

OSHA Data Clean-Up – MA 415 Midterm

William Dean

Contents

Objective	1
Problems with the OSHA Data – OSHA Clean-Up	1
Accident Table	1
OSHA Table	3
Complete Accidents Table	4
Exploring the Data	5
Bar Graphs	5
Time Series – Time Occurances	6
Histograms	7
Map Graphics	8
Final Remarks	9

Objective

The OSHA data presents Health & Safety Regulations and all workplace injuries and accidents. This data is ready and available, however it is very messy and is not able to be interpreted easily. The goal of this project is to clean this OSHA dataframes given with the intent of looking for the most dangerous place to work in the state of Massachusetts and provide insite of the data through Data Exploration.

Problems with the OSHA Data – OSHA Clean-Up

Accident Table

The Accident Table from OSHA is a record of the accidents and the situation behind each accident. There is a lot of information as can be seen in the display of it below but it needs to be cleaned up to before exploring it and anything more.

ACTIVITY	SYNCD	STATE	TIME	RELINSP	SEX	DEGREE	NATURE	BODY	SOURCE	EVENT	ENVIRON	INMANS	SKHAZSU	OCC_CODE	OFF
10096592MA	NA	10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	
10096592MA	NA	10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	
10096592MA	NA	10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	

Unnecessary Atributes

At first glance, the dataframe that represents the accidents in MA is very messy and has many attributes that would not be helpful to looking into the most dangerous places to work in Massachusetts, so these were the first to be eliminated. A Person's name or the RELINSP (An inspections key) doesn't give much use to look at dangerous locations and are taken out to simplify the data.

ACTIVITY	SEX	DEGREE	NATURE	BODYPART	SOURCE	EVENT	ENVIRON	NUM	TASK	HAZ	SUB	OCC_CODE	AGE
10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	
10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	
10096592	NA	3	21	04	16	08	07	06	2	NA	000	0	

Getting Rid of Incorrect Information (NA Values)

Although this is better, there is much information that is stored as placeholder values when they should be NA. For example, an OCC_CODE of “000” does not refer to an Occupation and should be NA. Taking out all instances of bad values cleans the data more.

ACTIVITY	SEX	DEGREE	NATURE	BODYPART	SOURCE	EVENT	ENVIRON	NUM	TASK	HAZ	SUB	OCC_CODE	AGE
10096592	NA	3	21	04	16	08	07	06	2	NA	NA	NA	
10096592	NA	3	21	04	16	08	07	06	2	NA	NA	NA	
10096592	NA	3	21	04	16	08	07	06	2	NA	NA	NA	

Give Appropriate Labels

There is a lot of information in the table that is not stored in this table but in other tables that would be better suited to have all in one location. Bringing this information to one table makes it easier to understand what the data means and will make exploration and analysis easier. The resulting table is much cleaner and can be seen below.

ACTIVITY	SEX	DEGREE	TASK	AGE	OCCUPATION	HAZ	SUB	FACT	BODYPART	SOURCE	ENVIRON	EVENT
10096592	NA	Non-Hospitalized	Irregular	NA	NA	NA	EQUIP.	IN-APPROPR FOR OPERATION	BODYSYSTEM	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT
10096592	NA	Non-Hospitalized	Irregular	NA	NA	NA	EQUIP.	IN-APPROPR FOR OPERATION	BODYSYSTEM	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT
10096592	NA	Non-Hospitalized	Irregular	NA	NA	NA	EQUIP.	IN-APPROPR FOR OPERATION	BODYSYSTEM	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT

Final Version

This provides much a much cleaner version of the original accident table but there exists rows which have the same values everywhere that are redundant. Taking them out gives a final version of the table we started with.

