

Northeastern University



College of Professional Studies, Northeastern University, Boston, MA 02215

ALY 6110 Data Management & Big Data

CRN: 72083

Project name: U.S. Traffic Accident Analysis

Instructor: Valeriy Shevchenko

Team:

Anju Yeh: yeh.an@northeastern.edu

Binyu Zhai: zhai.b@northeastern.edu

Liuzhao Tang: tang.liu@northeastern.edu

Jincheng Wang: wang.jinc@northeastern.edu

Zhongying Lan: lan.z@northeastern.edu

Summary

In this project, we use the U.S. Traffic Accidents dataset from Kaggle. There are 3.5 million records and 49 columns. This is a countrywide car accident dataset, which covers 49 states of the USA. It was collected from February 2016 to June 2020. The columns contain the Location of accidents, the Severity of accidents and the Weather condition of accidents, etc.

We will use Databricks as our cloud platform for analyzing. Besides, we choose Python language to do the whole process. Include cleaning up the raw data, performing the exploratory data analysis (EDA) and drawing a series of traffic accidents maps to figure out what is the trend of the U.S. traffic accidents. And then, we build a model to do the feature selection. At last, we would provide some relevant insights and suggestions.

This project aims to reveal the factors that might lead the accident and even its severity changes. We would like to find out as well as to show the concepts such as what time scale is the peak, which city takes the lead by the highlights and visualizations to the audience. Also, we will make modeling and predictions to indicate what factors are significant to the accidents' severities and try to predict the result. In general, we expect the result of the project can give an overall view to the audiences. As well as we hope that people can avoid tragic accident by paying attention to the key factors that are pointed out by the project.

Content

Step 1. Data clean

We conduct a preliminary clean-up of the dataset in order to perform further analysis. At very first, we quickly peek through the data shows we have all sorts of data types, like string, datetime, float, boolean, and integers. Then, we check the missing values, the outcome is showing as follow:

	Total	Percent			
End_Lat	2478818	70.548896	Wind_Direction	58874	1.675595
End_Lng	2478818	70.548896	Pressure(in)	55882	1.590441
Number	2262864	64.402694	Weather_Timestamp	43323	1.233003
Precipitation(in)	2025874	57.657793	Airport_Code	6758	0.192337
Wind_Chill(F)	1868249	53.171675	Timezone	3880	0.110428
TMC	1034799	29.451104	Zipcode	1069	0.030424
Wind_Speed(mph)	454609	12.938490	Nautical_Twilight	115	0.003273
Weather_Condition	76138	2.166941	Astronomical_Twilight	115	0.003273
Visibility(mi)	75856	2.158915	Civil_Twilight	115	0.003273
Humidity(%)	69687	1.983341	Sunrise_Sunset	115	0.003273
Temperature(F)	65732	1.870779	City	112	0.003188
			Description	1	0.000028

Then, fill 'Humidity (%)', 'Precipitation(in)', 'Wind_Chill(F)', 'Wind_Speed(mph)', 'Visibility(mi)' with 0 because it is possible to have no records for these columns. For example, it's possible to have zero rain if rain didn't fall that day. Besides, we fill the 'Temperature(F)' and 'Pressure(in)' with the average value. Then, we Fill the 'City' column. Since we have a 'State' column. So, we fill the city column with the most occurring city of the state it belongs to. At last, we drop useless columns and empty rows. Following are the reasons for deleting these columns:

- 1.It's hard to impute End-lat and End-lng, because it can only be collected when the accident affected a huge area of road. And it's 77% the same as start-Lat and start-lng. So, dropping is the logical choice.
2. TMC stands for Traffic Message Channel. Its purpose is to deliver information about traffic distortions or warnings to mobile receivers such as navigation devices. Which we don't need.
3. Weather_Timestamp: Shows the time-stamp of the weather observation record (in local time).
4. Drop some useless columns like ID, Source, TMC, etc.

There is zero missing value in our dataset now.

```

Cmd 11
1 df.isnull().any().sum()

Out[31]: 0

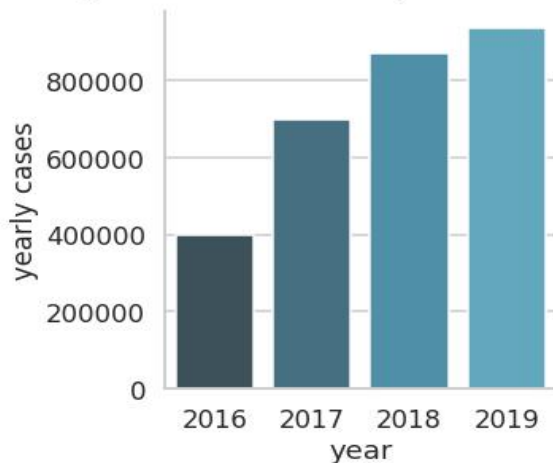
```

Step 2. Exploratory Data Analysis

2.1 Analyze number of accidents by Time

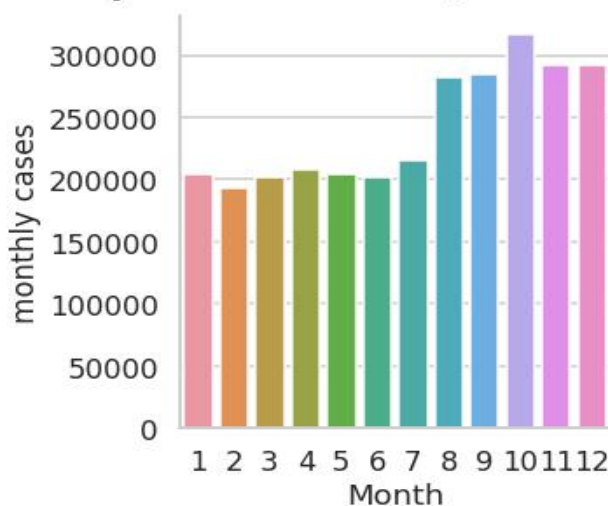
In this section, first we have made statistics on traffic accidents that occurred in the four years from 2016 to 2019. From the figure, we can find that in these four years, traffic accidents have been increasing year by year. In 2019, there were approximately 900,000 traffic accidents in the United States.

