

COVID-19 Analysis and Visualization

ABSTRACT

In today's age of big data, an enormous volume of information has been produced and gathered from a diverse range of rich data sources. Within this vast sea of data lies valuable insights and essential knowledge. A prime example can be found in healthcare and epidemiological data, particularly concerning patients impacted by epidemic diseases such as the coronavirus disease 2019 (COVID-19). The insights gleaned from this epidemiological data equip researchers, epidemiologists, and policymakers with a deeper understanding of the disease, potentially leading to innovative methods for its detection, control, and management. The adage “a picture is worth a thousand words” underscores the importance of employing visualization techniques to make this extensive data and its insights more accessible. In this paper, we introduce a visualization and visual analytics tool specifically designed for the analysis of COVID-19 epidemiological data. This tool enables users to gain a clearer understanding of the information regarding confirmed COVID-19 cases. While tailored for this particular dataset, the tool can also be utilized for visualizing and analyzing big data across a variety of other real-world applications and services.

Keywords – COVID-19, Data Visualization, Analysis, Pandemics, Decision-making

I. INTRODUCTION

As we are living an era of big data, big data are everywhere. To elaborate, with advances in technology, a huge amount of data has been easily generated and collected from a wide variety of rich data sources at a rapid rate. These big data can be of different levels of veracity (e.g., precise data, imprecise and uncertain data [1]-[3]). Examples of big data include social network data [4]-[8], financial time series [9]-[11], omic data (e.g., genomic data) [1], [12], disease reports [13], [14], as well as epidemic data and statistics.

Embedded in these big data are useful information and valuable knowledge. This calls for data science, which aims to discover knowledge from these big data via data mining algorithms, machine learning tools, mathematical and statistical models, data analytics, and visual analytics. The discovered knowledge is useful. For instance, knowledge discovered from these epidemiological data helps researchers, epidemiologists and policy makers to get a better understanding of the disease, which may inspire them to come up ways to detect, prevent, and/or control diseases such as viral diseases. Examples of viral diseases include:

- severe acute respiratory syndrome (SARS), with outbreak in 2002–2004;
- Middle East respiratory syndrome (MERS), with outbreak in 2012–2015; and
- coronavirus disease 2019 (COVID-19), with outbreak started in 2019 and became pandemic in 2020.

knowledge [18]. A majority of the existing visualizers on the COVID-19 epidemiological data focused on showing the numbers of confirmed cases and mortality spatially and/or temporally. In other words, they show

- spatial differences among different continents, countries, regions, or sovereignties; and/or
- temporal differences among weeks or days along the timeline—e.g., to show the effects of public health strategies and mitigation techniques such as social/physical distancing or stay-at-home orders in “flattening the (epidemic) curve”.

As the numbers of inhabitants and tests both play roles in the data and their analyses, they help in the computation of figures like (a) the numbers of confirmed cases and mortality per thousand/million inhabitants and (b) the number of tests per thousand inhabitants.

However, in addition to the number of cases or mortality, there are other important knowledge that can be discovered from the epidemiological data via data mining. For instance, frequent pattern mining finds relationships among attributes (or features) associated with confirmed COVID-19 cases. Moreover, a visual representation of this discovered knowledge gives a more comprehensive representation, which in turn leads to a better insight and understanding of the data and discovered knowledge. Hence, in this paper, we present a tool for big data visualization and visual analytics of COVID-19 epidemiological data. Due to the nature of these data, it is not unusual to have NULL values for some of the attributes (e.g., unstated transmission methods of disease).

II. LITERATURE SURVEY

Due to the COVID-19 pandemic, many visualizers and dashboards have been developed over the past few months. Some of them [19]-[21] visualized literatures related to COVID-19 research, and some others [22] visualized economic impact of COVID-19. However, a majority of them [23] focused on the actual COVID-19 cases. Globally notable visualizers include (a) World Health Organization (WHO) Coronavirus Disease 2019 (COVID-19) Dashboard [24], (b) COVID-19 Dashboard by the Center for Systems Science and Engineering (CSSE) at Johns Hopkins University (JHU)¹, and (c) COVID-19 dashboard by European Center for Disease Prevention and Control (ECDC)². In addition, Government of Canada³, provincial and territorial governments, major news channels/media/networks (e.g., TV⁴, newspaper), and Wikipedia⁵ capture data and provide dashboards on information about the COVID-19 pandemic situations in Canada. There are several commonality among these visualizers and dashboards. For instance, they mostly focused on the total numbers of new cases, confirmed cases, and deaths.

