# Analysis of Novel Coronavirus (COVID-19) data for United States
# ATDM

JULY 19

**St Clair College of Applied Arts and Technology**
**Authored by: Data Wizard- Group 01**

ST. CLAIR
COLLEGE

# Contents

## Project Guide

**Yingge Wang ([YWANG@stclaircollege.ca](mailto:YWANG@stclaircollege.ca))**

## Project Group Members

| Student Name | Student ID | Student Email |
|---|---|---|
| Harsh Parmar | 779093 | hp97@myscc.ca |
| Nishil Patel | 782837 | np85@myscc.ca |
| Raman Keshari | 783376 | rk409@myscc.ca |
| Shreyas Mahendra | 770049 | sm208@myscc.ca |
| Priyanka Harsukhbhai Ghetiya | 779258 | pg48@myscc.ca |

## Section Number- 002

## Academic Integrity

We, Harsh Parmar, Nishil Patel, Raman Keshari, Shreyas Mahendra and Priyanka Harsukhbhai Ghetiya, hereby state that we have not communicated with or gained information in any way from any person or resource that would violate the College's academic integrity policies, and that all work presented is our own. In addition, we also agree not to share our work in any way, before or after submission, that would violate the College's academic integrity policies.

# Analysis of Novel Coronavirus (COVID-19) data for United States

Start Date- 1st June 2021                                       End Date- 28th July 2021

**Objectives- Analysis of the dataset using EDA and visualization**

This project study aims to describe the health outcomes of people diagnosed with COVID-19 in the US, over time and in relation to the characteristics of the virus, by combining COVID-19 information from different states, hospital, general practice and death registry data.

**Background**

The novel coronavirus disease, named COVID-19 on 11 February 2020, is caused by SARS-CoV-2 virus. The outbreak was declared a Public Health Emergency of International Concern on 30 January 2020. While the number of confirmed cases worldwide and in the US is reported daily, detailed data on the outcomes of people who test positive for SARS-CoV-2, and predictors of outcomes, are still scarce.

Outcomes are likely to vary with context, including according to extensiveness of surveillance and testing, health systems functioning and population characteristics. Evidence to date has come primarily from various states of the US that are further along in the pandemic.

**Goal**

The aim of the project is to provide a list of priority areas for work and health research, that addresses evidence gaps and emerging evidence needs within the context of global pandemics generally and COVID19 specifically.

**Requirements**

The dataset is present on GitHub, which has been downloaded and loaded for our data analysis. Following are the technological requirement and tools used for this project.

| Technological Requirements |
| --- |
| Python 3.9 |
| Anaconda- Jupyter Notebook |
| Tableau |
| Excel |
| GitHub |
| MS office |

# Project- Submission

**DataSet**- **USA daily state reports for** **Novel Coronavirus (COVID-19)**
**Source**- https://github.com/CSSEGISandData/COVID-19/tree/master/csse_covid_19_data
**GitHub Link**- https://github.com/ShreyasM6/ATDM-Project-Repo

# Proposal

**We will be working with the publicly available dataset from JHU CSSE's COVID-19 data repository. The data repository consists of different datasets related to covid-19.**

## What is the problem you are solving?

Data in this dataset is a collection of the covid-19 cases in the US, in which various locations are taken into consideration. The aim of the project is to provide data analysis of covid-19 pandemic, various cases have been studied like most affected areas due to this pandemic. Study of data from various states is combined to show the growth of cases and recovery. Through this project we are trying to achieve a sustainable solution and understanding of the behavioral trend of the crisis and, a step towards helping people to understand the spread and predict the raise in cases for the US.
        The data analysis will provide a complete understanding of the critical information which will help address the business question that we are trying to answer. EDA and visualization are the most prominent way to perform analysis and obtain the required or desired details.

## Who benefits from this project?

By the analysis of the dataset the results help any Government or national body to examine the current situation and establish clear protocols to help the stop of the virus. The measures collected in the dataset will provide an accurate statistic to take critical actions in order to contain the spread of the virus. The analysis of the tends in the dataset will be beneficial to have a better idea in the behavioral patterns of the cases and help to mitigate the spike. With the help of the results determined by the analysis, the local bodies can help to suppress any serious conditions that may arise, as the result would give us clear directions of what precautionary measures to be implemented. Hence by carrying out any kind of analytics to provide insights about the dataset will help a number of different sectors and avoid similar situation in future.
        Analysis of the dataset using EDA and visualizing, will provide a way to convey information for any desired audience who has no prior knowledge about the dataset. The analysis using EDA and visualization can be taken advantage by the Government or national body to squeeze important details as much as possible.

## Why is the project important?

The pandemic has already taken grip over life of people. Since the start of the pandemic, some countries are facing problem of ever-increasing cases. Through the data analysis of cases one can analyze how the US is doing in terms of controlling the pandemic. Analyzing data leads to adapt the prevention model that are doing great in terms of lowering the cases, Predictions can be made with the dataset available to the individual/country/organizations, thus helping to decide how far we are able to control the pandemic or up to what extent we should guide preventive measures.

Analysis is more important than ever during these unprecedented times. As a society, we've seen how important even basic line graphs, bar charts, and heat maps are to understanding the spread of the virus. We've heard a lot about various models in terms of predicting deaths associated with COVID-19. Many people want to see the data and understand the facts in this rapidly changing environment.

## What are the business questions you are looking to answer or objectives you are looking to achieve from this project?

### Business questions-

Which state was the most affected from covid-19?
Which state has the highest deaths from covid-19?
When did the COVID-19 cases had the highest peak?
When did the COVID-19 cases had the highest death rate?
When did the COVID-19 cases started to drop?
When did the death rate started to drop?
With the EDA and visualization performed, did you observe any pattern in the dataset?
Top five states confirmed and deaths statistics in the US?
Did you observe any unusual behavior of the dataset?
Why do you think EDA and visualization important?

### Objectives-

With the help of the data we can estimate the cases and requirements that are essential for the treatment of patients effected by COVID 19. Improve on distribution and control the tools that are necessary in this critical times for example:- If a country is suffering from high amount of covid cases then the essentials are more required there and with the project we can understand how much tools are required for the stability of the situation.

The main objective is to describe the health outcomes of people diagnosed with COVID-19 in the US, over time and in relation to the characteristics of the virus, by combining COVID-19 information from different states, hospital, general practice and death registry data.

## Describe the data at a high level, explain the data collection process, source of data, etc

The dataset represents the detailed characteristics of the 2019 Novel Coronavirus, which provides information about the US and provinces with the number of confirmed, deaths and recovered cases. This list includes a complete list of all sources ever used in the data set since January 2020. Some sources listed here (e.g. ECDC, US CDC, BNO News) are not currently relied upon as a source of data. The data is collected from multiple sources such as World Health Organization (WHO) and different US data sources at the state (Admin1) or county/city (Admin2) level.

        The dataset which is used for analysis has categorical, continuous, and discrete data, the dataset has different data type like int, float, and object.

## What insights does this data give and how can it be used in future?

The analyzed data provides information of different states in the United states which consists of the confirmed, death and recovered cases based on a timeframe and the location. This information can be used by the state government to analyze the current situation and help to suppress the severity.

The dataset provides the insights of the day-to-day measures of increase and decrease of the cases, deaths, and recovery for the US. This kind of data is very useful in the current situation as it gives the government, the Ngo's, and other medical and social organization to analyze the covid affected areas for a better decision making.

This data is also future proof as it shows the pattern of transmission of the covid 19 and this kind of dataset can help the countries more resilient to any kind of pandemic in future.

## What type of problems are solved with this approach?

With analytics, we are trying to acknowledge the general public about the pandemic and providing solutions to overcome this critical situation. In this project we are approaching with an open mindset and providing solution for general public rather than businesses.

Any problem related to COVID-19 or information can be fetched by the following preformed analysis, the insights gained by this can hugely incorporated in any real word model to address the needs of the situation in any given time period.

## Can the performed analysis be implemented to any real word scenario?

Since the considered dataset is live data recorded for the US from 1st January 2020 to 13th July 2021, the performed analysis is a best fit for any real word scenario related to COVID-19. With the help of visuals provided below, decisions can be taken faster and act on the situation quicker than expected.

# Descriptive Statistics and Exploratory Data Analysis (EDA)

**Step 1**: **Confirm the data is correctly loaded**

Before loading the dataset to a dataframe, we need to verify and make sure if the dataset has consistent column names across all the files.

Importing the dataset can be achieved with the pandas packages in python, following is the Code to import a single file to dataframe-

```
import pandas as pd
covid=pd.read_csv("https://raw.githubusercontent.com/CSSEGISandData/COVID-19/master/csse_covid_19_data/csse_covid_19_daily_reports/01-01-2021.csv")
```

**Step 2**: **Describe the data**

Describe the data: The dataframe has 18 columns which gives information about the Province name and timeframe, the number of Confirmed cases, death toll and Recovered cases. It also provides statistical information which are derived columns (calculated columns) like Active cases, Incident_Rate and Case_Fatality_Ratio which will be helpful to perform more EDA on the dataset.

Summary statistics: Summary statistics can be performed in python using dataframe.describe() function. The describe function will provide the measures to describe the dataframe like- count, mean, std, min, max and much more.

Breakdown of variables (numeric, ordinal, categorical)

| Variable Type | Column Names |
|---|---|
| **Categorical** | **Province_State, Country_Region, ISO3** |
| **Numerical** | **Lat, Long, Confirmed, Deaths, Recovered, Active, Incident_Rate, Total_Test_Results, People_Hospitalized,Case_Fatality_Ratio, UID, Testing_Rate, Hospitalization_Rate** |
| **Date** | **Last_Update** |
| **Ordinal** | **FIPS** |

**Step 3**: **Check the validity of data**

Define schema: With the help of df.dtype function from the pandas package, we can determine the schema of the dataframe. Most of the variables are objects, float, or integer types.

| Variable Type | Column Names |
|---|---|
| Object | Province_State, Country_Region, ISO3 |
| Integer | Confirmed, Deaths, Recovered, Active, FIPS, Total_Test_Results, People_Hospitalized, UID |
| Float | Last_Update, Lat, Long, Incident_Rate, Case_Fatality_Ratio, Testing_Rate, Hospitalization_Rate |
| Date | Last_Update |

Understand the data: The primary information from the dataset is the daily cases of COVID-19 from US. Which provides statistics about the number of confirmed, death and recovered based on the province and timeframe.

**Step 4**:  Answer the following questions:

**1. Does the data include missing, incomplete, or invalid records?**
Yes, our dataset has missing and incomplete records. The missing and incomplete data is handled based on the requirement of the analysis to be performed. In our cases we are replacing the missing values with 0.

**2. Does your data include outliers?**
Yes, our dataset has outliers. Since the recorded data is live data, we can not eliminate the outliers as they also play an important role in providing the critical information.

**3. Is the data segmented into groups?**
Yes, the dataset is divided in groups based on the province, and timeframe. Grouped data will help us to provide information based on states with number of confirmed/death/recovered, and month/year.

**4. Is the data imbalanced (a large number of the records represent a majority class and very few records represent the minority class)?**
No, the dataset is not largely concentrated or representing a majority class. The data is equally divided for all the stated in the US.

**5. Are some data elements highly correlated with each other?**
Yes, the dataset has correlation with each other, as few columns are dependent and helps us to determine the trend of another column and also derive a calculated column.

**6. How was the data collected?**

The data is collected from multiple sources such as World Health Organization (WHO), European Centre for Disease Prevention and Control (ECDC) and different US data sources at the state or county/city level. The data was loaded using pandas library into a dataframe from github.

**7. What are the inclusion criteria for your data?**

Inclusion criteria are defined on the key target feature which includes analysis on demographic, and geographic characteristics.

**8. Can you generate preliminary visualizations for individual features?**

Using the dataset we can induce visualization to check which province has the highest number of deaths or highest number of confirmed cases. By the visualization created we can have more data insights about the dataset.

**Step 5**: Use visualization to understand and explore, but not to explain

By using Matplotlib package we can produce visualization to understand the behavior of the dataset. The following plot is a scatter plot displaying the relation between the number of confirmed and deaths. As we can see the number of confirmed and deaths are directly proportional.

```
plt.scatter(covid.Confirmed, covid.Deaths)
```

`<matplotlib.collections.PathCollection at 0x2472b6d1bb0>`

# Data Cleaning and Transformation

The COVID-19 data files considered for EDA is from 1$^{st}$ January 2020 to 13$^{th}$ July 2021, about 460 files close to 27 thousand rows and 20+ columns.

**Transformation using Python-**

1)  Load all the required libraries for analysis-
    Python provides a wide range of libraries for analysis which must be imported prior performing our analysis. In our EDA process we are importing the following libraries.

    1)numpy package for arithmetic operation.
    2)pandas package for dataframe handling.
    3)OS package for interacting with the operating system.
    4)glob package for retrieve files/pathnames.
    5)date package to work with date and time.

