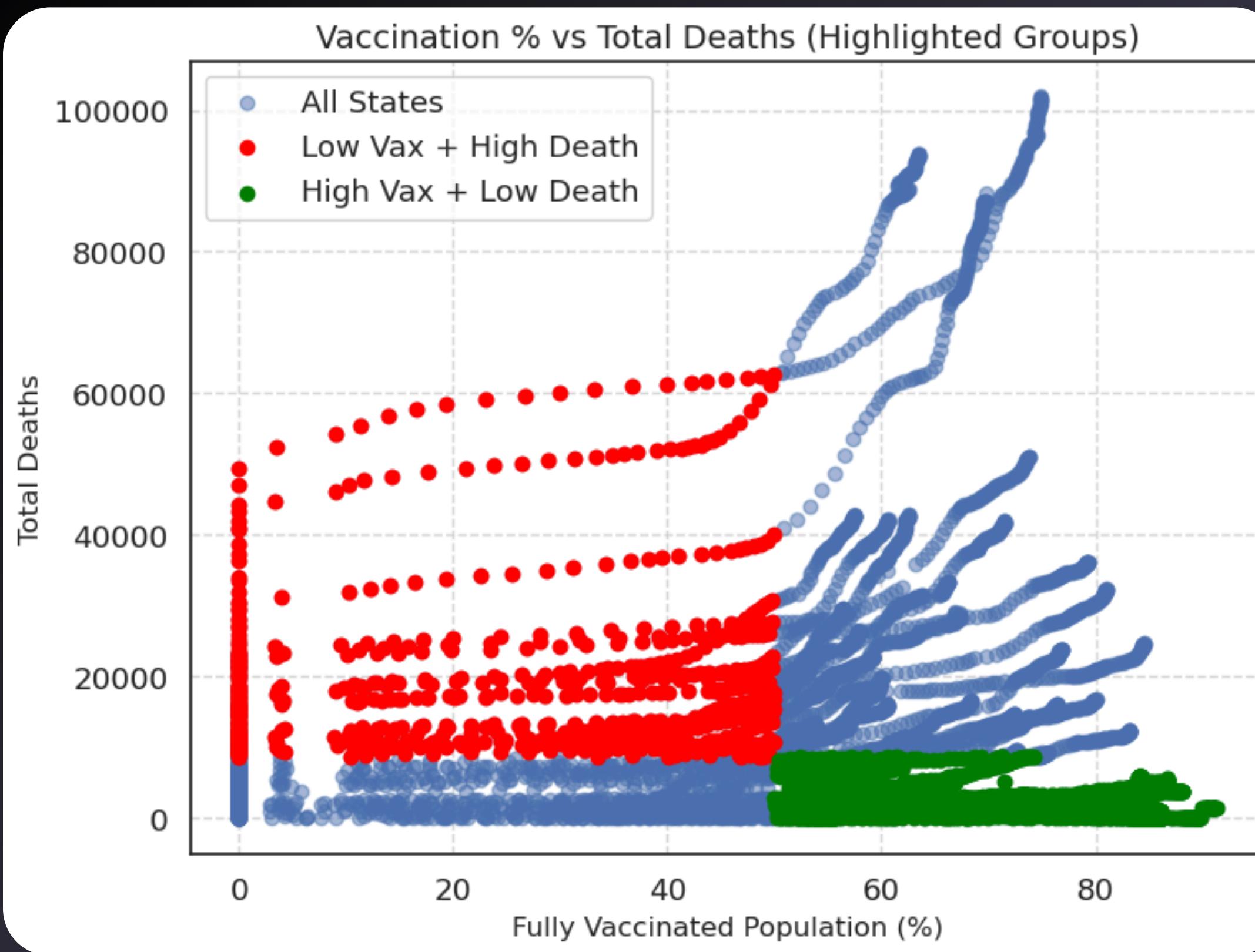


```
1 # Set Up Figure and Style
2 plt.figure(figsize=(9, 6))
3 sns.set(style="whitegrid", font_scale=1.15)
4
5 # Plot Scatter Points by Vaccination Category
6 sns.scatterplot(
7     data=df_merged,
8     x='fully_vaccinated_pct',
9     y='case_fatality_rate',
10    hue='vaccination_category',
11    palette='Set1',
12    alpha=0.6,
13    s=70,
14    edgecolor=None
15 )
16
17 # Add Regression Line
18 sns.regplot(
19     data=df_merged,
20     x='fully_vaccinated_pct',
21     y='case_fatality_rate',
22     scatter=False,
23     color='blue',
24     line_kws={'linewidth': 2, 'alpha': 0.7}
25 )
26
27 # Add Title, Labels, Legend, and Grid
28 plt.title('Fully Vaccinated % vs COVID-19 Case Fatality Rate', fontsize=15, fontweight='bold')
29 plt.xlabel('Fully Vaccinated Population (%)', fontsize=13)
30 plt.ylabel('Case Fatality Rate (%)', fontsize=13)
31 plt.legend(title='Vaccination Category')
32 plt.grid(True, linestyle='--', alpha=0.6)
33 plt.tight_layout()
34
35 # Show Plot
36 plt.show()
```

- Use **sns.scatterplot()** to visualize the relationship between vaccination % and fatality rate, with points colored by vaccination category.
- Add a trend line using **sns.regplot()** to show the overall correlation, and customize the plot with titles, labels, and grid for clarity.



