

PROJECT REPORT ON BIG DATA VISUALIZATION MADE EASY

Project Report

Presented to

Prof. Kevin Smith Department of

Computer Science San José State

University

In Partial Fulfillment

Of the Requirement for the Class

CS 235

By

Rashmeet Kaur Khanuja (012409982)

December, 2017

Abstract

“A picture is worth a thousand words” – especially when you are trying to find relationships and understand your data – the size of which may know no bounds, and the organization no structure. To create meaningful visuals of such data, there are a variety of design patterns to consider to make the visualization user-friendly. Data size and composition play an important role in this process. Our project aims to visualize humongous raw data which consists of US road traffic volume for the year 2015. This

~2gb dataset contains daily volumes of traffic, binned by the hour. The data used has the following fields: daily observation of traffic volume divided into 24 hourly bins, station_id, traffic flow direction, and type of road. Some fields we used as is, while the others we performed pre-processing and aggregation on to derive meaningful inferences. For visualization, we majorly used D3.js along with other libraries like datamaps, d3-svg-legend, d3-tooltip, and bootstrap (CDN). For data pre-processing, we used python. Using these technologies, we aim to create an interface for multiple user scenarios that is backed by strong analytics while providing a friendly-yet-informative dashboard to the users to identify trends to suit their needs.

Keywords : big data, road traffic, big data visualization, d3.js, datamaps, heatmaps, preattentive variables, bootstrap, road traffic data US

Table of Contents

1.0 Introduction	4
2.0 Functional Design	5
3.0 System Design	8
4.0 Implementation	9
5.0 System Testing	13
6.0 Conclusion and Future Work	15
References	16

1.0 Introduction

ROAD TRAFFIC is the flux or passage of motorized and non-motorized vehicles on the road. With the rapid development of transportation systems and integration of newer technologies with vehicles, traffic has become a vital part of human life and significantly influenced the quality of life. For instance, American drivers spend an average of more than 17,600 minutes behind the wheel each year, according to a new survey from the AAA Foundation for Traffic Safety [1]. With the growing population and large number of vehicles problems, like congestion, accidents and air pollution are increasing. A number of efforts to address these problems have been proven effective, including intelligent transportation systems (ITSs), public transportation systems, safety seat belts, and air-bags. Among these solutions, ITSs are deemed attractive because they enhance the efficiency and functionalities of transportation systems with advanced information technology. In particular, the role of data in ITS has become essential because of the size of the collected data. Traffic data refers to datasets generated and collected on moving vehicles and objects. Data visualization is an efficient means to represent distributions and structures of datasets and reveal hidden patterns in the data. Huge amount of numbers in big spreadsheets, no matter how meaningful, are not only tedious to go through but also get very time consuming and boring for the users. Presenting these numbers in a form that is visually appealing and informative as well provides a platform for the users to retrieve relevant results while cutting out the amount of time and efforts spent on manual, labor-intensive experimentation. Data visualization is not a new concept; a lot of work has been done in this field before. Visualization-based data discovery platforms allow business users to mash up disparate data sources to create custom analytical views with flexibility and ease of use that simply didn't exist with the raw data. "Watch 24 Hours of Traffic, Visualized as a Living Circulatory System" by Max Galka uses similar technologies namely D3.js and Three.js to visualize traffic on the highways as living circulatory system.[2] We also found the work done by Matteo Picozzi, Nervo Verdezoto, Matti Pouke, Jarkko Vajus-Anttila and Aaron Quigley in their paper - "Traffic Visualization - Applying Information Visualization Techniques to Enhance Traffic Planning" relevant in terms of technologies they used namely D3.js along with highcharts and Google Maps API in order to represent geospatial data for exploratory visualizations. Other related work in this field has also been done using real-time data that we shall elaborate upon in the future scope for our project section.

The dataset we chose here consists of road traffic volumes in US throughout the year 2015. Since our data file was in textual format, we used python to convert it into CSV (since it is compatible with D3.js). Furthermore, it converted the raw data into a more structured format to perform aggregations and other calculations for various derivations. Our objectives behind visualizing this data were -

- *Getting answers from the visualization of data such as, where are the heaviest traffic volumes? By time of day? Any interesting seasonal patterns to traffic volumes? Analyze traffic volumes during the holiday season? Traffic trends in rural and urban areas?*
- *Study of traffic routes between various cities within states*

- *Plot interactive maps of trends obtained from the data analysis, e.g. Busiest states during the year? Traffic in areas within states, e.g. - station ids? Direction of travel?*
- *Provide a dashboard suitable for multiple user scenarios so as to serve various user segments*
- *Learning new interactive data visualization technologies such as D3.js and datamaps.*

Our ultimate motivation for this project was to help users make data driven decisions through traffic visualization.