2)  Set the path from where the CSV files are picked-
    Using the glob package we can set the path where the COVID-19 data files are present to load them into a single dataframe. Also list all the files present in the directory using the OS package.

3)  Create a dataframe-
    Using the pandas package, create an empty dataframe. The created empty dataframe will be used to load all the CSV files dynamically in a single run.

4)  Load the files into a single DataFrame-
    Loading multiple CSV files into a dataframe is achieved in a single piece of code, which iterates between all the files and append the data to the dataframe. The for-loop operator helps us to iterate each file and load all the files to a temporary dataframe present in the directory. The append operator consolidates all the files loaded in the temporary dataframe to the newly created dataframe. To avoid ambiguity while performing EDA- we have created a new column which holds the filename for each row, which allows us an ability to trace back for verification and validity.

5)  Retrieve the information of the loaded dataframe-
    After loading the data into the dataframe, we need to check the dimension of the dataframe and the data type of each columns, which helps us to understand the loaded data and perform any modifications if required. The .info operator retrieves all the information of a dataframe describing the categorical and continuous variables.

6) Find the sum of NaN values for each column-
   It is important to know the number of NaN values present in the dataframe for each column before performing EDA. As the result will determine the way to handle the NaN values. In general, any real time dataset will contain many NaN values which requires us to eliminate handle them in an efficient manner which would not effect the original dataset.

7) Replace the NaN values-
   As part of EDA, the key aspect is the way of handling the NaN values. Part of this solution is to drop the NaN values, however dropping the NaN values in our dataset will wipe out half of the data, resulting in losing valuable insights from the dataset. Although dropping NaN values might be a solution in some cases but replacing them with 0 is the best practice.

8) Data Transformation-
   Substring the filename column to retrieve only the date information and exclude the file type. This transformation will help to create new columns and perform analysis based on a timeframe. It is very important to have a column which represent the timeframe of the data, which will let us to represent the data in a specific period of time.

9) Creating new columns from the existing columns-
   Derive columns from already existing columns which provide information about the timestamp (Month, Year).

10) Rename the columns for readability-
    Renaming the columns will help us to understand the column information better.

11) Write the transformed data into a CSV file-
    To avoid performing transformation repeatedly, we are writing the transformed data from the dataframe to a CSV file.

*All the mentioned above EDA and transformation is performed using python on Jupyter Notebook, the notebook has the required code along with headings and comments which has been uploaded on GitHub.*

# Data Analysis

The COVID-19 dataset is visualized in Tableau after transformation, to identify the key patterns and to answer the business questions.

## Key patterns and visualization-
**Load the transformed data to Tableau for visualization**

### *Confirmed cases in US*

Display of Cases in different states of US



Map based on Longitude (generated) and Latitude (generated). Size shows sum of Confirmed. Details are shown for Country and State.

The above geographical distribution of the cases in US displays the spread of COVID-19 virus. The larger the size of the circle, larger is the intensity.

**Statistical information about the confirmed cases in all the states of the US for each month starting from January 2020 to July 2021**