Yearly accidents cases(2016-2019)



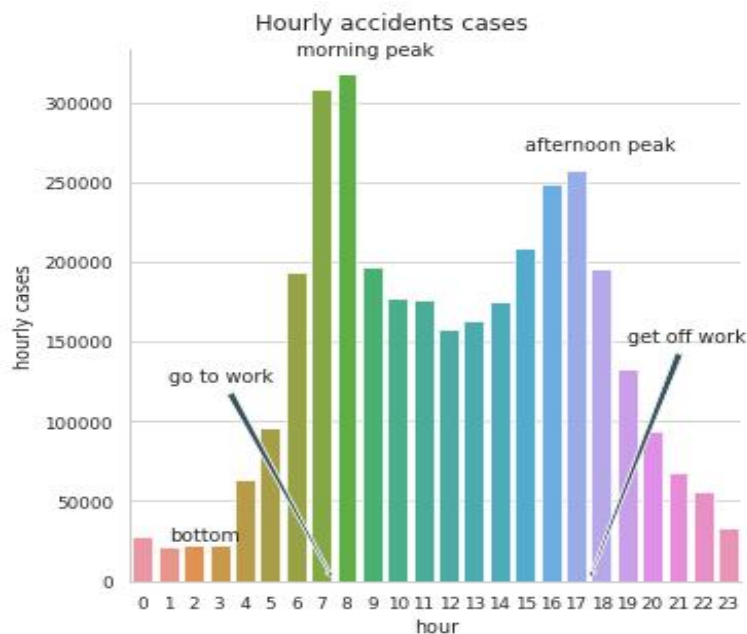
Next, we counted the average traffic accidents that occurred every month for the period from 2016 to 2019. Through analysis, we can know that the number of traffic accidents in the five months from August to December is significantly higher than other months. We guess there are worse weather conditions in the summer and winter. Weather factors may be the cause of traffic accidents

monthly accidents cases(2016-2019)



Finally, we counted the average number of traffic accidents that occurred per hour in a day. From the figure, we can know that most accidents occur during the day, and peaks occur at 8 am and 5 pm. Eight in the morning is when people go to work and five in the afternoon is when people get off work to get home. These two time periods are

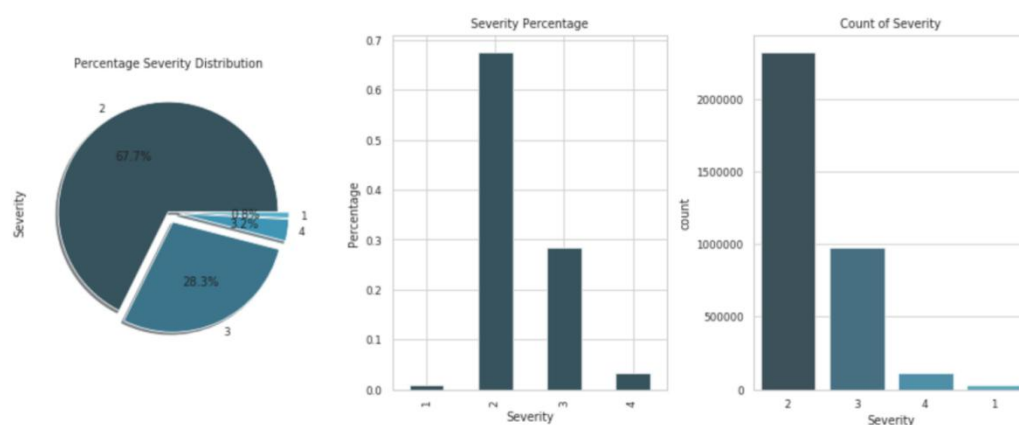
peak times for people to travel, so the possibility of traffic accidents is greater. Compared with these two times, at other times, people tend to rest at home instead of traveling, so the number of traffic accidents is much lower.



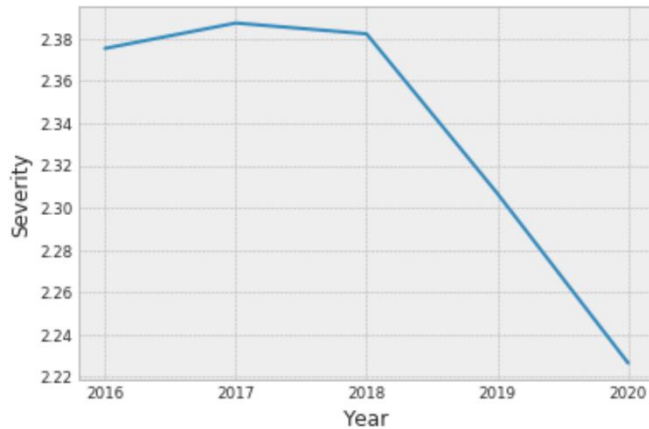
2.2 Analyze number of accidents by Severity

Column 'Severity' shows the severity of the accident, the number between 1 and 4, where 1 indicates the least impact on traffic (i.e., short delay as a result of the accident) and 4 indicates a significant impact on traffic (i.e., long delay).

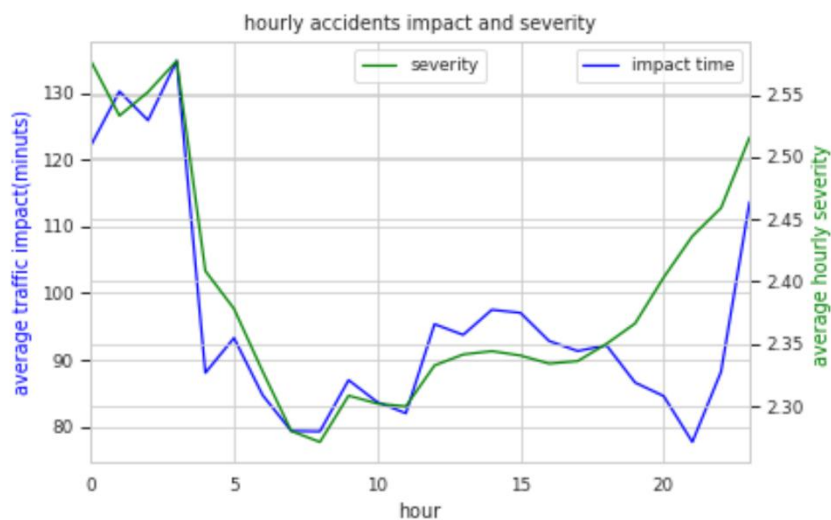
First, we take a look at this feature. Obviously, 67.7 % Accidents fall in the Severity class 2 followed by Severity class 3,4 and 1.



Then, we try to figure out the 'severity' changes in years from 2016 to 2020. It's easy to find that it has declined year by year from 2017 to 2020.

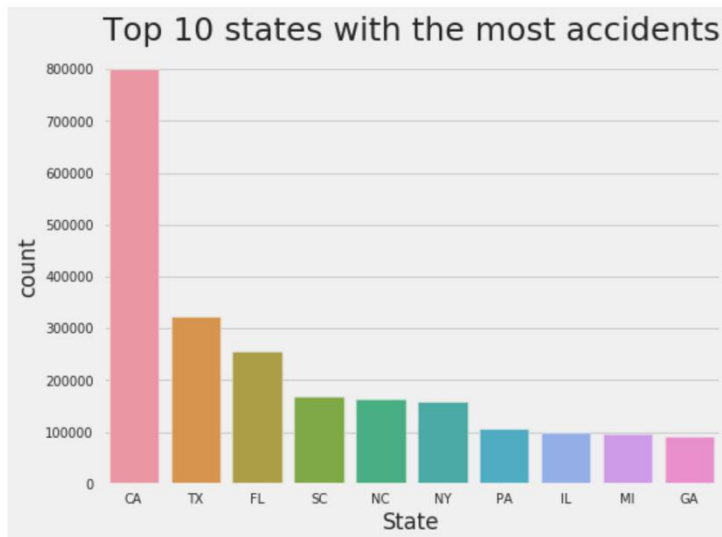


At last, we create the following chart to discuss the severity changes hourly. Serious car accidents always happen between 5pm and 3am.