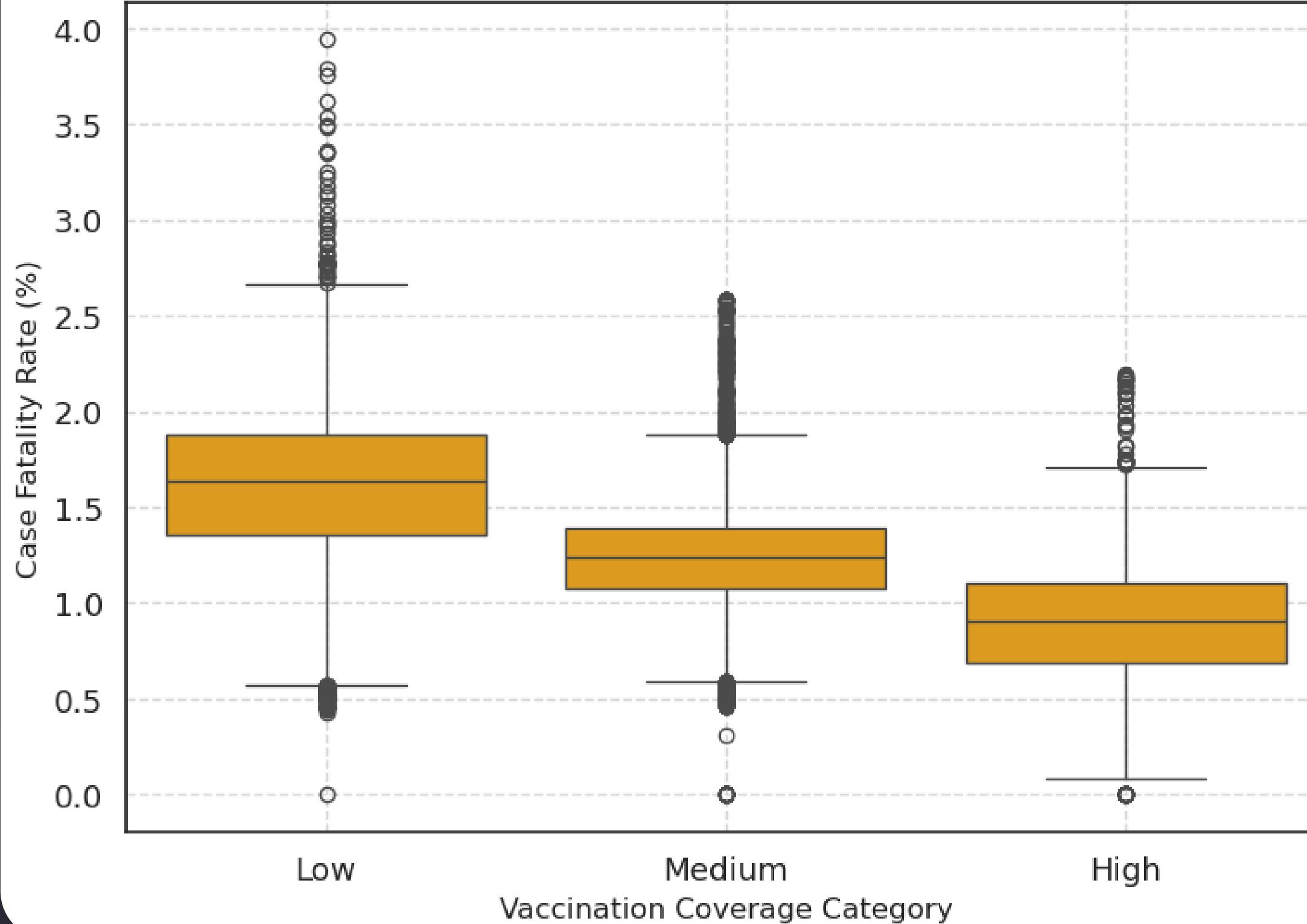
- Higher booster coverage tends to link with lower case fatality rates, but the relationship is weak and scattered.
- When booster coverage exceeds ~30–40%, CFR values often cluster below 1%, indicating improved protection.
- Low booster coverage (<20%) shows high CFR variability, suggesting other factors also impact fatality rates.



- Low vaccination + high death states (red) cluster at lower vaccination rates and higher total deaths.
- High vaccination + low death states (green) cluster at higher vaccination rates and lower death counts.
- All other states (blue) scatter between the two extremes, indicating that vaccination was important but other factors (population size, health access) also influenced deaths.



Case Fatality Rate by Vaccination Coverage Category

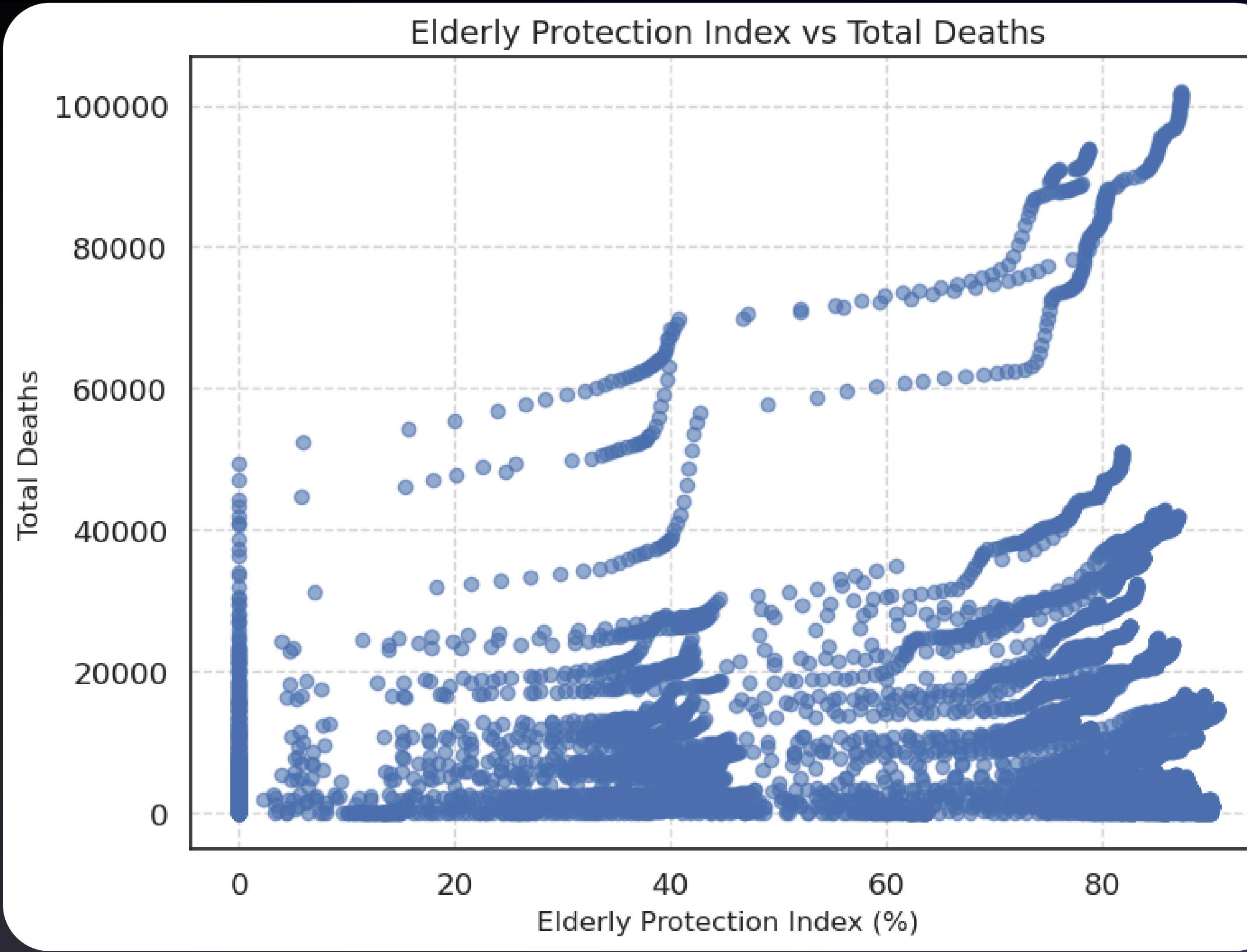
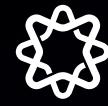


- States with **high** vaccination coverage show lower median CFRs and fewer extreme outliers.
- **Low** coverage states have higher median CFRs and more variability, reflecting worse outcomes.
- Overall, CFR tends to **decrease** as vaccination coverage **improves** from Low to High.