Confirmed cases across US in each month

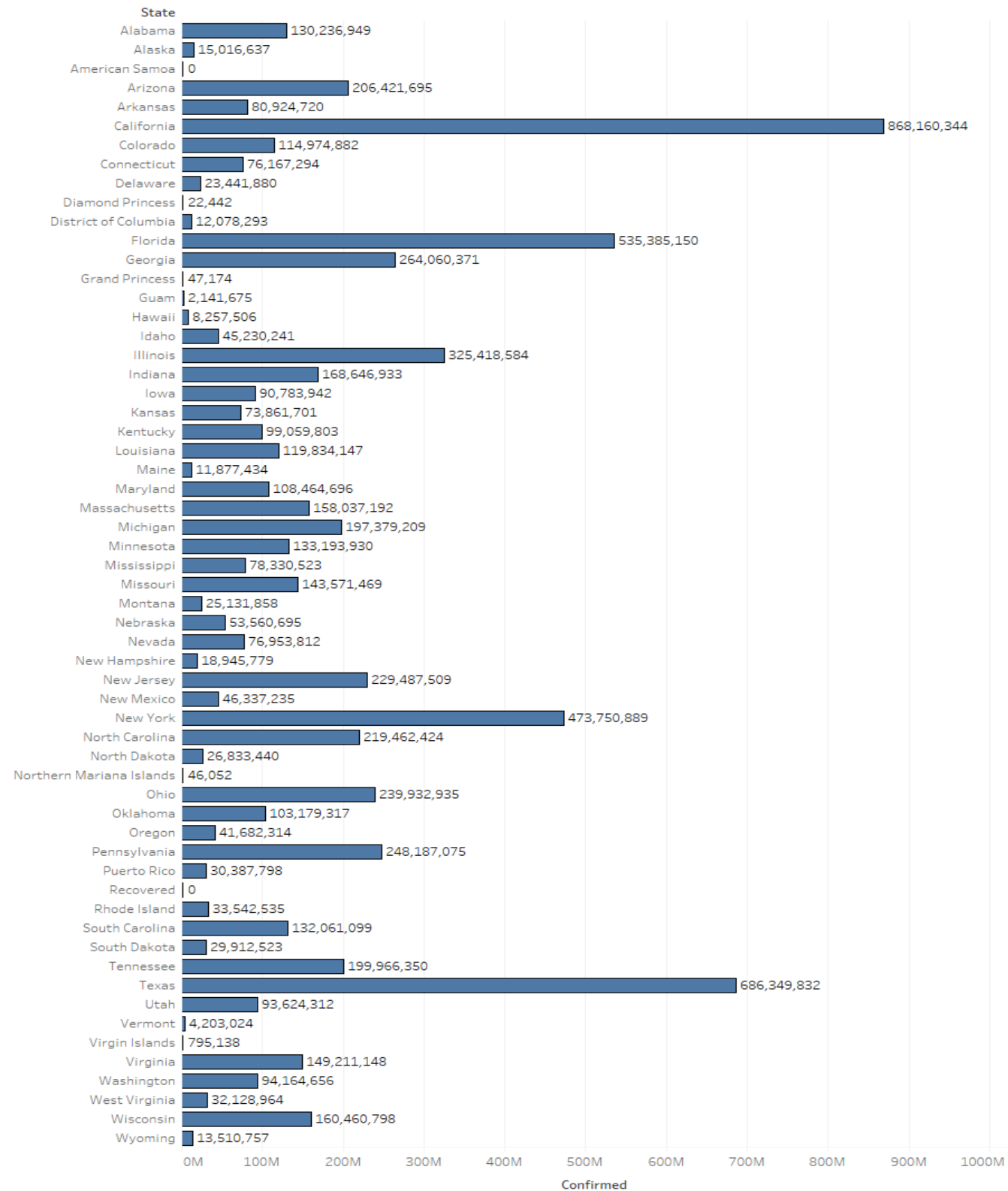| State | January 20... | February 2... | March 2020 | April 2020 | May 2020 | June 2020 | July 2020 | August 2020 | September ... | October 20... | November ... | December ... | January 20... | February 2... | March 2021 | April 2021 | May 2021 | June 2021 | July 2021 | August 2021 | September ... | October 20... | November ... | December ... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Alabama | 886,230 | 897,394 | 906,957 | 1,015,966 | 1,196,092 | 1,497,968 | 2,324,610 | 3,143,495 | 3,596,806 | 4,394,927 | 5,154,037 | 7,251,141 | 11,736,966 | 11,253,636 | 13,188,890 | 12,920,071 | 13,774,828 | 13,383,391 | 4,063,699 | 3,515,838 | 3,524,570 | 3,530,787 | 3,536,064 | 3,542,586 |
| Alaska | 68,198 | 69,635 | 71,174 | 78,611 | 82,159 | 90,082 | 119,427 | 168,035 | 206,441 | 320,204 | 579,507 | 926,784 | 1,436,307 | 1,353,786 | 1,605,192 | 1,637,563 | 1,765,155 | 1,715,208 | 511,057 | 440,632 | 441,617 | 442,473 | 443,167 | 444,223 |
| American S... | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Arizona | 1,295,371 | 1,307,282 | 1,322,956 | 1,433,912 | 1,651,367 | 2,348,339 | 4,238,727 | 5,144,883 | 5,265,166 | 5,899,430 | 6,804,576 | 10,326,188 | 18,986,348 | 18,512,668 | 21,530,395 | 21,057,843 | 22,307,215 | 21,702,256 | 6,589,741 | 5,705,174 | 5,724,316 | 5,740,962 | 5,756,291 | 5,770,289 |
| Arkansas | 493,834 | 500,128 | 507,158 | 556,190 | 627,120 | 818,267 | 1,209,648 | 1,605,763 | 1,929,982 | 2,510,110 | 3,155,568 | 4,471,794 | 7,525,616 | 7,266,199 | 8,462,852 | 8,235,398 | 8,688,740 | 8,467,809 | 2,600,684 | 2,247,334 | 2,253,756 | 2,259,031 | 2,262,468 | 2,269,271 |
| California | 4,674,958 | 4,726,594 | 4,775,967 | 5,494,215 | 6,636,697 | 8,369,348 | 13,094,395 | 17,708,149 | 19,430,427 | 22,293,271 | 25,535,120 | 42,989,730 | 83,778,142 | 80,812,574 | 93,697,308 | 91,560,130 | 96,338,998 | 93,215,694 | 28,538,207 | 24,763,483 | 24,837,548 | 24,897,808 | 24,960,069 | 25,031,512 |
| Colorado | 598,879 | 607,328 | 618,035 | 826,385 | 1,084,134 | 1,201,002 | 1,454,425 | 1,709,575 | 1,862,117 | 2,427,470 | 4,201,316 | 6,659,001 | 10,509,904 | 9,972,713 | 11,840,508 | 12,196,758 | 13,472,760 | 13,266,181 | 3,859,101 | 3,304,659 | 3,313,880 | 3,321,492 | 3,329,657 | 3,337,602 |
| Connecticut | 468,595 | 475,411 | 481,683 | 861,427 | 1,234,718 | 1,314,074 | 1,419,922 | 1,486,127 | 1,518,724 | 1,765,493 | 2,364,446 | 3,730,409 | 6,563,748 | 6,467,619 | 7,781,080 | 8,082,014 | 8,676,326 | 8,380,861 | 2,457,021 | 2,117,578 | 2,121,282 | 2,123,566 | 2,133,108 | 2,142,062 |
| Delaware | 140,701 | 141,898 | 143,558 | 201,617 | 306,244 | 342,703 | 413,519 | 471,311 | 508,934 | 603,392 | 725,975 | 1,143,003 | 2,029,550 | 1,998,191 | 2,401,912 | 2,484,088 | 2,700,763 | 2,623,218 | 763,893 | 655,468 | 657,507 | 659,357 | 661,663 | 663,415 |
| Diamond Pr... | 392 | 392 | 392 | 1,274 | 1,323 | 1,274 | 1,323 | 1,274 | 1,323 | 1,274 | 1,372 | 1,274 | 1,274 | 1,127 | 1,274 | 1,225 | 1,274 | 1,225 | 392 | 343 | 343 | 343 | 343 | 343 |
| District of C... | 104,398 | 104,910 | 106,076 | 163,391 | 255,666 | 289,128 | 327,773 | 367,247 | 380,995 | 427,417 | 473,977 | 627,483 | 959,347 | 932,293 | 1,123,832 | 1,144,948 | 1,225,913 | 1,188,275 | 351,737 | 303,037 | 304,044 | 304,708 | 305,379 | 306,319 |
| Florida | 3,886,345 | 3,924,617 | 3,962,654 | 4,504,186 | 4,976,666 | 5,952,994 | 11,313,284 | 15,440,749 | 16,535,168 | 18,831,453 | 21,073,603 | 27,539,349 | 44,808,692 | 43,874,907 | 52,327,764 | 53,463,443 | 57,955,756 | 56,580,773 | 16,630,955 | 14,266,555 | 14,327,548 | 14,361,453 | 14,406,856 | 14,439,380 |
| Georgia | 1,780,088 | 1,797,731 | 1,845,448 | 2,228,907 | 2,665,863 | 3,099,292 | 4,827,423 | 6,703,748 | 7,479,415 | 8,553,647 | 10,024,010 | 13,453,203 | 23,177,413 | 22,784,347 | 26,937,051 | 26,669,722 | 28,354,923 | 27,476,810 | 8,270,849 | 7,149,852 | 7,170,644 | 7,186,555 | 7,202,386 | 7,221,044 |
| Grand Princ... | 824 | 824 | 824 | 2,678 | 2,781 | 2,678 | 2,781 | 2,678 | 2,781 | 2,678 | 2,884 | 2,678 | 2,678 | 2,575 | 2,678 | 2,575 | 2,678 | 2,575 | 824 | 721 | 721 | 721 | 721 | 721 |
| Guam | 16,395 | 16,498 | 16,908 | 19,573 | 20,415 | 21,429 | 23,900 | 34,053 | 57,263 | 94,362 | 132,972 | 156,448 | 197,987 | 178,476 | 203,007 | 197,303 | 209,539 | 204,512 | 63,718 | 55,316 | 55,346 | 55,365 | 55,411 | 55,479 |
| Hawaii | 58,982 | 59,626 | 60,434 | 71,602 | 74,152 | 77,164 | 92,436 | 184,965 | 272,127 | 337,903 | 379,337 | 464,115 | 702,737 | 663,930 | 785,763 | 812,257 | 897,662 | 892,355 | 258,791 | 220,948 | 221,584 | 222,161 | 222,887 | 223,588 |
| Idaho | 275,641 | 279,178 | 283,766 | 319,870 | 341,563 | 374,939 | 608,516 | 867,339 | 1,002,485 | 1,393,751 | 1,961,032 | 2,830,222 | 4,239,700 | 3,942,338 | 4,615,134 | 4,582,095 | 4,870,509 | 4,740,824 | 1,441,524 | 1,247,707 | 1,250,187 | 1,251,975 | 1,253,585 | 1,256,361 |
| Illinois | 2,201,529 | 2,227,621 | 2,255,914 | 2,928,475 | 4,299,506 | 4,815,153 | 5,552,304 | 6,572,547 | 7,461,479 | 9,433,604 | 14,123,670 | 19,826,819 | 29,354,737 | 27,436,438 | 31,945,923 | 32,291,918 | 34,804,901 | 33,740,829 | 10,147,071 | 8,768,055 | 8,785,162 | 8,799,800 | 8,814,019 | 8,831,110 |
| Indiana | 909,082 | 923,321 | 939,402 | 1,182,727 | 1,545,656 | 1,753,711 | 2,121,789 | 2,644,065 | 3,058,586 | 4,019,805 | 6,292,587 | 10,032,332 | 16,182,220 | 15,203,577 | 17,631,165 | 17,540,627 | 18,794,293 | 18,293,745 | 5,536,320 | 4,788,793 | 4,799,494 | 4,809,282 | 4,817,795 | 4,826,559 |
| Iowa | 622,701 | 631,264 | 638,885 | 719,243 | 971,086 | 1,136,552 | 1,439,274 | 1,764,010 | 2,165,906 | 2,851,235 | 4,461,099 | 5,845,991 | 8,309,459 | 7,729,617 | 8,998,506 | 8,920,303 | 9,454,317 | 9,140,144 | 2,801,599 | 2,429,664 | 2,433,877 | 2,436,814 | 2,439,320 | 2,443,076 |
| Kansas | 412,383 | 419,288 | 424,436 | 479,385 | 606,197 | 669,087 | 897,869 | 1,179,725 | 1,460,266 | 1,950,750 | 2,974,515 | 4,434,499 | 7,123,682 | 6,741,314 | 7,795,936 | 7,614,753 | 8,039,980 | 7,795,418 | 2,397,722 | 2,082,977 | 2,084,623 | 2,087,576 | 2,093,023 | 2,096,297 |
| Kentucky | 472,212 | 479,481 | 488,289 | 555,584 | 664,458 | 763,130 | 990,120 | 1,342,758 | 1,662,195 | 2,315,389 | 3,358,275 | 5,231,222 | 9,287,454 | 9,183,845 | 10,859,288 | 10,765,160 | 11,515,936 | 11,223,954 | 3,353,811 | 2,893,907 | 2,902,937 | 2,910,008 | 2,916,520 | 2,923,870 |
| Louisiana | 982,057 | 993,384 | 1,001,507 | 1,458,530 | 1,708,721 | 1,947,590 | 2,905,889 | 3,737,752 | 3,964,154 | 4,444,093 | 4,992,750 | 6,610,394 | 10,217,106 | 9,793,122 | 11,396,940 | 11,203,524 | 11,922,689 | 11,652,389 | 3,543,563 | 3,059,054 | 3,063,533 | 3,071,420 | 3,076,851 | 3,087,135 |
| Maine | 39,372 | 39,856 | 40,527 | 57,398 | 78,207 | 95,914 | 113,893 | 125,399 | 136,406 | 161,118 | 229,621 | 425,953 | 1,038,047 | 1,055,377 | 1,285,620 | 1,414,844 | 1,631,681 | 1,605,606 | 437,485 | 370,031 | 371,737 | 373,081 | 374,248 | 376,013 |
| Maryland | 815,735 | 823,236 | 830,545 | 1,109,051 | 1,684,310 | 2,015,743 | 2,393,529 | 2,835,506 | 3,050,743 | 3,516,141 | 4,140,882 | 5,768,984 | 9,189,911 | 8,824,169 | 10,437,975 | 10,724,514 | 11,523,718 | 11,147,611 | 3,306,182 | 2,850,280 | 2,858,190 | 2,865,645 | 2,872,854 | 2,879,242 |
| Massachus... | 1,042,708 | 1,052,816 | 1,056,305 | 1,857,306 | 2,788,438 | 3,016,525 | 3,270,687 | 3,488,640 | 3,422,196 | 3,939,830 | 4,826,214 | 7,469,468 | 13,529,498 | 13,338,408 | 15,977,281 | 16,454,168 | 17,661,384 | 17,066,136 | 5,023,369 | 4,325,819 | 4,339,603 | 4,351,860 | 4,363,047 | 4,375,486 |
| Michigan | 1,108,000 | 1,126,050 | 1,141,504 | 1,764,395 | 2,195,885 | 2,410,979 | 2,801,774 | 3,240,536 | 3,582,749 | 4,521,380 | 7,179,191 | 10,847,166 | 16,752,880 | 15,686,477 | 18,839,459 | 21,629,947 | 24,256,566 | 23,540,080 | 6,586,168 | 5,603,477 | 5,618,805 | 5,633,614 | 5,645,816 | 5,666,311 |
| Minnesota | 772,381 | 784,445 | 798,349 | 863,256 | 1,186,751 | 1,435,336 | 1,772,633 | 2,179,377 | 2,523,029 | 3,352,420 | 5,670,498 | 8,503,878 | 12,254,026 | 11,332,270 | 13,317,185 | 13,797,635 | 15,052,491 | 14,590,302 | 4,321,818 | 3,730,806 | 3,738,321 | 3,743,804 | 3,750,506 | — |
| Mississippi | 568,126 | 574,461 | 580,509 | 674,740 | 835,004 | 1,004,961 | 1,499,208 | 2,070,485 | 2,310,000 | 2,773,722 | 3,187,992 | 4,393,888 | 6,993,536 | 6,699,141 | 7,810,209 | 7,635,857 | 8,069,289 | 7,831,596 | 2,394,078 | 2,075,865 | 2,080,847 | 2,085,223 | 2,088,660 | 2,093,126 |
| Missouri | 836,665 | 851,654 | 864,269 | 990,723 | 1,118,581 | 1,242,624 | 1,654,782 | 2,374,300 | 3,087,908 | 4,223,497 | 6,124,768 | 8,776,725 | 13,702,279 | 12,873,941 | 14,912,283 | 14,563,741 | 15,437,825 | 15,132,332 | 4,646,146 | 4,014,509 | 4,024,502 | 4,032,106 | 4,039,250 | 4,046,059 |
| Montana | 123,778 | 126,392 | 128,849 | 139,307 | 142,784 | 149,534 | 193,494 | 262,295 | 337,335 | 636,590 | 1,126,630 | 1,626,870 | 2,442,703 | 2,293,885 | 2,678,615 | 2,656,172 | 2,832,492 | 2,760,788 | 837,615 | 724,787 | 726,212 | 727,169 | 728,210 | 729,352 |
| Nebraska | 347,673 | 353,209 | 358,804 | 401,741 | 590,675 | 697,355 | 828,833 | 990,443 | 1,141,113 | 1,566,309 | 2,428,511 | 3,407,216 | 4,960,172 | 4,621,936 | 5,369,267 | 5,359,720 | 5,681,714 | 5,483,970 | 1,678,275 | 1,454,756 | 1,456,962 | 1,458,675 | 1,460,647 | 1,462,719 |
| Nevada | 488,487 | 494,350 | 500,828 | 580,514 | 658,026 | 769,777 | 1,260,593 | 1,762,084 | 1,918,010 | 2,319,811 | 2,965,369 | 4,495,580 | 7,160,703 | 6,735,268 | 7,814,773 | 7,707,453 | 8,218,515 | 8,042,077 | 2,448,118 | 2,112,564 | 2,117,854 | 2,122,945 | 2,126,687 | 2,133,430 |
| New Hamps... | 67,940 | 69,193 | 70,757 | 99,956 | 148,878 | 174,431 | 195,675 | 212,061 | 222,075 | 271,150 | 392,372 | 811,844 | 1,725,521 | 1,738,044 | 2,109,315 | 2,237,147 | 2,438,537 | 2,363,353 | 677,460 | 579,817 | 582,055 | 583,957 | 586,169 | 588,072 |
| New Jersey | 1,672,817 | 1,685,032 | 1,698,794 | 3,390,937 | 4,613,023 | 4,796,857 | 5,170,721 | 5,420,013 | 5,470,939 | 6,243,428 | 7,728,094 | 11,104,943 | 18,356,945 | 18,249,217 | 22,555,162 | 23,742,683 | 25,303,461 | 24,457,668 | 7,116,801 | 6,105,588 | 6,124,962 | 6,143,089 | 6,160,399 | 6,175,936 |
| New Mexico | 246,157 | 249,971 | 254,106 | 298,028 | 387,440 | 459,181 | 610,651 | 732,716 | 780,688 | 1,037,642 | 1,719,028 | 2,799,156 | 4,480,993 | 4,226,766 | 4,915,506 | 4,835,116 | 5,147,150 | 5,003,706 | 1,524,967 | 1,320,818 | 1,323,248 | 1,325,874 | 1,327,831 | 1,330,496 |
| New York | 3,567,872 | 3,587,682 | 3,608,449 | 8,257,831 | 10,439,927 | 10,658,092 | 11,459,029 | 11,864,836 | 11,896,649 | 13,148,292 | 14,617,654 | 20,679,828 | 37,033,918 | 37,664,106 | 46,489,016 | 48,571,575 | 52,148,233 | 50,447,830 | 14,594,538 | 12,519,660 | 12,564,827 | 12,604,247 | 12,645,169 | 12,681,629 |
| North Carol... | 1,257,828 | 1,270,016 | 1,287,633 | 1,436,877 | 1,736,180 | 2,297,156 | 3,350,733 | 4,282,199 | 4,922,354 | 6,219,091 | 7,430,900 | 10,692,679 | 19,425,656 | 19,530,141 | 23,174,603 | 23,285,187 | 25,068,654 | 24,367,053 | 7,208,491 | 6,207,511 | 6,226,471 | 6,245,215 | 6,258,904 | 6,280,892 |
| North Dako... | 172,865 | 175,227 | 178,432 | 193,462 | 226,088 | 247,107 | 293,090 | 378,544 | 531,514 | 871,888 | 1,484,159 | 1,922,661 | 2,555,299 | 2,309,932 | 2,658,615 | 2,636,599 | 2,805,066 | 2,715,891 | 837,103 | 726,906 | 727,336 | 727,971 | 728,549 | 729,136 |
| Northern M... | 442 | 445 | 447 | 692 | 858 | 999 | 1,192 | 1,484 | 1,670 | 2,127 | 2,339 | 2,727 | 3,558 | 3,321 | 4,077 | 4,022 | 4,473 | 4,390 | 1,281 | 1,099 | 1,100 | 1,100 | 1,104 | 1,105 |
| Ohio | 1,128,692 | 1,143,994 | 1,163,991 | 1,415,574 | 1,788,033 | 2,043,413 | 2,739,728 | 3,445,773 | 3,921,040 | 4,964,557 | 7,652,342 | 13,415,953 | 23,064,621 | 22,143,289 | 25,954,172 | 25,988,109 | 27,772,103 | 26,931,826 | 8,096,882 | 6,999,604 | 7,016,563 | 7,031,204 | 7,047,194 | 7,064,278 |
| Oklahoma | 533,531 | 540,512 | 546,797 | 606,616 | 675,220 | 766,195 | 1,113,691 | 1,583,976 | 2,009,412 | 2,739,723 | 3,722,551 | 5,654,231 | 9,880,282 | 9,628,543 | 11,222,067 | 10,963,366 | 11,525,883 | 11,156,971 | 3,417,627 | 2,963,237 | 2,969,991 | 2,979,782 | 2,986,351 | 2,992,762 |
| Oregon | 217,520 | 220,737 | 223,681 | 264,457 | 304,584 | 360,327 | 531,351 | 709,734 | 808,220 | 1,028,197 | 1,411,624 | 2,238,792 | 3,768,296 | 3,623,284 | 4,260,497 | 4,388,101 | 4,953,268 | 4,918,674 | 1,407,690 | 1,201,025 | 1,205,321 | 1,208,829 | 1,212,134 | 1,215,852 |
| Pennsylvan... | 1,228,051 | 1,244,633 | 1,262,873 | 1,934,678 | 2,627,659 | 2,878,602 | 3,397,426 | 3,880,159 | 4,187,349 | 5,093,174 | 7,030,382 | 12,386,603 | 22,243,440 | 21,750,683 | 26,039,261 | 27,361,376 | 29,908,293 | 29,039,938 | 8,410,638 | 7,214,396 | 7,237,786 | 7,257,038 | 7,276,346 | 7,296,291 |
| Puerto Rico | 233,326 | 237,341 | 239,926 | 264,170 | 302,610 | 366,339 | 504,682 | 809,845 | 990,091 | 1,370,053 | 1,071,496 | 1,553,787 | 2,493,008 | 2,376,611 | 2,786,897 | 3,056,908 | 3,414,022 | 3,316,134 | 943,171 | 806,239 | 809,476 | 811,485 | 813,706 | 816,475 |
| Recovered | | | | 0 | | | | | | | | | 0 | | | | | | | | | | | |
| Rhode Island | 197,524 | 200,643 | 203,417 | 308,463 | 464,503 | 504,046 | 557,669 | 615,403 | 648,665 | 784,572 | 1,099,589 | 1,755,546 | 2,991,792 | 2,897,214 | 3,467,909 | 3,550,592 | 3,804,284 | 3,678,896 | 1,089,859 | 940,664 | 941,881 | 943,063 | 947,244 | 949,097 |
| South Carol... | 812,028 | 819,833 | 829,126 | 923,502 | 1,032,612 | 1,321,485 | 2,279,452 | 2,986,233 | 3,400,856 | 4,075,762 | 4,593,194 | 6,214,353 | 11,330,757 | 11,623,497 | 13,919,458 | 13,919,297 | 14,858,986 | 14,377,679 | 4,265,091 | 3,675,744 | 3,686,514 | 3,696,420 | 3,704,966 | 3,714,254 |
| South Dako... | 188,373 | 191,081 | 194,316 | 229,180 | 282,938 | 316,503 | 360,667 | 424,272 | 560,743 | 914,200 | 1,518,443 | 2,034,427 | 2,832,798 | 2,590,588 | 3,016,591 | 3,000,260 | 3,167,943 | 3,054,816 | 940,735 | 816,843 | 817,725 | 818,862 | 819,656 | 820,563 |
| Tennessee | 1,183,329 | 1,198,626 | 1,211,108 | 1,368,067 | 1,609,448 | 1,895,558 | 2,885,424 | 3,986,291 | 4,634,860 | 5,872,038 | 7,413,347 | 11,394,444 | 18,778,721 | 17,743,356 | 20,735,902 | 20,597,237 | 21,876,401 | 21,144,158 | 6,439,146 | 5,577,717 | 5,588,415 | 5,602,017 | 5,609,378 | 5,621,362 |
| Texas | 4,420,876 | 4,471,737 | 4,527,314 | 4,966,274 | 5,654,862 | 6,837,305 | 11,560,901 | 16,032,596 | 18,038,039 | 21,707,507 | 26,459,118 | 35,888,780 | 60,647,134 | 60,791,637 | 71,036,327 | 70,214,615 | 74,586,616 | 72,537,976 | 21,827,195 | 18,848,549 | 18,900,849 | 18,944,074 | 18,989,957 | 19,043,510 |
| Utah | 519,424 | 527,537 | 536,158 | 606,651 | 707,257 | 882,562 | 1,243,938 | 1,505,123 | 1,758,428 | 2,502,223 | 3,719,817 | 5,455,280 | 8,909,186 | 8,465,834 | 9,870,110 | 9,712,302 | 10,303,031 | 10,051,547 | 3,061,115 | 2,646,760 | 2,652,541 | 2,657,701 | 2,661,907 | 2,667,880 |
| Vermont | 14,318 | 14,612 | 14,844 | 29,730 | 33,395 | 36,404 | 41,541 | 45,290 | 46,902 | 54,393 | 79,535 | 142,486 | 329,715 | 357,220 | 467,341 | 528,350 | 585,452 | 567,914 | 154,440 | 130,718 | 131,239 | 131,737 | 132,412 | 133,036 |
| Virgin Islan... | 6,042 | 6,084 | 6,124 | 7,167 | 7,550 | 7,644 | 12,294 | 23,819 | 29,372 | 32,030 | 33,418 | 43,323 | 64,249 | 61,880 | 74,226 | 75,670 | 84,219 | 88,739 | 24,980 | 21,112 | 21,210 | 21,244 | 21,308 | 21,421 |
| Virginia | 909,095 | 917,597 | 926,698 | 1,121,437 | 1,616,290 | 2,003,837 | 2,498,924 | 3,107,399 | 3,533,732 | 4,237,114 | 4,928,048 | 7,009,430 | 12,865,567 | 13,084,728 | 15,618,111 | 15,807,424 | 16,900,950 | 16,360,126 | 4,833,623 | 4,161,431 | 4,174,592 | 4,186,719 | 4,198,370 | 4,209,906 |
| Washington | 582,281 | 588,504 | 594,550 | 824,946 | 976,581 | 1,145,768 | 1,577,115 | 1,984,876 | 2,169,932 | 2,616,749 | 3,372,206 | 5,003,010 | 8,209,573 | 7,905,303 | 9,349,916 | 9,625,013 | 10,729,745 | 10,666,686 | 3,066,114 | 2,620,985 | 2,627,998 | 2,633,255 | 2,642,446 | 2,651,074 |
| West Virgin... | 113,137 | 115,260 | 117,333 | 135,756 | 153,849 | 170,653 | 229,038 | 303,435 | 389,362 | 537,522 | 853,007 | 1,557,040 | 3,114,743 | 3,038,066 | 3,596,607 | 3,672,668 | 4,010,331 | 3,922,017 | 1,145,780 | 984,487 | 988,052 | 991,002 | 993,646 | 996,173 |
| Wisconsin | 968,590 | 982,131 | 999,228 | 1,105,647 | 1,308,715 | 1,501,065 | 1,917,439 | 2,401,284 | 2,969,468 | 4,806,756 | 7,781,688 | 10,571,906 | 15,341,968 | 14,178,315 | 16,368,157 | 16,171,909 | 17,172,195 | 16,599,712 | 5,102,420 | 4,430,559 | 4,437,349 | 4,442,542 | 4,447,593 | 4,454,162 |
| Wyoming | 63,244 | 64,648 | 65,936 | 74,635 | 83,368 | 92,013 | 114,506 | 139,305 | 165,097 | 272,072 | 574,545 | 878,020 | 1,345,085 | 1,248,688 | 1,448,396 | 1,426,462 | 1,523,926 | 1,501,987 | 455,364 | 393,473 | 393,924 | 394,335 | 395,274 | 396,454 |