ACTIVITY	SEX	DEGREE	TASK	AGE	OCCUPATION	HAZ	SUB	FACT	BODYPART	SOURCE	ENVIRON	EVENT
10096592	NA	Non-Hospitalized	Irregular	NA	NA	NA	EQUIP.	INAP-PROPR FOR OPERATION	BODYSYSTEM	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT	SMOKE/VAPOR/INSTANT

ACTIVITY	SYMBOL	REASON	TASK	OCCUPATION	HAZARD	SUBFACT	BODY PART	SOURCE	ENVIRON	EVENT
3055487	145	Fatal	Regular	47	ELECTRICIAN	NA	INSUF/LACK/WRK CLTHG/EQUIP	PROBLEM	OTHER APPARAT/WIRING	FALL(FROM ELEVATION)
3068112	28	Fatal	Regular	39	ELECTRICAL POWER IN- STALLERS AND REPAIRERS	NA	OTHER	NECK	OTHER MOV- ING/FALLING OBJ AC	OVERHEAD CAUGHT IN OR BETWEEN

OSHA Table

The OSHA Table contains the information for all of the Inspections that have happened which is able to give us more information about the accidents in Massachusetts. Our goal is to combine the important information from the OSHA table to the Accident table we have already made to how a condensed table with all the accident information. However, we run into a lot of the same problems as the Accident table so the OSHA table has to be cleaned and preped before joining. The given OSHA table is too large and messy and can be seen below:

CONT	TYPE	OSHA	STOR	REV	PRE	ACT	NO	REN	CS	DO	JOB	TITLE	FIRST	LAST	NAME	SITE	ADD
NA	H	198402	NA	NA	0	1023677	0	111100	NA	C	000000	00	DUBE	RT 1 MAIN	ST		
NA	M	199105	NA	NA	0	1033936	3	111100	NA	I	000000	00	KNOWLTON	NEW			
NA	H	198806	NA	NA	0	1875003	4	111140	NA	NA	000000	00	RENTAL	NA			
													& FROST				

Cleaning the Table

After taking only the information we want, getting rid of incorrect values which should be NAs, and relabeling the keys stored as integers with more descriptive and accurate labels, we come of with a version of the OSHA table that will provide important information about what happened at each inspection. A final version of the OSHA table before combining with the Accident table can be seen below and is easier to understand the information compared to how it was originally presented.

ACTIVITY	SYMBOL	NAME	HEAD	OWNER	OPEN DATE	CLOSE DATE	INDUSTRY	city	latitude	longitude	Decade	
1023677	0	DUBE	RT 1 MAIN ST	NA	1983- 12-15	NA	NA	PLASTERING, DRYWALL, AND INSULATION	NA	NA	NA	80s
1033936	3	KNOWLTON	NEW MA- CHINE CO.	Private	1990- 07-17	1990- 07-20	NA	INDUSTRIAL MACHINERY, NEC	Salem	42.51685	- 70.8985	90s

ACTIVITY	SYMBOL	NA	HEADD	OWNER	OPER	DATE	INDUSTRY	city	latitude	longitude	Decade
1875003	RENTAL	NA	NA	1979-05-14	1979-05-14	NA	SHEET METAL WORK	NA	NA	NA	70s
	&										
	FROST										

Complete Accidents Table

We have two tables with information about the accidents in MA, but now we want to store it all in one table. Using a left join, we can keep all the information about the accidents from the accident table, but we can also incorporate any additional information that the OSHA table will give about each accident. Joining the tables results in one single table that has all the information from the data frames which would be helpful in exploring the most dangerous places in Massachusetts.

ACTIVITY	SYMBOL	NA	HEADD	OWNER	OPER	DATE	INDUSTRY	city	latitude	longitude	Decade
10096392	Non-Irregular	NA	NA	NA	EQUIPMENT	1980-05-05	CONSTRUCTION	MASSACHUSETTS	41.88	-71.16	70s
	Hospitalized										
30554875	Fatal	REG	47	ELECTRICAL	CONSTRUCTION	2003-04-04	ELECTRICAL	MASSACHUSETTS	41.88	-71.16	2000s
30681192	Fatal	REG	30	ELECTRICAL	CONSTRUCTION	2004-05-07	POWER	MASSACHUSETTS	41.88	-71.16	2000s

