CS590 Big Data and Cloud Computing

# Data Analysis
# Chicago Transportation Condition Awareness

Min Namgung & Shehryar Shahid & Tahmid Alam

### Overview of project

**1**

- The project problem description
- Applications
- Motivation of project problem
- The goal and objectives

### Approach & Implementation

**2**

- Challenges
- Methodology
- General methods and technologies
- Approaches

### Interpretation & Result

**3**

- Interpretation of results
- Results
- Recommendation

### Summary & Conclusion

**4**

- A summary of project
- A description of limitations of current work
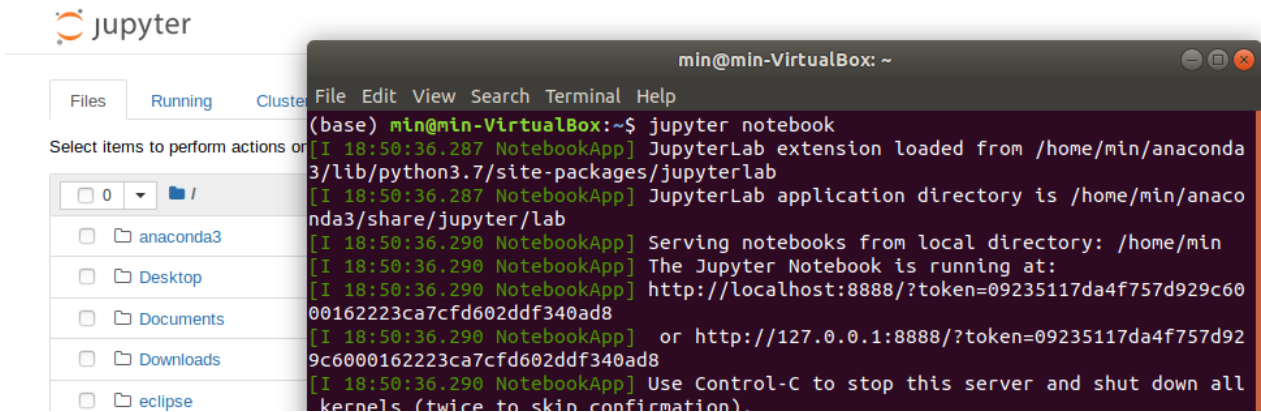- A description of possible future work directions

# Overview of project

- **Chicago car crash data**

  Public open dataset, big data



- **Tools to use**

  Jupyter notebook – python && Pyspark



Tools and dataset we used for the project.

# 1.1
Overview of project

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
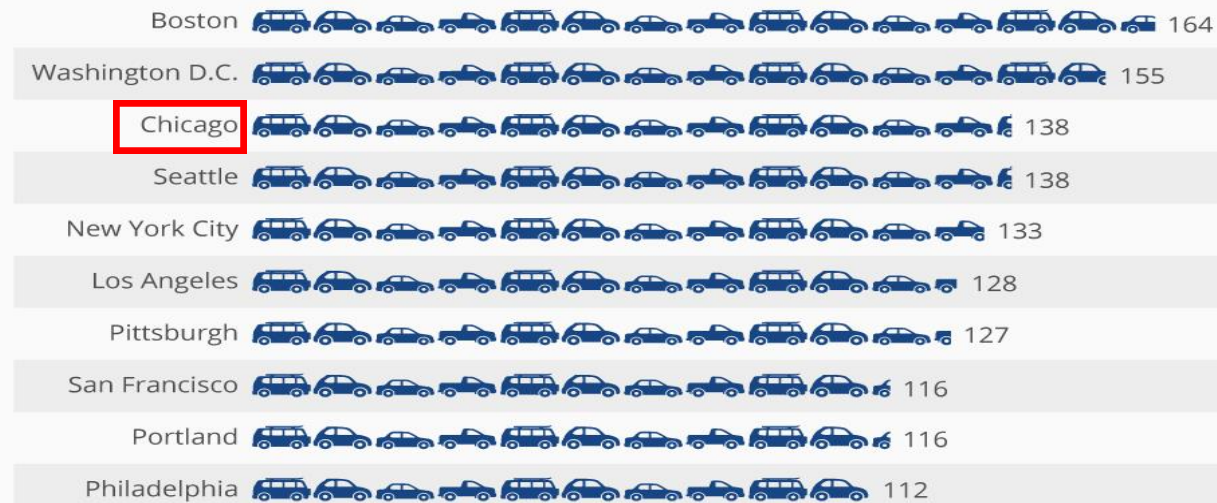condition awareness