The *spatial information* about the total of numbers of confirmed cases and deaths in different countries (or regions/sovereignties) is usually represented by (a) the bubble map or (b) the choropleth map. To elaborate, in a *bubble map*, the total number of confirmed cases for each country is indicated by the radius of the bubble representing the country. See Fig. 1. Similarly, the bubble map can also show the number of new cases or deaths, in an absolute value (e.g., N_1 newly reported cases) or a relative figure with respect to population (e.g., N_2 deaths per one million population). See Fig. 2.

While the severity of COVID-19 in many countries can be representing by the sizes or radii of the bubbles representing these countries, many bubbles overlap. The overlapping and/or containment of bubbles makes it challenging to visualize the severity of the disease in countries in dense regions such as Eastern Caribbean and Southeastern Europe as shown in Fig. 2.

Fig. 1.

A snapshot of a bubble map [24] showing the total number of confirmed cases among different countries in the world as of August 07, 2020.



III. Our Visualizer for Analyzing COVID-19 Epidemiological Data

To help users and researchers to get a better understanding of COVID-19 disease, we design and develop a big data visualization and visual analytics tools for COVID-19 epidemiological data. In this section, we illustrate our idea on COVID-19 epidemiological data for Canada [36].

A. Data Collection

Our tool first collects different categories of big COVID-19 epidemiological data from different sources (e.g., federal and provincial/territorial governments). These include:

- administrative information, which includes a unique privacy-preserving identifier for each case, its region and episode week (i.e., symptom onset week or its closest week);
- case details, which include gender, age group, and occupation of the cases (e.g., health care workers, school or daycare workers, long-term care residents);
- symptom-related data, which include additional information for the case who is not asymptomatic (i.e., symptomatic case) such as onset week of symptoms, as well as a collection of symptoms (e.g., cough, fever, chills, sore throat, runny nose, shortness of breath, nausea, headache, weakness, pain, irritability, diarrhea, and other symptoms);
- clinical course and outcomes, which include hospital status (e.g., hospitalized in the intensive care unit (ICU), non-ICU hospitalized, not hospitalized, unstated). For recovered case, it also includes additional information such as the recovery week. For the case who has not recovered, it indicates that the case died while infected by COVID-19; as well as
- exposures, which include transmission methods (e.g., domestic acquisition via contact of COVID-19 case or contact with traveller; international travel).

B. Data Preprocessing

After gathering related data from heterogeneous sources, our tool preprocesses the data. Given the nature of these cases, it is not unusual to have missing, unstated or unknown information (i.e., NULL values). For example, for some “Boolean” attributes (e.g., attribute “asymptomatic”), we observed three values: asymptomatic, symptomatic, and unstated (i.e., NULL value). Our tool links all values (including NULL) for each attribute. Moreover, our tool also detects and flags any data inconsistency.

C. Visualization of Frequent Patterns of Cardinality 1

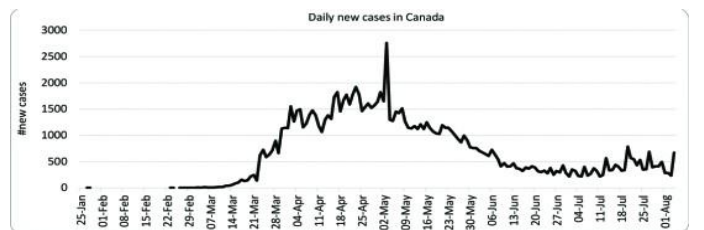
Once the data are preprocessed, we integrate the datamining capability into our tool. Specifically, via frequent pattern mining, our tool first discovers frequently occurring domain attributes (i.e., 1-itemsets). Our tool uses a *pie chart* or *sunburst diagram* to represent values for each frequently occurring domain attribute.