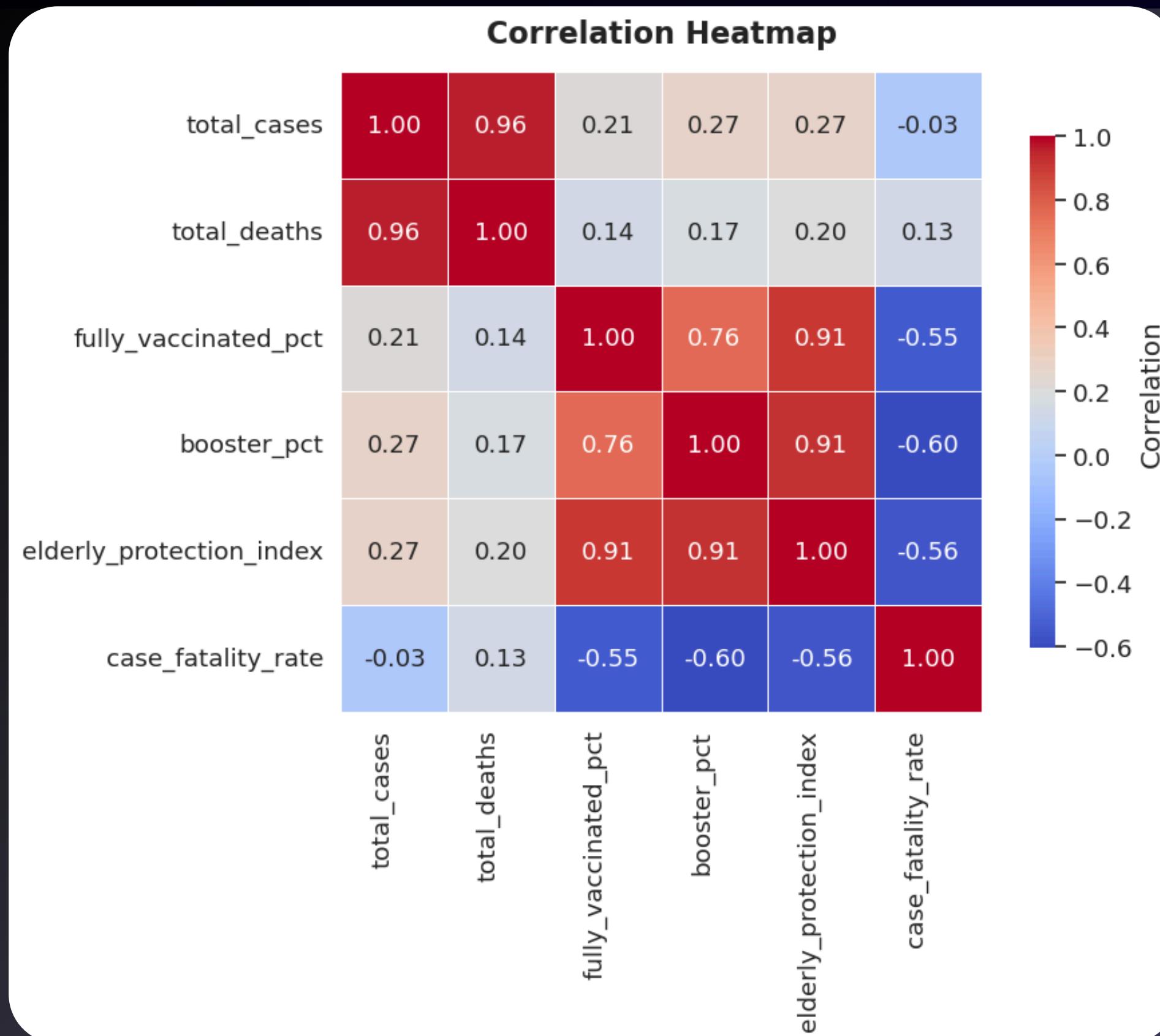


```
● ● ●  
1 # Create Figure  
2 plt.figure(figsize=(8,6))  
3  
4 # Draw Boxplot  
5 sns.boxplot(  
6     x='vaccination_category',  
7     y='case_fatality_rate',  
8     data=df_merged,  
9     order=['Low', 'Medium', 'High'],  
10    color='orange' # <-- manually specifying order  
11 )  
12  
13 # Add Title and Labels  
14 plt.title('Case Fatality Rate by Vaccination Coverage Category', fontweight='bold')  
15 plt.xlabel('Vaccination Coverage Category', fontsize=12)  
16 plt.ylabel('Case Fatality Rate (%)', fontsize=12)  
17  
18 # Add Grid and Layout  
19 plt.grid(True, linestyle='--', alpha=0.7)  
20 plt.tight_layout()  
21  
22 # Show Plot  
23 plt.show()
```

- Use **sns.boxplot()** to compare the distribution of case fatality rates across vaccination coverage levels (Low, Medium, High).
- The **order** parameter ensures consistent category sorting on the x-axis.
- Add titles, labels, and gridlines to clearly present the relationship between vaccination and fatality outcomes.



- States with higher elderly protection (**60%–90%**) still experienced a spread in total deaths, implying that other factors like outbreak size, timing, or healthcare access mattered.



- Vaccination metrics (full vax %, booster %, elderly protection) are strongly correlated (**0.76–0.91**), showing states with strong initial rollout continued strong coverage.
- CFR has moderate negative correlation (~ **-0.55** to **-0.60**) with vaccination, supporting its role in reducing death severity.
- Total cases and total deaths are highly correlated (**0.96**), reflecting the overall impact of widespread infections.

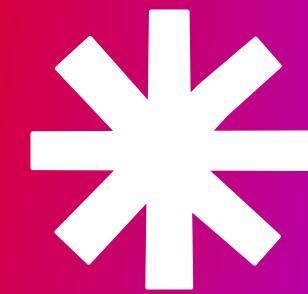


```
 1 # Pick Numeric Columns Only
 2 corr_columns = [
 3     'total_cases',
 4     'total_deaths',
 5     'fully_vaccinated_pct',
 6     'booster_pct',
 7     'elderly_protection_index',
 8     'case_fatality_rate'
 9 ]
10
11 # Calculate Correlation Matrix
12 corr_matrix = df_merged[corr_columns].corr()
13
14 # Create Heatmap with Seaborn Styling
15 plt.figure(figsize=(10, 8))
16 sns.set(style="white", font_scale=1.2)
17
18 ax = sns.heatmap(
19     corr_matrix,
20     annot=True,
21     cmap='coolwarm',
22     fmt='.2f',
23     linewidths=0.5,
24     square=True,
25     cbar_kws={'shrink': 0.8, 'label': 'Correlation'},
26     annot_kws={"size": 13}
27 )
28
29 # Add Title and Ticks
30 ax.set_title('Correlation Heatmap', fontsize=16, fontweight='bold', pad=16)
31 plt.yticks(rotation=0)
32 plt.xticks(rotation=90, ha='right')
33 plt.tight_layout()
34
35 # Show Plot
36 plt.show()
```

- Select relevant numeric columns and use **.corr()** to compute the correlation matrix.
- Visualize the matrix with **sns.heatmap()**, showing relationships between variables with color and annotated values.
- Customize the plot (title, labels, color scale) for readability and professional presentation.