Sum of Confirmed broken down by Month Year (MY) vs. State.

While delivering the results found during statistics, the report has to be precise and concise, with the help of the above tabular representation of the data for the total confirmed cases for each state in a month, provides complete information in a single report. Anybody looking at this report can acquire knowledge about the number of cases confirmed in the US, in addition this report provides specific information with states and timeframe.

**Bar Plot representation of the confirmed cases in US**

### Confirmed cases count for each states

| State | Confirmed |
|---|---|
| Alabama | 130,236,949 |
| Alaska | 15,016,637 |
| American Samoa | 0 |
| Arizona | 206,421,695 |
| Arkansas | 80,924,720 |
| California | 868,160,344 |
| Colorado | 114,974,882 |
| Connecticut | 76,167,294 |
| Delaware | 23,441,880 |
| Diamond Princess | 22,442 |
| District of Columbia | 12,078,293 |
| Florida | 535,385,150 |
| Georgia | 264,060,371 |
| Grand Princess | 47,174 |
| Guam | 2,141,675 |
| Hawaii | 8,257,506 |
| Idaho | 45,230,241 |
| Illinois | 325,418,584 |
| Indiana | 168,646,933 |
| Iowa | 90,783,942 |
| Kansas | 73,861,701 |
| Kentucky | 99,059,803 |
| Louisiana | 119,834,147 |
| Maine | 11,877,434 |
| Maryland | 108,464,696 |
| Massachusetts | 158,037,192 |
| Michigan | 197,379,209 |
| Minnesota | 133,193,930 |
| Mississippi | 78,330,523 |
| Missouri | 143,571,469 |
| Montana | 25,131,858 |
| Nebraska | 53,560,695 |
| Nevada | 76,953,812 |
| New Hampshire | 18,945,779 |
| New Jersey | 229,487,509 |
| New Mexico | 46,337,235 |
| New York | 473,750,889 |
| North Carolina | 219,462,424 |
| North Dakota | 26,833,440 |
| Northern Mariana Islands | 46,052 |
| Ohio | 239,932,935 |
| Oklahoma | 103,179,317 |
| Oregon | 41,682,314 |
| Pennsylvania | 248,187,075 |
| Puerto Rico | 30,387,798 |
| Recovered | 0 |
| Rhode Island | 33,542,535 |
| South Carolina | 132,061,099 |
| South Dakota | 29,912,523 |
| Tennessee | 199,966,350 |
| Texas | 686,349,832 |
| Utah | 93,624,312 |
| Vermont | 4,203,024 |
| Virgin Islands | 795,138 |
| Virginia | 149,211,148 |
| Washington | 94,164,656 |
| West Virginia | 32,128,964 |
| Wisconsin | 160,460,798 |
| Wyoming | 13,510,757 |

0M    100M   200M   300M   400M   500M   600M   700M   800M   900M   1000M

Confirmed

Sum of Confirmed for each State.

**Line plot representing the trend of confirmed cases over the timeframe**

## Confirmed Trend



The trend of sum of Confirmed for Month Year Month.

The above line graph displays the information of the rise and fall of the COVID-19 cases in US. Line graph helps us understand the key insights of the cases confirmed in the given timeframe, information like when was the spike and decrease in cases can be fetched without investigating the actual numbers.

**Dashboard for Confirmed cases**

## Display of Cases in different states of US



© 2021 Mapbox © OpenStreetMap

## Confirmed cases across US in each month
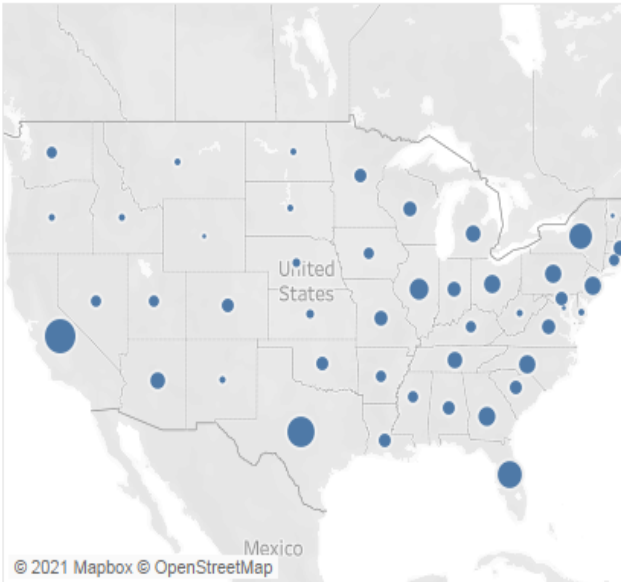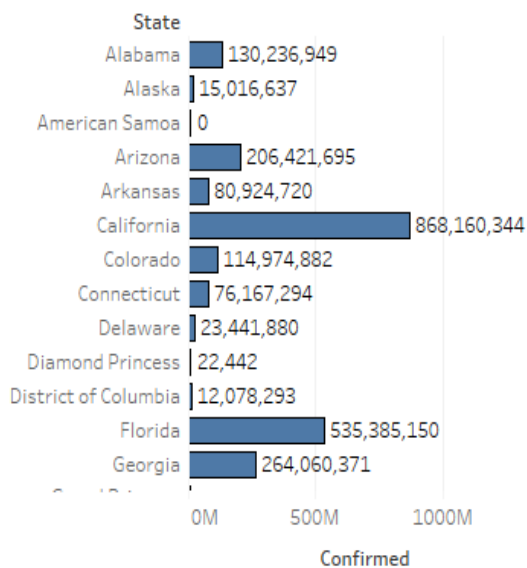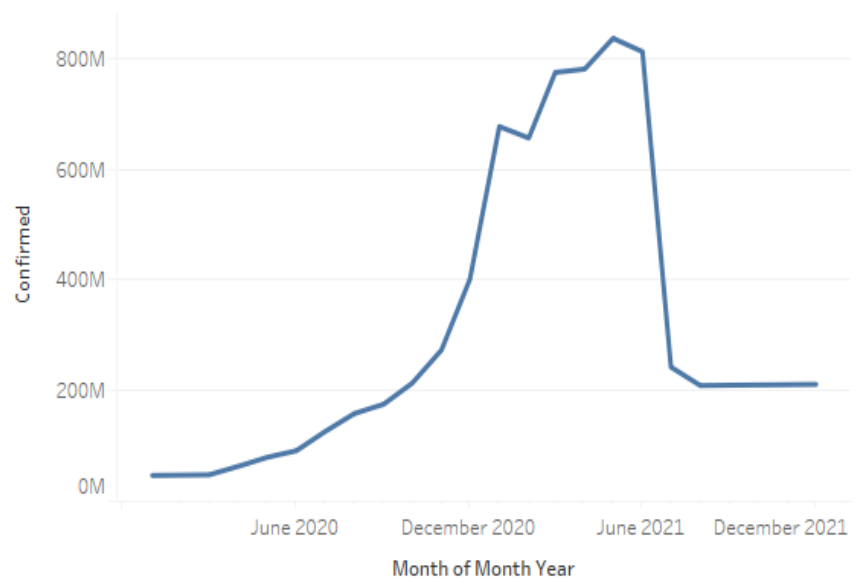
| | Month Year | | | | | |
| State | January 20.. | February 2.. | March 2020 | April 2020 | May 2020 | June 2020 |
|---|---|---|---|---|---|---|
| Alabama | 886,230 | 897,394 | 906,957 | 1,015,966 | 1,196,092 | 1,497,9 |
| Alaska | 68,198 | 69,635 | 71,174 | 78,611 | 82,159 | 90,0 |
| American S.. | 0 | 0 | 0 | 0 | 0 | |
| Arizona | 1,295,371 | 1,307,282 | 1,322,956 | 1,433,912 | 1,651,367 | 2,348,3 |
| Arkansas | 493,834 | 500,128 | 507,158 | 556,190 | 627,120 | 818,2 |
| California | 4,674,958 | 4,726,594 | 4,775,967 | 5,494,215 | 6,636,697 | 8,369,3 |
| Colorado | 598,879 | 607,328 | 618,035 | 826,385 | 1,084,134 | 1,201,0 |
| Connecticut | 468,595 | 475,411 | 481,683 | 861,427 | 1,234,718 | 1,314,0 |
| Delaware | 140,701 | 141,898 | 143,558 | 201,617 | 306,244 | 342,7 |
| Diamond Pr.. | 392 | 392 | 392 | 1,274 | 1,323 | 1,2 |
| District of C.. | 104,398 | 104,910 | 106,076 | 163,391 | 255,666 | 289,1 |
| Florida | 3,886,345 | 3,924,617 | 3,962,654 | 4,504,186 | 4,976,666 | 5,952,9 |
| Georgia | 1,780,088 | 1,797,731 | 1,845,448 | 2,228,907 | 2,665,863 | 3,099,2 |
| Grand Princ.. | 824 | 824 | 824 | 2,678 | 2,781 | 2,6 |

