

The Map-based Airlines' Visualization with Twitter Sentiment Analysis

Yuncong Ma, Weixi Liu

Show all airports

@SouthwestAir for the win as always- saved my day and got me on a direct to Orlando. 🙌🙌🙌

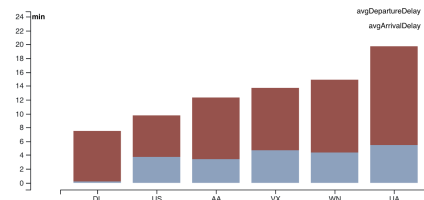
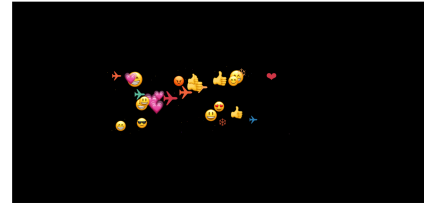
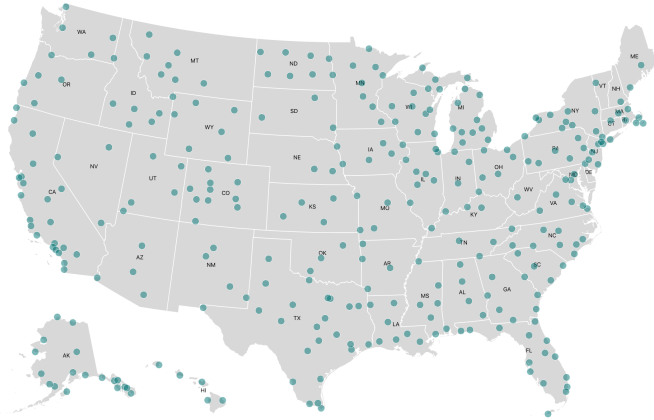


Fig. 1: Project Overview

Abstract— The goal of this web-based consulting interactive visualization is mainly to display statistic inform for the past airlines by combining three kinds of visualization patterns, and it provides a chance for audiences to explore the potential useful inform throughout the vis patterns. At the meantime, we combined the advantage of the React.js with D3 to creative interactive visualizations and using the Amazon Website Service to deal with big data problem (instead of uploading the whole large dataset along with website, we store the data into MongoDB and transport information via AWS). This vis is consisted of a US national map, of which circles represent different airports that are distributed across the entire country. The position of each airport is located according to the longitude and latitude data coming from us.json file. While the users move the mouse over the map, the chosen state or airport will be highlight and the more corresponding detail inform will pop up. As a extend practice of assignment 4, there are another two coordinated views sited beside the main map view—a stack bar chart and a word cloud. The stack bar chart is used to display the statistical data of the degree of flight delays among six different airline companies. Since the passengers' feedback is the second half of the vis topic, the word cloud will present those key words of which the appearance has highest frequency.

Index Terms—Sentiment Analysis, Interactive Visualization, React, Amazon Website Service, MongoDB, Word Cloud.

1 INTRODUCTION

As a major public transportation, airline market has been ushering a booming development. On a global scale, a continuous world-wide growth of air traffic could be observed, and according to several market researches, the growth is expected to maintain positive rates up to 2030.

However, there are many factors that affect the performance of the commercial aviation system, which can lead to annoying results to their passengers sometimes. Given the uncertain factors of the whole aviation system, passengers usually have to plan their travel many days or even months before the departure date. Meanwhile, in order to de-

crease the trip costs, avoid the rush traffic hours, and then obtain a relaxed travel experience, travelers also hope to gain as more detailed information as possible.

Converting the traditional numeric information into a more vivid visualization form, could help the viewers gain their desired information efficiently and easily. So we intend to build a map based interactive consulting visualization, which combines two date sets coming from the US Department of Transportation Bureau of Transportation Statistics. We hope this application can reveal some potential patterns under the flight records and display them to the viewers.

We apply D3.js library to build the whole data visualization including one map view which is used to depict the airports and the airlines. While users move the mouse over an area belonged to a specific state, the area will highlight of which the color change and the name of that state will pop up. Once users click the state, it will filter out other states, that only displays the airports sited inside the chosen state. For sure, we also provide a button placed at the top middle of the web page, to reset the map view to its original status. Correspondingly, the right-side bar chart will change once users click a specific state.

- Yuncong Ma is with Worcester Polytechnic Institute. E-mail: yma10@wpi.edu.
- Weixi Liu is with Worcester Polytechnic Institute. E-mail: wliu6@wpi.edu.

Manuscript received 31 March 2014; accepted 1 August 2014; posted online 13 October 2014; mailed on 4 October 2014.

For information on obtaining reprints of this article, please send e-mail to: tvcg@computer.org.