2.0 Functional Design

In this section, we shall discuss what our prototype does and how. We shall also cover the requirements and design approach.

Before we decided on the design approach, it was important to first understand the nature of data. We answered questions like how our data was organized, how the elements related to each other, how the user could operate the data (filtering, sorting, etc.), and how the details of our data could be focused upon. Answering these help us take steps towards deciding the functional design of our software. Our dataset was highly unstructured and needed cleaning, matching, and organizing; that called for a lot of data pre-processing. We chose python for the entire pre-processing step. First we used it to write scripts in order to change the data format from text (Fig. 1) to Comma Separated Values (Fig. 2). Furthermore, we added more understandable columns by mapping them with the available fields like state names corresponding to the fips code provided in the data set (Fig. 3), month names corresponding to numeric values (Fig. 4) day names like Sunday, Monday, and so on for values 1, 2, 3...given in our data set (Fig.

5) , and more . We further performed aggregation on data and added a column for the sum of traffic (Fig. 6). Also we used it to create a separate data file for the information of traffic volume during holidays(Fig. 7). Once our data organization was clear, the next step was to decide upon the design patterns to be used. In order to pick design patterns to represent our data, we also needed a clear picture of functional, non-functional, and user experience requirements.

Fig. 1

Fig. 2

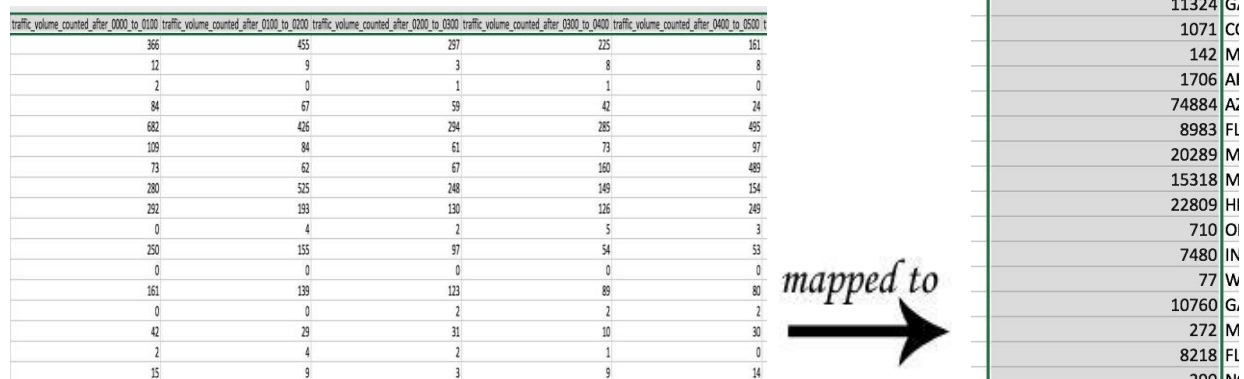


Fig. 6

Holiday Name	Traffic Volume	Traffic Volume on day
New Year's Day: Jan1	16560444	30-Dec
Martin Luther King Jr. Day: Jan 19	16560444	16-Jan
Valentine's Day: Feb 14	145000	13-Feb
Memorial Day: May 25	130000	22-May
Independence Day: July 4	100000	2-Jul
Labor Day: Sep 7	80000	4-Sep
Halloween: Oct 31	100000	29-Oct
Veterans Day: Nov 11	100000	12-Nov
Thanksgiving Day: Nov 26	100000	25-Nov
Christmas Eve: Dec 24	100000	23-Dec

Fig. 7

Requirements - In order to visualize past traffic data, decision makers need to be able to navigate through various fields of the data and be able to view different elements. They must be able to make sense of the visuals by looking at different views of the data such as traffic in states throughout the year, traffic distribution based on months, days, and hours, options to view traffic in rural and urban regions, and on holidays. One of the major functional requirements of our project was for the visualization to be interactive and flexible so as to provide the users with functionality to refine the views based on their requirements. The visualization needed to be informative enough to fulfill all these requirements, yet not too cluttered to prevent overwhelming the user.

