



SPARK INDIVIDUAL ASSIGNMENT

Final Report

ABSTRACT

In this final report, I discuss the background, goal of analysis, the analysis, and conclusion of a dataset regarding roughly 3.5 million car accidents that took place in the United States between February 2016 and June 2020.

Nelson, Ckalib

Spark – Professor Raul Marin Perez

Table of Contents

<i>Background/Scenario Description.....</i>	<i>2</i>
Background Information	2
Data Description	3
<i>CRISP-DM</i>	<i>4</i>
<i>Goal of the Analysis.....</i>	<i>5</i>
<i>Completed Steps.....</i>	<i>6</i>
<i>Deep Dive Analysis</i>	<i>7</i>
<i>Future Analysis</i>	<i>8</i>
<i>Conclusions/Insights.....</i>	<i>9</i>
<i>Video Presentation</i>	<i>9</i>

Background/Scenario Description

Background Information

During the early 2010s, many Americans benefited from a better economy and significantly cheaper oil prices. Both of these trends greatly incentivize individuals to travel via automobile. Additionally, Americans tend to drive much more than others around the world, especially on rural roads with insufficient lightening and barriers.

Unfortunately, the United States had one of the worst road-safety records in the world during this time period. Per *The Economist*, America's road deaths rose by the highest annual percentage increase in 50 years in 2013.¹ Furthermore, the rate of death per 100,000 people per year was roughly 10.9, which is twice as high as Belgium, the country with the second worst record. Surely, these horrifying incidents prompted many public institutions to collect data on accidents in an effort to prevent deaths on the road.



¹ The Economist: America's road-safety record is the worst in the rich world. <https://www.economist.com/graphic-detail/2016/09/05/americas-road-safety-record-is-the-worst-in-the-rich-world>

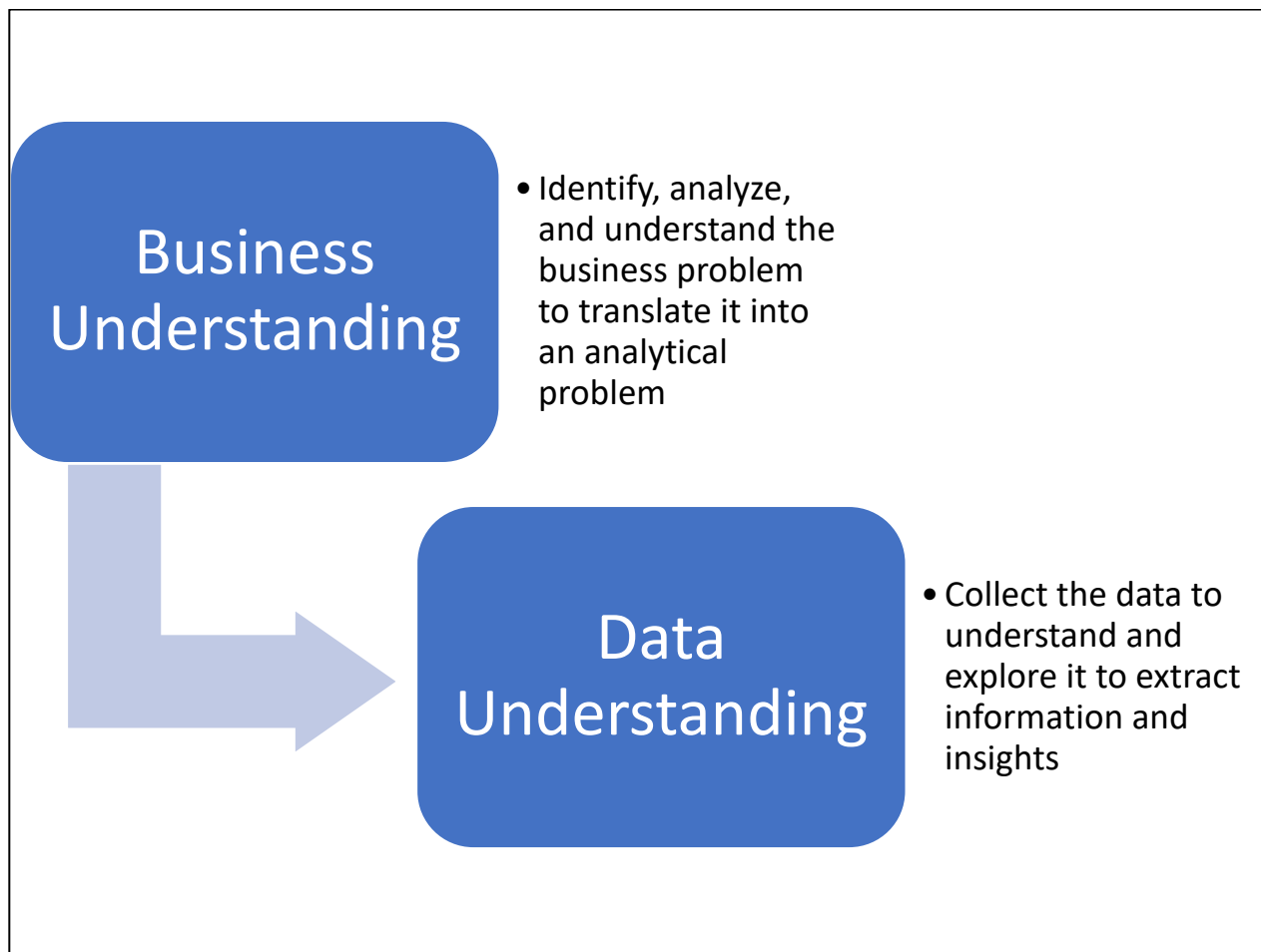
To complete the individual assignment, I decided to use the 3.5 million records of US Accidents per [Kaggle](#). The data was collected in real-time via APIs from various local and government institutions and represents traffic accidents in 49 states from February 2016 to June 2020. Per the notebook, here is an illustration of the data on accidents in the United States:

[illegible]

Ease of Use: *Accessible APIs*

CRISP-DM

To complete this project in an organized and methodological manner, I utilized CRISP-DM. CRISP-DM is a widely used and accepted methodology for data related projects that was developed by a consortium of over 200 organizations in 1996. Although CRISP-DM is a six-phase methodology, my project will not require the last four phases given the nature of the project. Thus, I only completed two phases: Business Understanding and Data Understanding:



Once I completed the Business Understanding and Data Understanding phase, I built tables to answer questions about the dataset, which will be discussed in the following section.

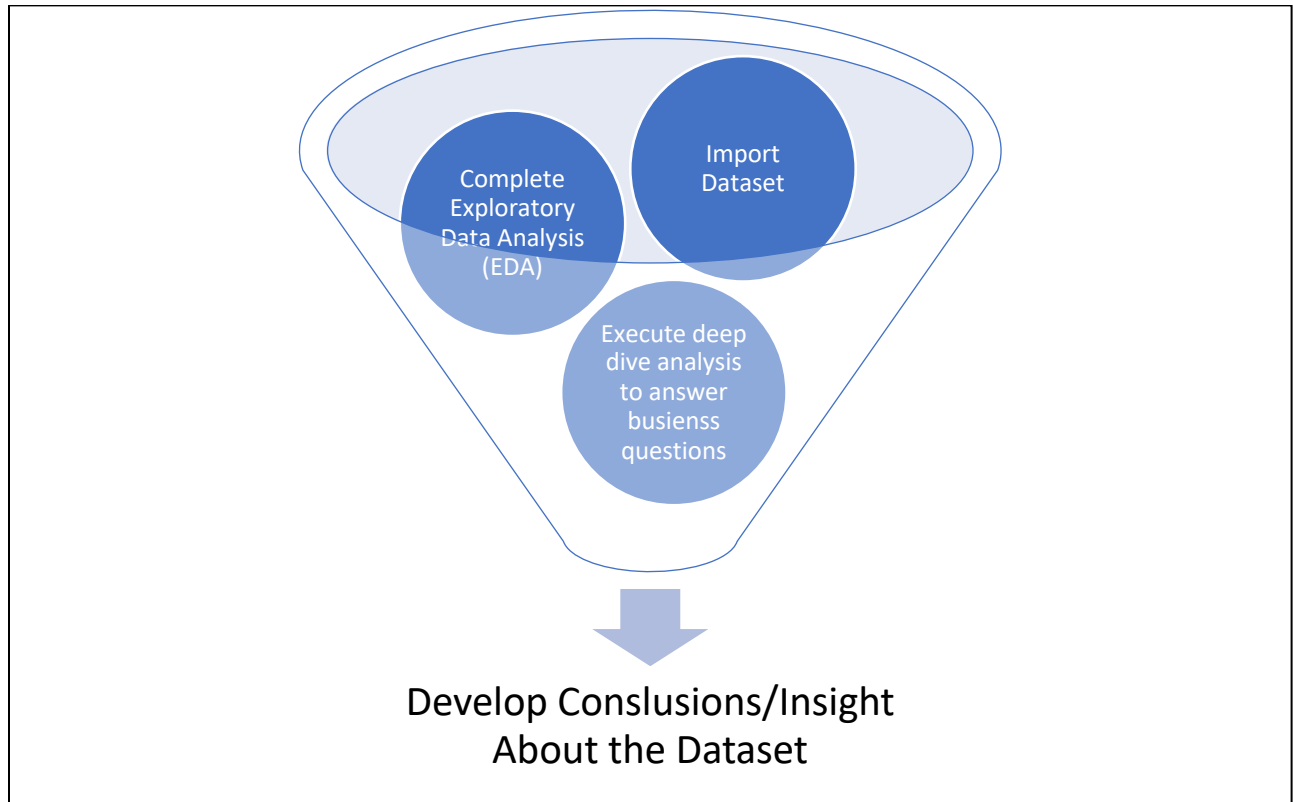
Goal of the Analysis

The goal of my analysis is to answer the following ten questions about my data:

Business Question 1	• In which states do most accidents occur?
Business Question 2	• Which one of the aforementioned states have the most severe accidents?
Business Question 3	• What are the top descriptions of the accidents?
Business Question 4	• Under what conditions do most of these accidents occur?
Business Question 5	• What is the landscape of accidents in the area in which I live?
Business Question 6	• What is the average wind speed for each of the severity levels?
Business Question 7	• What is the average precipitation of each of the severity levels?
Business Question 8	• What is the average temperature for each of the severity levels?
Business Question 9	• How severe are accidents at crossings?
Business Question 10	• How severe are accidents at railways?

Completed Steps

After determining which business questions, I completed the following steps to answer those questions:



Deep Dive Analysis

Leveraging the lectures and notes from the professor, I completed a deep dive analysis of the dataset to answer the aforementioned questions. As a result of my analysis, here are the answers to those questions:

In which states do most accidents occur?