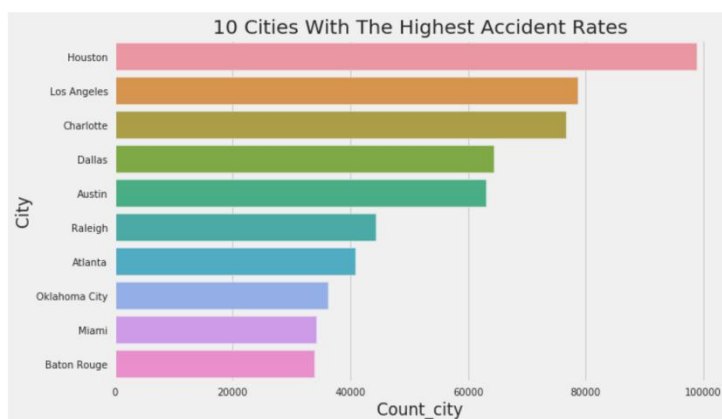


2.3 Analyze number of accidents by Weather condition

In this section, our goal is trying to identify what kind of weather conditions cause the most car accidents. There are a total five weather conditions that can be found from the following pie chart such as partly cloudy, overcast, mostly cloudy, fair and clear. As we have excluded the minimal cause data such as light rain, rain, light snow, thunderstorm and scattered clouds etc. in order to extract the most significant factors that caused the car accidents. Surprisingly, most car accidents happened during the clear weather with 31.4%. Therefore, we could infer that drivers are usually more cautious driving during poor weather conditions.



Taking a closer look into the city bar chart, we found out that Houston has the most accidents followed by Los Angeles and Charlotte among all US cities. In this case, we came up with a point that Texas is the 2nd biggest state in the USA, thus this can lead to more drivers speeding to reach their destination faster.



Step 3. Draw maps

Figure X displays maps of traffic accidents from 2016 to 2020. The first finding from the map is that the cases of traffic accidents are increasing. Second, we can find that most of the cases occur along the coastline area and Great Lakes Region. Finally, the traffic accidents in the central region are increasing in number and scope year after year.

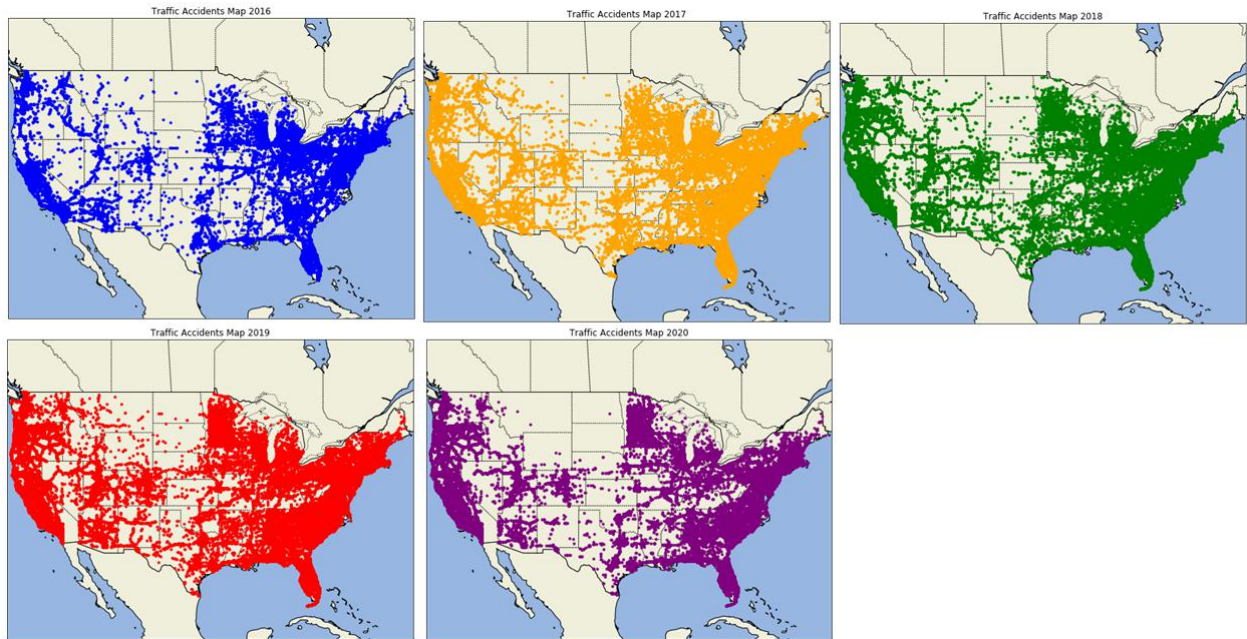


Figure X: Maps of traffic accidents from 2016 to 2020

Step 4. Feature Selection and Modeling Prediction

In this step we mainly perform two processes to make some determinations, which are feature selections by random forest and forecasting severity of accident by multi-logistic regression modeling. By the feature selection, we would like to introduce the key factors that lead the severity, as well as predict the severity of accidents by model.

We believe that this is an important and interesting step because data can usually give the audience some more surprising information than the imagination. Here, we can logically assume some key factors such as the humidity, wind level, and weather might affect the severity of accidents.

First, we need to prepare more on reshaping the data set before doing the work. If we directly use “get_dummies” to transform the categorical information, we will get more than 10000 columns, which makes the work to be harsh. So, we drop some not essential columns like “Country”, “City”, “Descriptions”, etc. and then the data size has been cut down to the 203 columns as shown below.