Design Approach - Taking into account the objectives of traffic data analysis, the tasks of traffic data visualization can be classified as follows -

2.1 Traffic Situations Monitoring - In road traffic data, interesting trends may be hidden such as days with maximum traffic, routes with the busiest roads, stations that need more attention to avoid jams, pollution information, etc.

2.2 Trend Search - One of the most important goals of data visualization is to discover patterns. These patterns/trends reflect useful information such as “On Fridays, between 5pm-11pm, the route between San Francisco and San Jose shows maximum heat”. This, if noticed on a repeated basis, can be termed as a trend and such trends can help various users in varied segments. It also helps discover the relationship between different parameters such as Day, Time, and Station Id with respect to Traffic Volume in the above example.

2.3 Contingent exploration and forecasting - Having data visualized to convey meaningful information, it can be used to study past trends and predict the future. For eg., taxi trips are found to be most frequent in the downtown on Friday nights after 10 pm. Such data can, for instance, be used by Uber drivers to plan their driving schedule for the upcoming weekend accordingly and make most profit.

2.4 Route Selection Mechanism - Human intelligence coupled with data vis insights can be used to plan and recommend routes in order to facilitate better analytics and improve traffic control efficiency .

We based our design approach on the above requirements and tasks that must be taken care of by an efficient visualization system. We brainstormed, surveyed, and analyzed various user-scenarios in order to incorporate elements that would benefit multiple user segments and provide a platform to visualize big data in an appealing way. We incorporated meta overview, interactive filtering, details on demand, macro and micro views, preattentive variables to focus user’s subconscious attention, textual representation of relevant highlights, and layering and positioning to prevent overloading the user with information. We shall discuss about the use of various design patterns and their implementation in the upcoming sections.

3.0 System Design

Our system architecture provides the aforementioned interactive visualization dashboard. The high level architecture is composed of two components - the Data Processor and the View Generator (Fig. 8).

Input Data - Input data here refers to the raw data file that we obtained in text format. This dataset was unstructured and unorganized which is why it needed to be passed into the data processor in order to clean, map, and organize it. The original data consisted of 38 columns.

Data Processor - This component takes the raw input data and runs python scripts on that data to provide CSV files in proper structure. It also does the job of performing the necessary cleaning, rearranging, and sorting of data in order to make it usable according to the needed visualization input.

Online Processing - This step consists of the various components where the front end work is done. The task of fetching the processed data and creating various design patterns based on it is accomplished here. It facilitates analysis and exploration of traffic variations flow using javascript libraries such as D3, datamaps, d3-svg-legend, and d3-tooltip along with html, css, and bootstrap to support the illustrations.

In this setup, some of the system requirements are as follows -

- Versions of languages used -

- D3 V4
- Datamaps 0.4.4 or higher
- Python 2.7 or higher
- Memory requirement - minimum 8GB RAM
- IDE used for coding - Visual Studio Code 1.18, pyCharm 2017.3