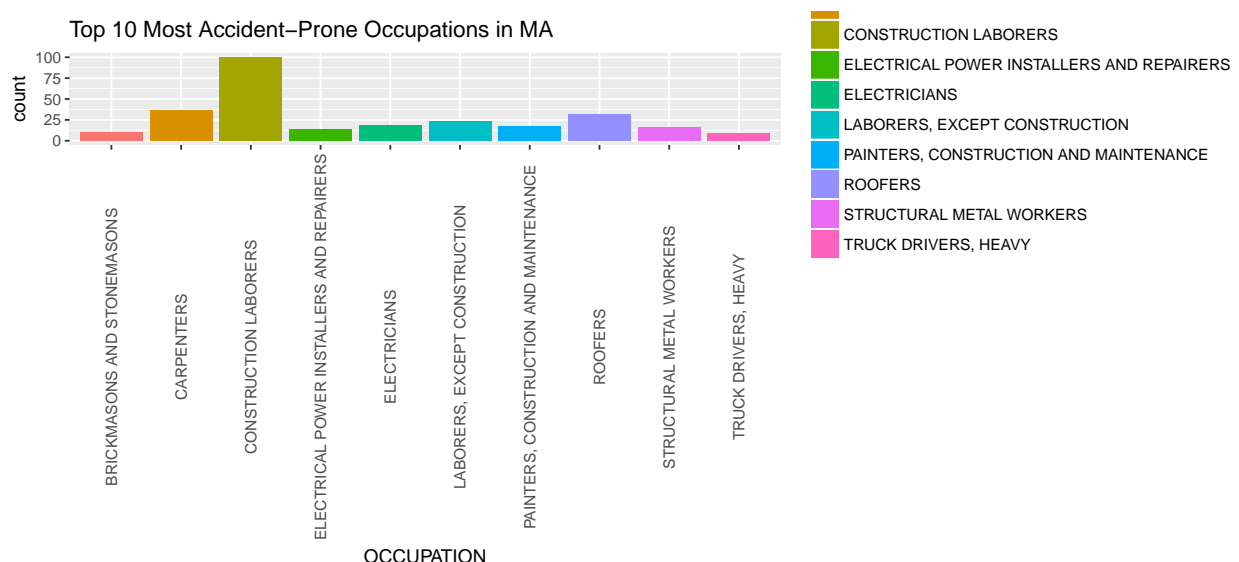
ACCIDENT ID	DATE	TIME	LOCATION	OCCUPATION	DESCRIPTION	CAUSE	INVESTIGATION	STATUS	INDUSTRY	DECADE	
3068	11-27-2004	12:00	Regis	Bar	NA	IN-APRO-PRI-ATE FOR TASK	TOO DB-(POWERED) ACTION	CON-COR-STRUC-TION, LLC	NEW MUL-TI-FAM-ILY HOUS-ING CON-STRUC-TION (EX-CEPT OPER-ATIVE BUILDERS)	70.61	1990s
3055	11-27-2004	12:00	Regis	Bar	NA	NA MISJUDGMENT	OTHER HAN-IN DLG OR EQUIPMENT	VIEW FIRE WROCK MAN'S BENT	WOOD PROD-UCTS, NEC	70.07	1970s

Exploring the Data

Now that we have a Table with all the useful information about all the accidents that happened and have it organized, we can begin to explore the accidents in Massachusetts. Because of the our tidy data which brings in many types of attributes for each accidents, we can produce many visualizations from the data. Here are a few types:

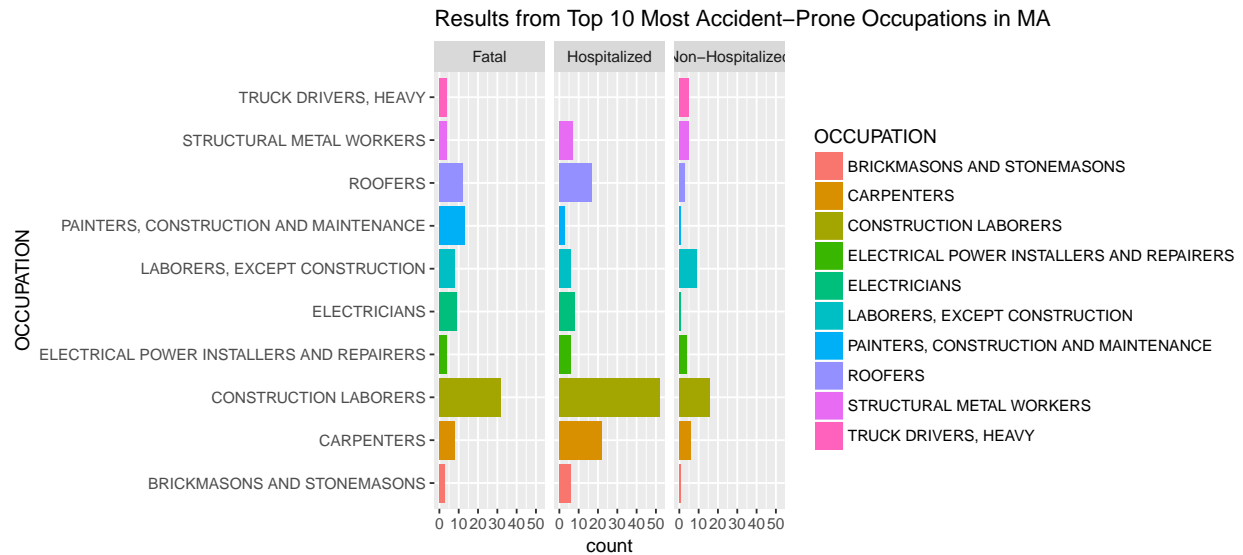
Bar Graphs

Since there are so many discrete variables, we can explore them with the use of bar graphs. Let's first consider the 10 most frequent **Occupations** in our data:



The plot above gives incite to which occupations lead to the most amount of injuries, but it does not tell us the severity of those injuries.

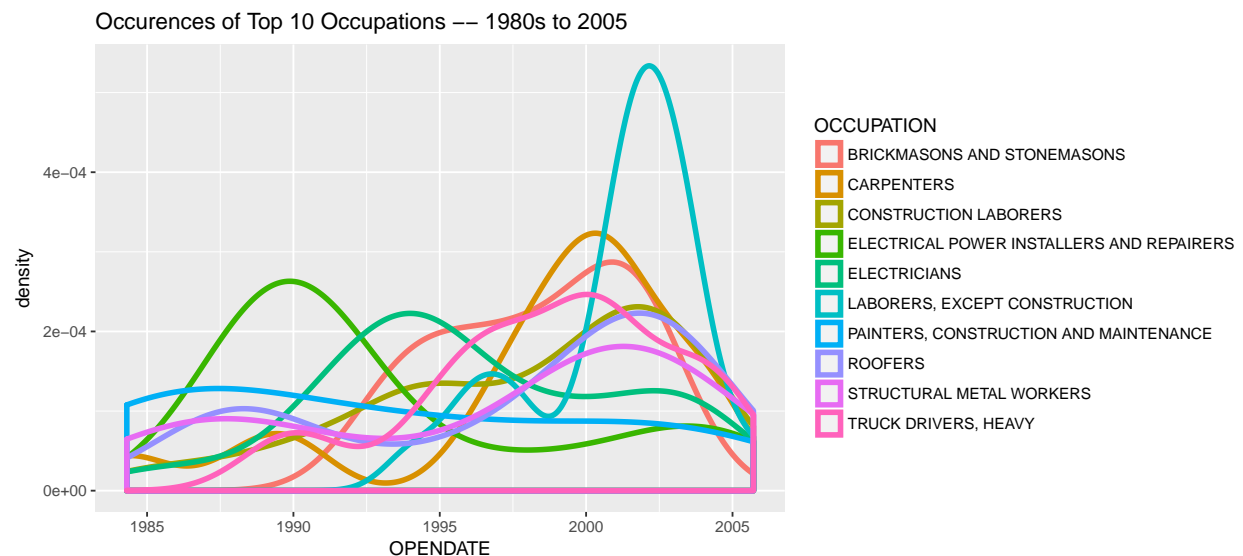
Because of that, consider the same set of data, but grouped by another variable, **Degree** of Accidents. Here we can consider the outcome of the 10 most frequent Occupations in the data:



We are able to consider if the most accident-prone jobs tend to have serious injuries or just many none serious injuries.