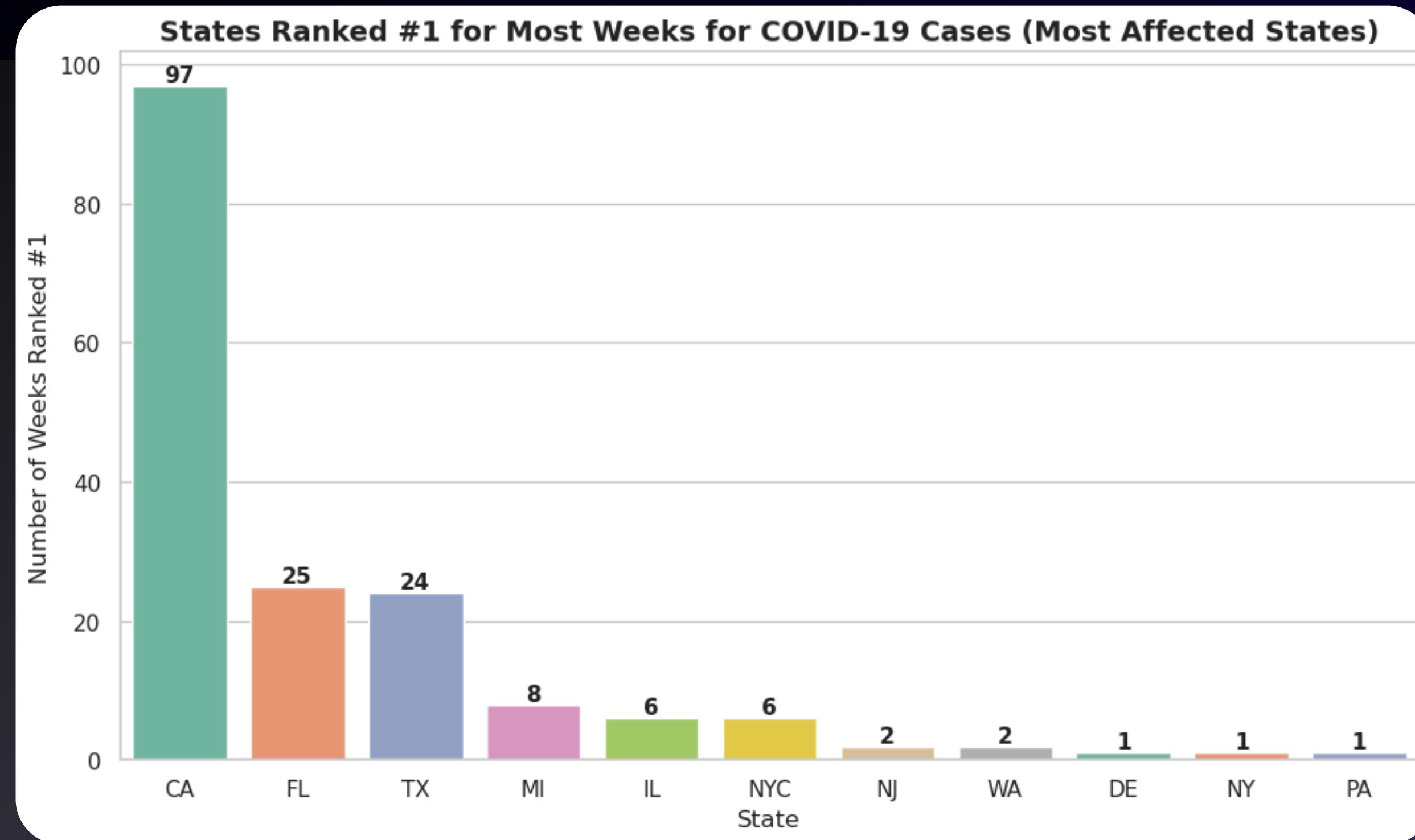
[Intro](#)[Data Preprocessing](#)[Feature Engineering](#)[EDA](#)[SQL](#)[ML Model](#)[Conclusion](#)

SQL Analysis



Most Important Insights

```
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
SELECT empCode, empName, empSalary
FROM Employee
WHERE empName in
(SELECT DISTINCT empName
FROM population
WHERE Country = "TH")
AND empSalary >=
(SELECT AVG(salary)
FROM Salary
WHERE gender = "M")
```