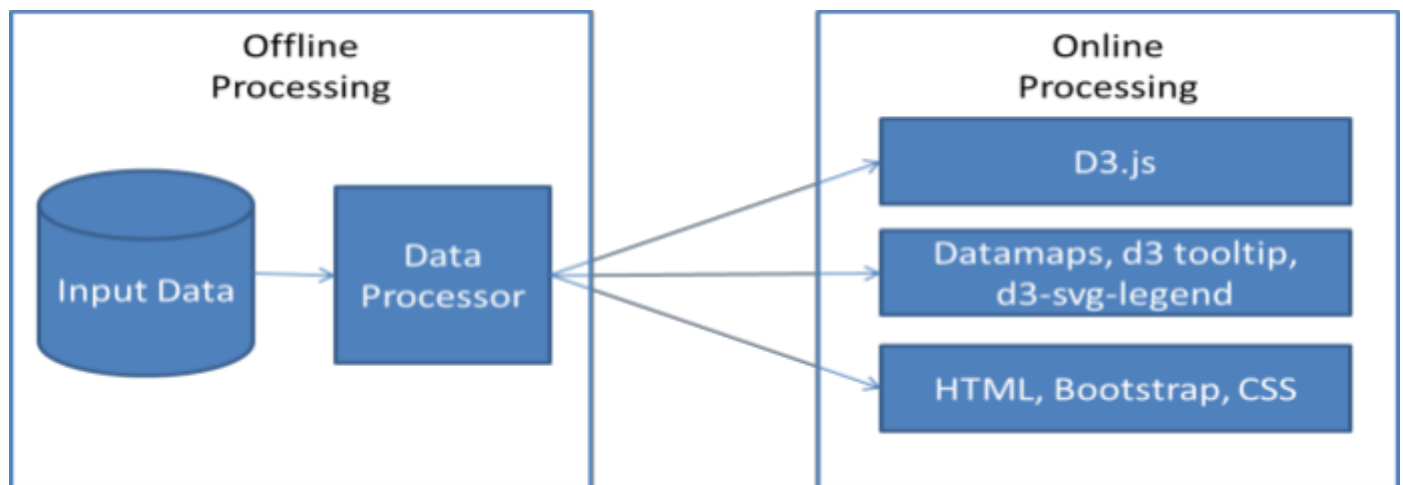


Fig. 8

4.0 Implementation

To show the different attributes of our data in a way that facilitates analysis and appeals the user, we decided to incorporate multiple **design patterns** to serve different needs. We needed to make sure that the information conveyed to the user reduces their cognitive load, hence we adopted the following approach.

4.1. For the primary view (Fig. 9), we chose **country heatmap** to display the distribution of traffic across different states in the US for the entire year. It was best suited for our data because displaying state-wise information along with additional details was required. The depth of colors on the map corresponds to the volume of traffic in that state; on mouse hover the state color changes to show it is in focus. Alongside the heatmap, we provided a drop down menu for data filtering. It enables the user to select the month for which they wish to see the traffic volumes. Data legend is used to represent data values based on different colors. To show points of interest, labels are used to display the traffic information of the state on which mouse is hovered. The preattentive variable used in heatmap is **color saturation**. There's a table below the dropdown menu to display useful highlights such as total traffic during the selected month, state with the highest traffic, and month with the highest traffic.

