

6.1: Sourcing Open Data

Source

In early 2020, a worldwide pandemic was triggered by a novel coronavirus referred to as COVID-19. The United States stood out as the nation with the highest number of documented COVID-19 cases. Two datasets, namely `covid_us_county.csv` and `us_county.csv`, were made available by [Johns Hopkins University](#) and can be accessed for download on [Kaggle](#). These datasets encompass data regarding the cumulative counts of COVID-19 cases and fatalities categorized by county and gender.

Collection

Information was gathered and consolidated from the public health departments of numerous counties and states across the United States. This comprehensive dataset represents a wide geographical range and encompasses a substantial volume of data collected from various regions and administrative divisions within the country.

Limitations

The datasets in question are individual datasets that must be integrated into a unified dataset. However, it's important to note that these datasets lack information concerning COVID vaccinations and the age of patients. Additionally, they do not encompass data related to hospitalizations associated with COVID-19 cases. It's crucial to recognize that this data is constrained to individuals who have access to healthcare facilities, are willing to seek medical attention at a hospital, and have the means to do so. Consequently, the reported cases and deaths may potentially underestimate the true figures due to the exclusion of individuals who may not have these privileges or resources.

Ethics

The dataset doesn't include any HIPAA-related data. Nevertheless, it's worth noting that certain counties within the dataset may have relatively small populations. In such cases, even non-HIPAA data could potentially be utilized to indirectly identify patients. Therefore, it's crucial to perform thorough data analysis and implement robust privacy safeguards to ensure that patient identities remain protected and confidential. This proactive approach to data handling and analysis is essential for maintaining data privacy and security, especially when working with datasets that may inadvertently reveal sensitive information about individuals.

Data Profile

Dataset: `_covid_us_county.csv`

Variable	Description	Structured/Unstructured	Qualitative/Quantitative	Discrete/Continuous	Nominal/Ordinary
<i>fips</i>	County code in numerical format	Structured	Qualitative		Nominal
<i>county</i>	United States's County	Structured	Qualitative		Nominal
<i>state</i>	United States's State	Structured	Qualitative		Nominal
<i>state code</i>	Two letter abbreviation of State	Structured	Qualitative		Nominal
<i>lat</i>	Latitude of the County	Structured	Quantitative	Continuous	
<i>long</i>	Longitude of the County	Structured	Quantitative	Continuous	
<i>cases</i>	Total COVID-19 cases	Structured	Quantitative	Discrete	
<i>deaths</i>	Total COVID-19 deaths	Structured	Quantitative	Discrete	
<i>date</i>	Date	Structured	Qualitative		Ordinal

Dataset : us_county.csv

Variable	Description	Structured/Unstructured	Qualitative/Quantitative	Discrete/Continuous	Nominal/Ordinary
<i>fips</i>	County code in numerical format	Structured	Qualitative		Nominal
<i>county</i>	United States's County	Structured	Qualitative		Nominal
<i>state</i>	United States's State	Structured	Qualitative		Nominal
<i>state code</i>	Two letter abbreviation of State	Structured	Qualitative		Nominal
<i>female</i>	Female Population	Structured	Quantitative	Discrete	
<i>male</i>	Male Population	Structured	Quantitative	Discrete	
<i>population</i>	Total Population	Structured	Quantitative	Discrete	
<i>female_percentage</i>	Percentage of Female population	Structured	Quantitative	Continuous	
<i>median_age</i>	County's median Age	Structured	Quantitative	Continuous	
<i>lat</i>	Latitude of the County	Structured	Quantitative	Continuous	
<i>long</i>	Longitude of the County	Structured	Quantitative	Continuous	

Data Cleaning

Dataset: covid_us_county.csv

Column	Type of Inconsistency	Actions
<i>fips</i>	Rename column	Change name to better fit data. New name is "county_code"
<i>county</i>	Mixed data types	county set to string
<i>county</i>	2% of records were unassigned according to Kaggle	Couldn't find the records, so I didn't delete them. The records could have been deleted from deleting the 10,810 records in the country_code column.
<i>county_code</i>	Missing values	Dropped missing values, because only makes up p <1% of data. (Total 10810 records)
<i>state code</i>	Mixed data types	state_code set to string
<i>Multiple dated column names</i>	Data inconsistency	Dropped columns because they were not informative

Dataset : us_county.csv

Column	Type of Inconsistency	Actions
<i>fips</i>	Rename column	Change name to better fit data. New name is "county_code"
<i>state code</i>	Mixed data types	state_code set to string

The two datasets were combined to create a new dataset which we will work with `deaths_cases_gender.csv`

Variable	Description
<i>fips</i>	County code in numerical format
<i>county</i>	United States's County
<i>state</i>	United States's State
<i>state code</i>	Two letter abbreviation of State
<i>female</i>	Female Population
<i>male</i>	Male Population
<i>population</i>	Total Population
<i>female_percentage</i>	Percentage of Female population
<i>median_age</i>	County's median Age
<i>lat</i>	Latitude of the County
<i>long</i>	Longitude of the County
<i>cases</i>	Total COVID-19 cases
<i>deaths</i>	Total COVID-19 deaths
<i>date</i>	Date
<i>region</i>	the region column groups states into U.S. regions using the information from Wikipedia

Key Questions

- Which counties had the most COVID-19 impact?
- Which states saw the most COVID-19 impact?
- When did COVID-19 cases reach their highest point?
- When did COVID-19 deaths reach their highest point?
- Does the size of the population affect the percentage of cases and deaths?
- Does the median age of people affect the percentage of cases and deaths?