- It seems that accidents occur the most in California (CA), Texas (TX), and Florida (FL).

Which one of the aforementioned states have the most severe accidents?

- Between CA, TX, and FL, it seems as if FL has the most severe accidents.

What are the top descriptions of the accidents?

- It seems as if the descriptions of the accidents include information regarding congestion on interstates, which may align with the intuition we have about safety of interstates. Maybe public officials in these areas should reconsider these dangerous areas!

Under what conditions do most of these accidents occur?

- Contrary to intuition, it seems that most accidents occur while the weather conditions are clear, fair, or mostly cloudy. This most likely has to do with the fact that climates differ depending on which state you live in. For example, it may tend to rain more in Florida than in California.

What is the landscape of accidents in the area in which I live?

- I have lived in Charlottesville, VA since August of 2015. Although I am not familiar with the first accident listed above, I am very familiar with the roads/interstates mentioned in the description of these accidents. Fortunately, there haven't been very severe accidents in the area nor have I been involved in them!

What is the average wind speed for each of the severity levels?

- It does seem as if wind and severity are positively correlated albeit marginally.

What is the average precipitation for each of the severity levels?

- It does seem as if rain and severity are positively correlated albeit marginally.

What is the average temperature for each of the severity levels?

- It does seem as if temp and severity are negatively correlated. The colder it is, or the lower the temperature, the more likely there is to be an accident.

How severe are accidents at crossings?

- It seems as if accidents at crossings are quite severe. Most accidents at crossings are either a 3 or 4 category.

How severe are accidents at railways?