Time Series – Time Occurances

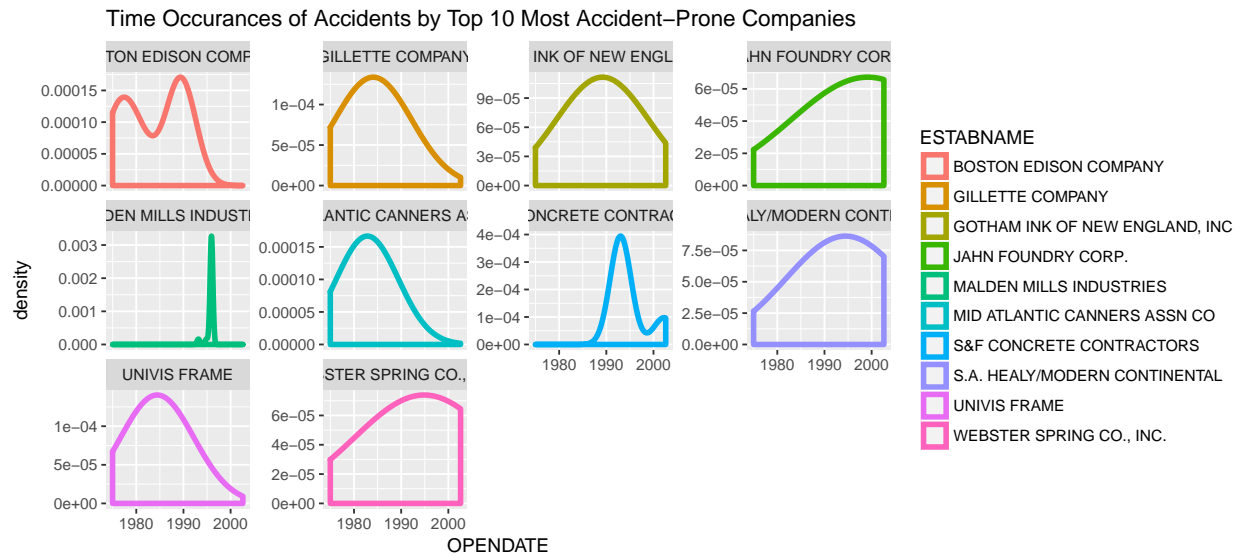
Using the Accident's Date, we can explore when and how many accidents happen. Consider still the Top 10 Accident-Prone Occupations. We can see when theses accidents happened:



Many of these results skew either left or right, so it may be considered if these are disappearing career paths or rising, more dangerous careers. Perhaps the job are changing.

Let's now explore another Attribute over time, Company names.

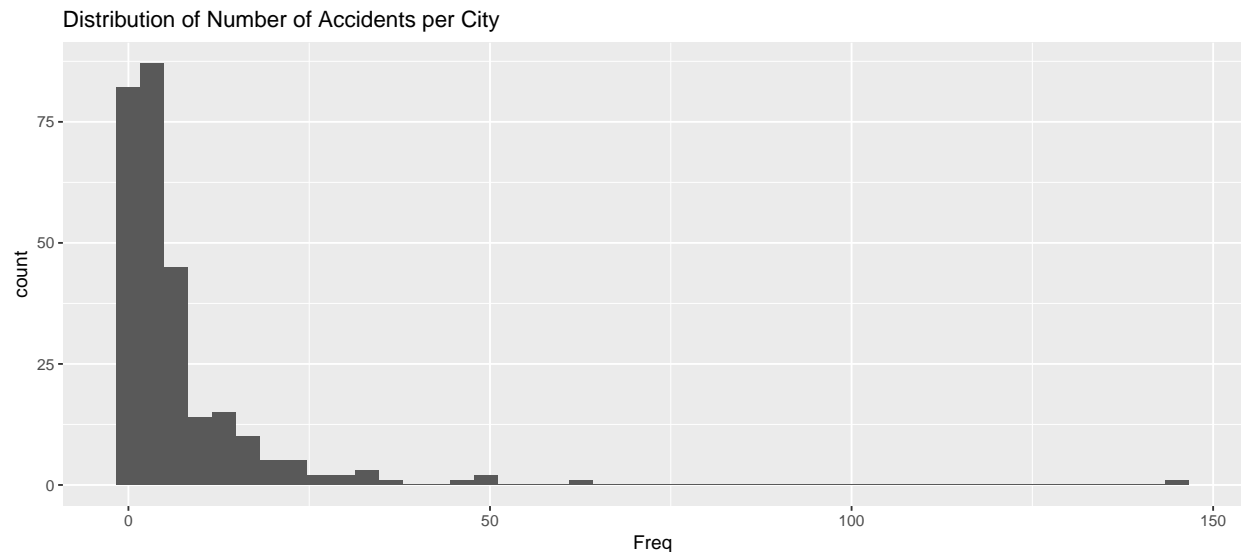
Sub-setting by the 10 Most Occuring Companies, we can see when in time these Accident-Prone companies had these Accidents:



All of these companies had many accidents but the time over which they had these accidents differ and can be visualized with the plots. Some have smaller accidents throughout time while some have large accidents that occur not so often.