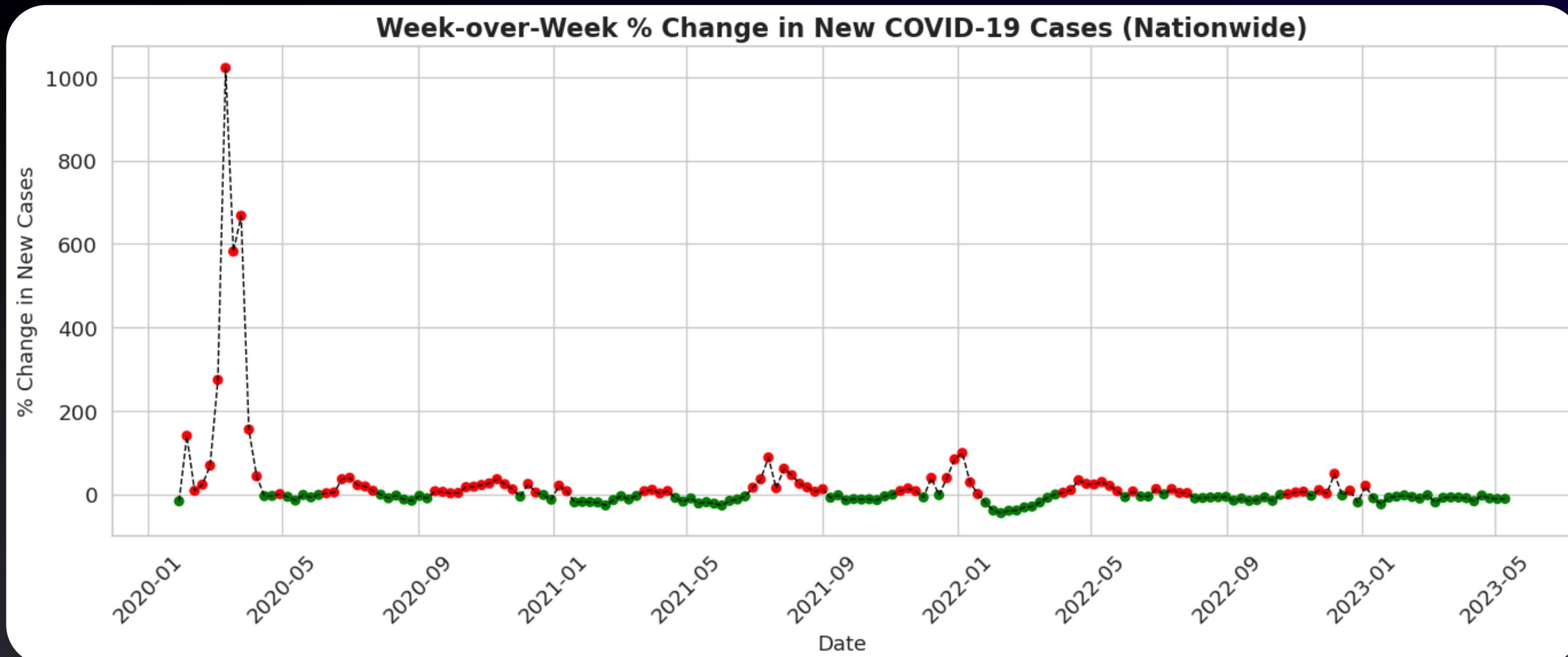


- California (CA), Florida (FL), and Texas (TX) were the most affected states.
- They ranked #1 in weekly death counts for 97, 25, and 24 weeks respectively.



```
● ● ●  
1 query = """  
2 -- QUERY 4: Rank states by new cases per date using window function,  
3 -- only extract rank 1  
4 SELECT  
5     state,  
6     COUNT(date) AS num_weeks_rank1  
7 FROM (  
8     SELECT  
9         state,  
10        date,  
11        new_cases,  
12        RANK() OVER (PARTITION BY date ORDER BY new_cases DESC) AS  
13        rank_by_cases  
14    FROM cases_deaths  
15 )  
16 WHERE rank_by_cases = 1  
17 GROUP BY state  
18 ORDER BY num_weeks_rank1 DESC, state;  
19 """  
20  
21 ranked_cases = pd.read_sql(query, conn) # Check unique states with  
# rank 1  
22 ranked_cases.to_csv("sql_results/query_4.csv", index=False) # Save  
# to CSV for further analysis  
23 ranked_cases
```

- Use a SQL window function (**RANK() OVER**) to rank states by **new_cases** for each date.
- Filter for **rank = 1** to find which state had the most cases per week.
- Group and count results in Python with **pd.read_sql()** and save to **CSV** for analysis.



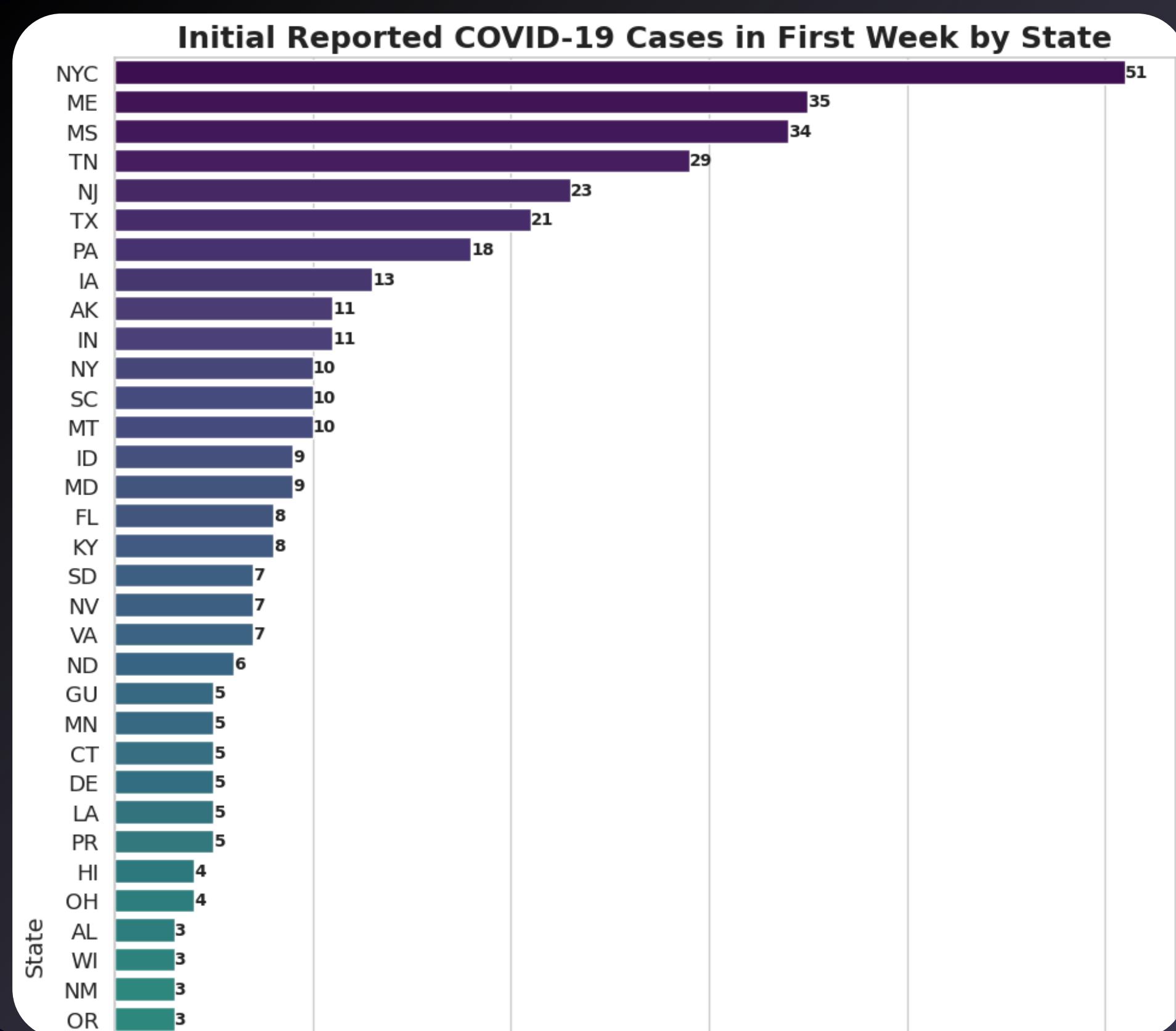
- Early 2020 saw extreme spikes in weekly case growth (1000%+), driven by rapid initial spread.
- After mid-2020, % changes stabilized to -30% to +30%, except during variant waves (Delta, Omicron).
- By 2022–23, WoW changes flattened, reflecting controlled spread due to vaccines and prior exposure.



```
● ● ●

1 query = """
2 -- QUERY 5: National week-over-week % change in total new cases
3
4 WITH weekly_totals AS (
5     SELECT
6         date,
7         SUM(new_cases) AS total_new_cases
8     FROM cases_deaths
9     GROUP BY date
10 ),
11 weekly_change AS (
12     SELECT
13         date,
14         total_new_cases,
15         LAG(total_new_cases) OVER (ORDER BY date) AS previous_total_cases
16     FROM weekly_totals
17 )
18 SELECT
19     date,
20     total_new_cases,
21     previous_total_cases,
22     ROUND(
23         CASE
24             WHEN previous_total_cases = 0 THEN NULL
25             ELSE ((total_new_cases - previous_total_cases) * 100.0 / previous_total_cases)
26         END, 2
27     ) AS percent_change
28 FROM weekly_change
29 ORDER BY date;
30 """
31
32 national Wow df = pd.read_sql(query, conn)
33 national Wow df.to_csv("sql_results/query_5.csv", index=False) # Save to CSV for further analysis
34 national Wow df.head(20)
35
```

- Use CTEs to first sum new cases by date, creating weekly totals (**weekly_totals**).
- Calculate the previous week's total using SQL's window function **LAG()** to compare week-over-week case numbers.
- Compute the percentage change between weeks, handling division by zero safely, then load the results into pandas and save for further analysis.