Out[15]:

	Severity	Start_Lat	Start_Long	Distance(mi)	Temperature(F)	Wind_Cat(0F)	Humidity(%)	Pressure(in)	Visibility(mi)	Wind_Speed(mph)	Precipitation(in)	Alcohol	Bump	Crossing	One_Way	Junction	No_Sall
ID																	
2218289	3	44.981962	-81.236412	0.00	10.0	-5.9	87.0	30.14	10.0	13.0	0.00	False	False	False	False	False	False
2482082	3	30.409206	-81.857265	0.00	73.9	0.0	97.0	29.91	10.0	4.0	0.00	False	False	False	False	True	False
192924	5	42.000784	-83.648857	-0.01	32.0	24.8	96.0	30.17	0.0	8.1	0.04	False	False	False	False	False	False
2192046	3	35.520986	-87.514491	0.00	34.0	23.0	44.0	30.16	10.0	16.1	0.00	False	False	False	False	False	False
826302	2	38.241489	-122.269399	0.00	58.0	38.0	86.0	29.99	0.0	0.0	0.00	False	False	False	False	False	False

8 rows × 203 columns

Control panel: 8/22 | 10/22 | 11/22 | 12/22 | 1/23 | 2/23 | 3/23 | 4/23 | 5/23 | 6/23 | 7/23 | 8/23 | 9/23 | 10/23 | 11/23 | 12/23 | 1/24 | 2/24 | 3/24 | 4/24 | 5/24 | 6/24 | 7/24 | 8/24 | 9/24 | 10/24 | 11/24 | 12/24 | 1/25 | 2/25 | 3/25 | 4/25 | 5/25 | 6/25 | 7/25 | 8/25 | 9/25 | 10/25 | 11/25 | 12/25 | 1/26 | 2/26 | 3/26 | 4/26 | 5/26 | 6/26 | 7/26 | 8/26 | 9/26 | 10/26 | 11/26 | 12/26 | 1/27 | 2/27 | 3/27 | 4/27 | 5/27 | 6/27 | 7/27 | 8/27 | 9/27 | 10/27 | 11/27 | 12/27 | 1/28 | 2/28 | 3/28 | 4/28 | 5/28 | 6/28 | 7/28 | 8/28 | 9/28 | 10/28 | 11/28 | 12/28 | 1/29 | 2/29 | 3/29 | 4/29 | 5/29 | 6/29 | 7/29 | 8/29 | 9/29 | 10/29 | 11/29 | 12/29 | 1/30 | 2/30 | 3/30 | 4/30 | 5/30 | 6/30 | 7/30 | 8/30 | 9/30 | 10/30 | 11/30 | 12/30 | 1/31 | 2/31 | 3/31 | 4/31 | 5/31 | 6/31 | 7/31 | 8/31 | 9/31 | 10/31 | 11/31 | 12/31 | 1/32 | 2/32 | 3/32 | 4/32 | 5/32 | 6/32 | 7/32 | 8/32 | 9/32 | 10/32 | 11/32 | 12/32 | 1/33 | 2/33 | 3/33 | 4/33 | 5/33 | 6/33 | 7/33 | 8/33 | 9/33 | 10/33 | 11/33 | 12/33 | 1/34 | 2/34 | 3/34 | 4/34 | 5/34 | 6/34 | 7/34 | 8/34 | 9/34 | 10/34 | 11/34 | 12/34 | 1/35 | 2/35 | 3/35 | 4/35 | 5/35 | 6/35 | 7/35 | 8/35 | 9/35 | 10/35 | 11/35 | 12/35 | 1/36 | 2/36 | 3/36 | 4/36 | 5/36 | 6/36 | 7/36 | 8/36 | 9/36 | 10/36 | 11/36 | 12/36 | 1/37 | 2/37 | 3/37 | 4/37 | 5/37 | 6/37 | 7/37 | 8/37 | 9/37 | 10/37 | 11/37 | 12/37 | 1/38 | 2/38 | 3/38 | 4/38 | 5/38 | 6/38 | 7/38 | 8/38 | 9/38 | 10/38 | 11/38 | 12/38 | 1/39 | 2/39 | 3/39 | 4/39 | 5/39 | 6/39 | 7/39 | 8/39 | 9/39 | 10/39 | 11/39 | 12/39 | 1/40 | 2/40 | 3/40 | 4/40 | 5/40 | 6/40 | 7/40 | 8/40 | 9/40 | 10/40 | 11/40 | 12/40 | 1/41 | 2/41 | 3/41 | 4/41 | 5/41 | 6/41 | 7/41 | 8/41 | 9/41 | 10/41 | 11/41 | 12/41 | 1/42 | 2/42 | 3/42 | 4/42 | 5/42 | 6/42 | 7/42 | 8/42 | 9/42 | 10/42 | 11/42 | 12/42 | 1/43 | 2/43 | 3/43 | 4/43 | 5/43 | 6/43 | 7/43 | 8/43 | 9/43 | 10/43 | 11/43 | 12/43 | 1/44 | 2/44 | 3/44 | 4/44 | 5/44 | 6/44 | 7/44 | 8/44 | 9/44 | 10/44 | 11/44 | 12/44 | 1/45 | 2/45 | 3/45 | 4/45 | 5/45 | 6/45 | 7/45 | 8/45 | 9/45 | 10/45 | 11/45 | 12/45 | 1/46 | 2/46 | 3/46 | 4/46 | 5/46 | 6/46 | 7/46 | 8/46 | 9/46 | 10/46 | 11/46 | 12/46 | 1/47 | 2/47 | 3/47 | 4/47 | 5/47 | 6/47 | 7/47 | 8/47 | 9/47 | 10/47 | 11/47 | 12/47 | 1/48 | 2/48 | 3/48 | 4/48 | 5/48 | 6/48 | 7/48 | 8/48 | 9/48 | 10/48 | 11/48 | 12/48 | 1/49 | 2/49 | 3/49 | 4/49 | 5/49 | 6/49 | 7/49 | 8/49 | 9/49 | 10/49 | 11/49 | 12/49 | 1/50 | 2/50 | 3/50 | 4/50 | 5/50 | 6/50 | 7/50 | 8/50 | 9/50 | 10/50 | 11/50 | 12/50 | 1/51 | 2/51 | 3/51 | 4/51 | 5/51 | 6/51 | 7/51 | 8/51 | 9/51 | 10/51 | 11/51 | 12/51 | 1/52 | 2/52 | 3/52 | 4/52 | 5/52 | 6/52 | 7/52 | 8/52 | 9/52 | 10/52 | 11/52 | 12/52 | 1/53 | 2/53 | 3/53 | 4/53 | 5/53 | 6/53 | 7/53 | 8/53 | 9/53 | 10/53 | 11/53 | 12/53 | 1/54 | 2/54 | 3/54 | 4/54 | 5/54 | 6/54 | 7/54 | 8/54 | 9/54 | 10/54 | 11/54 | 12/54 | 1/55 | 2/55 | 3/55 | 4/55 | 5/55 | 6/55 | 7/55 | 8/55 | 9/55 | 10/55 | 11/55 | 12/55 | 1/56 | 2/56 | 3/56 | 4/56 | 5/56 | 6/56 | 7/56 | 8/56 | 9/56 | 10/56 | 11/56 | 12/56 | 1/57 | 2/57 | 3/57 | 4/57 | 5/57 | 6/57 | 7/57 | 8/57 | 9/57 | 10/57 | 11/57 | 12/57 | 1/58 | 2/58 | 3/58 | 4/58 | 5/58 | 6/58 | 7/58 | 8/58 | 9/58 | 10/58 | 11/58 | 12/58 | 1/59 | 2/59 | 3/59 | 4/59 | 5/59 | 6/59 | 7/59 | 8/59 | 9/59 | 10/59 | 11/59 | 12/59 | 1/60 | 2/60 | 3/60 | 4/60 | 5/60 | 6/60 | 7/60 | 8/60 | 9/60 | 10/60 | 11/60 | 12/60 | 1/61 | 2/61 | 3/61 | 4/61 | 5/61 | 6/61 | 7/61 | 8/61 | 9/61 | 10/61 | 11/61 | 12/61 | 1/62 | 2/62 | 3/62 | 4/62 | 5/62 | 6/62 | 7/62 | 8/62 | 9/62 | 10/62 | 11/62 | 12/62 | 1/63 | 2/63 | 3/63 | 4/63 | 5/63 | 6/63 | 7/63 | 8/63 | 9/63 | 10/63 | 11/63 | 12/63 | 1/64 | 2/64 | 3/64 | 4/64 | 5/64 | 6/64 | 7/64 | 8/64 | 9/64 | 10/64 | 11/64 | 12/64 | 1/65 | 2/65 | 3/65 | 4/65 | 5/65 | 6/65 | 7/65 | 8/65 | 9/65 | 10/65 | 11/65 | 12/65 | 1/66 | 2/66 | 3/66 | 4/66 | 5/66 | 6/66 | 7/66 | 8/66 | 9/66 | 10/66 | 11/66 | 12/66 | 1/67 | 2/67 | 3/67 | 4/67 | 5/67 | 6/67 | 7/67 | 8/67 | 9/67 | 10/67 | 11/67 | 12/67 | 1/68 | 2/68 | 3/68 | 4/68 | 5/68 | 6/68 | 7/68 | 8/68 | 9/68 | 10/68 | 11/68 | 12/68 | 1/69 | 2/69 | 3/69 | 4/69 | 5/69 | 6/69 | 7/69 | 8/69 | 9/69 | 10/69 | 11/69 | 12/69 | 1/70 | 2/70 | 3/70 | 4/70 | 5/70 | 6/70 | 7/70 | 8/70 | 9/70 | 10/70 | 11/70 | 12/70 | 1/71 | 2/71 | 3/71 | 4/71 | 5/71 | 6/71 | 7/71 | 8/71 | 9/71 | 10/71 | 11/71 | 12/71 | 1/72 | 2/72 | 3/72 | 4/72 | 5/72 | 6/72 | 7/72 | 8/72 | 9/72 | 10/72 | 11/72 | 12/72 | 1/73 | 2/73 | 3/73 | 4/73 | 5/73 | 6/73 | 7/73 | 8/73 | 9/73 | 10/73 | 11/73 | 12/73 | 1/74 | 2/74 | 3/74 | 4/74 | 5/74 | 6/74 | 7/74 | 8/74 | 9/74 | 10/74 | 11/74 | 12/74 | 1/75 | 2/75 | 3/75 | 4/75 | 5/75 | 6/75 | 7/75 | 8/75 | 9/75 | 10/75 | 11/75 | 12/75 | 1/76 | 2/76 | 3/76 | 4/76 | 5/76 | 6/76 | 7/76 | 8/76 | 9/76 | 10/76 | 11/76 | 12/76 | 1/77 | 2/77 | 3/77 | 4/77 | 5/77 | 6/77 | 7/77 | 8/77 | 9/77 | 10/77 | 11/77 | 12/77 | 1/78 | 2/78 | 3/78 | 4/78 | 5/78 | 6/78 | 7/78 | 8/78 | 9/78 | 10/78 | 11/78 | 12/78 | 1/79 | 2/79 | 3/79 | 4/79 | 5/79 | 6/79 | 7/79 | 8/79 | 9/79 | 10/79 | 11/79 | 12/79 | 1/80 | 2/80 | 3/80 | 4/80 | 5/80 | 6/80 | 7/80 | 8/80 | 9/80 | 10/80 | 11/80 | 12/80 | 1/81 | 2/81 | 3/81 | 4/81 | 5/81 | 6/81 | 7/81 | 8/81 | 9/81 | 10/81 | 11/81 | 12/81 | 1/82 | 2/82 | 3/82 | 4/82 | 5/82 | 6/82 | 7/82 | 8/82 | 9/82 | 10/82 | 11/82 | 12/82 | 1/83 | 2/83 | 3/83 | 4/83 | 5/83 | 6/83 | 7/83 | 8/83 | 9/83 | 10/83 | 11/83 | 12/83 | 1/84 | 2/84 | 3/84 | 4/84 | 5/84 | 6/84 | 7/84 | 8/84 | 9/84 | 10/84 | 11/84 | 12/84 | 1/85 | 2/85 | 3/85 | 4/85 | 5/85 | 6/85 | 7/85 | 8/85 | 9/85 | 10/85 | 11/85 | 12/85 | 1/86 | 2/86 | 3/86 | 4/86 | 5/86 | 6/86 | 7/86 | 8/86 | 9/86 | 10/86 | 11/86 | 12/86 | 1/87 | 2/87 | 3/87 | 4/87 | 5/87 | 6/87 | 7/87 | 8/87 | 9/87 | 10/87 | 11/87 | 12/87 | 1/88 | 2/88 | 3/88 | 4/88 | 5/88 | 6/88 | 7/88 | 8/88 | 9/88 | 10/88 | 11/88 | 12/88 | 1/89 | 2/89 | 3/89 | 4/89 | 5/89 | 6/89 | 7/89 | 8/89 | 9/89 | 10/89 | 11/89 | 12/89 | 1/90 | 2/90 | 3/90 | 4/90 | 5/90 | 6/90 | 7/90 | 8/90 | 9/90 | 10/90 | 11/90 | 12/90 | 1/91 | 2/91 | 3/91 | 4/91 | 5/91 | 6/91 | 7/91 | 8/91 | 9/91 | 10/91 | 11/91 | 12/91 | 1/92 | 2/92 | 3/92 | 4/92 | 5/92 | 6/92 | 7/92 | 8/92 | 9/92 | 10/92 | 11/92 | 12/92 | 1/93 | 2/93 | 3/93 | 4/93 | 5/93 | 6/93 | 7/93 | 8/93 | 9/93 | 10/93 | 11/93 | 12/93 | 1/94 | 2/94 | 3/94 | 4/94 | 5/94 | 6/94 | 7/94 | 8/94 | 9/94 | 10/94 | 11/94 | 12/94 | 1/95 | 2/95 | 3/95 | 4/95 | 5/95 | 6/95 | 7/95 | 8/95 | 9/95 | 10/95 | 11/95 | 12/95 | 1/96 | 2/96 | 3/96 | 4/96 | 5/96 | 6/96 | 7/96 | 8/96 | 9/96 | 10/96 | 11/96 | 12/96 | 1/97 | 2/97 | 3/97 | 4/97 | 5/97 | 6/97 | 7/97 | 8/97 | 9/97 | 10/97 | 11/97 | 12/97 | 1/98 | 2/98 | 3/98 | 4/98 | 5/98 | 6/98 | 7/98 | 8/98 | 9/98 | 10/98 | 11/98 | 12/98 | 1/99 | 2/99 | 3/99 | 4/99 | 5/99 | 6/99 | 7/99 | 8/99 | 9/99 | 10/99 | 11/99 | 12/99 | 1/100 | 2/100 | 3/100 | 4/100 | 5/100 | 6/100 | 7/100 | 8/100 | 9/100 | 10/100 | 11/100 | 12/100 | 1/101 | 2/101 | 3/101 | 4/101 | 5/101 | 6/101 | 7/101 | 8/101 | 9/101 | 10/101 | 11/101 | 12/101 | 1/102 | 2/102 | 3/102 | 4/102 | 5/102 | 6/102 | 7/102 | 8/102 | 9/102 | 10/102 | 11/102 | 12/102 | 1/103 | 2/103 | 3/103 | 4/103 | 5/103 | 6/103 | 7/103 | 8/103 | 9/103 | 10/103 | 11/103 | 12/103 | 1/104 | 2/104 | 3/104 | 4/104 | 5/104 | 6/104 | 7/104 | 8/104 | 9/104 | 10/104 | 11/104 | 12/104 | 1/105 | 2/105 | 3/105 | 4/105 | 5/105 | 6/105 | 7/105 | 8/105 | 9/105 | 10/105 | 11/105 | 12/105 | 1/106 | 2/106 | 3/106 | 4/106 | 5/106 | 6/106 | 7/106 | 8/106 | 9/106 | 10/106 | 11/106 | 12/106 | 1/107 | 2/107 | 3/107 | 4/107 | 5/107 | 6/107 | 7/107 | 8/107 | 9/107 | 10/107 | 11/107 | 12/107 | 1/108 | 2/108 | 3/108 | 4/108 | 5/108 | 6/108 | 7/108 | 8/108 | 9/108 | 10/108 | 11/108 | 12/108 | 1/109 | 2/109 | 3/109 | 4/109 | 5/109 | 6/109 | 7/109 | 8/109 | 9/109 | 10/109 | 11/109 | 12/109 | 1/110 | 2/110 | 3/110 | 4/110 | 5/110 | 6/110 | 7/110 | 8/110 | 9/110 | 10/110 | 11/110 | 12/110 | 1/111 | 2/111 | 3/111 | 4/111 | 5/111 | 6/111 | 7/111 | 8/111 | 9/111 | 10/111 | 11/111 | 12/111 | 1/112 | 2/112 | 3/112 | 4/112 | 5/112 | 6/112 | 7/112 | 8/112 | 