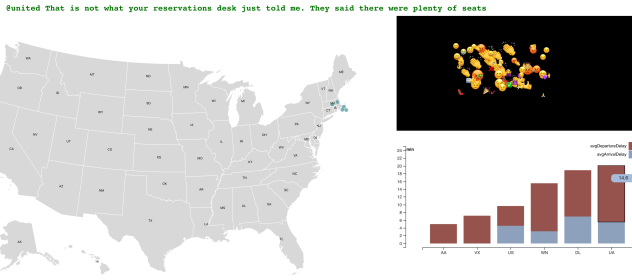


Fig. 2: A snapshot of an interaction.

The bar chart we use here is to display a comparison of flight delay period among at most six airline companies. This bar chart is a variety of the plain bar, which is divided horizontally into upper and lower part, each of them present the degree of departure and arrival delay respectively. Initially, this stack bar chart will display average delay time across all the flight data. Of course, to make the bar chart more practical, we also append the numeric value of the flight delay time while users move the mouse over a specific part of the bar chart.[9] Meanwhile, the word cloud view will change as a reaction of the click action on the stack bar chart, the key words it shows every time will change correspond with the content of passengers' tweets. Most of the words are the reflection of the feel of the passengers, either is positive or negative.

2 BACKGROUND

As the time goes by, the use of coordinated multiple views has been changing and expanding a lot, in addition it also becomes part of larger sense making environments where the techniques are being used to analyze large datasets, integrate alternate viewpoints, and generate nuggets of information.[16] Nowadays, D3.js library is one of the most popular tools to implement the coordinated multiple views and then analysis large data set. It is worth and quite practical to apply these ideas and tools while building the final project. Meanwhile, data visualization is not just a way that simply transforms the data into several tables or charts, instead, it also involves pre-processing on data such as clearing, filtering, mapping or other aggregation operation. The process that chose appropriate visualization pattern with the dataset is challenging, but on the other hand, a good vis always provides its audiences an intuitive and logical experience. Utilizing all the handy techniques we have to develop an extension of the previous project, is a good study path for our future work.

There are three important parts: designing the whole visualization, fixing the big data problem (our original dataset contains 300M+ rows), achieving the interactions between visualizations with the help of d3 and react instead of using dispatch.

3 METHOD

Basically we used D3.js to build each data visualization including US Map, Stacked Bar Chart, Word Cloud.

3.1 Interaction

People usually use dispatch which is a tool to create interaction between data visualization views to create interactions, but we decide to use React.js to do the same thing instead of dispatch.

React is a JavaScript library for building user interfaces, the core idea behind it is integrating separate parts into a whole thing (class) and controlling each part's state. For example, we can create a simple spinning button by combining a button element, a spinning figure, and a state variable to control spinning or not. We embed the three things into a class SpinningButton, we set the spinning state with false as default, when user click the button, setting the state with true, showing the spinning figure and starting up a timer function as the same time, when time running out, set the state with false back, hidden the spinning figure.

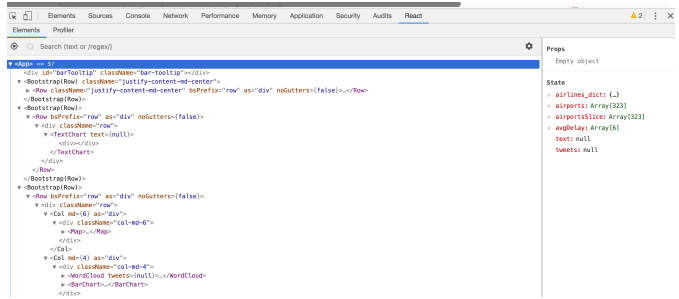


Fig. 3: React states example.

Collection Name *	Documents	Avg. Document Size	Total Document Size	Num. Indexes	Total Index Size
airlines	14	74.4 B	1.0 KB	1	32.0 KB
airports	223	186.6 B	58.9 KB	1	16.0 KB
flights	3,589,846	354.5 B	1.2 GB	6	116.6 MB
tweets	14,873	542.9 B	77 MB	2	232.0 KB

Fig. 4: Dataset in MongoDB.

As you can see, we can use this kind of idea to control our visualization information and transport the information by react states, the figure 3 shows our state.

3.2 Big Data Problem

As we mentioned before, our original dataset contains 300M+ rows, the data size is above 600MB, it is definitely not a good idea to store the data along with the website. Instead of uploading the dataset, we create a MongoDB (a NoSQL database) database to store our data. There are two reasons we choose MongoDB[1], the first reason is it is NoSQL data structure as JSON format, it's convenient to I/O data. It has a good compatible with Amazon Website Service which is also good for transporting data. Of course, MongoDB is also good for aggregation operation which help us to retrieve different kind of query combinations. [14]

Amazon Website Service could help us to build our APIs to retrieve MongoDB data in cloud, it has a high computing performance and compatible with MongoDB[13].