## ■ Problem statement

Busy Chicago's traffic leads to many problems such as **injuries** or **car crashes** which cause another collision. The car crash can cause other **traffic jams** and it brings about **more busy traffic** which possibly generates another car accident



The U.S. Cities With The Worst Traffic Problems
Average hours lost to congestion per driver in major U.S. cities in 2018

| City | Hours |
| --- | --- |
| Boston | 164 |
| Washington D.C. | 155 |
| Chicago | 138 |
| Seattle | 138 |
| New York City | 133 |
| Los Angeles | 128 |
| Pittsburgh | 127 |
| San Francisco | 116 |
| Portland | 116 |
| Philadelphia | 112 |

@StatistaCharts    Source: INRIX

statista

# 1.2
Overview of project

■ Application

# 1.3
Overview of project

## ■ Motivation of project



- The third-largest city in U.S.

- A large amount of daily commuters and population

- Easy access to public data

- Car crash data has collected for the past 6 years -> a huge data collection

# 1.3
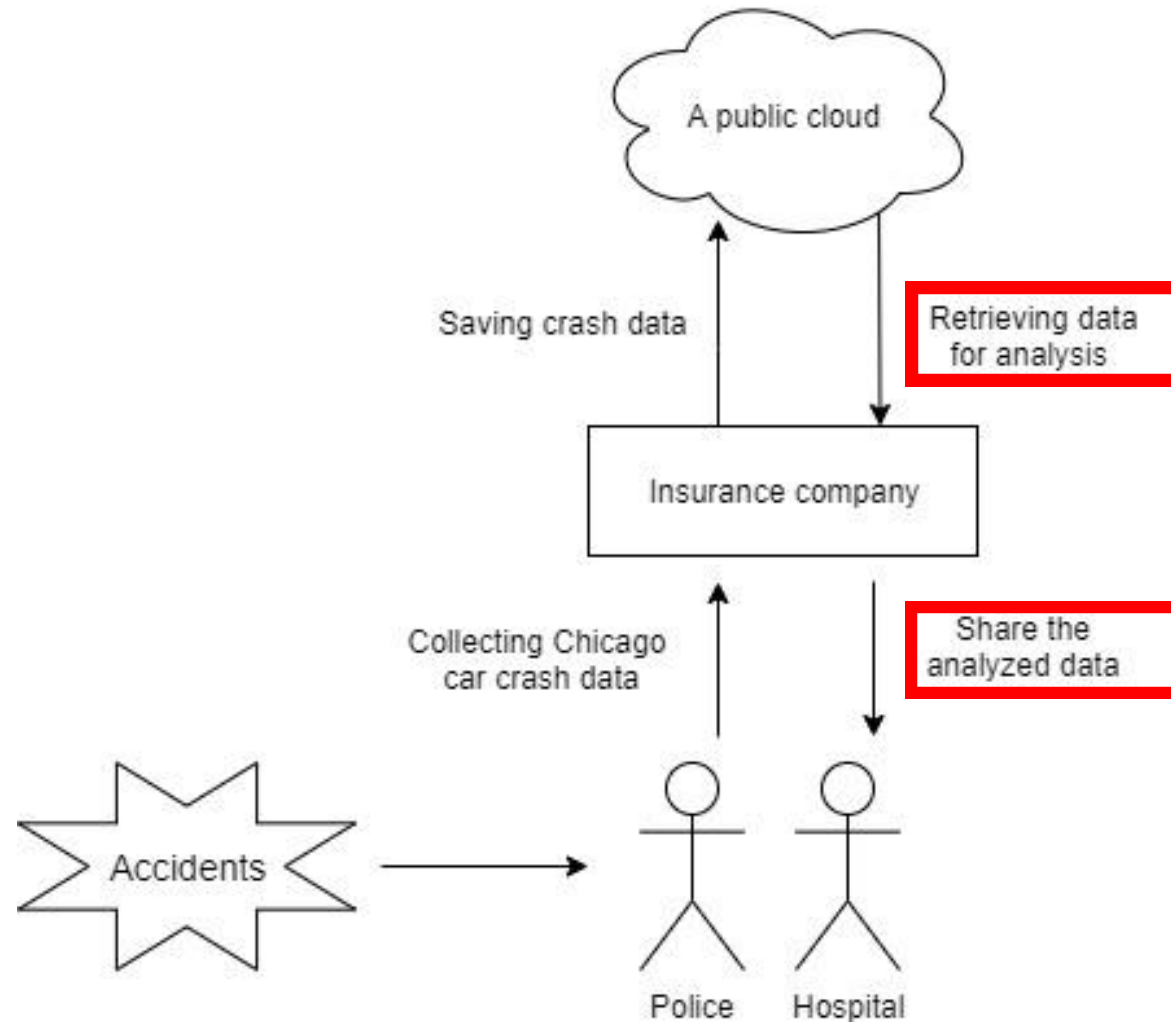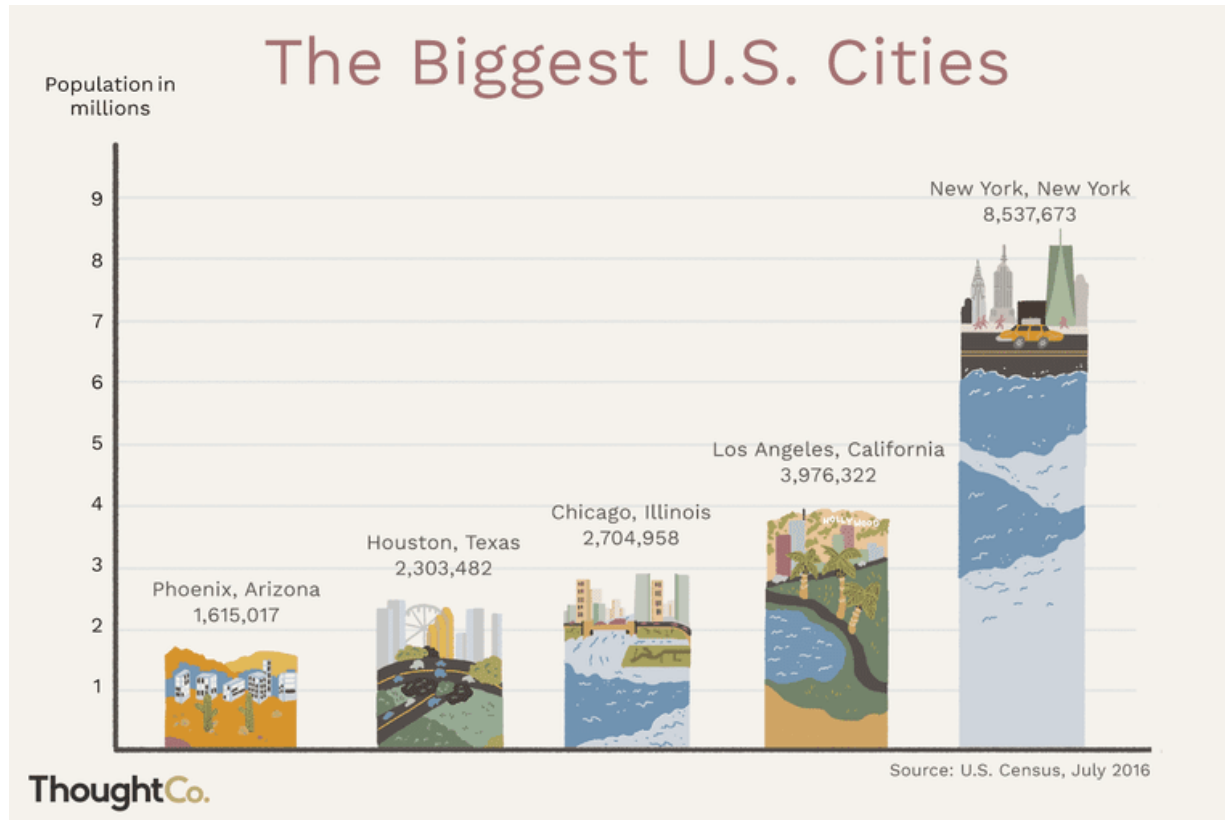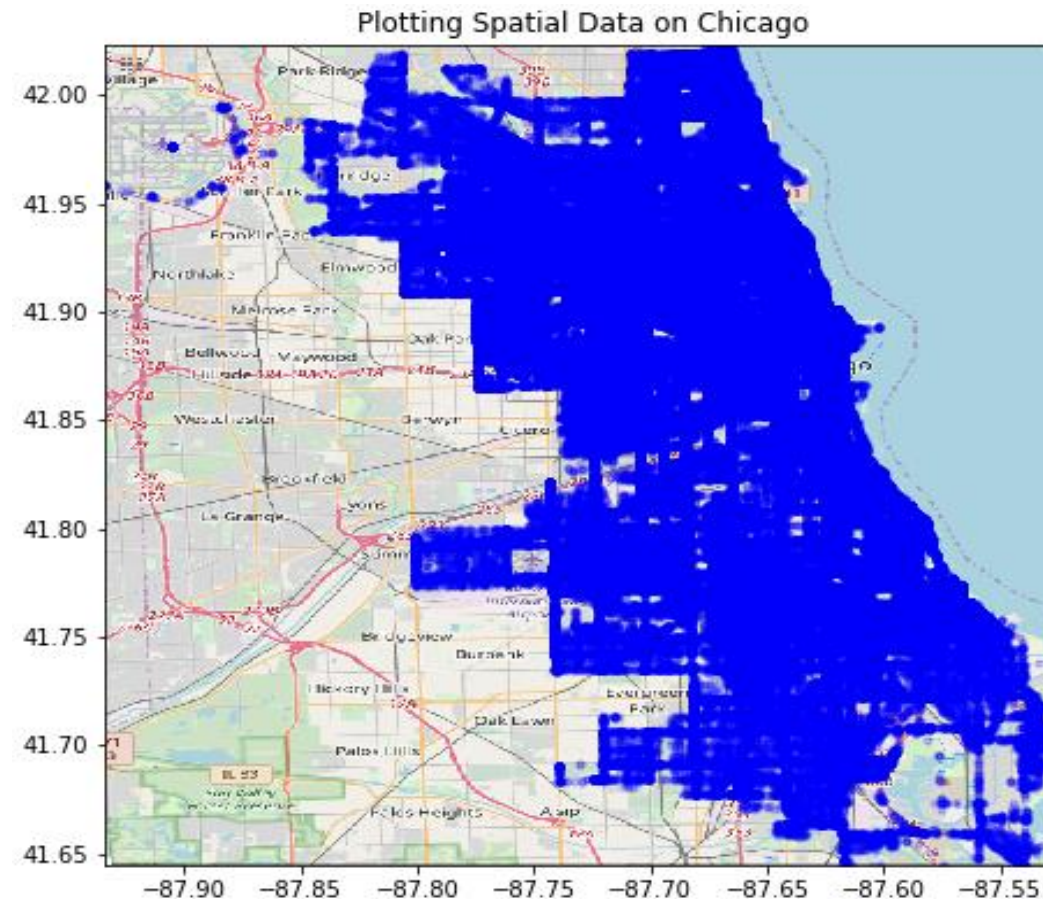
Overview of project