9/112 | 10/112 | 11/112 | 12/112 | 1/113 | 2/113 | 3/113 | 4/113 | 5/113 | 6/113 | 7/113 | 8/113 | 9/113 | 10/113 | 11/113 | 12/113 | 1/114 | 2/114 | 3/114 | 4/114 | 5/114 | 6/114 | 7/114 | 8/114 | 9/114 | 10/114 | 11/114 | 12/114 | 1/115 | 2/115 | 3/115 | 4/115 | 5/115 | 6/115 | 7/115 | 8/115 | 9/115 | 10/115 | 11/115 | 12/115 | 1/116 | 2/116 | 3/116 | 4/116 | 5/116 | 6/116 | 7/116 | 8/116 | 9/116 | 10/116 | 11/116 | 12/116 | 1/117 | 2/117 | 3/117 | 4/117 | 5/117 | 6/117 | 7/117 | 8/117 | 9/117 | 10/117 | 11/117 | 12/117 | 1/118 | 2/118 | 3/118 | 4/118 | 5/118 | 6/118 | 7/118 | 8/118 | 9/118 | 10/118 | 11/118 | 12/118 | 1/119 | 2/119 | 3/119 | 4/119 | 5/119 | 6/119 | 7/119 | 8/119 | 9/119 | 10/119 | 11/119 | 12/119 | 1/120 | 2/120 | 3/120 | 4/120 | 5/120 | 6/120 | 7/120 | 8/120 | 9/120 | 10/120 | 11/120 | 12/120 | 1/121 | 2/121 | 3/121 | 4/121 | 5/121 | 6/121 | 7/121 | 8/121 | 9/121 | 10/121 | 11/121 | 12/121 | 1/122 | 2/122 | 3/122 | 4/122 | 5/122 | 6/122 | 7/122 | 8/122 | 9/122 | 10/122 | 11/122 | 12/122 | 1/123 | 2/123 | 3/123 | 4/123 | 5/123 | 6/123 | 7/123 | 8/123 | 9/123 | 10/123 | 11/123 | 12/123 | 1/124 | 2/124 | 3/124 | 4/124 | 5/124 | 6/124 | 7/124 | 8/124 | 9/124 | 10/124 | 11/124 | 12/124 | 1/125 | 2/125 | 3/125 | 4/125 | 5/125 | 6/125 | 7/125 | 8/125 | 9/125 | 10/125 | 11/125 | 12/125 | 1/126 | 2/126 | 3/126 | 4/126 | 5/126 | 6/126 | 7/126 | 8/126 | 9/126 | 10/126 | 11/126 | 12/126 | 1/127 | 2/127 | 3/127 | 4/127 | 5/127 | 6/127 | 7/127 | 8/127 | 9/127 | 10/127 | 11/127 | 12/127 | 1/128 | 2/128 | 3/128 | 4/128 | 5/128 | 6/128 | 7/128 | 8/128 | 9/128 | 10/128 | 11/128 | 12/128 | 1/129 | 2/129 | 3/129 | 4/129 | 5/129 | 6/129 | 7/129 | 8/129 | 9/129 | 10/129 | 11/129 | 12/129 | 1/130 | 2/130 | 3/130 | 4/130 | 5/130 | 6/130 | 7/130 | 8/130 | 9/130 | 10/130 | 11/130 | 12/130 | 1/131 | 2/131 | 3/131 | 4/131 | 5/131 | 6/131 | 7/131 | 8/131 | 9/131 | 10/131 | 11/131 | 12/131 | 1/132 | 2/132 | 3/132 | 4/132 | 5/132 | 6/132 | 7/132 | 8/132 | 9/132 | 10/132 | 11/132 | 12/132 | 1/133 | 2/133 | 3/133 | 4/133 | 5/133 | 6/133 | 7/133 | 8/133 | 9/133 | 10/133 | 11/133 | 12/133 | 1/134 | 2/134 | 3/134 | 4/134 | 5/134 | 6/134 | 7/134 | 8/134 | 9/134 | 10/134 | 11/134 | 12/134 | 1/135 | 2/135 | 3/135 | 4/135 | 5/135 | 6/135 | 7/135 | 8/135 | 9/135 | 10/135 | 11/135 | 12/135 | 1/136 | 2/136 | 3/136 | 4/136 | 5/136 | 6/136 | 7/136 | 8/136 | 9/136 | 10/136 | 11/136 | 12/136 | 1/137 | 2/137 | 3/137 | 4/137 | 5/137 | 6/137 | 7/137 | 8/137 | 9/137 | 10/137 | 11/137 | 12/137 | 1/138 | 2/138 | 3/138 | 4/138 | 5/138 | 6/138 | 7/138 | 8/138 | 9/138 | 10/138 | 11/138 | 12/138 | 1/139 | 2/139 | 3/139 | 4/139 | 5/139 | 6/139 | 7/139 | 8/139 | 9/139 | 10/139 | 11/139 | 12/139 | 1/140 | 2/140 | 3/140 | 4/140 | 5/140 | 6/140 | 7/140 | 8/140 | 9/140 | 10/140 | 11/140 | 12/140 | 1/141 | 2/141 | 3/141 | 4/141 | 5/141 | 6/141 | 7/141 | 8/141 | 9/141 | 10/141 | 11/141 | 12/141 | 1/142 | 2/142 | 3/142 | 4/142 | 5/142 | 6/142 | 7/142 | 8/142 | 9/142 | 10/142 | 11/142 | 12/142 | 1/143 | 2/143 | 3/143 | 4/143 | 5/143 | 6/143 | 7/143 | 8/143 | 9/143 | 10/143 | 11/143 | 12/143 | 1/144 | 2/144 | 3/144 | 4/144 | 5/144 | 6/144 | 7/144 | 8/144 | 9/144 | 10/144 | 11/144 | 12/144 | 1/145 | 2/145 | 3/145 | 4/145 | 5/145 | 6/145 | 7/145 | 8/145 | 9/145 | 10/145 | 11/145 | 12/145 | 1/146 | 2/146 | 3/146 | 4/146 | 5/146 | 6/146 | 7/146 | 8/146 | 9/146 | 10/146 | 11/146 | 12/146 | 1/147 | 2/147 | 3/147 | 4/147 | 5/147 | 6/147 | 7/147 | 8/147 | 9/147 | 10/147 | 11/147 | 12/147 | 1/148 | 2/148 | 3/148 | 4/148 | 5/148 | 6/148 | 7/148 | 8/148 | 9/148 | 10/148 | 11/148 | 12/148 | 1/149 | 2/149 | 3/149 | 4/149 | 5/149 | 6/149 | 7/149 | 8/149 | 9/149 | 10/149 | 11/149 | 12/149 | 1/150 | 2/150 | 3/150 | 4/150 | 5/150 | 6/150 | 7/150 | 8/150 | 9/150 | 10/150 | 11/150 | 12/150 | 1/151 | 2/151 | 3/151 | 4/151 | 5/151 | 6/151 | 7/151 | 8/151 | 9/151 | 10/151 | 11/151 | 12/151 | 1/152 | 2/152 | 3/152 | 4/152 | 5/152 | 6/152 | 7/152 | 8/152 | 9/152 | 10/152 | 11/152 | 12/152 | 1/153 | 2/153 | 3/153 | 4/153 | 5/153 | 6/153 | 7/153 | 8/153 | 9/153 | 10/153 | 11/153 | 12/153 | 1/154 | 2/154 | 3/154 | 4/154 | 5/154 | 6/154 | 7/154 | 8/154 | 9/154 | 10/154 | 11/154 | 12/154 | 1/155 | 2/155 | 3/155 | 4/155 | 5/155 | 6/155 | 7/155 | 8/155 | 9/155 | 10/155 | 11/155 | 12/155 | 1/156 | 2/156 | 3/156 | 4/156 | 5/156 | 6/156 | 7/156 | 8/156 | 9/156 | 10/156 | 11/156 | 12/156 |