With the help of MongoDB and AWS, we don't need to uploading the whole large dataset which could make the website run slowly. On the other hand, we can keep the whole dataset without any pruning and modification which could lead information loss.

We can see in figure 4, our flights collection is up to 1.2G which is super large.

3.3 D3 Transition

When the data changes in a d3 views, what is the interesting thing we may notice? That is transition. It's boring if a data visualization just change their view by refreshing the whole view, that is deleting the whole dataset, and replacing with new dataset. D3 provides a really good way to deal with this problem, it's called, enter, update and exit as shown in figure 5. It's super useful and interesting part of D3.js and also it's hard to understand, as lease we spent a lot of time to understand these kind of concepts. Generally, the processes are, supposing our view has been implemented. Now, we need to update the view by the new dataset. First, the D3 collects the previous dataset and new coming dataset, and secondly, the D3 finds the common part and different part (we could or we better specify which data field we are going to compare), and lastly, the D3 update the different part and remove the useless data. From these procedures, we could efficiently and perfectly control the transitions.[3]

Fig. 5: D3 Transition example.

Fig. 6: Emoji Word Cloud example.

3.4 Word Cloud

Word Cloud is a powerful tool for sentiment analyzing[10]. We firstly tried to tokenize Tweeter user's whole text and count each word's frequency for showing in the word cloud. But it turns out not a good choice since there are so many stop words and useless characters, even though we preprocessed the user's text, eliminate all possible stops words and special characters[15], our results are not as good as what we think (see in figure 7). Thus we choose an alternate way to present user attitude, that is extract all emoji characters! We can see although it's not good enough, we can still tell the overall attitude from this emoji word cloud (see in figure 6).

4 RESULTS

With the help of D3.js, React.js, MongoDB and Amazon Website Service, we successfully create our whole interactive data visualization. Our visualization shows the average airline departure delay time, the average airline arrival delay time and Tweeter user's sentiment. User could click any state to show the information corresponding the all airports in this specific state. User also could click the stack bar chart to show each specific airline's corresponding information.

5 OTHER USEFUL ARTICLES INFLUENCE US

It gives an idea that how to depict a visualization of adjacency relations in hierarchical data, especially with a huge dataset. [12]

They presented a visual analysis of Twitter time-series, which combines sentiment and stream analysis with geoand time-based interactive visualizations for the exploration of real-world Twitter data streams. [6]

It introduced a novel visualization called NodeTrix.[11]

This paper mainly introduces how and why the authors visualize cyber security session data by an interactive behavior map. They encoded an action as city and user sessions as trajectories going through the cities. They illustrated how they explore relationships between

Fig. 7: Bad Word Cloud example.

actions, identity patterns of the typical session and detect anomaly behaviors.[4]

The D3.js library has become a super popular tool for developing visualization on the website. More and more people use this kind of technology to improve their website outlook as well as the actual needs. However, it is time-consuming when people only want to change their visualization's appearance. This article tells how to use a pair of tools for deconstructing and restyling existing D3 visualization without understanding the underlying code and logic.[8]

This article presents their approach to combine sentiment analysis with a new term association technique as well as a geo-temporal visualization for an effective analysis of large customer feedback streams. They introduce a feature and geo-based stream analysis technique that automatically detects which attributes (features) are frequently commented on, which attributes have interesting sentiments patterns, which attributes cluster significantly in certain geo-locations, and what terms often occur together. They also introduce two new geo-temporal visualizations (pixel sentiment and key term geo maps) that help users analyze large volumes of web surveys and Twitter data.[7]

Nowadays data visualization exists in everywhere our lives. We can see data visualization in newspaper, TV, Internet, commercial advertisements, etc. With the widely use of data visualization, there is a growing awareness of the uncertainty problem within the visualization community, and many traditional techniques are being extended to represent not only the data, but also the uncertainty information associated with the data. We call this visualization of uncertainty. In addition, it is important to realise that, even if there is certainty about the data, errors can occur in the process of turning the data into a picture. We call this uncertainty of visualization. This paper aims to review the current state of the art in uncertainty in scientific visualization, looking at both of these aspects.[2]

Automatic exploration of semantic classification of documents is becoming more and more popular, there are many tasks measuring the documents or chunks of text such as automatic sentiment analysis, which measures the text based on emotive categories such as positive and negative. However, there is limited progress on how to visualize the affective content of documents and describes an interactive capability for exploring emotion in a large document collection.

This paper described an approach of how these measurements might be presented to the users for exploration and added value. This approach includes identifying the affective content of documents, as well as possible ways of visualizing it and presenting those results

analytically.[5]

5.1 Conclusion

This paper is mainly trying to find a way to create a interactive visualization by combining React.js and D3.js. The visualization is also helpful for people who want to know what is the average delay of each airlines and what is the people's attitude by presenting Tweeter's sentiment information.