## Confirmed cases count for each states



State
Alabama 130,236,949
Alaska 15,016,637
American Samoa 0
Arizona 206,421,695
Arkansas 80,924,720
California 868,160,344
Colorado 114,974,882
Connecticut 76,167,294
Delaware 23,441,880
Diamond Princess 22,442
District of Columbia 12,078,293
Florida 535,385,150
Georgia 264,060,371

## Confirmed Trend



Dashboard helps us to consolidate many reports into a single frame, providing multiple aspect of information at a time. Dashboards can provide us quick information hazel free, in the above dashboard we can fetch any can kind of information with respect to the confirmed cases for any state in the US.

*Death cases in US*

## Display of Deaths in different states of US



Map based on Longitude (generated) and Latitude (generated). Size shows sum of Deaths. Details are shown for Country and State.

The above geographical distribution of deaths in US displays the severity of COVID-19 virus. The larger the size of the circle, larger is the intensity.
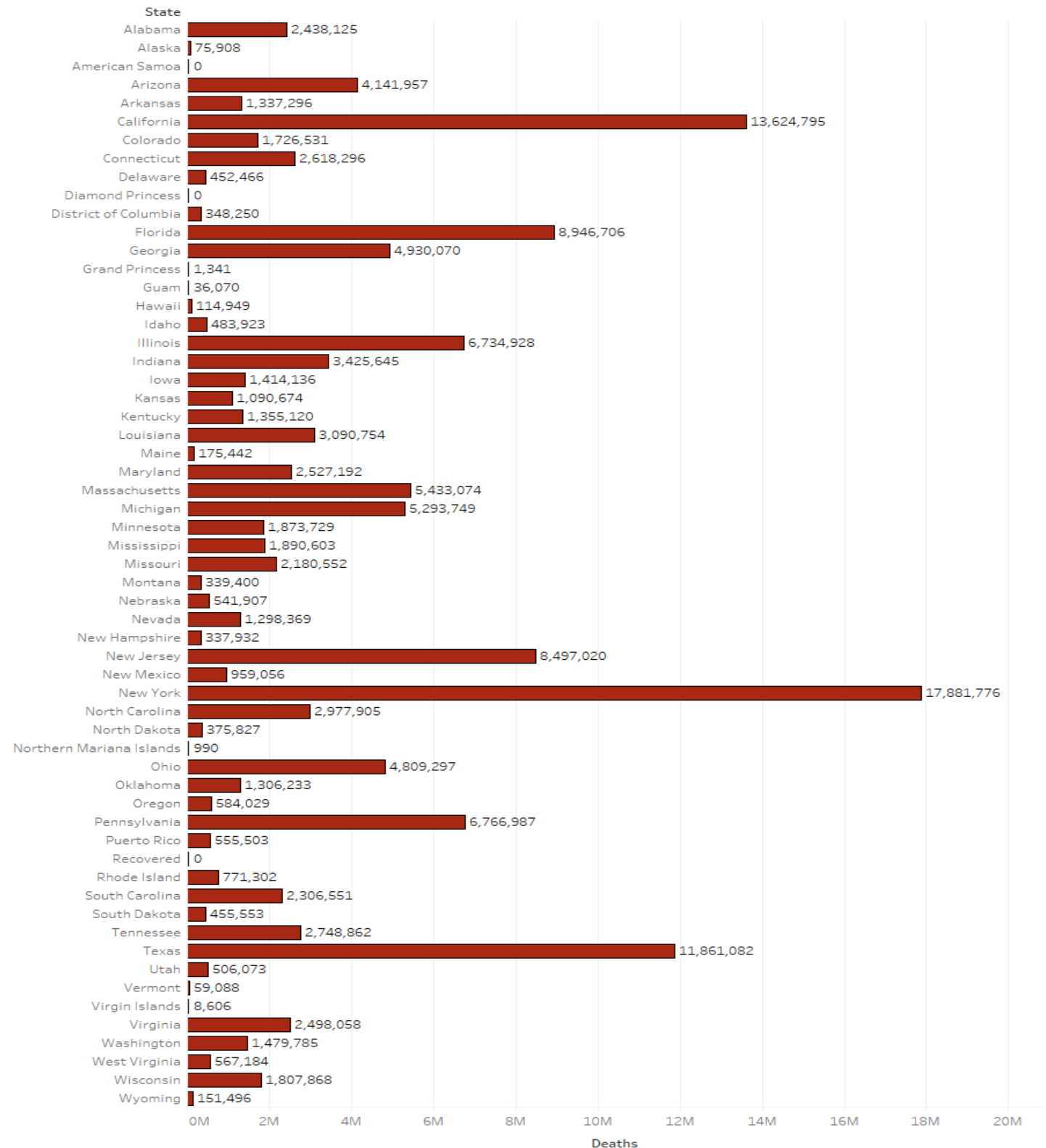
## Statistical information about the death cases in all the states of the US for each month starting from January 2020 to July 2021