We use randomforest classifier and import feature selection function from sklearn package to get the result as the figure shown below.

[illegible]

(In the figure, True represents that the columns are selected as the feature. False represent that the columns are not as significant as the feature)

In general, among the more than 200 factors, only 15 of them are constantly significant to the severity of accidents. They are 'Start_Lat', 'Start_Lng', "Distance(mi)", 'Temperature(F)', 'Wind_Chill(F)', 'Humidity(%)', 'Pressure(in)', 'Visibility(mi)', 'Wind_Speed(mph)', 'Precipitation(in)', 'Crossing', 'Junction', 'Traffic_Signal', 'Side_L', and 'Side_R'. The result is a little amazing that the state and weather are not the factors that lead to the severity of the accident.

Prediction by Multi-Logistic Regression:

Because the Severity is from 1 to 4 and each number represents its accident severity situation, we prefer to deal with categorical data, which need to use logistic regression as the model. (The result of the Prediction V.S. Reality is briefly shown in the figure below)

Out[21]:

	Actually Severity	Predicted Severity	Start_Lat	Start_Lng	Distance(mi)	Temperature(F)	Wind_Chill(F)	Humidity(%)	Pressure(in)	Visibility(mi)	Wind_Speed(mph)	Precipitation(in)
ID												
A-927574	2	2	26.608829	-80.065338	0.000	82.0	82.0	69.0	29.95	10.0	13.0	0.00
A-823441	3	2	37.751369	-79.414055	0.000	42.0	33.0	86.0	25.82	10.0	23.0	0.00
A-625908	2	2	40.126068	-75.331192	0.000	55.0	55.0	40.0	29.30	10.0	7.0	0.00
A-1537866	3	2	37.203964	-76.572983	0.000	34.9	0.0	86.0	30.20	10.0	0.0	0.00
A-63005	2	2	33.187656	-117.279549	0.010	60.1	0.0	75.0	30.02	10.0	9.2	0.00
A-2571089	3	2	41.943270	-87.716210	0.364	79.0	0.0	77.0	29.87	10.0	9.2	0.02
A-2841848	2	2	46.061200	-95.830190	0.000	58.0	58.0	99.0	28.49	10.0	8.0	0.00
A-2395069	3	2	38.216862	-122.137840	0.000	80.1	0.0	18.0	29.83	10.0	15.0	0.00
A-202006	3	2	41.819302	-71.389595	0.010	44.1	37.4	68.0	29.29	10.0	13.8	0.00
A-1367723	3	3	28.175329	-82.391678	1.610	73.4	0.0	83.0	30.16	10.0	0.0	0.00

687496 rows x 204 columns

Command took 0.88 seconds -- by wang.jinc@northeastern.edu at 12/4/2020, 7:03:53 PM on 6110

The accuracy of the predictive model is around 67% to 70%, which means that this model is appropriate for forecasting, but it is better to do more improvement.

```
1 acc = metrics.accuracy_score(test_y, a)
2 print(acc)
```

0.6776126697464422

Command took 2.23 seconds -- by wang.jinc@northeastern.edu at 12/4/2020, 7:03:55 PM on 6110

Comments

The first difficulty that we encountered was when we were importing the data into Databricks. The error specified that it failed to find a data source which mainly because that the data was loaded in pyspark format which required us to convert the data frame from Spark to Pandas. In addition to this, Databricks sometimes shut down without any signal when we were running the multiple commands. We believe that this issue can be improved by having Databricks create a checklist to detect the performance or pop up any alerts which allows us to maximize the work efficiency. What's more, the challenge of data visualization is to present a clear and concise accident map. At first, we put all data into one map. However, the map is very hard to understand due to the overlap area, even though we use different colors and legend to classify these points as figure Y shows. Then, we try to adjust the transparency of the map by setting a lower alpha value. But it

does not work because the overlap part will become new colors and the reader may be confused about what these colors mean. Finally, we create five plots to demonstrate the trend of traffic accidents from 2016 to 2020. Audiences can easily understand the meaning of each map and find the trend of accident cases.