- It seems as if accidents at railways are not as severe as they could be with most severe accidents at railways being a category 2. Relative to crossings, railways are a bit safer.

Future Analysis

Although my analysis answers ten interesting business questions about accidents in the United States, there is so much more analysis that can be done! Here are five additional questions I would like to answer once I develop my Spark skills to a greater level of expertise:

Can explanatory variables predict the severity level of accidents?

Do accidents occur during a specific time of day, such as morning, afternoon, evening, or night?

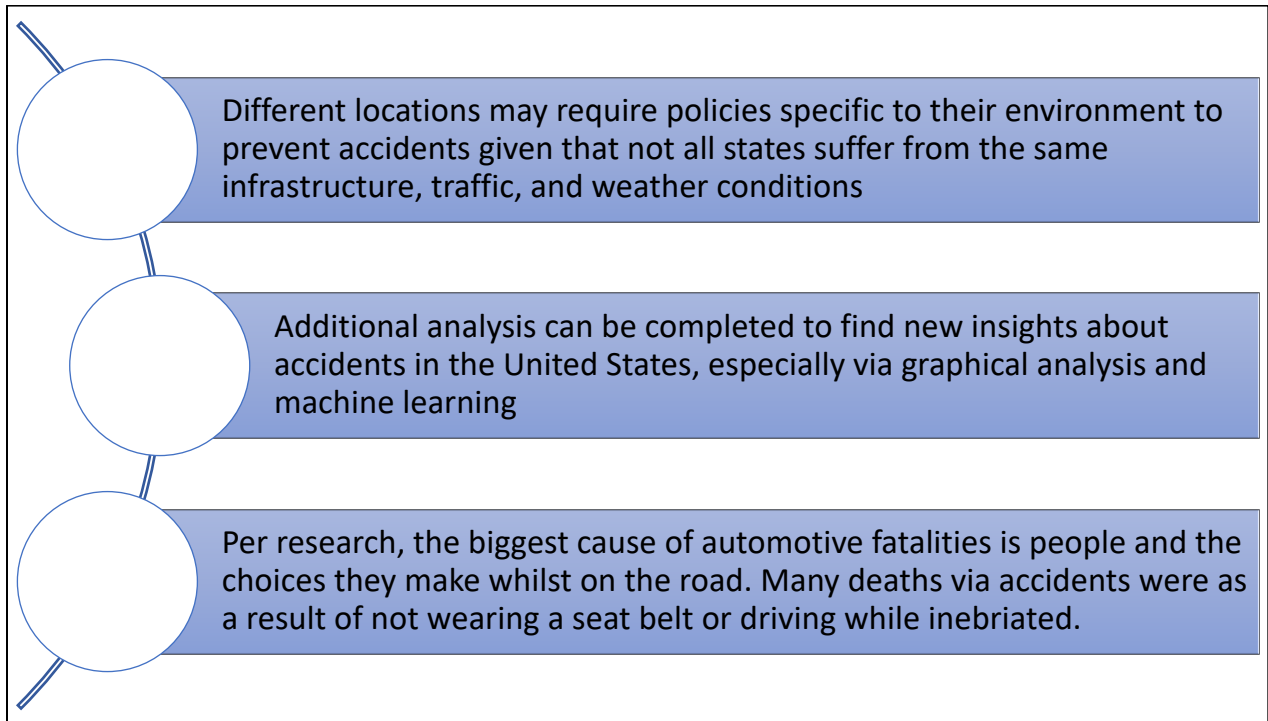
Are accidents with a greater distance more severe than accidents with a smaller distance?

Can additional variables be added to the dataset for analysis, such as approximate number of cars on the road at the time of the accident?

Can analysis of this dataset support infrastructure laws, such as construction for safer roads and regulations?

Conclusions/Insights

Although I have answered several business questions and could answer additional questions in the future, there will still be several elements about accidents in the United States that the data won't be able to answer. Nonetheless, there needs to be policies implemented in the United States intended to reduce the accident and death rate on the road. Thus, what else might we be able to take away from this analysis and what additional information can we learn about accidents in the United States?



Video Presentation

Here is the link to the video presentation that discusses this project in detail: [Spark Individual Assignment \(C. Nelson\)](#)