Deaths across US in each month

| State | January 20.. | February 2.. | March 2020 | April 2020 | May 2020 | June 2020 | July 2020 | August 2020 | September .. | October 20.. | November .. | December .. | January 20.. | February 2.. | March 2021 | April 2021 | May 2021 | June 2021 | July 2021 | August 2021 | September .. | October 20.. | November .. | December .. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Alabama | 14,869 | 15,004 | 15,136 | 18,586 | 25,672 | 30,867 | 41,634 | 54,015 | 60,091 | 69,984 | 78,143 | 101,694 | 192,041 | 219,867 | 265,631 | 262,077 | 278,287 | 271,439 | 79,367 | 68,122 | 68,375 | 68,788 | 69,007 | 69,429 |
| Alaska | 362 | 363 | 375 | 548 | 580 | 617 | 761 | 1,006 | 1,267 | 1,730 | 2,389 | 4,082 | 6,881 | 6,808 | 8,100 | 8,244 | 9,150 | 8,937 | 2,596 | 2,215 | 2,219 | 2,219 | 2,225 | 2,234 |
| American S.. | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Arizona | 30,104 | 30,281 | 30,598 | 34,793 | 45,610 | 56,314 | 88,100 | 121,575 | 130,999 | 143,526 | 148,673 | 186,645 | 334,711 | 355,558 | 427,128 | 418,779 | 442,736 | 431,212 | 128,290 | 110,499 | 110,843 | 111,279 | 111,557 | 112,147 |
| Arkansas | 7,600 | 7,690 | 7,800 | 8,704 | 10,147 | 12,172 | 15,493 | 21,213 | 30,223 | 42,131 | 50,730 | 71,734 | 123,044 | 122,399 | 142,497 | 139,973 | 147,780 | 143,319 | 43,492 | 37,633 | 37,739 | 37,828 | 37,928 | 38,027 |
| California | 87,971 | 88,656 | 89,298 | 114,926 | 159,061 | 191,345 | 245,448 | 321,140 | 366,254 | 421,149 | 432,006 | 539,048 | 1,057,784 | 1,164,826 | 1,463,196 | 1,468,400 | 1,559,219 | 1,514,387 | 439,932 | 377,217 | 378,846 | 380,096 | 381,520 | 383,070 |
| Colorado | 15,217 | 15,345 | 15,527 | 25,096 | 40,329 | 45,443 | 49,582 | 52,273 | 52,591 | 58,353 | 66,568 | 100,713 | 146,037 | 136,464 | 157,760 | 154,519 | 166,076 | 164,023 | 49,625 | 42,825 | 42,929 | 43,007 | 43,073 | 43,156 |
| Connecticut | 33,698 | 33,864 | 33,924 | 60,486 | 102,509 | 111,050 | 118,228 | 119,383 | 115,712 | 121,766 | 122,356 | 143,793 | 182,444 | 173,825 | 202,368 | 198,160 | 209,658 | 202,602 | 62,087 | 53,942 | 53,978 | 54,039 | 54,171 | 54,253 |
| Delaware | 4,729 | 4,754 | 4,781 | 7,050 | 13,088 | 14,828 | 16,082 | 16,545 | 16,446 | 18,322 | 19,225 | 23,703 | 33,947 | 32,454 | 39,469 | 39,341 | 41,897 | 40,825 | 12,183 | 10,506 | 10,534 | 10,567 | 10,586 | 10,604 |
| Diamond Pr.. | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| District of C.. | 4,402 | 4,422 | 4,442 | 6,787 | 12,364 | 14,106 | 15,515 | 15,979 | 15,747 | 16,794 | 16,660 | 18,858 | 23,728 | 23,092 | 27,155 | 26,951 | 28,616 | 27,771 | 8,399 | 7,267 | 7,279 | 7,290 | 7,306 | 7,320 |
| Florida | 75,632 | 76,260 | 76,828 | 93,518 | 119,067 | 135,980 | 182,026 | 273,703 | 323,094 | 389,271 | 405,871 | 477,776 | 697,255 | 700,662 | 847,527 | 848,342 | 917,299 | 901,677 | 264,444 | 226,821 | 227,582 | 228,007 | 228,767 | 229,297 |
| Georgia | 40,207 | 40,435 | 41,134 | 55,940 | 75,662 | 89,810 | 105,802 | 138,985 | 163,176 | 189,999 | 208,911 | 243,661 | 371,466 | 391,016 | 478,499 | 482,991 | 517,959 | 509,784 | 147,848 | 126,628 | 127,062 | 127,351 | 127,598 | 128,146 |
| Grand Princ.. | 24 | 24 | 24 | 48 | 81 | 78 | 81 | 81 | 78 | 81 | 78 | 81 | 78 | 69 | 78 | 75 | 78 | 75 | 24 | 21 | 21 | 21 | 21 | 21 |
| Guam | 273 | 273 | 274 | 369 | 381 | 380 | 387 | 428 | 928 | 1,615 | 2,132 | 2,592 | 3,350 | 3,013 | 3,477 | 3,382 | 3,581 | 3,444 | 1,082 | 941 | 942 | 942 | 942 | 942 |
| Hawaii | 753 | 757 | 776 | 1,004 | 1,124 | 1,118 | 1,281 | 1,745 | 2,974 | 4,645 | 4,993 | 6,215 | 9,778 | 9,998 | 11,771 | 11,691 | 12,517 | 12,348 | 3,659 | 3,147 | 3,153 | 3,165 | 3,167 | 3,170 |
| Idaho | 2,863 | 2,905 | 2,953 | 3,878 | 4,497 | 4,662 | 5,686 | 8,998 | 11,238 | 13,887 | 18,526 | 28,230 | 44,777 | 42,598 | 50,251 | 49,834 | 52,906 | 51,819 | 15,637 | 13,507 | 13,535 | 13,561 | 13,573 | 13,602 |
| Illinois | 62,939 | 63,531 | 64,100 | 92,270 | 153,446 | 185,229 | 209,596 | 220,247 | 224,292 | 251,178 | 283,856 | 386,349 | 550,198 | 519,949 | 605,702 | 594,695 | 636,392 | 623,168 | 188,772 | 163,302 | 163,534 | 163,757 | 164,063 | 164,363 |
| Indiana | 26,806 | 27,057 | 27,299 | 42,244 | 64,057 | 72,580 | 81,337 | 88,575 | 92,150 | 108,608 | 133,016 | 192,566 | 293,679 | 285,260 | 333,650 | 326,457 | 345,460 | 336,760 | 102,520 | 88,764 | 88,964 | 89,104 | 89,256 | 89,476 |
| Iowa | 8,989 | 9,117 | 9,263 | 11,065 | 17,512 | 21,908 | 25,059 | 29,413 | 32,917 | 40,337 | 49,930 | 78,860 | 122,647 | 123,917 | 146,461 | 144,792 | 153,163 | 148,884 | 44,921 | 38,847 | 38,913 | 38,999 | 39,048 | 39,174 |
| Kansas | 4,722 | 4,897 | 4,934 | 6,966 | 8,986 | 9,846 | 11,438 | 13,437 | 16,431 | 23,224 | 30,992 | 51,855 | 100,187 | 104,945 | 124,021 | 120,991 | 128,173 | 124,557 | 37,359 | 32,390 | 32,409 | 32,548 | 32,659 | 32,707 |
| Kentucky | 7,568 | 7,611 | 7,730 | 10,920 | 15,051 | 17,491 | 20,957 | 24,514 | 28,267 | 34,294 | 40,221 | 54,558 | 101,564 | 109,168 | 145,670 | 153,398 | 165,980 | 168,126 | 46,094 | 38,930 | 39,037 | 39,191 | 39,326 | 39,454 |
| Louisiana | 34,842 | 35,074 | 35,225 | 61,535 | 85,400 | 91,880 | 105,017 | 125,322 | 133,204 | 146,744 | 149,387 | 173,442 | 228,412 | 219,195 | 258,445 | 253,957 | 268,385 | 261,009 | 79,395 | 68,759 | 68,825 | 68,970 | 69,081 | 69,249 |
| Maine | 1,006 | 1,018 | 1,024 | 1,722 | 2,494 | 2,889 | 3,287 | 3,523 | 3,579 | 3,845 | 4,285 | 6,832 | 14,987 | 15,578 | 18,800 | 18,768 | 20,425 | 20,323 | 5,829 | 5,012 | 5,044 | 5,046 | 5,055 | 5,071 |
| Maryland | 26,878 | 27,161 | 27,329 | 39,312 | 69,240 | 82,969 | 92,634 | 98,119 | 98,352 | 105,830 | 108,525 | 133,817 | 186,227 | 180,836 | 212,889 | 212,723 | 231,367 | 232,867 | 67,868 | 58,215 | 58,341 | 58,454 | 58,546 | 58,693 |
| Massachus.. | 66,765 | 67,098 | 67,446 | 105,478 | 185,190 | 209,786 | 229,542 | 238,485 | 237,152 | 255,971 | 258,893 | 295,759 | 379,397 | 366,589 | 435,497 | 429,720 | 453,184 | 438,529 | 133,286 | 115,444 | 115,640 | 115,884 | 116,059 | 116,280 |
| Michigan | 53,364 | 53,692 | 54,016 | 103,279 | 152,234 | 164,718 | 176,653 | 182,038 | 181,953 | 198,166 | 216,680 | 290,772 | 407,381 | 383,866 | 444,814 | 451,075 | 505,199 | 500,507 | 145,734 | 124,954 | 125,344 | 125,529 | 125,788 | 125,993 |
| Minnesota | 14,734 | 14,908 | 15,085 | 18,593 | 31,361 | 40,873 | 45,991 | 50,060 | 52,460 | 60,889 | 75,093 | 110,805 | 163,237 | 151,517 | 178,179 | 176,179 | 189,449 | 185,695 | 55,908 | 48,271 | 48,490 | 48,573 | 48,654 | 48,725 |
| Mississippi | 16,451 | 16,577 | 16,709 | 20,141 | 28,349 | 34,617 | 44,317 | 60,000 | 68,610 | 79,044 | 83,866 | 103,184 | 154,795 | 151,606 | 179,082 | 175,791 | 185,355 | 180,010 | 54,615 | 47,269 | 47,422 | 47,495 | 47,568 | 47,730 |
| Missouri | 14,261 | 14,459 | 14,600 | 19,065 | 27,595 | 32,440 | 37,708 | 43,341 | 50,245 | 67,578 | 81,518 | 115,843 | 186,946 | 187,378 | 225,155 | 222,059 | 237,973 | 233,308 | 69,386 | 59,610 | 59,794 | 59,842 | 60,099 | 60,349 |
| Montana | 1,476 | 1,510 | 1,538 | 1,761 | 1,879 | 1,955 | 2,440 | 3,397 | 4,544 | 7,027 | 12,458 | 18,654 | 31,470 | 31,365 | 36,908 | 37,941 | 40,517 | 39,850 | 11,794 | 10,136 | 10,176 | 10,189 | 10,198 | 10,217 |
| Nebraska | 3,477 | 3,533 | 3,580 | 4,340 | 6,306 | 8,156 | 9,642 | 10,964 | 12,031 | 15,007 | 20,093 | 32,903 | 50,186 | 47,265 | 55,469 | 54,785 | 57,537 | 55,384 | 17,015 | 14,789 | 14,813 | 14,850 | 14,869 | 14,913 |
| Nevada | 8,867 | 8,960 | 9,104 | 12,325 | 16,435 | 18,069 | 22,614 | 31,917 | 37,280 | 42,599 | 46,009 | 63,630 | 109,400 | 111,952 | 133,085 | 131,879 | 140,328 | 136,666 | 40,745 | 35,084 | 35,221 | 35,295 | 35,400 | 35,505 |
| New Hamps.. | 2,999 | 3,026 | 3,045 | 3,893 | 6,838 | 9,290 | 10,787 | 11,273 | 11,050 | 12,146 | 12,403 | 16,014 | 26,681 | 26,700 | 31,471 | 31,324 | 33,837 | 33,001 | 9,788 | 8,432 | 8,456 | 8,478 | 8,487 | 8,513 |
| New Jersey | 115,696 | 116,053 | 116,447 | 204,725 | 323,388 | 361,666 | 416,819 | 421,013 | 408,276 | 428,162 | 421,210 | 471,000 | 562,375 | 532,465 | 627,162 | 622,382 | 663,018 | 643,317 | 195,018 | 168,709 | 169,031 | 169,364 | 169,537 | 169,897 |
| New Mexico | 5,928 | 6,004 | 6,091 | 7,422 | 11,897 | 14,711 | 17,577 | 20,430 | 21,813 | 24,742 | 31,298 | 48,842 | 84,372 | 84,029 | 100,292 | 98,712 | 105,802 | 104,304 | 30,940 | 26,645 | 26,712 | 26,780 | 26,823 | 26,890 |
| New York | 253,429 | 253,792 | 254,519 | 591,971 | 802,767 | 817,172 | 875,360 | 881,981 | 854,975 | 896,595 | 875,854 | 966,776 | 1,136,408 | 1,087,389 | 1,281,920 | 1,273,281 | 1,350,731 | 1,307,476 | 396,678 | 343,407 | 343,975 | 344,520 | 345,121 | 345,679 |
| North Carol.. | 20,731 | 20,927 | 21,194 | 25,851 | 36,153 | 44,365 | 54,862 | 69,606 | 81,472 | 99,877 | 112,804 | 142,368 | 240,010 | 252,725 | 304,979 | 304,092 | 326,154 | 320,612 | 93,852 | 80,545 | 80,764 | 81,110 | 81,291 | 81,561 |
| North Dako.. | 2,159 | 2,198 | 2,235 | 2,500 | 3,279 | 3,738 | 4,221 | 4,956 | 6,148 | 10,845 | 17,646 | 26,147 | 37,345 | 33,836 | 38,744 | 37,616 | 39,609 | 38,410 | 11,977 | 10,428 | 10,437 | 10,442 | 10,453 | 10,458 |
| Northern M.. | 16 | 16 | 30 | 66 | 68 | 67 | 71 | 54 | 52 | 54 | 52 | 56 | 52 | 46 | 52 | 50 | 52 | 50 | 16 | 14 | 14 | 14 | 14 | 14 |
| Ohio | 31,134 | 31,482 | 31,826 | 45,289 | 68,163 | 77,792 | 89,839 | 103,929 | 111,709 | 129,354 | 160,400 | 264,874 | 441,952 | 424,024 | 493,031 | 483,431 | 512,564 | 497,464 | 151,538 | 131,395 | 131,650 | 131,897 | 132,152 | 132,408 |
| Oklahoma | 6,448 | 6,553 | 6,625 | 9,641 | 12,576 | 13,515 | 15,894 | 20,638 | 24,304 | 30,234 | 36,345 | 50,219 | 97,447 | 104,020 | 128,490 | 158,136 | 170,744 | 169,946 | 46,648 | 39,277 | 39,541 | 39,594 | 39,656 | 39,742 |
| Oregon | 3,442 | 3,477 | 3,520 | 4,989 | 6,359 | 7,076 | 8,908 | 11,689 | 13,394 | 16,033 | 18,747 | 29,661 | 51,184 | 50,678 | 60,930 | 60,634 | 66,011 | 65,860 | 19,125 | 16,354 | 16,415 | 16,444 | 16,492 | 16,607 |
| Pennsylvan.. | 57,357 | 57,773 | 58,163 | 86,086 | 153,064 | 175,668 | 194,461 | 203,891 | 204,692 | 224,224 | 237,725 | 333,089 | 552,145 | 545,126 | 640,511 | 634,530 | 682,247 | 667,024 | 198,658 | 171,267 | 171,789 | 172,144 | 172,461 | 172,892 |
| Puerto Rico | 3,662 | 3,715 | 3,747 | 5,021 | 6,237 | 6,565 | 7,505 | 11,114 | 14,961 | 18,928 | 22,300 | 30,553 | 48,103 | 46,431 | 54,824 | 55,333 | 61,958 | 60,826 | 17,668 | 15,134 | 15,179 | 15,210 | 15,241 | 15,288 |
| Recovered | | | | 0 | | | | | | | | 0 | | | | | | | | | | | | |
| Rhode Island | 8,070 | 8,133 | 8,187 | 12,276 | 22,069 | 25,347 | 27,711 | 28,656 | 28,547 | 31,128 | 32,786 | 42,854 | 60,092 | 57,347 | 66,998 | 65,483 | 69,077 | 66,788 | 20,488 | 17,789 | 17,822 | 17,850 | 17,889 | 17,915 |
| South Carol.. | 17,770 | 17,922 | 18,045 | 20,853 | 26,447 | 30,412 | 43,502 | 65,603 | 77,068 | 90,229 | 95,766 | 114,341 | 181,952 | 191,661 | 231,004 | 228,831 | 244,031 | 236,702 | 70,298 | 60,566 | 60,661 | 60,810 | 60,979 | 61,098 |
| South Dako.. | 2,098 | 2,154 | 2,211 | 2,422 | 3,255 | 3,861 | 4,673 | 5,446 | 6,149 | 9,171 | 16,496 | 28,823 | 45,239 | 42,881 | 49,691 | 48,371 | 51,188 | 49,708 | 15,271 | 13,247 | 13,268 | 13,296 | 13,311 | 13,323 |
| Tennessee | 14,517 | 14,662 | 14,933 | 17,996 | 21,422 | 25,093 | 32,718 | 45,112 | 56,750 | 74,196 | 92,857 | 134,633 | 244,800 | 255,653 | 300,872 | 296,007 | 313,565 | 304,176 | 91,417 | 79,076 | 79,269 | 79,491 | 79,687 | 79,960 |
| Texas | 83,518 | 84,168 | 85,025 | 96,466 | 116,042 | 130,045 | 183,563 | 313,759 | 369,920 | 434,171 | 474,538 | 593,544 | 962,821 | 988,799 | 1,214,955 | 1,210,867 | 1,288,800 | 1,255,666 | 370,593 | 319,008 | 319,878 | 320,746 | 321,576 | 322,614 |
| Utah | 3,014 | 3,053 | 3,089 | 3,683 | 4,924 | 6,020 | 8,176 | 10,579 | 11,308 | 14,040 | 17,568 | 25,678 | 43,493 | 43,528 | 53,309 | 53,190 | 57,402 | 56,200 | 16,535 | 14,173 | 14,220 | 14,245 | 14,292 | 14,354 |
| Vermont | 464 | 467 | 470 | 1,188 | 1,499 | 1,477 | 1,544 | 1,584 | 1,527 | 1,588 | 1,619 | 2,703 | 4,684 | 4,651 | 5,696 | 5,887 | 6,357 | 6,142 | 1,797 | 1,540 | 1,547 | 1,548 | 1,551 | 1,558 |
| Virgin Islan.. | 103 | 103 | 105 | 156 | 220 | 216 | 235 | 319 | 454 | 506 | 524 | 550 | 638 | 583 | 664 | 664 | 697 | 719 | 218 | 186 | 186 | 186 | 187 | 187 |
| Virginia | 19,562 | 19,720 | 19,886 | 26,135 | 41,871 | 49,635 | 59,368 | 67,230 | 74,974 | 87,572 | 91,601 | 110,603 | 177,845 | 184,173 | 257,372 | 256,408 | 275,880 | 270,004 | 77,078 | 65,830 | 66,123 | 66,244 | 66,377 | 66,567 |
| Washington | 14,070 | 14,148 | 14,223 | 27,288 | 34,661 | 37,035 | 42,998 | 49,700 | 51,510 | 57,725 | 61,778 | 75,357 | 112,473 | 112,504 | 133,926 | 133,035 | 143,870 | 141,146 | 41,845 | 35,864 | 36,009 | 36,109 | 36,166 | 36,345 |
| West Virgin.. | 2,132 | 2,169 | 2,199 | 2,673 | 3,594 | 3,921 | 4,279 | 5,775 | 8,129 | 10,441 | 14,271 | 24,664 | 52,103 | 52,841 | 66,210 | 66,646 | 69,682 | 68,685 | 20,044 | 17,197 | 17,263 | 17,314 | 17,355 | 17,597 |
| Wisconsin | 10,791 | 10,930 | 11,093 | 15,579 | 20,910 | 24,885 | 28,169 | 32,208 | 34,545 | 44,855 | 68,493 | 104,301 | 165,482 | 159,690 | 187,329 | 184,071 | 197,486 | 195,449 | 58,481 | 50,410 | 50,568 | 50,645 | 50,699 | 50,799 |
| Wyoming | 486 | 500 | 524 | 614 | 753 | 885 | 1,036 | 1,227 | 1,452 | 1,881 | 3,886 | 7,563 | 15,221 | 15,219 | 17,802 | 17,305 | 18,231 | 17,908 | 5,428 | 4,682 | 4,714 | 4,714 | 4,716 | 4,749 |