Recall that existing frequent pattern visualizers mentioned in Section II-B display all frequent patterns. In contrast, our tool displays frequent patterns one-by-one. For a frequent 1-itemset, we show not only the frequent $\langle \text{attribute}, \text{value} \rangle$ -pair but also pairs associated with other values for that *attribute*. This enables users to get insight about the portion of that *value* for the *attribute*.

Example 1 As a preview, when applying our tool to the COVID-19 epidemiological data for Canada up to August 06, 2020, our tool finds a frequent 1-itemset $\langle \text{transmission}, \text{domestic acquisition} \rangle$ -pair with a frequency of 97,052 out of 107,916 cases (90%). Our tool also displays other pairs for attribute “transmission” in the pie chart: $\langle \text{transmission}, \text{international travel} \rangle$ (4%) and $\langle \text{transmission}, \text{NULL} \rangle$ (6%).

In addition, our tool gives users an option to ignore the NULL value for any attribute. If the user select this option, our tool displays not only the frequent $\langle \text{attribute}, \text{value} \rangle$ -pair but also pairs associated with other non-NULL values for that *attribute*. This enables users to get insight about the portions of all non-NULL values for the *attribute*.

Example 2 Continue with our preview, without NULL values, our tool finds a frequent 1-itemset $\langle \text{transmission}, \text{domestic acquisition} \rangle$ -pair with a frequency of 96% of cases with stated transmission methods. Our tool also finds and displays another pair with a non-NULL value for attribute “transmission”: $\langle \text{transmission}, \text{international travel} \rangle$ (4%).

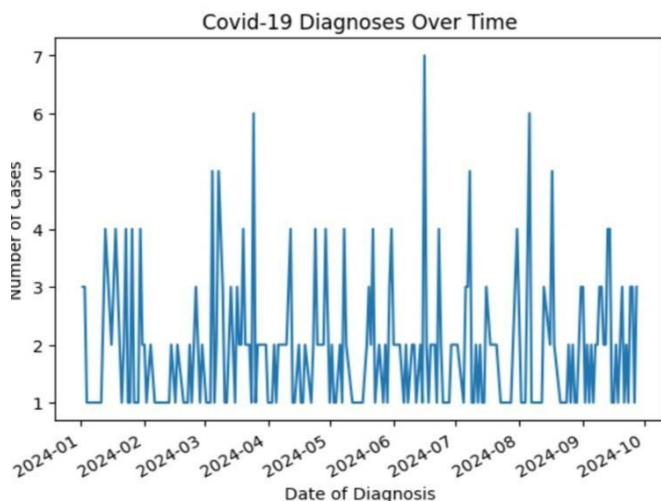


IV.RESULT AND DISCUSSION

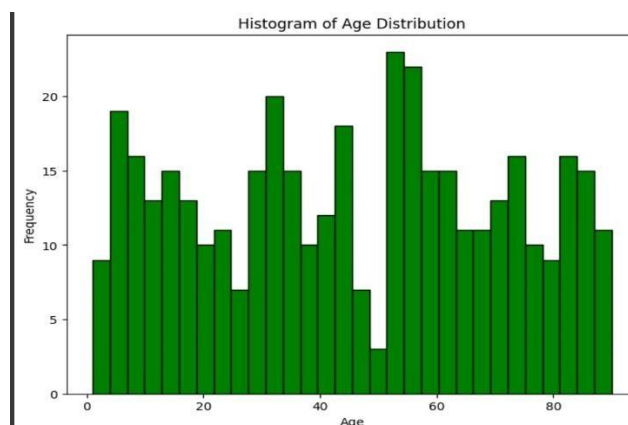
Trends in COVID-19 Cases and Recoveries

Using **line plots**, the dataset revealed clear trends in the progression of COVID-19:

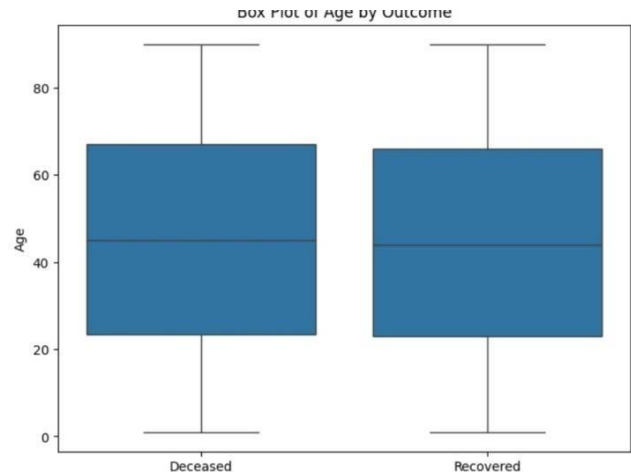
- A **sharp increase in cases** was observed during specific time periods, likely corresponding to waves or outbreaks.
 - Recovery trends** often lagged behind the spikes in cases, highlighting the time required for treatment and recovery.
 - Lockdown periods** and vaccination campaigns were reflected in reduced case numbers during certain phases.
- Line Plot:** Visualizing trends over time for cases, recoveries, and deaths.



A **histogram plot** is a graphical representation of the distribution of a dataset. It groups continuous data into intervals, called **bins**, and shows how many data points fall into each bin.



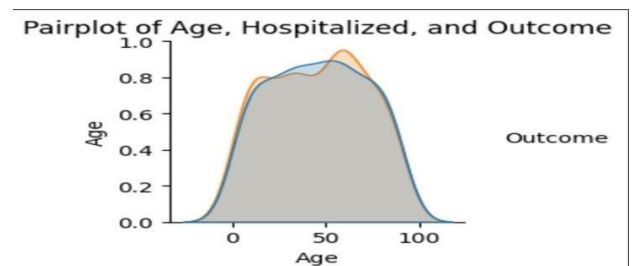
Box Plot: Analyzing the spread of data, e.g., recovery times across age groups.



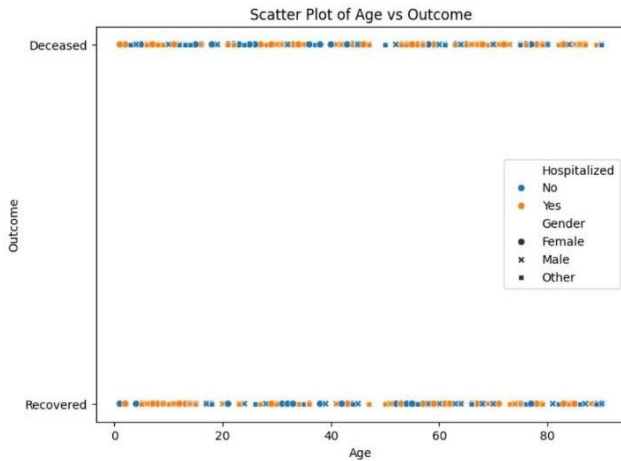
Violin Plot: Distribution of recovery status by demographics.



A **pair plot** is a visualization that shows pairwise relationships between multiple numerical features in a dataset. It plots a grid of scatterplots.



A **scatter plot** is a type of data visualization used to display the relationship between two numerical variables. Each data point is represented as a dot, plotted along two axes, with one variable determining the x-axis and the other determining the y-axis.



I. CONCLUSION

Over the past few months, there have been works on visualizing and analyzing different aspects of big data related to COVID-19. However, a majority of existing visualizers focus on showing the temporal and/or spatial trends on the numbers of cases and mortality. In this paper, we focus on epidemiological aspects of COVID-19 data. Our key contributions include our design and development of a big data visualization and visual analytics tool for COVID-19 epidemiological data. By incorporating big data mining and data analytics into our tool, we discover frequent patterns, together with their related patterns. We provide users with flexibility to include or exclude NULL values. Evaluation results show benefits of our tool in providing a comprehensive and effective visual representation of these important patterns, which in turn helps researchers to get a better understanding of the COVID-19 disease and thus enables them to combat the disease. Moreover, our tool can be applicable to other real-life applications with NULL values. As ongoing and future work, we explore possibility to incorporate other techniques [37]-[40] into our tool.