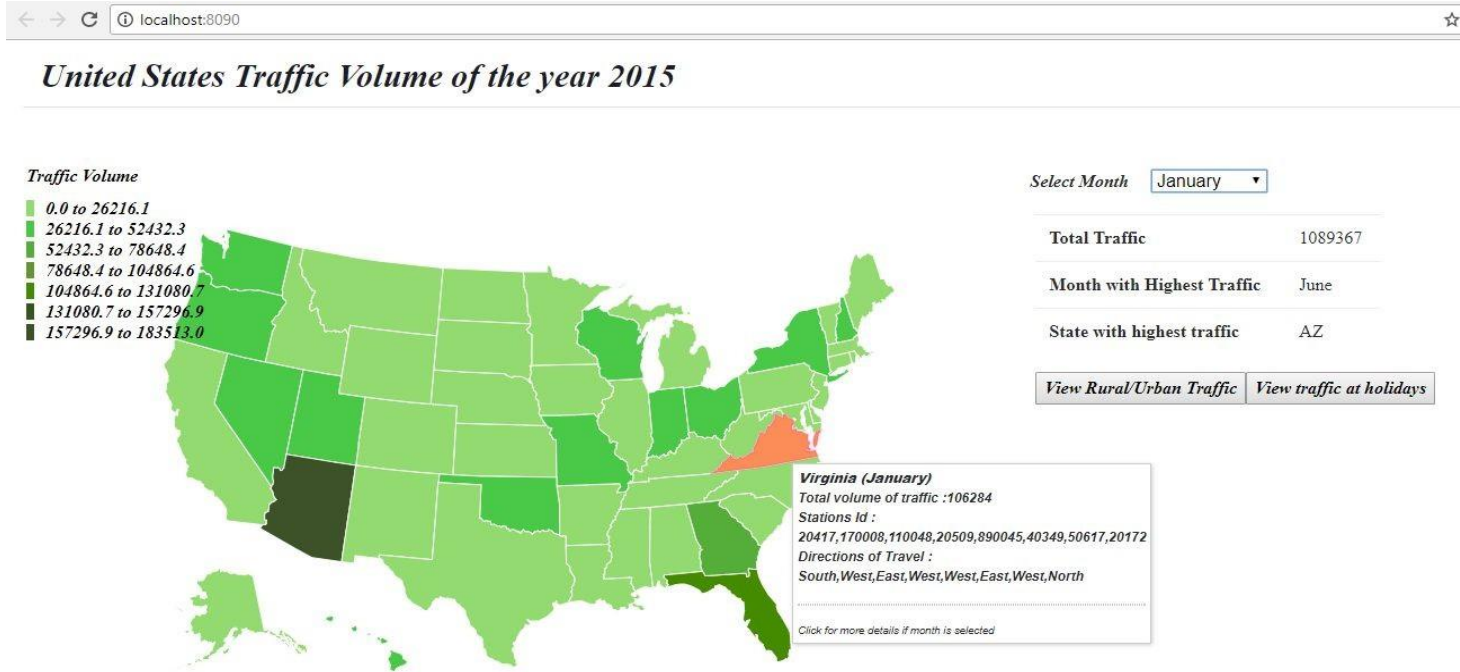


Fig. 9

4.2 When the user clicks on a particular state in the country heatmap, a modal opens that shows a **histogram** to represent the day-wise traffic of the selected state in the selected month (Fig. 10). To represent data values, we used numeric framing where X-axis shows the date and Y-axis the traffic volume. The height of the bars corresponds to the volume of traffic, thus the preattentive variable here being **size**. Upon mouse hover, the selected bar changes color and a tooltip appears that displays more information about the selected day's traffic count.



Fig. 10

4.3 If a user wishes to go deeper and view hourly traffic for a selected day in a particular state in the selected month, they can click on the day's bar and a **line graph** appears (Fig. 11). X-axis represents the 24 hours and Y-axis the traffic volume. Since the data points were more in number, line graph was the best fit to show peaks and valleys so as to make it easy for the user to just take a glance and infer peak traffic hours.

Fig. 11

Hourly Traffic in Florida for January 2

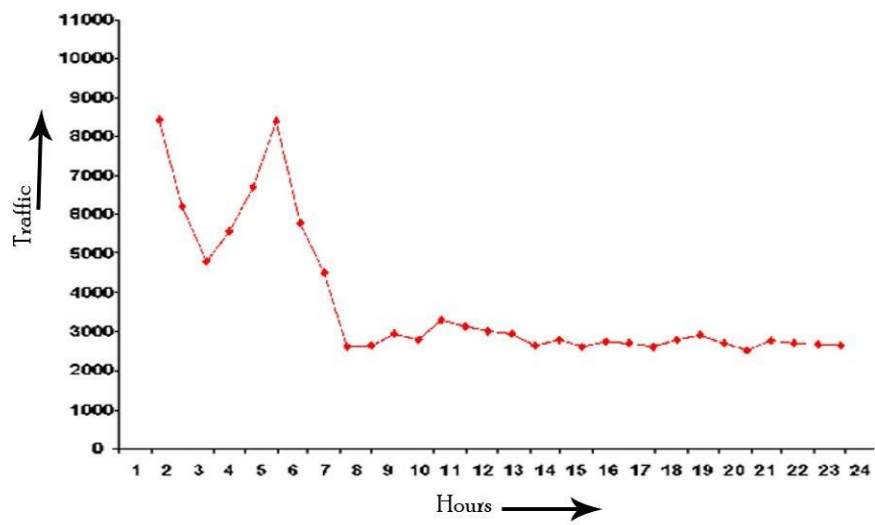


Fig. 11

While the modal for the histogram or line graph is open, rest of the information such as the heatmap and other buttons are defocused in the background so as to let the user isolate relevant information.

4.4 To view traffic volumes in rural and urban, **pie chart** was our choice since it conveniently shows the comparison of traffic volumes between the two (Fig. 12). Both **color** and **size** were used as preattentive variables here.

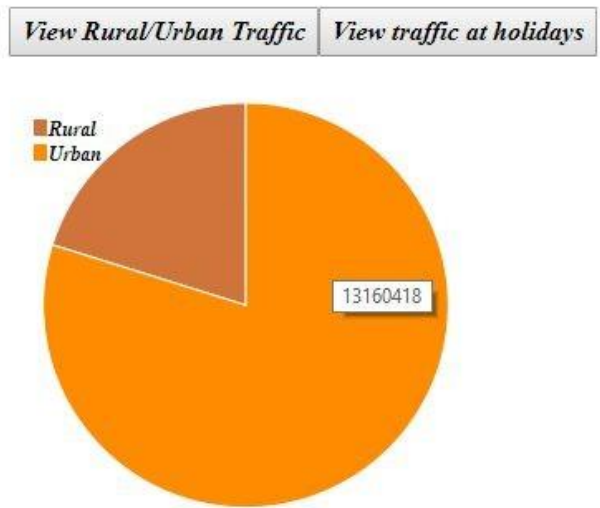


Fig. 12

4.5 We chose a **dynamic bubble chart** to represent the traffic data on holidays (Fig. 13). The size of bubbles corresponds to the volume of traffic, thus **size** being the preattentive variable. On mouse hover, titles appear that display the traffic volume and name of the holiday so as to let the users explore traffic trends during holidays.