- New York City had the highest initial cases (**51**), showing an early intense outbreak.
- Maine, Mississippi, Tennessee, and New Jersey also had high first-week cases, indicating early spread.
- Many states like California, Michigan, Georgia, and District of Columbia started with low cases, while Washington, Arizona, Illinois, Massachusetts, and North Carolina reported just **1** case, suggesting isolated introductions.



```
1 query = """
2 -- QUERY 6: First week with > 0 new cases per state
3
4
5 WITH filtered_data AS (
6     SELECT *
7     FROM cases_deaths
8     WHERE new_cases > 0
9 ),
10 ranked_first_case AS (
11     SELECT
12         state,
13         date,
14         new_cases,
15         ROW_NUMBER() OVER (PARTITION BY state ORDER BY date) AS rn
16     FROM filtered_data
17 )
18 SELECT
19     state,
20     date AS first_reported_week,
21     new_cases AS initial_reported_cases
22 FROM ranked_first_case
23 WHERE rn = 1
24 ORDER BY first_reported_week;
25 """
26
27 first_nonzero_df = pd.read_sql(query, conn)
28 first_nonzero_df.to_csv("sql_results/query_6.csv", index=False) # Save
29 to CSV for further analysis
30 first_nonzero_df.head()
```

- Filter the dataset to only include weeks with new cases (**new_cases > 0**).
- Use **ROW_NUMBER()** window function, partitioned by state, ordered by date, to find each state's first week with reported cases.
- **Select** and **order** results by the first reported week, then load into **pandas** and save as **CSV** for further analysis.



States with Below-Average Vaccination and Above-Average Deaths

	state	avg_vax	avg_deaths_per_week
0	Texas	48.9%	555
1	Georgia	43.9%	260
2	Ohio	47.9%	258
3	Michigan	49.4%	255
4	Arizona	50.0%	209
5	Tennessee	43.7%	193
6	North Carolina	49.9%	180
7	Indiana	45.4%	152
8	Alabama	40.9%	136
9	Kentucky	47.3%	130
10	Missouri	45.6%	129
11	South Carolina	46.0%	121
12	Oklahoma	46.4%	110

- States like Texas, Georgia, Alabama, and Tennessee show low vaccination rates combined with high death rates.
- Texas has the lowest vaccination rate (49%) and highest weekly deaths (555).
- These areas are high-risk and may need stronger public health interventions.



```
● ● ●

1 query = """
2 -- QUERY 6: First week with > 0 new cases per state
3
4
5 WITH filtered_data AS (
6     SELECT *
7     FROM cases_deaths
8     WHERE new_cases > 0
9 ),
10 ranked_first_case AS (
11     SELECT
12         state,
13         date,
14         new_cases,
15         ROW_NUMBER() OVER (PARTITION BY state ORDER BY date) AS rn
16     FROM filtered_data
17 )
18 SELECT
19     state,
20     date AS first_reported_week,
21     new_cases AS initial_reported_cases
22 FROM ranked_first_case
23 WHERE rn = 1
24 ORDER BY first_reported_week;
25 """
26
27 first_nonzero_df = pd.read_sql(query, conn)
28 first_nonzero_df.to_csv("sql_results/query_6.csv", index=False) # Save
29 to CSV for further analysis
30 first_nonzero_df.head()
```

- Join vaccination and case data on state and date, then compute **average vaccination %** and deaths per state using **AVG()**
- Filter states with below-average vaccination and above-average deaths using subqueries.
- Rank by deaths (**ORDER BY**), load with **pd.read_sql()**, and save to CSV for analysis.



Weeks with $\geq 50\%$ Spike in New Cases

	date	total_new_cases	prev_week_cases	percent_change
0	2020-02-05	12	5	140%
1	2020-02-26	27	16	69%
2	2020-03-04	101	27	274%
3	2020-03-11	1133	101	1022%
4	2020-03-18	7725	1133	582%
5	2020-03-25	59286	7725	667%
6	2020-04-01	151264	59286	155%
7	2021-07-14	251109	133430	88%
8	2021-07-28	464744	288157	61%
9	2021-12-29	2141043	1165714	84%
10	2022-01-05	4258203	2141043	99%

- Weeks with over 50% increase in new cases highlight significant COVID-19 transmission rebounds.
- The largest spike occurred on March 11, 2020, with a 1021% increase.



```
● ● ●  
1 query = ""  
2 -- QUERY 11: COVID rebound alert - spikes of 50%+ in national new cases  
3  
4 WITH daily_totals AS (  
5     SELECT  
6         date,  
7         SUM(new_cases) AS total_new_cases  
8     FROM cases_deaths  
9     GROUP BY date  
10 ),  
11 with_lagged AS (  
12     SELECT  
13         date,  
14         total_new_cases,  
15         LAG(total_new_cases) OVER (ORDER BY date) AS prev_week_cases  
16     FROM daily_totals  
17 ),  
18 spike_alerts AS (  
19     SELECT  
20         date,  
21         total_new_cases,  
22         prev_week_cases,  
23         ROUND((total_new_cases - prev_week_cases) * 100.0 / prev_week_cases, 2) AS percent_change  
24     FROM with_lagged  
25     WHERE prev_week_cases IS NOT NULL  
26 )  
27  
28 SELECT  
29     SUBSTR(date, 1, 10) AS date, -- Only keep YYYY-MM-DD  
30     total_new_cases,  
31     prev_week_cases,  
32     percent_change  
33 FROM spike_alerts  
34 WHERE percent_change >= 50  
35 ORDER BY date;  
36 """  
37  
38 spike_alerts_df = pd.read_sql(query, conn)  
39 spike_alerts_df.to_csv("sql_results/query_11.csv", index=False) # Save to CSV for further analysis  
40 spike_alerts_df
```

- Uses **CTEs** to first sum daily new cases (**daily_totals**), then add previous day's cases with the **LAG()** window function (**with_lagged**).
- Calculates percentage change in new cases and filters for spikes of **50%** or more (**spike_alerts**).
- Selects and formats key columns, ordering results by date to identify significant COVID-19 rebounds.

[Intro](#)[Data Preprocessing](#)[Feature Engineering](#)[EDA](#)[SQL](#)[ML Model](#)[Conclusion](#)

Machine Learning Model

Linear Regression

**Predict Weekly New COVID-19 Cases
Based on Vaccination Coverage**

Goal

To estimate how different levels of vaccination coverage impact the number of new COVID-19 cases reported in a given week across U.S. states.