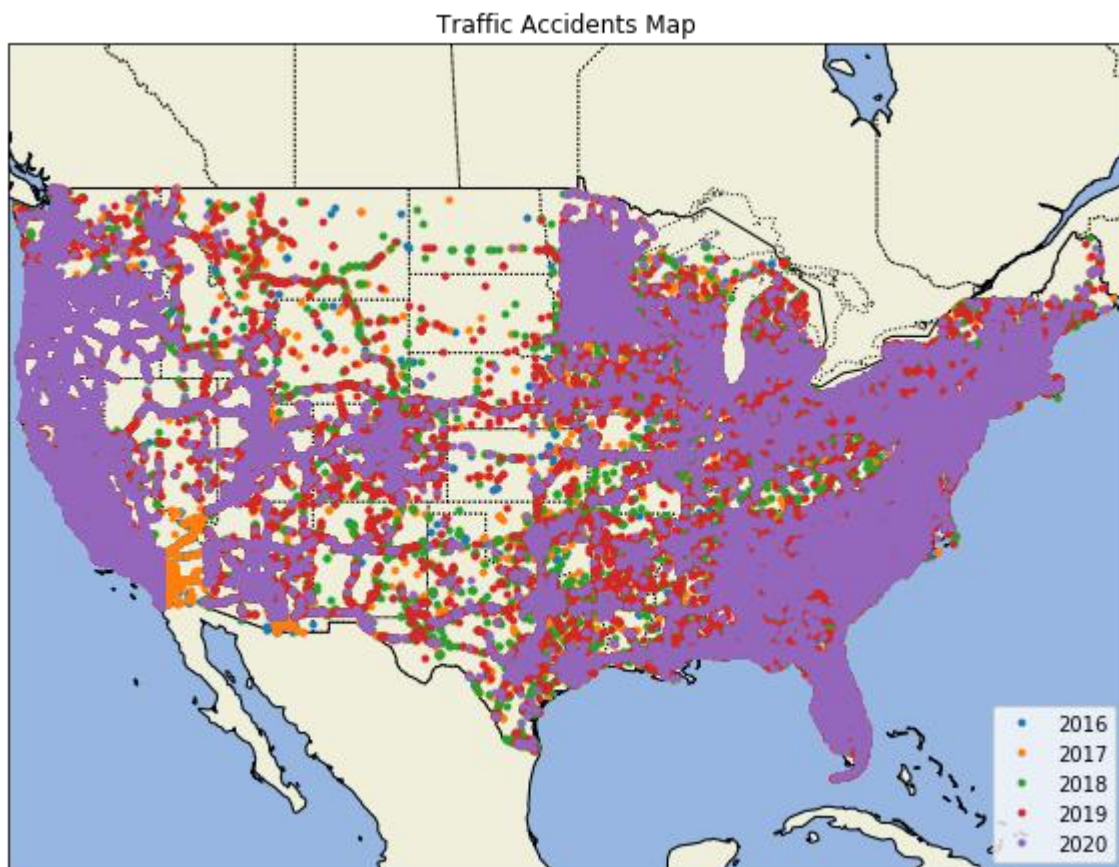


Figure Y: the traffic accidents map (old version)

For the model part, we fell hard to reshape the dataset. Because the dataset contains both numerical and categorical information, we must use 'get_dummies' to transform the categorical data. This will cause the volume of the data set to be too large to do any analysis. Here, the task asks us to determine which information should be included in the model.

Conclusions summary

We have a generally conclusion for the total four parts as the following:

- ◆ There were more cases during 8-12 compared to other months, excluding the data from 2020, guess there are worse weather conditions in the winter.

- ◆ Most accidents happened during the daytime, and there are two peaks on 7-8 and 16-17 when people are on commute between workplace and home.
- ◆ During 23 to 3 o'clock before dawn. Cases numbers are relatively at the bottom level as most people are in sleep
- ◆ Accidents are categorized into four groups (1-4) based on their Severity. 1 being least and 4 being more severe. 67.7 % Accidents fall in the Severity class 2 followed by Severity class 3,4 and 1.
- ◆ Our assumptions generally are that bad weather could lead to more accidents. But here we can see that more accidents occur when the weather is clear. This may be because people drive more carefully when the weather is bad.
- ◆ Blocked, Accident, Rd Accident, due, Lane Blocked are the most frequent words in the description of traffic accident.
- ◆ CA has the most accidents happened. Since CA has highest population in the US states with logical assumption. Houston has the most accidents among all US cities. Being that Texas is the 2nd biggest state in USA, this can lead to more drivers Speeding to reach their destination faster.
- ◆ For the map section, in the time scale, traffic accident cases are increasing by the dot's density. In terms of scope, these cases are spreading from the coastline to the central region.
- ◆ Among all the factors (over than 200 counts), only 15 of them are constant significant in each run of the model. They are 'Start_Lat', 'Start_Lng', 'Distance(mi)', 'Temperature(F)', 'Wind_Chill(F)', 'Humidity (%)', 'Pressure(in)', 'Visibility(mi)', 'Wind_Speed(mph)', 'Precipitation(in)', 'Crossing', 'Junction', 'Traffic_Signal', 'Side_L', and 'Side_R'.
- ◆ The forecasting of severity of accidents is 67.76%, which means the model is appropriate. We think that in the reshaping process, we probably need to make some improvements to reaching the higher accuracy.

References

Nikam, Swapnil Kisan, "ANALYSIS OF US ACCIDENTS AND SOLUTIONS" (2020).

Electronic Theses, Projects, and Dissertations. 979.

<https://scholarworks.lib.csusb.edu/etd/979>

Binu. (November, 2020). Road Accidents in US. Retrieved December 1, 2020

from <https://www.kaggle.com/biphili/road-accidents-in-us>.

William Roosevelt. (July, 2020). US Accidents data cleaning + eda. Retrieved December

1, 2020 from <https://www.kaggle.com/williamrroosevelt/us-accidents-data-cleaning-eda>.