Sum of Deaths broken down by Month Year (MY) vs. State.

The above tabular representation of the data for the total death cases for each state in a month, provides complete information in a single report. Having a look at this report can acquire knowledge about the number of deaths confirmed in the US, in addition this report provides specific information with states and timeframe.

**Bar Plot representation of the death cases in US**

Death counts for each states

| State | Deaths |
|-------|--------|
| Alabama | 2,438,125 |
| Alaska | 75,908 |
| American Samoa | 0 |
| Arizona | 4,141,957 |
| Arkansas | 1,337,296 |
| California | 13,624,795 |
| Colorado | 1,726,531 |
| Connecticut | 2,618,296 |
| Delaware | 452,466 |
| Diamond Princess | 0 |
| District of Columbia | 348,250 |
| Florida | 8,946,706 |
| Georgia | 4,930,070 |
| Grand Princess | 1,341 |
| Guam | 36,070 |
| Hawaii | 114,949 |
| Idaho | 483,923 |
| Illinois | 6,734,928 |
| Indiana | 3,425,645 |
| Iowa | 1,414,136 |
| Kansas | 1,090,674 |
| Kentucky | 1,355,120 |
| Louisiana | 3,090,754 |
| Maine | 175,442 |
| Maryland | 2,527,192 |
| Massachusetts | 5,433,074 |
| Michigan | 5,293,749 |
| Minnesota | 1,873,729 |
| Mississippi | 1,890,603 |
| Missouri | 2,180,552 |
| Montana | 339,400 |
| Nebraska | 541,907 |
| Nevada | 1,298,369 |
| New Hampshire | 337,932 |
| New Jersey | 8,497,020 |
| New Mexico | 959,056 |
| New York | 17,881,776 |
| North Carolina | 2,977,905 |
| North Dakota | 375,827 |
| Northern Mariana Islands | 990 |
| Ohio | 4,809,297 |
| Oklahoma | 1,306,233 |
| Oregon | 584,029 |
| Pennsylvania | 6,766,987 |
| Puerto Rico | 555,503 |
| Recovered | 0 |
| Rhode Island | 771,302 |
| South Carolina | 2,306,551 |
| South Dakota | 455,553 |
| Tennessee | 2,748,862 |
| Texas | 11,861,082 |
| Utah | 506,073 |
| Vermont | 59,088 |
| Virgin Islands | 8,606 |
| Virginia | 2,498,058 |
| Washington | 1,479,785 |
| West Virginia | 567,184 |
| Wisconsin | 1,807,868 |
| Wyoming | 151,496 |

Deaths

Sum of Deaths for each State.

**Line plot representing the trend of death cases over the timeframe**

## Deaths Trend



The trend of sum of Deaths for Month Year Month.

The above line graph displays the information of the rise and fall of the COVID-19 deaths in US.

**Dashboard for Death cases**

## Display of Deaths in different states of US



© 2021 Mapbox © OpenStreetMap

## Deaths across US in each month

| State | January 20.. | February 2.. | March 2020 | April 2020 | May 2020 |
|---|---|---|---|---|---|
| Alabama | 14,869 | 15,004 | 15,136 | 18,586 | 25,67 |
| Alaska | 362 | 363 | 375 | 548 | 58 |
| American S.. | 0 | 0 | 0 | 0 | |
| Arizona | 30,104 | 30,281 | 30,598 | 34,793 | 45,61 |
| Arkansas | 7,600 | 7,690 | 7,800 | 8,704 | 10,14 |
| California | 87,971 | 88,656 | 89,298 | 114,926 | 159,06 |
| Colorado | 15,217 | 15,345 | 15,527 | 25,096 | 40,32 |
| Connecticut | 33,698 | 33,864 | 33,924 | 60,486 | 102,50 |
| Delaware | 4,729 | 4,754 | 4,781 | 7,050 | 13,08 |
| Diamond Pr.. | 0 | 0 | 0 | 0 | |
| District of C.. | 4,402 | 4,422 | 4,442 | 6,787 | 12,36 |
| Florida | 75,632 | 76,260 | 76,828 | 93,518 | 119,06 |
| Georgia | 40,207 | 40,435 | 41,134 | 55,940 | 75,66 |
| Grand Princ.. | 24 | 24 | 24 | 48 | 8 |

## Death counts for each states



## Deaths Trend



The above dashboard providing wide range of information for the death cases in the US.

**Severity of Deaths in the US**

Severity of Deaths



Map based on Longitude (generated) and Latitude (generated). Color shows sum of Deaths. Size shows sum of Confirmed. Details are shown for Country and State.

The above geographical distribution of the deaths in US displays the severity of COVID-19 virus. The darker and larger the size of the circle, larger is the intensity.

**Quick statistics of COVID-19 dataset**

## Quick Monthly Stats

| Month, Year of .. | Confirmed | Deaths | Recovered |
|---|---|---|---|
| January 2020 | 46,798,022 | 1,351,445 | 16,734,598 |
| February 2020 | 47,333,350 | 1,360,381 | 16,992,144 |
| March 2020 | 47,914,950 | 1,370,284 | 17,236,687 |
| April 2020 | 63,406,912 | 2,197,532 | 20,195,591 |
| May 2020 | 79,413,113 | 3,204,163 | 24,109,082 |
| June 2020 | 91,488,194 | 3,574,633 | 29,301,608 |
| July 2020 | 125,815,453 | 4,138,614 | 41,317,857 |
| August 2020 | 158,486,131 | 4,758,193 | 56,027,565 |
| September 2020 | 175,419,534 | 5,039,464 | 66,156,287 |
| October 2020 | 213,284,888 | 5,682,291 | 82,587,710 |
| November 2020 | 273,340,896 | 6,091,050 | 101,571,431 |
| December 2020 | 401,481,360 | 7,703,867 | 154,628,148 |
| January 2021 | 676,947,647 | 11,652,951 | 217,033,957 |
| February 2021 | 656,212,305 | 11,670,034 | 211,853,236 |
| March 2021 | 774,713,304 | 14,006,968 | 29,784,582 |
| April 2021 | 780,699,757 | 13,944,255 | 30,126,980 |
| May 2021 | 836,234,052 | 14,895,929 | 30,290,799 |
| June 2021 | 812,164,631 | 14,546,888 | 30,466,868 |
| July 2021 | 242,504,163 | 4,321,305 | 19,732,363 |
| August 2021 | 209,166,538 | 3,722,242 | 19,890,001 |
| September 2021 | 209,751,462 | 3,732,495 | 20,066,622 |
| October 2021 | 210,245,088 | 3,741,081 | 20,224,281 |
| November 2021 | 210,739,019 | 3,749,623 | 20,416,423 |
| December 2021 | 211,276,315 | 3,760,332 | 20,557,708 |

Confirmed, Deaths and Recovered broken down by Month
Year (MY).

The above table displays the information of confirmed, deaths and recovered cases for each month in
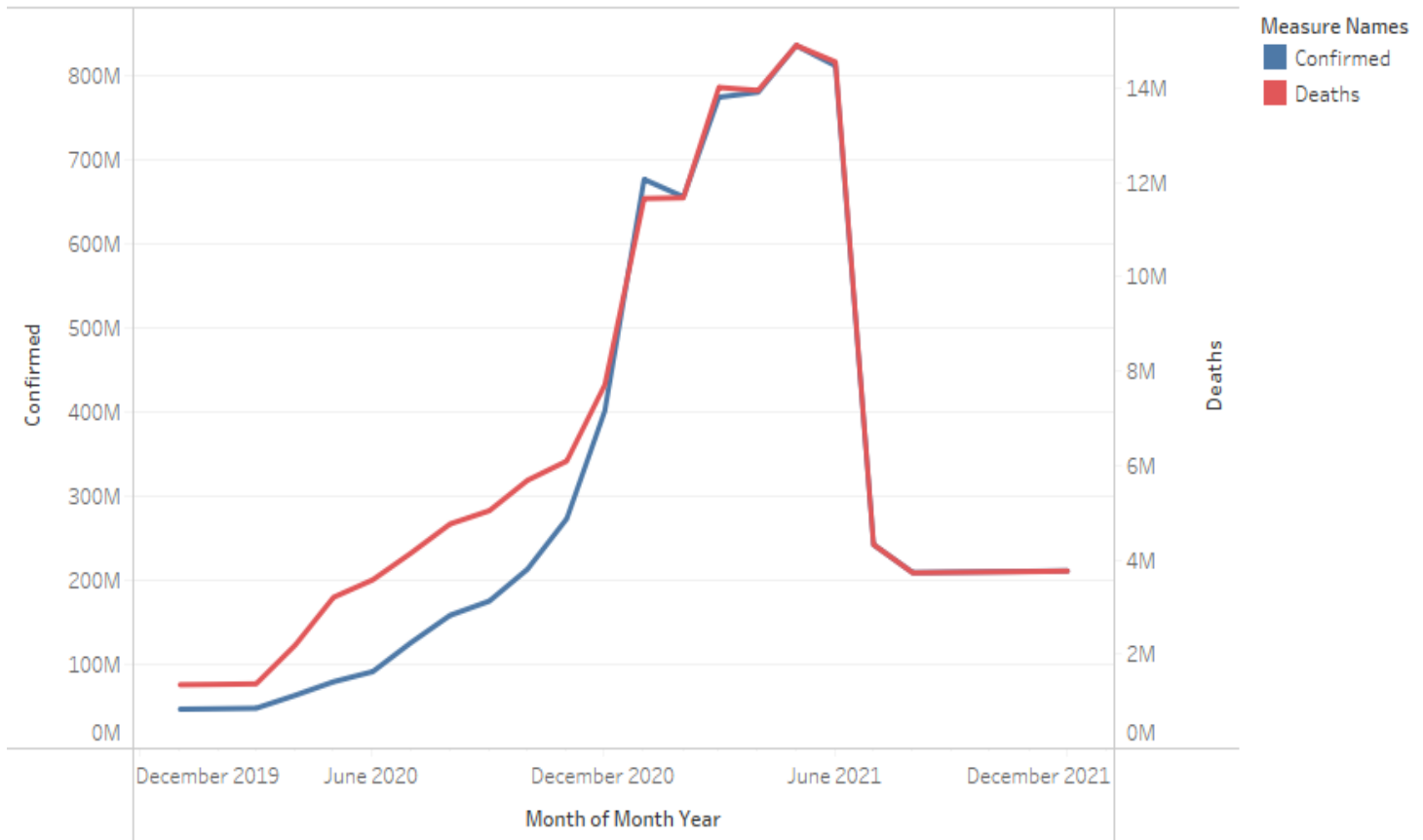the US.