Importing scikit-learn, Data Preparation and Feature Engineering

```
● ● ●  
1 from sklearn.model_selection import train_test_split  
2 from sklearn.linear_model import LinearRegression  
3 from sklearn.metrics import mean_squared_error, r2_score  
4 import numpy as np  
5  
6 # Select Relevant Features and Drop Nulls  
7 model_df = df_merged[  
8     ["new_cases", "new_deaths", "fully_vaccinated_pct",  
9      "booster_pct", "one_dose_pct", "date", "elderly_protection_index"]]  
10 ].dropna()  
11  
12 # Create a 'Week' Feature to Encode Time Progression  
13 model_df["week"] = model_df["date"].dt.isocalendar().week  
14  
15 # Define X and y  
16 X = model_df.drop(columns=["new_cases", "date"])  
17 y = model_df["new_cases"]  
18  
19 # Train/Test Split  
20 x_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, ran  
dom_state=42, shuffle=True)
```

- Selected features include vaccination rates (**fully_vaccinated_pct**, **booster_pct**, **one_dose_pct**), **elderly_protection_index**, **week** (from date), and encoded state data.
- These variables were chosen to reflect vaccine-driven immunity and regional variation over time.
- Defined **X** by dropping new_cases and date; **y** was set as new_cases.
- Used **train_test_split()** to split the data into training and test sets (**80/20**) with **shuffling**.



Model Training, Predictions and Evaluating Metrics

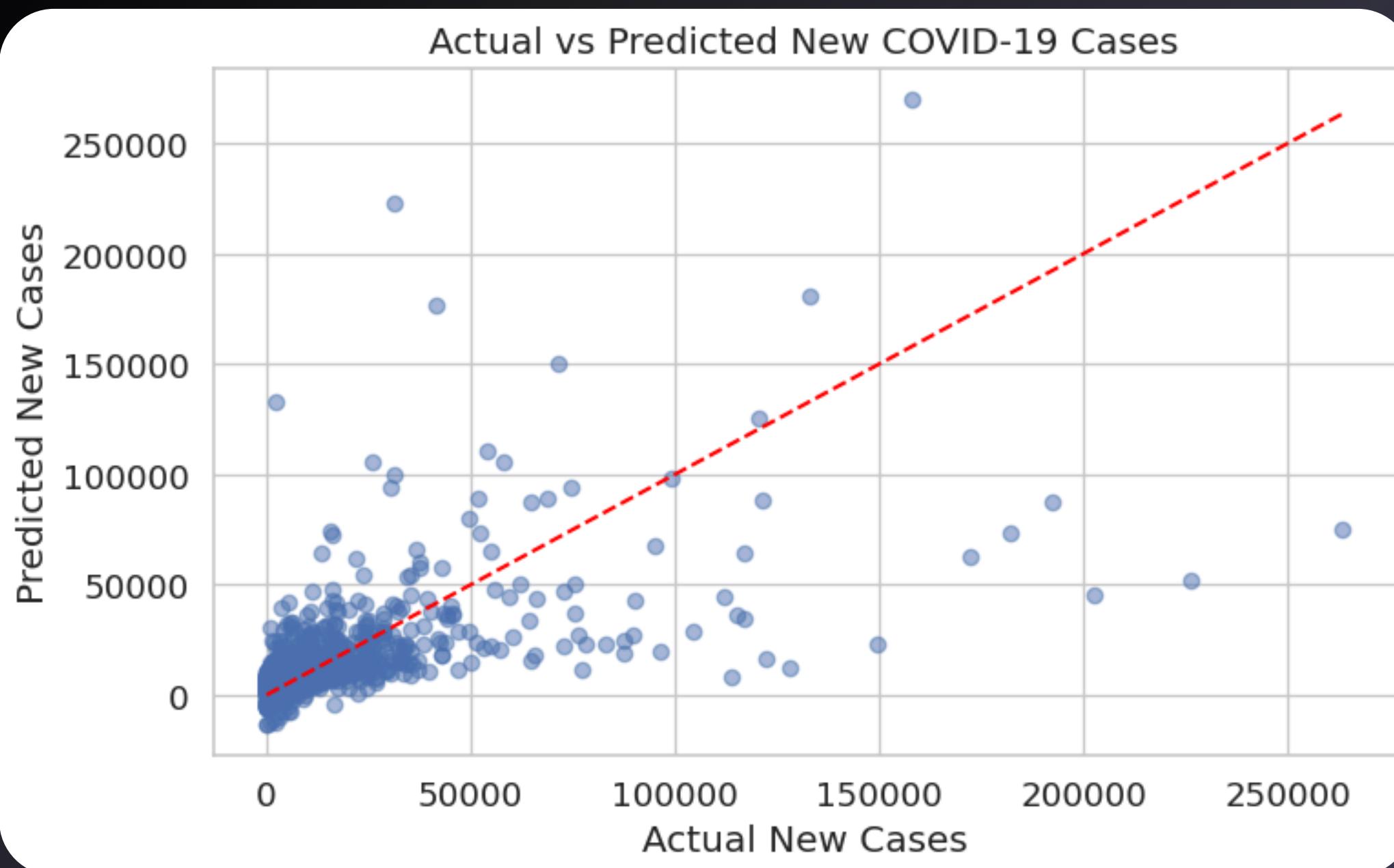


```
1 # Instantiate and Train Linear Regression Model
2 lr = LinearRegression()
3 lr.fit(X_train, y_train)
4
5 # Predict on Test Data
6 y_pred = lr.predict(X_test)
7
8 # Evaluate Model using R^2 and RMSE
9 r2 = r2_score(y_test, y_pred)
10 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
11
12 print(f"R2 Score: {r2:.3f}")
13 print(f"RMSE: {rmse:.2f}")
```

- Used **LinearRegression()** from scikit-learn to model the relationship between vaccine features, time (week), and new COVID-19 cases.
- Trained the model with **.fit()** and made predictions on test data using **.predict()**.
- Evaluated performance using **r2_score** and **mean_squared_error** (converted to RMSE with **np.sqrt**).
- Output metrics (R^2 , RMSE) indicate how accurately the model captures case trends.



Actual vs Predictions Scatter Plot



- The model explains about **29%** of case variation ($R^2 = 0.291$) with a moderate average error of **~18,283** cases (RMSE), reflecting limited predictive power due to real-world data variability.
- Most predictions are clustered near the bottom left (**0–50,000 range**), suggesting the model captures low/moderate case weeks better.
- However, it underperforms on high-case weeks, which drags down R^2 . This is common in linear models when data has outliers or a nonlinear pattern.

