# Analyzing the NEISS Data

Moneka Bommasani

## Introduction:
This is an analysis report on NEISS data.

## Source:
United States Consumer Product Safety Commission's National Electronic Injury Surveillance System (NEISS) - http://www.cpsc.gov/en/Research--Statistics/NEISS-Injury-Data/

## Tools Used:

| | |
|---|---|
| RStudio | For Q1 – Q4 |
| JMP | For Q5 |

## Questions:

First, let's load the data into RStudio using "read.csv".
Now import the required libraries:

- library(dplyr)

## 1)

### a) What are the top three body parts most frequently represented in this dataset?

**Code:**
```
bodyparts <- left_join(NEISS2014, BodyParts, by = c("body_part" = "Code"))
tail(names(sort(table(bodyparts$BodyPart))),3)
```

**Result:**

                "Finger"        "Face"          "Head"

**Explanation:**
First, let's join the two datasets (NEISS2014, BodyParts), so as to display the name of the body parts. Now we can calculate the frequency of each body part with the table() function; then sort them to know the top three body parts using tail() function. This is because, by default sort gives an ASCENDING order. So we need the least 3 values to get the top three frequent body parts.

## b) What are the top three body parts that are least frequently represented?

**Code:**
- bodyparts <- left_join(NEISS2014, BodyParts, by = c("body_part" = "Code"))
- head(names(sort(table(bodyparts$BodyPart))),3)

**Result:**

"25-50% of body"    "Pubic region"    "Not Recorded"

**Explanation:**
Same as above, we join the two data sets to display the name of the body part. Now we can calculate the frequency of each body part with the table() function; then sort them to know the least three body parts using head() function. This is because, by default sort gives an ASCENDING order. So we need the top 3 values to get the bottom three frequent body parts.

**Note:**
If we need the exact body parts, then the answer will be as follows:

- head(names(sort(table(bodyparts$BodyPart))),5)

"Pubic region"    "Internal"    "Arm, upper"


## 2)

### a) How many injuries in this dataset involve a skateboard?
### b) Of those injuries, what percentage were male and what percentage were female?
### c) What was the average age of someone injured in an incident involving a skateboard?

**Approach 1:**

**Code:**
- nrow(data.frame(grep("skateboard", NEISS2014$narrative, ignore.case = TRUE )))

- subset <- NEISS2014[grep("skateboard", NEISS2014$narrative, ignore.case = TRUE), ]
- options(digits=4)
- transform(as.data.frame(table(subset$sex)),percentage_column=Freq/nrow(subset)*100)

- +(subset$age > 200)
- mean(subset$age)
- summary(subset$age)

Result:

Explanation:
We can search the narrative section of the data for the term "skateboard" using grep() function and count the number of rows to find out the number of injuries involving skateboard.
Then, we can subset the above data to find out the number and percentage of males and females involved in skateboard injuries.
Now, let's check if there are any kids less than 2 years involved in the skateboard injuries. The output shows that there are none. So, the mean or summary of the age column can be seen to find out the average age of a person injured in the skateboard incident.

Approach 2:

If we see the manual, the code for skateboard is 1333. Using the product code:

Code:

- skate <- subset(NEISS2014 , NEISS2014$prod1 == 1333 | NEISS2014$prod2 == 1333)
- nrow(skate)

- transform(as.data.frame(table(skate$sex)),percentage_column=Freq/nrow(skate)*100)

- +(subset$age > 200)
- mean(skate$age)
- summary(skate$age)

Result:

**Explanation:**

We can now use the product code 1333 (using logical OR for prod1 and prod2) to get the subset of data involving skateboard injuries and count the number of rows.

Then, we can calculate the frequency and percentage of males and females involved in skateboard injuries.

Now, let's check if there are any kids less than 2 years involved in the skateboard injuries. The output shows that there are none. So, the mean or summary of the age column can be seen to find out the average age of a person injured in the skateboard incident.

## Note:

There is a difference in the results for both the approaches. This is because of the misspelling of the word "skateboard" in the narrative section of the original dataset in some instances.

So, the **approach 2** result is more appropriate.

Also, the manual includes "skateboard" injuries under some other products (5042, 5035, 3215, 5036, 5044) which are not considered in this result.