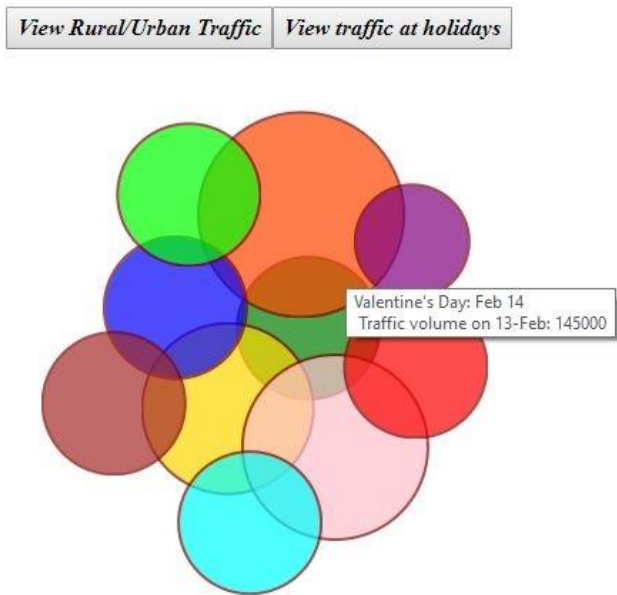


Fig. 13

4.6 Finally we used **dynamic queries** for filtering and sorting at various levels to provide atomic visualization options to the users.

All these views enable various levels of details to provide adequate information and their layering and positioning helps in decluttering so as not to make it difficult for the user by providing too much information at once.

This functional design is aimed to serve multiple user segments. Some examples of **user scenarios** that could benefit by using our dashboard are -

- common citizens - who want to view traffic details
- data analysts - to study road traffic trends
- cab companies - to plan their strategies
- traffic control department - to understand the traffic distribution and likelihood of congestion, and to manage assigning force accordingly

5.0 System Testing

To validate that the requirements we had set initially were being met and that our user-centered design framework goal was achieved, we performed continuous **usability testing** on our system. We approached the design process in a way to make it user friendly and keep the visible options at a single time limited while maintaining the delivery of relevant information.

Rather than testing the system in entirety upon completion, we adopted **unit testing** approach wherein we kept testing the working of our dashboard at various stages to check how various elements perform. The first phase was to test it ourselves in order to check each functionality. While testing the system for proper working of the country heatmap, we figured that it would be a good idea if we incorporated similar heatmaps for each state on mouse click. Since we had already implemented a bar chart on mouse click that appeared in a modal, and it was getting really difficult to open two modals at the same time, we prepared a wireframe mock-up of the state heatmaps. The idea was to show not only the traffic volume heat in different counties and stations within a state, but also the routes connecting these places and differentiating between the busy and light ones (Fig. 14).



Fig. 14

Upon testing more functionality, we discovered that even though we had performed data pre-processing at the onset itself, we still needed to use python for more data processing during development stages as well in order to add more columns, drop the columns that we found were irrelevant during the course of development, and to perform further aggregations.

Further we performed **integration testing** to ensure that the combination of various elements worked fine. In order to do this we selected a month from the drop down menu to check if the heatmap changes colors accordingly, the values in legend update on the basis of that month's data, the label values updated, and the overview table showed the correct highlights. In this process we discovered that the label values show too much information when no month is selected. While we tried to remove this feature for the yearly data (when none option is selected in the drop down menu), it was getting removed from the rest of the options as well (when a month was selected). Therefore we chose to let the labels be the same for all options.

Moreover, we performed **system testing** by running our dashboard on different operating systems namely Windows 10 and Mac OS. It turned out that the dashboard worked consistently on all of these. The load-time and refresh rate varied slightly depending upon the configuration of the systems.

One area where our system did not perform as expected was during **regression testing**. When the data size was increased above 0.5 GB, javascript would run out of heap memory and the dashboard wouldn't load. We tried fixing this issue but our attempts didn't seem to work. Finally we worked with a small data size for now since this project was mainly focused on the User Interface Design.