REFERENCES

- [1] F. Jeanquartier, C. Jean-Quartier and A. Holzinger, "Visualizing uncertainty for comparing genomic pediatric brain cancer data", *IV 2019 Part I*, pp. 388-391.
Show in Context [View Article](#) [Google Scholar](#)
- [2] C.K. Leung, "Uncertain frequent pattern mining", *Frequent Pattern Mining*, pp. 417-453, 2014.
Show in Context [CrossRef](#) [Google Scholar](#)
- [3] C.K. Leung, M.A.F. Mateo and D.A. Brajczuk, "A tree-based approach for frequent pattern mining from uncertain data", *PAKDD*, pp. 653-661, 2008.
Show in Context [CrossRef](#) [Google Scholar](#)
- [4] S. Fernando et al., "Compositional microservices for immersive social visual analytics", *IV 2019 Part I*, pp. 216-223.
Show in Context [View Article](#) [Google Scholar](#)
- [5] F. Jiang, C.K. Leung and S.K. Tanbeer, "Finding popular friends in social networks", *CGC*, pp. 501-508, 2012.
Show in Context [View Article](#) [Google Scholar](#)
- [6] C.K. Leung, S.K. Tanbeer and J.J. Cameron, "Interactive discovery of influential friends from social networks", *SNAM*, vol. 4, no. 1, pp. 154:1-154:13, 2014.
Show in Context [CrossRef](#) [Google Scholar](#)
- [7] M. Mai et al., "Big data analytics of Twitter data and its application for physician assistants: who is talking about your profession in Twitter?", *Data Management and Analysis*, pp. 17-32, 2020.
Show in Context [CrossRef](#) [Google Scholar](#)
- [8] S.K. Tanbeer, C.K. Leung and J.J. Cameron, "Interactive mining of strong friends from social networks and its applications in e-commerce", *JOCEC*, vol. 24, no. 2 - 3, pp. 157-173, 2014.
Show in Context [CrossRef](#) [Google Scholar](#)

- [9] R.C. Camara et al., "Fuzzy logic-based data analytics on predicting the effect of hurricanes on the stock market", *FUZZ-IEEE*, pp. 576-583, 2018.
Show in Context [View Article](#) [Google Scholar](#)
- [10] D. Jonker, R. Brath and S. Langevin, "Industry-driven visual analytics for understanding financial timeseries models", *IV 2019 Part I*, pp. 210-215.
Show in Context [View Article](#) [Google Scholar](#)
- [11] C.K. Leung, R.K. MacKinnon and Y. Wang, "A machine learning approach for stock price prediction", *IDEAS*, pp. 274-277, 2014.
Show in Context [CrossRef](#) [Google Scholar](#)
- [12] O.A. Sarumi and C.K. Leung, "Scalable datascience and machine learning algorithm for gene prediction", *BigDAS*, pp. 118-126, 2019.
Show in Context [Google Scholar](#)
- [13] A.Diallo et al., "Proportional visualization of genotypes and phenotypes with rainbow boxes: methods and application to sickle cell disease", *IV 2019 Part I*, pp. 1-6.
Show in Context [View Article](#) [Google Scholar](#)
- [14] J. Souza, C.K. Leung and A. Cuzzocrea, "An innovative big data predictive analytics framework over hybrid big data sources with an application for disease analytics", *AINA 2020*, pp. 669-680.
Show in Context [CrossRef](#) [Google Scholar](#)
- [15] A.A. Ardakani et al., "Application of deep learning technique to manage COVID-19 in routine clinical practice using CT images: results of 10 convolutional neural networks", *Comp. Bio. Med.*, vol. 121, pp. 103795:1-103795:9, 2020.