It also find a way to deal with big data problem when the dataset it's too large to upload to website.

Moreover, it gives a good example to use D3 transition for improving the visualization's quality.

At the end, we find that combining React and D3 is a good way to develop interactive visualization, it's efficient, easy to control and extensible. Although our current visualization's interaction is not fast, that is because the transporting problem between the AWS and MongoDB (it may be resulted by some bad query sentences), and the dataset is too large to do complicate query operation like finding all given airline in some airports. That is nothing to do with the React and D3.

In the future, we could add more views to allow user to discover more information. For example, we could add a timeline to find the delay information according the time range. We also could add a functionality that user could click any two of airports to see the corresponding information between the two airports, or click any two of us states to show corresponding information. Furthermore, we could find a more efficient and useful way to preprocess Tweeter User's tweet and present them in Word Cloud.

ACKNOWLEDGMENTS

We would like to acknowledge Professor Lane Harrison for his insight in helping come up with ideas on how to build the visualization and the guidance provided throughout the project.

REFERENCES

- [1] A. Boicea, F. Radulescu, and L. I. Agapin. Mongodb vs oracle-database comparison. In *2012 third international conference on emerging intelligent data and web technologies*, pages 330–335. IEEE, 2012.
- [2] K. Brodlie, R. A. Osorio, and A. Lopes. A review of uncertainty in data visualization. In *Expanding the frontiers of visual analytics and visualization*, pages 81–109. Springer, 2012.
- [3] A. Buja, J. A. McDonald, J. Michalak, and W. Stuetzle. Interactive data visualization using focusing and linking. In *Proceeding Visualization'91*, pages 156–163. IEEE, 1991.
- [4] M. F. De Oliveira and H. Levkowitz. From visual data exploration to visual data mining: a survey. *IEEE Transactions on Visualization and Computer Graphics*, 9(3):378–394, 2003.
- [5] M. L. Gregory, N. Chinchor, P. Whitney, R. Carter, E. Hetzler, and A. Turner. User-directed sentiment analysis: Visualizing the affective content of documents. In *Proceedings of the Workshop on Sentiment and Subjectivity in Text*, pages 23–30. Association for Computational Linguistics, 2006.
- [6] M. Hao, C. Rohrdantz, H. Janetzko, U. Dayal, D. A. Keim, L.-E. Haug, and M.-C. Hsu. Visual sentiment analysis on twitter data streams. In *2011 IEEE Conference on Visual Analytics Science and Technology (VAST)*, pages 277–278. IEEE, 2011.
- [7] M. C. Hao, C. Rohrdantz, H. Janetzko, D. A. Keim, U. Dayal, L. E. Haug, M. Hsu, and F. Stoffel. Visual sentiment analysis of customer feedback streams using geo-temporal term associations. *Information Visualization*, 12(3-4):273–290, 2013.
- [8] J. Harper and M. Agrawala. Deconstructing and restyling d3 visualizations. In *Proceedings of the 27th annual ACM symposium on User interface software and technology*, pages 253–262. ACM, 2014.
- [9] R. M. Heiberger, N. B. Robbins, et al. Design of diverging stacked bar charts for likert scales and other applications. *Journal of Statistical Software*, 57(5):1–32, 2014.
- [10] F. Heimerl, S. Lohmann, S. Lange, and T. Ertl. Word cloud explorer: Text analytics based on word clouds. In *2014 47th Hawaii International Conference on System Sciences*, pages 1833–1842. IEEE, 2014.
- [11] N. Henry, J.-D. Fekete, and M. J. McGuffin. Nodetrix: a hybrid visualization of social networks. *IEEE transactions on visualization and computer graphics*, 13(6):1302–1309, 2007.
- [12] D. Holten. Hierarchical edge bundles: Visualization of adjacency relations in hierarchical data. *IEEE Transactions on visualization and computer graphics*, 12(5):741–748, 2006.
- [13] K. R. Jackson, L. Ramakrishnan, K. Muriki, S. Canon, S. Cholia, J. Shalf, H. J. Wasserman, and N. J. Wright. Performance analysis of high performance computing applications on the amazon web services cloud. In *2010 IEEE second international conference on cloud computing technology and science*, pages 159–168. IEEE, 2010.
- [14] D. Keim, H. Qu, and K.-L. Ma. Big-data visualization. *IEEE Computer Graphics and Applications*, 33(4):20–21, 2013.
- [15] G. Miner, J. Elder IV, A. Fast, T. Hill, R. Nisbet, and D. Delen. *Practical text mining and statistical analysis for non-structured text data applications*. Academic Press, 2012.
- [16] J. C. Roberts. State of the art: Coordinated & multiple views in exploratory visualization. In *Fifth International Conference on Coordinated and Multiple Views in Exploratory Visualization (CMV 2007)*, pages 61–71. IEEE, 2007.