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ▪ Motivation of project



Plotting Spatial Data on Chicago

# 1.4
Overview of project

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Goal of project

- To analyze Chicago's traffic accident data for the past 6 years **to offer valuable insight to both the general public, and the police**

- The car crashes dataset include some **external factors** such as weather conditions, traffic way types, road defection, or the location of accidents.

- With this additional information, **Chicago drivers** will be able to **avoid the most accident-like conditions**.

- Stakeholders such as **police or insurance company** will be able to **understand** how to help **for reducing** these situations, and **how to prevent accidents** before they occur

# 2.1

Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

# ▪ Challenge

- Collecting the data which will help to give us better analysis result.

- Some of the important data which could have made more impact on the analysis are missing.

- Lack of domain knowledge to choose features which will give us more robust result.

- Finding some important library in spark environment to support our approach.

- Not enough compatibility with numpy library.

**2.2**

Approach & Implementation

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ■ Methodology

- Collection of data

- Exploration of collected data

- Data cleaning

- Feature selection

- Elbow method

- Feature indexing

- Scaling of data

- Kmeans Algorithm

# 2.3
Approach & Implementation

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ■ Collection of data

- We collected our data from city of Chicago
- There are plenty of car crash data to work
- There are many outcomes and analysis with our approach
- The data we collected consists of almost 350K+ rows and 48 columns.
- Data range from 2003-2019.

# 2.3
## Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ■ Exploration of Collected data

- **Simple query:**
  - Counting the number of accidents occurred in a month, weather type, road type, time of the day and week

- **Getting the average on numeric data. E.g:** posted speed limit

- **Checking for null values for each variables**

- **Collection of unique values in categorical data**
  - **For weather:** Clear, Rain, Cloudy/Overcast, Snow etc.
  - **Lighting Condition:** Daylight, Darkness, Darkness-Lighted condition etc.
  - **First Crash:** Rear end, Parked, Sideswipe, Turning, etc.

# 2.3

Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Data cleaning

- Very crucial step before any kind of analysis

- Removing the corrupted values

- Removing the rows with null or no values:
  - If missing value is less.

# 2.3
Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Feature selection

- For **Kmeans algorithm** we have to find features

- After doing primary data exploration we selected **17 columns as features**

- Data can help **to create clusters and insight** for our purpose

- Some of the selected features are:
  1. Weather type
  2. Time of the day
  3. Day of the week
  4. Month
  5. Traffic way type etc.

# 2.3

Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Elbow method

- Used to find out optimum number of K values for the data

- Uses wssse value as a metric

- Kmeans is run on from number of clusters from 2 to 20 or more

- Then using the wssse data a graph is generated

- After plotting this graph it looks like an elbow, hence the name

# 2.3
### Approach & Implementation

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

■ Elbow method



**K = 9**

# 2.3
Approach & Implementation

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ▪ Feature Indexing

- Spark Kmeans cannot work on categorical data.
- Have to convert the categorical data into numerical value.
- Using spark mlib library we can achieve that.
- "StringIndexer encodes a string column of labels to a column of label indices. The indices are in [0, numLabels), and four ordering options are supported."

```
id | category
----|----------
0  | a
1  | b
2  | c
3  | a
4  | a
5  | c
```

```
id | category | categoryIndex
----|----------|---------------
0  | a        | 0.0
1  | b        | 2.0
2  | c        | 1.0
3  | a        | 0.0
4  | a        | 0.0
5  | c        | 1.0
```

Source:https://spark.apache.org/docs/latest/ml-features

## 2.3
Approach & Implementation

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Feature Scaling

- Necessity of feature scaling

- Types of scaling (In spark mLib):
  - **Standard deviation**
  - **Mean**

- We used standard deviation method for our purpose

# ▪ K-means

- Unsupervised Machine Learning Algoritm

- Works with Eucledian distance

- Creates cluster based on closest centroid

- K = 9

**3.1**
Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ▪ Cluster Identification

We divided our clusters into 5 regions:
- Low Region
- Low-Mid Region
- Mid Region
- High Region
- Very High Region