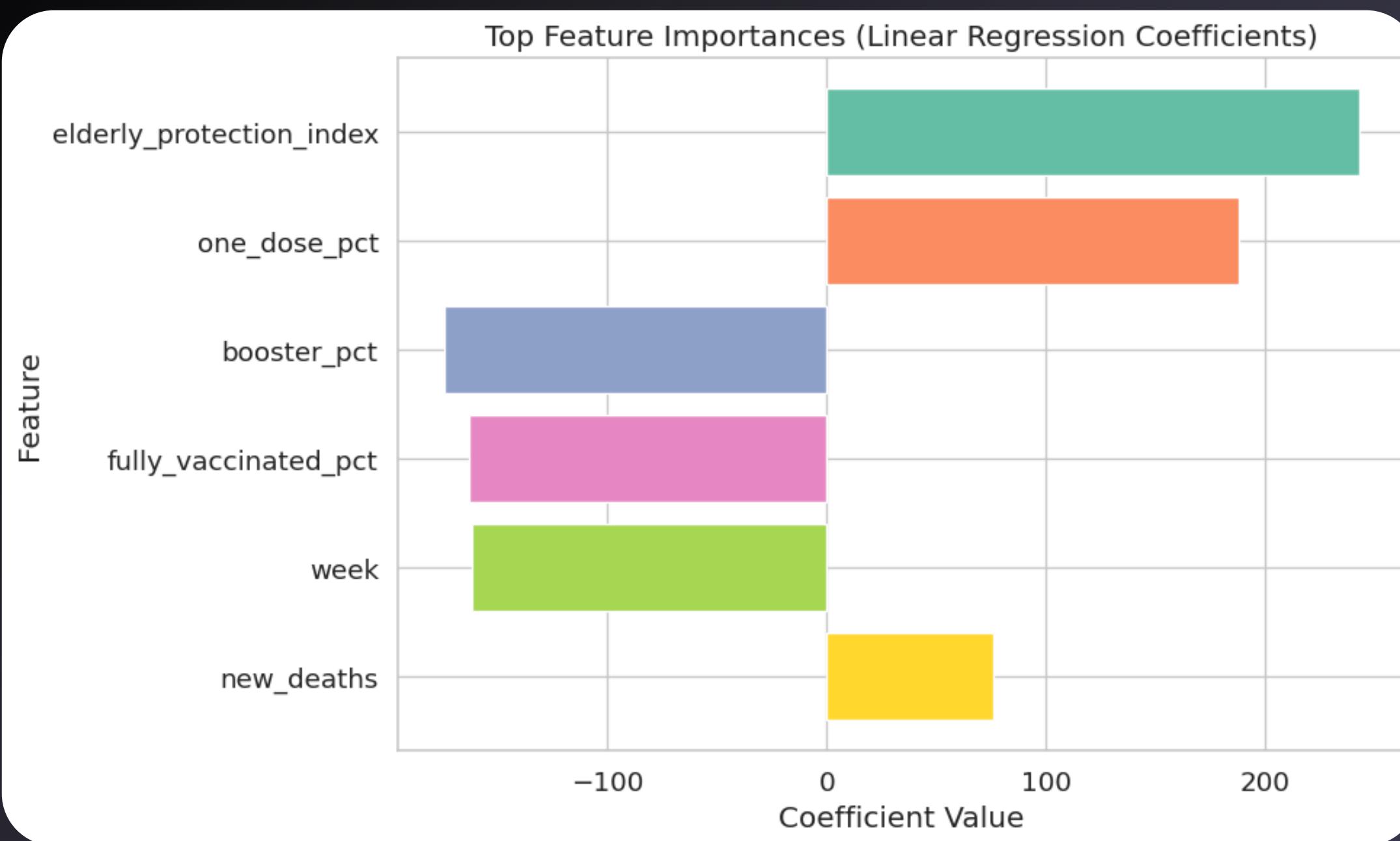


```
● ● ●  
1 plt.figure(figsize=(8, 5))  
2  
3 plt.scatter(y_test, y_pred, alpha=0.5)  
4  
5 plt.plot([0, max(y_test)], [0, max(y_test)], color='red', linestyle='--')  
6  
7 plt.xlabel("Actual New Cases")  
8 plt.ylabel("Predicted New Cases")  
9 plt.title("Actual vs Predicted New COVID-19 Cases")  
10 plt.tight_layout()  
11  
12 plt.show()
```

- **plt.scatter()** plots actual vs predicted cases with transparency (`alpha=0.5`) for better visibility of overlapping points.
- **plt.plot()** line adds a red dashed line representing perfect predictions where values match exactly.
- `plt.xlabel()`, `plt.ylabel()`, and `plt.title()` add descriptive labels and a title; `plt.tight_layout()` adjusts spacing; `plt.show()` displays the final plot.



Feature Importance (Coefficient Plot)



- **Booster** and **full vaccination** rates show **negative coefficients**, indicating they help reduce predicted COVID-19 cases.
- The **negative** week coefficient suggests cases tend to **decline** over time.
- **Positive** coefficients for **elderly protection**, **one-dose percentage**, and new deaths may reflect reactive trends or lagging indicators.



```
● ● ●  
1 # Create Coefficient DataFrame  
2 coef_df = pd.DataFrame({  
3     "Feature": X.columns,  
4     "Coefficient": lr.coef_  
5 }).sort_values(by="Coefficient", key=abs, ascending=False)  
6  
7 # Define Top N  
8 top_n = 10  
9 top_features = coef_df.head(top_n)  
10  
11 # Manually Assign Colors (length = top_n)  
12 bar_colors = sns.color_palette("Set2", n_colors=top_n)  
13  
14 # Plot  
15 plt.figure(figsize=(10, 6))  
16 bars = plt.barh(  
17     y=top_features["Feature"],  
18     width=top_features["Coefficient"],  
19     color=bar_colors # assign directly to matplotlib's barh  
20 )  
21  
22 plt.title("Top Feature Importances (Linear Regression Coefficients)")  
23 plt.xlabel("Coefficient Value")  
24 plt.ylabel("Feature")  
25 plt.gca().invert_yaxis() # Highest at top  
26 plt.tight_layout()  
27 plt.show()
```

- Creates a DataFrame **coef_df** to pair features with their linear regression **coefficients**, sorting by absolute value for importance.
- Selects the top features (**top_features**) for visualization and assigns distinct colors using **Seaborn's "Set2" palette**.
- Uses **plt.barh()** to draw a horizontal bar chart of feature coefficients, inverts the y-axis to show highest importance on top, and adds labels and title for clarity.



Conclusion

- Higher vaccination linked to fewer cases; booster rates varied.
- Elderly protection reduces deaths.
- CA, FL, TX had highest deaths; CA led 97 weeks.
- Early spikes included 1000% surge in March 2020.
- NYC saw early outbreaks; some states had minimal initial cases.
- TX, GA, AL, TN showed low vaccination and high deaths.

- Model predicts low/moderate case weeks well but misses large surges. Booster and full vaccination rates lower predicted cases. Cases decline over time.
- Vaccination and time features alone aren't sufficient for a robust model, as many other factors (e.g., variants, mobility, policies) influence case trends.

The analysis used state-level data without population normalization or socioeconomic factors, limiting depth; future work could include classification, time series, or causal methods.



Intro

Data Preprocessing

Feature Engineering

EDA

SQL

ML Model

Conclusion

Thank You!