**Bar plot for Confirmed, Death and Recovered cases**

## Confirmed, Deaths and Recovered

| State | Confirmed | Deaths | Recovered |
|---|---|---|---|
| Alabama | 130,236,949 | 2,438,125 | 31,385,579 |
| Alaska | 15,016,637 | 75,908 | 947,789 |
| American Samoa | 0 | 0 | 0 |
| Arizona | 206,421,695 | 4,141,957 | 12,439,312 |
| Arkansas | 80,924,720 | 1,337,296 | 33,515,473 |
| California | 868,160,344 | 13,624,795 | 0 |
| Colorado | 114,974,882 | 1,726,531 | 2,844,365 |
| Connecticut | 76,167,294 | 2,618,296 | 2,628,636 |
| Delaware | 23,441,880 | 452,466 | 2,881,647 |
| Diamond Princess | 22,442 | 0 | 0 |
| District of Columbia | 12,078,293 | 348,250 | 3,935,693 |
| Florida | 535,385,150 | 8,946,706 | 0 |
| Georgia | 264,060,371 | 4,930,070 | 0 |
| Grand Princess | 47,174 | 1,341 | 0 |
| Guam | 2,141,675 | 36,070 | 942,305 |
| Hawaii | 8,257,506 | 114,949 | 1,553,287 |
| Idaho | 45,230,241 | 483,923 | 9,744,844 |
| Illinois | 325,418,584 | 6,734,928 | 0 |
| Indiana | 168,646,933 | 3,425,645 | 43,799,780 |
| Iowa | 90,783,942 | 1,414,136 | 34,871,232 |
| Kansas | 73,861,701 | 1,090,674 | 743,554 |
| Kentucky | 99,059,803 | 1,355,120 | 6,054,657 |
| Louisiana | 119,834,147 | 3,090,754 | 51,046,120 |
| Maine | 11,877,434 | 175,442 | 1,940,905 |
| Maryland | 108,464,696 | 2,527,192 | 2,196,903 |
| Massachusetts | 158,037,192 | 5,433,074 | 48,605,219 |
| Michigan | 197,379,209 | 5,293,749 | 54,250,330 |
| Minnesota | 133,193,930 | 1,873,729 | 55,052,487 |
| Mississippi | 78,330,523 | 1,890,603 | 31,269,115 |
| Missouri | 143,571,469 | 2,180,552 | 397,080 |
| Montana | 25,131,858 | 339,400 | 9,463,264 |
| Nebraska | 53,560,695 | 541,907 | 16,530,878 |
| Nevada | 76,953,812 | 1,298,369 | 275,924 |
| New Hampshire | 18,945,779 | 337,932 | 5,972,572 |
| New Jersey | 229,487,509 | 8,497,020 | 11,916,136 |
| New Mexico | 46,337,235 | 959,056 | 11,019,177 |
| New York | 473,750,889 | 17,881,776 | 25,612,871 |
| North Carolina | 219,462,424 | 2,977,905 | 70,445,365 |
| North Dakota | 26,833,440 | 375,827 | 12,000,646 |
| Northern Mariana Islands | 46,052 | 990 | 8,164 |
| Ohio | 239,932,935 | 4,809,297 | 81,935,840 |
| Oklahoma | 103,179,317 | 1,306,233 | 40,071,963 |
| Oregon | 41,682,314 | 584,029 | 1,154,916 |
| Pennsylvania | 248,187,075 | 6,766,987 | 75,314,534 |
| Puerto Rico | 30,387,798 | 555,503 | 8,692,236 |
| Recovered | 0 | 0 | 1,384,375 |
| Rhode Island | 33,542,535 | 771,302 | 931,508 |
| South Carolina | 132,061,099 | 2,306,551 | 26,808,475 |
| South Dakota | 29,912,523 | 455,553 | 12,496,121 |
| Tennessee | 199,966,350 | 2,748,862 | 80,889,238 |
| Texas | 686,349,832 | 11,861,082 | 252,311,763 |
| Utah | 93,624,312 | 506,073 | 34,250,078 |
| Vermont | 4,203,024 | 59,088 | 989,871 |
| Virgin Islands | 795,138 | 8,606 | 341,591 |
| Virginia | 149,211,148 | 2,498,058 | 5,712,661 |
| Washington | 94,164,656 | 1,479,785 | 0 |
| West Virginia | 32,128,964 | 567,184 | 9,566,900 |
| Wisconsin | 160,460,798 | 1,807,868 | 62,898,072 |
| Wyoming | 13,510,757 | 151,496 | 5,261,077 |

Sum of Confirmed, sum of Deaths and sum of Recovered for each State.

**Confirmed VS Deaths trend**

## Confirmed VS Deaths



The trends of Confirmed and Deaths for Month Year Month. Color shows details about Confirmed and Deaths.

The above line graph displays the relation between the confirmed and death cases. We can see that as the confirmed cases increased the death cases also increase.

**NOTE-** The Y axis intervals are different for confirmed and deaths, confirmed cases intervals are displayed on the left, and deaths on the right.

## Death VS Recovered

### Death VS Recovered



The trends of Deaths and Recovered for Month Year Month. Color shows details about Deaths and Recovered.

**NOTE-** The Y axis intervals are not the same for both the measures.

## Trends of different measures

### Trend of Confirmed, Death and Recovered cases



The trends of sum of Confirmed, sum of Deaths and sum of Recovered for Month Year Month.

**Overall insights of the COVID-19 dataset**

## Quick Monthly Stats

| Month, Yea.. | Confirmed | Deaths | Recovered |
|---|---|---|---|
| January 20.. | 46,798,022 | 1,351,445 | 16,734,598 |
| February 2.. | 47,333,350 | 1,360,381 | 16,992,144 |
| March 2020 | 47,914,950 | 1,370,284 | 17,236,687 |
| April 2020 | 63,406,912 | 2,197,532 | 20,195,591 |
| May 2020 | 79,413,113 | 3,204,163 | 24,109,082 |
| June 2020 | 91,488,194 | 3,574,633 | 29,301,608 |
| July 2020 | 125,815,453 | 4,138,614 | 41,317,857 |
| August 2020 | 158,486,131 | 4,758,193 | 56,027,565 |
| September .. | 175,419,534 | 5,039,464 | 66,156,287 |
| October 20.. | 213,284,888 | 5,682,291 | 82,587,710 |
| November .. | 273,340,896 | 6,091,050 | 101,571,431 |
| December 2.. | 401,481,360 | 7,703,867 | 154,628,148 |
| January 20.. | 676,947,647 | 11,652,951 | 217,033,957 |
| February 2.. | 656,212,305 | 11,670,034 | 211,853,236 |
| March 2021 | 774,713,304 | 14,006,968 | 29,784,582 |

## Severity of Deaths



© Mapbox © OSM

## Confirmed, Deaths and Recovered

| State | Confirmed | Deaths | Recovered |
|---|---|---|---|
| Alabama | 130,236,949 | 2,438,125 | 31,385,57 |
| Alaska | 15,016,637 | 75,908 | 947,789 |
| American S.. | 0 | 0 | 0 |
| Arizona | 206,421,695 | 4,141,957 | 12,439,31 |
| Arkansas | 80,924,720 | 1,337,296 | 33,515,47 |
| California | 868,160,344 | 13,624,795 | 0 |
| Colorado | 114,974,882 | 1,726,531 | 2,844,365 |
| Connecticut | 76,167,294 | 2,618,296 | 2,628,636 |
| Delaware | 23,441,880 | 452,466 | 2,881,647 |
| Diamond Pr.. | 22,442 | 0 | 0 |
| District of C.. | 12,078,293 | 348,250 | 3,935,693 |
| Florida | 535,385,150 | 8,946,706 | 0 |

0.5B  1B  1.5B   0M   20M        200M 4

## Confirmed VS Deaths



Month of Month Year

## Death VS Recovered



Month of Month Year

## Trend of Confirmed, Death and Recovered cases



Month of Month Year

The above dashboard provides the overall insights and the behavior of the dataset.

# Business Questions

**Which state was the most affected from covid-19?**

Confirmed cases count for each states

| State | Confirmed cases |
|---|---|
| California | 868,160,344 |
| Texas | 686,349,832 |
| Florida | 535,385,150 |
| New York | 473,750,889 |
| Illinois | 325,418,584 |
| Georgia | 264,060,371 |
| Pennsylvania | 248,187,075 |
| Ohio | 239,932,935 |

As displayed in the above visual, California state is the most affected from COVID-19

**Which state has the highest deaths from covid-19?**

Death counts for each states

| State | Death counts |
|---|---|
| New York | 17,881,776 |
| California | 13,624,795 |
| Texas | 11,861,082 |
| Florida | 8,946,706 |
| New Jersey | 8,497,020 |
| Pennsylvania | 6,766,987 |
| Illinois | 6,734,928 |

As indicated in the above visual, New York is the state with most deaths.

**When did the COVID-19 cases had the highest peak?**

Confirmed Trend

Month of Month Year: May 2021
Confirmed: 836,234,052

May 2021 was the month which had the highest peak of COVID-19 cases, with 836,234,052 cases.

## When did the COVID-19 cases had the highest death rate?



May 2021 was the month with highest death cases of COVID-19, with deaths of 14,895,929.

## When did the COVID-19 cases started to drop?



During the month June 2021, there is a huge drop of COVID-19 cases.
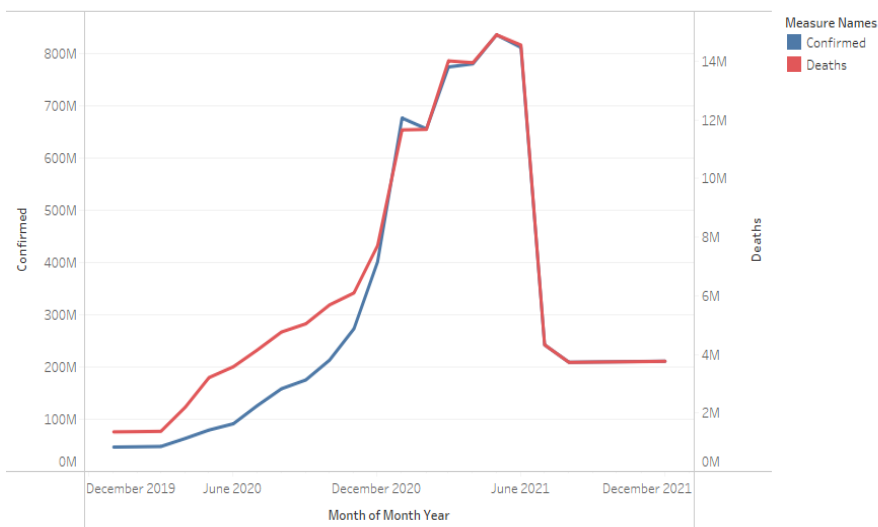
## When did the death rate started to drop?

**Deaths Trend**



During the month of June 2021, the death rate dropped drastically.

## With the EDA and visualization performed, did you observe any pattern in the dataset?

**Confirmed VS Deaths**



The trends of Confirmed and Deaths for Month Year Month. Color shows details about Confirmed and Deaths.

In the above trend, we can see there is a direct relationship between the confirmed and deaths column, as the number of confirmed cases increased, the death counts also increased.

**Top five states confirmed and deaths statistics in the US?**

| Confirmed | | Deaths | |
|---|---|---|---|
| State | Count | State | Count |
| California | 868160344 | New York | 17881776 |
| Florida | 535385150 | Texas | 11861082 |
| Illinois | 325418584 | New Jersey | 6766987 |
| New York | 473750889 | Florida | 8946706 |
| Texas | 686349832 | California | 13624795 |

**Did you observe any unusual behavior of the dataset?**

Yes, there are unusual behaviors, During the month of June 2021 the number of confirmed and death cases dropped thoroughly. The drop was caused by the massive vaccination drives conducted by the national governments to the general public to tackle the virus.

The number of confirmed and death cases had a linear raise until October 2020, however there is an immediate uplift of the cases until April 2021. This was due to, the US government had enforced strict lockdown all over the country until September 2020, after the relaxation of the rules the virus spread was escalated.

**Why do you think EDA and visualization is important?**

After performing all the above analysis using EDA and visualization, both helps us to provide accurate insights about the data before making any assumptions. EDA and visualization provide us a better understanding of the variables and the relationships between them, it helps to understand the data better to measure its impact on the business and communicates the insight visually to internal and external audiences. By performing EDA and visualization on our dataset provides precise information and helped us to answer all our business question effortlessly. EDA and visualization are two important aspect of data analytics when trying to get insights on a dataset.

*All the mentioned above screenshots and visualization is performed on Tableau the required code/file (.twb) has been uploaded on GitHub.*

**Thank you,**
**Data Wizard Team**