Cluster 1: High Region
Cluster 2 : Low-Mid Region
Cluster 3: Mid Region
Cluster 4 and 8: Low Region
Cluster 6: Very High Region

# 3.1
Insights and Recommendations

## ■ Cluster 1: High Accident Frequency Region

Accident Count: 84,407

Cluster 1 accounts for 23% of the data set

Key Features:
- Weather Condition: Clear (89%)
- Lightening Condition: Daylight (71%) and Darkness-Lighted Road(20%)
- Crash Type: Rear End(36%), Sideswipe(17%) and Turning(24%)
- Traffic Way Type: Not Divided(50%) and Divided–without raised median(24%)
- Traffic Control Device: Traffic Signal(73%) and Stop sign/Flasher(19%)
- Device Condition: Functioning Properly(87%) and Functioning Improperly(1.5%)

# 3.1
Insights and Recommendations

■ Cluster 1: High Accident Frequency Region

Insights:
• Most Accidents happen on Non-Divided Roads
• Under Clear Weather and properly functional Traffic Control devices

Recommendations:
• Speed limit control in Non-Divided Roads
• Caution Signs

## 3.2
Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

# ■ Cluster 2: Low-Mid Accident Frequency Region

Accident Count: 32,625

Cluster 2 accounts for only 9% of the data set

Key Features:
- Weather Condition: Snow(35%), Cloudy/Overcast(29%)
- Lightening Condition: Daylight (54%) and Darkness-Lighted Road(20%)
- Crash Type: Parked(33%) and Rear End(23%)
- Traffic Way Type: Not Divided(44%), Divided–w/o raised median(17%) and Oneway(18%)
- Traffic Control Device: No Control(66%) and Traffic Signal(23%) and Stop sign/Flasher(10%)
- Device Condition: No Control(66%) and Functioning Properly(30%)

# 3.2

Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ■ Cluster 2: Low-Mid Accident Frequency Region

Insights:
- Most Accidents Occurred in Winters
- Non-Divided and One-Way Roads
- No Control Device Present

Recommendations:
- Have traffic control devices/stop signs on these roads



Weather Conditions

■ Clear ■ Rain ■ Snow ■ Cloudy/Overcast



Cluster 2 Annual Accident Frequency

# 3.3

Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ■ Cluster 3: Mid Accident Frequency Region

Accident Count: 63,419

Cluster 3 accounts for 18% of the data set

Key Features:
- Weather Condition: Clear(90%) and Rain(9%)
- Lightening Condition: Daylight (67%) and Darkness-Lighted Road(21%)
- Crash Type: Parked(45%), Rear End(15%) and Side Swipe(15%)
- Traffic Way Type: Not Divided(46%), Oneway(25%) and Parking(4%)
- Traffic Control Device: No Control(99.4%)
- Device Condition: No Control(99.4%)

# 3.3
## Insights and Recommendations

## ■ Cluster 3: Mid Accident Frequency Region

Insights:
- Most Accidents happen on Non-Divided or One-Way Roads
- Under Clear Weather but with no Traffic Control device
- Major accident type is Parked Motor Vehicle but that does not happen in Parking area

Recommendations:
- Traffic Control Devices
- Caution Signs for Non-Divided and One-Way roads that have roadside parking

# 3.4
Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

## ■ Cluster 4: Low Accident Frequency Region

Insights:
- Most Accidents happen on Non-Divided
- Under Clear Weather, with and without Traffic Control device

Recommendations:
- More Traffic Control Devices on Non-Divided roads to prevent these kinds of accidents

# 3.5
Insights and Recommendations

## ■ Cluster 6: Very High Accident Frequency Region

Accident Count: 99,012 (Cluster with Maximum Accidents)

