# Train Wreck Analysis

*Long Nguyen*
*SJSU ID: 010806471*

## 1. Project abstract

Use train wreck datasets http://www.trainwreckdb.com/ with spark service in bluemix to figure out what are the 10 most dangerous places for accidents and why.
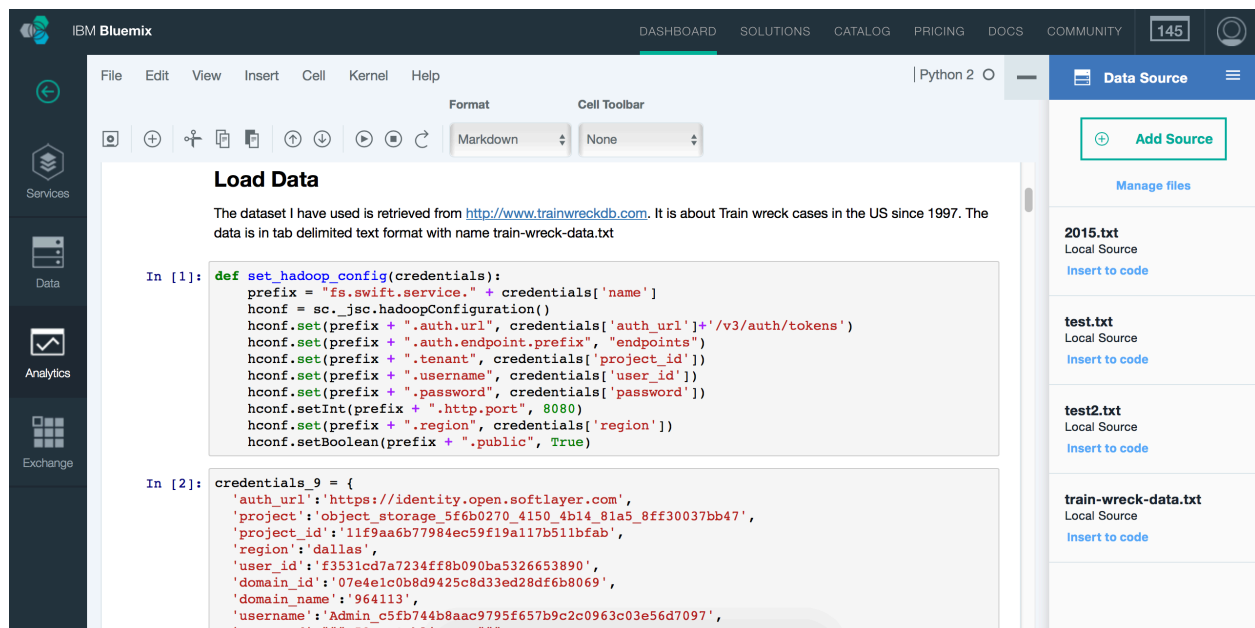
## 2. Project Scope

In this project, I will do analysis to point out the 10 most dangerous Cities and State for train accidents.

The dataset also contains information about the Street where accidents occurred. However, after do experimental analysis on the Street, the number of wreck cases mostly lay on the the streets with common name like Private or Private Rd (see Github link for more details). This result is not able to come to conclusion.

## 3. Github Link

https://github.com/longnguyen0708/TrainWreckAnalysis

## 4. Load Dataset

## 5. Point out 10 most dangerous Cities for train accidents

Perform Map and Reduce on data

Visualize the 10 Cities with the highest number of train accidents since 1997



**From the Graph above, we can conclude that top 10 dangerous cities for train accidents are**

```
City "HOUSTON, TEXAS" had 453 train wreck cases
City "CHICAGO, ILLINOIS" had 357 train wreck cases
City "MEMPHIS, TENNESSEE" had 187 train wreck cases
City "SAN ANTONIO, TEXAS" had 175 train wreck cases
City "LOUISVILLE, KENTUCKY" had 174 train wreck cases
City "BATON ROUGE, LOUISIANA" had 167 train wreck cases
City "PHOENIX, ARIZONA" had 159 train wreck cases
City "GARY, INDIANA" had 156 train wreck cases
City "JACKSONVILLE, FLORIDA" had 151 train wreck cases
City "LOS ANGELES, CALIFORNIA" had 139 train wreck cases
```

## 6. Point out 10 most dangerous States for train accidents

Perform Map and Reduce on data



Visualize the 10 States with the highest number of train accidents since 1997

**From the Graph above, we can conclude that top 10 dangerous cities for train accidents are**

```
State TEXAS had 5511 train wreck cases
State ILLINOIS had 3070 train wreck cases
State CALIFORNIA had 2880 train wreck cases
State INDIANA had 2853 train wreck cases
State LOUISIANA had 2510 train wreck cases
State GEORGIA had 2283 train wreck cases
State OHIO had 2187 train wreck cases
State ALABAMA had 1946 train wreck cases
State FLORIDA had 1625 train wreck cases
State MICHIGAN had 1600 train wreck cases
```