Histograms

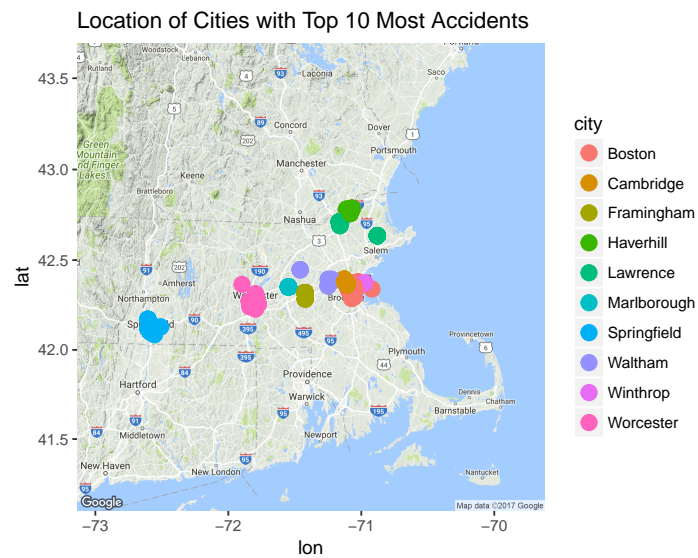
Histograms can be used to see the distributions of some of our discrete variables. Consider the accidents that occur per city in the whole time frame:



Over a time interval, we are able to see that lower amounts of accidents are more likely to occur. Cities with more accidents are less frequent and could be accounted by being larger, more industrial places.

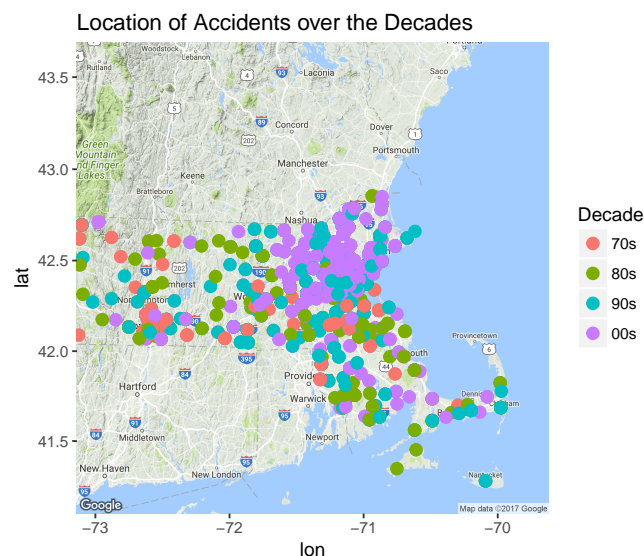
Map Graphics

With the use of location, we can visualize where in the state these accidents occurred. We are able to see where the cities with the most accidents are.



We can see that many of these locations are near Boston, but also near cities where there is a lot of industry compared to smaller places.

Another option is to see where in Massachusetts had the most accidents throughout time. To make it easier to see, let's group accidents by the decade in which they occurred. We can see whether accidents are moving to new locations throughout the state.



There could be a trend that is moving accidents toward Boston.

Final Remarks

By cleaning up the data but also bringing in the most information we can about each accident, we get one data frame that is easy to understand and ready to explore. At the beginning, our data of interest was scattered across many tables and had many values that were incorrect but, by using R, we are able to bring the data to one place.

The data is versatile in the many ways it can be visually explored. The ways to explore the data that are above are just a subset of all possible. The data can be seen through many of the discrete variables, the frequency of them, the time occurrence of accidents, and/or geographical locations of accidents.

The data that was prepared can be taken in many directions to discover the most dangerous place to work, and to uncover the multifaceted stories behind each.