Cluster 6 accounts for 28% of the data set

Key Features:
- Weather Condition: Clear(91%) and Rain(9%)
- Lightening Condition: Daylight (74%) and Darkness-Lighted Road(16%)
- Crash Type: Parked(26%), Side Swipe(21%) and Rear End(23%)
- Traffic Way Type: Not Divided(43%), Divided–w/o raised median(17%) and Parking(18%)
- Traffic Control Device: No Control(97%)
- Device Condition: No Control(98%)

3.5
Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

■ Cluster 6: Very High Accident Frequency Region

Insights:
- Most Parking Lot Accidents in this cluster
- Under Clear Weather and without Traffic Control device
- 13,586 out of 17,595 (77%) accidents happened in Parking lot during Daytime

Recommendations:
- People need to be extra cautious in Parking Lots
- Sign and signal inside Parking Lots



Traffic Way Type = Parking Lot



Traffic Control Device
■ No Control  ■ Traffic Signal  ■ Stop Sign/Flasher

## 3.4
Insights and Recommendations

CS590 Big data
Analysis
Chicago Transportation
condition awareness

■ Cluster 4: Low Accident Frequency Region

Accident Count: 20,673

Cluster 4 accounts for 6% of the data set

Key Features:
- Weather Condition: Clear(83%) and Rain(14%)
- Lightening Condition: Daylight (48%) and Darkness-Lighted Road(37%)
- Crash Type: Parked(17.5%), Angle(17.5%) and Rear End(14%)
- Traffic Way Type: Not Divided(42%) and Divided–w/o raised median(20%)
- Traffic Control Device: No Control(51%), Traffic Signal(33%) and Stop Sign(13%)
- Device Condition: No Control(53%) and Functioning Properly(44%)

**3.6**

Insights and Recommendations

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

- ## Cluster 8: Low Accident Frequency Region

Accident Count: 25,340

Cluster 8 accounts for just 7% of the data set

Key Features:
- Weather Condition: Clear(86%) and Rain(11%)
- Lightening Condition: Daylight (71%) and Darkness-Lighted Road(20%)
- Crash Type: Turning(15%), Side Swipe(15%) and Rear End(33%)
- Traffic Way Type: Not Divided(50%), Divided–w/o raised median(22%) and One-Way(14%)
- Traffic Control Device: Traffic Signal(45%) and Stop Sign/Flasher(35%)
- Device Condition: Functioning Properly(74%)

# 3.6
Insights and Recommendations

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ■ Cluster 8: Low Accident Frequency Region

Insights:
- Control Devices were functioning properly
- Turning, Side Swipe and Rear-end are major accident types
- Clearly shows human factor as major cause of accident

Recommendations:
- People need to be extra cautious when taking a turn or overtaking

# 4.1
Summary and Conclusion

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

## ■ Summary

- In the **highest** accident frequency region: People need to be extra cautious in **Parking Lots**

- In **high accident** frequency region: **Speed limit control in Non-Divided Roads**, Caution Signs

- In **low-mid accident** frequency region: Have **traffic control** devices/**stop signs**

- In **mid accident** frequency region: **Caution Signs for Non-Divided** and **One-Way roads** that have roadside parking

- In **low accident** frequency region: **More Traffic Control Devices** on Non-Divided roads

## ▪ Limitations & Difficulties

- Big data -> take lots of time for data sorting

- Hardware Limitation: running software with a large data in virtual box would often crash

- Jupyter kernal frequently crashed while running the program

- Missing data manipulation library

- Learning curve

**4.3**
Summary and Conclusion

C S 5 9 0   B i g   d a t a
Analysis
Chicago Transportation
condition awareness

### ▪ Future work

- With this analyzed data, the project can be improved as **a prediction of car crash**

- The project can add advanced **machine Learning algorithms**

- More functionalities

- Suggest the outcomes or summary to Stakeholders (Police department, hospital and insurance company etc)

CS590 Big Data and Cloud Computing

# Q & A

Min Namgung & Shehryar Shahid & Tahmid Alam