Finally we performed **acceptance testing** with the help of our colleagues. When the dashboard was presented to users, we got various types of feedback that helped us make some useful changes in order to make the GUI more user friendly. Some of the changes implemented at this stage were as follows -

- Placement of buttons to provide better clarity
- Removing the heatmap scrollbar - initially we had added a scrollbar to implement some functionality but that appeared to the users as if the drop down menu to select months and the heatmap were different sections, whereas they were supposed to be one section as a whole. To avoid this confusion, we removed the scrollbar and provided a complete page scroll functionality instead
- Levels of complexity - since there is a lot of information that can be inferred from our dataset, it was important to display views in a way that it doesn't overwhelm the user. Keeping this in mind, we used the "details on demand" approach wherein limited elements are visible at once and the user can click or hover to go into the deeper layers of functionality.
- Finally, we made some minor changes like pleasing fonts, vibrant colors, combination of colors to provide appeal yet avoid distraction, and so on.

6.0 Conclusion and Future Work

Big Data presents numerous opportunities to draw inferences and study trends but is also cumbersome to process and analyze in its raw form. Thus here we presented a working dashboard that visualizes US road traffic data for the year 2015. Our visualization addresses problems arising from large scale, multi-modal, unstructured data and provides comprehensible visuals to the users that are appealing to the eyes and make the job of analysis and forecasting easy and fun. We adopted various design patterns to show different aspects of data along with the corresponding pre-attentive variables to focus the user's attention on the highlights. The levels of complexity in our dashboard have been taken care of by displaying views on different layers. We performed various types and levels of testing including feedback from users from different segments and incorporated changes accordingly so as to make the GUI more user-friendly.

In future work, we will address augmentations that can be done to our current implementation - (i) Using large size dataset by storing it in a database like MongoDB or other and then writing web services to fetch the data and visualize it using dynamic queries. (ii) Using real-time data instead of static data so as to let the decision-makers view current road traffic conditions and plan their strategies accordingly. Kafka can be used for this purpose. (iii) Adding another vertical by creating an app for our dashboard so as to make this analytics platform mobile and provide the decision-makers on the go services (iv) Currently, the dashboard refers to road traffic only, further air-traffic and other data can also be used for similar visualizations that will help in horizontal growth.

References

1. AAA NewsRoom. (2017). *Americans Spend an Average of 17,600 Minutes Driving Each Year* / AAA NewsRoom. [online] Available at: <http://newsroom.aaa.com/2016/09/americans-spend-average-17600-minutes-driving-year/> [Accessed 3 Dec. 2017].
2. Metrocosm. (2017). *Watch 24 Hours of Traffic, Visualized as a Living Circulatory System* - Metrocosm. [online] Available at: <http://metrocosm.com/map-us-traffic/> [Accessed 3 Dec. 2017].
3. Chen, W., Guo, F. and Wang, F. (2015). A Survey of Traffic Data Visualization. [online] Available at: http://www.cad.zju.edu.cn/home/vagblog/VAG_Work/IEEEITS2015/TrafficVASurvey.pdf [Accessed 2 Dec. 2017]
4. Picozzi, M., Verdezoto, N. and Pouke, M. (2013). Traffic Visualization Applying Information Visualization Techniques to Enhance Traffic Planning. [online] Available at: <https://pure.au.dk/ws/files/70285046/TrafficVisualizationIVAPP2013.pdf> [Accessed 2 Dec. 2017]
5. Data link: <https://kaggle.com>
6. d3 library: <http://d3js.org/d3.v4.js>
7. Other libraries: <https://maxcdn.bootstrapcdn.com/bootstrap/4.0.0-alpha.6/css/bootstrap.min.css>
8. Other libraries: <https://code.jquery.com/jquery-3.1.1.slim.min.js>
9. Other libraries: <https://cdnjs.cloudflare.com/ajax/libs/tether/1.4.0/js/tether.min.js>
10. Other libraries: <https://maxcdn.bootstrapcdn.com/bootstrap/4.0.0-alpha.6/js/bootstrap.min.js>
11. Datamaps library: <http://datamaps.github.io/>
12. Python IDE: <https://www.jetbrains.com/pycharm/>
13. Visual Studio Code software: <https://code.visualstudio.com/download>
14. Code referred: <https://github.com/sakethsaxena/leading-causes-of-death-US-app>