## 3)

### a) What diagnosis had the highest hospitalization rate?

**Code:**

- Diagno <- left_join(NEISS2014, DiagnosisCodes, by = c("diag" = "Code"))
- hosp <- subset(Diagno , Diagno$disposition == 4)
- tail(names(sort(table(hosp$Diagnosis))),1)
- transform(as.data.frame(table(hosp$Diagnosis)),percentage_column=Freq/nrow(hosp)*100)

**Result:**

"Fracture"

**Explanation:**

First, let's join the two datasets (NEISS2014 and DiagnosisCodes) to display the diagnosis names. Since, the disposition code for hospitalization is 4, we can subset the dataset using the equal to operator. Then, we can group the types of diagnosis and calculate their frequency as well as the percentage rate to know the diagnosis corresponding to the highest number.

### b) What diagnosis most often concluded with the individual leaving without being seen?

**Code:**

- Diagno <- left_join(NEISS2014, DiagnosisCodes, by = c("diag" = "Code"))

- leave <- subset(Diagno , Diagno$disposition == 6)
- tail(names(sort(table(leave$Diagnosis))),2)
- tail(names(sort(table(leave$diag_other))),2)

**Result:**

"Laceration"        "Other/Not Stated"

"PAIN"              ""

**Explanation:**
First, let's join the two datasets (NEISS2014 and DiagnosisCodes) to display the diagnosis names. Since, the disposition code for individual leaving without being seen is 6, we can subset the dataset using the equal to operator. Then, we can group the types of diagnosis and calculate their frequency know the diagnosis corresponding to the highest number.

The most often diagnosis comes out to be "Other/Not Stated". So, with further investigation, it can be seen that **"Laceration"** and **"Pain"** both are the most often diagnosis where the individual leaves without being seen.

## c) Briefly discuss your findings and any caveats you'd mention when discussing this data

- People who had fractures, had to be hospitalized for the future treatment.
- People in general, who were suffering from pain or deep cuts(Laceration) left without much ado.
- The meaning for "hospitalization rate" seemed to be vague. It was not clear whether the disposition = 4 was the only case to be considered or not.

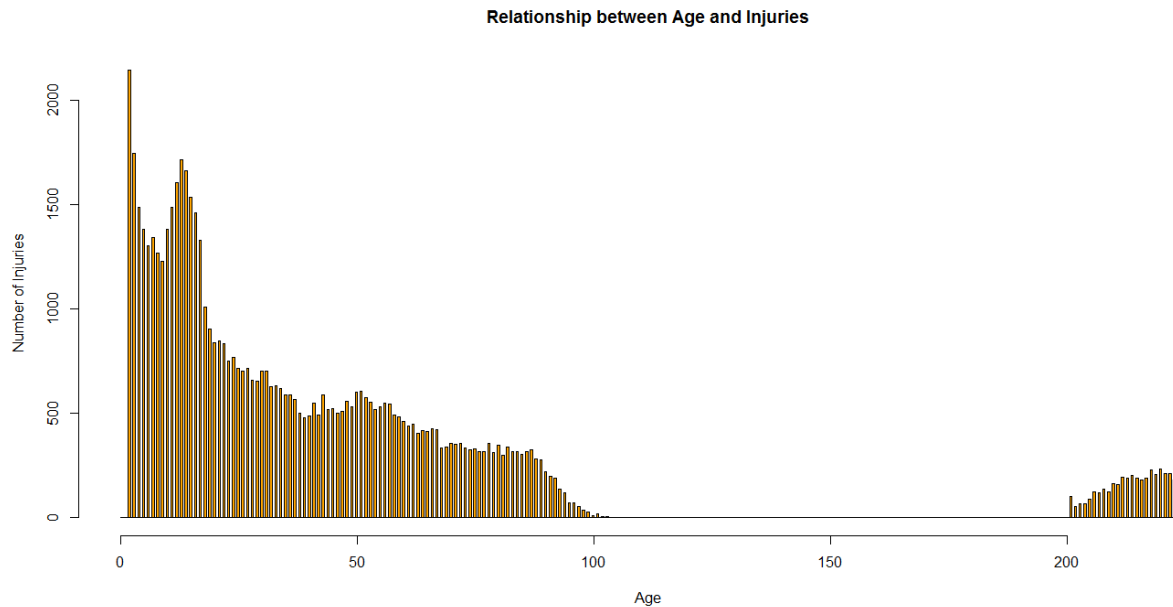## 4) Visualize any existing relationship between age and reported injuries

**Code:**
- hist(NEISS2014$age, col="orange", main = "Relationship between Age and Injuries" , breaks = 500, xlab = "Age", ylab = "Number of Injuries")
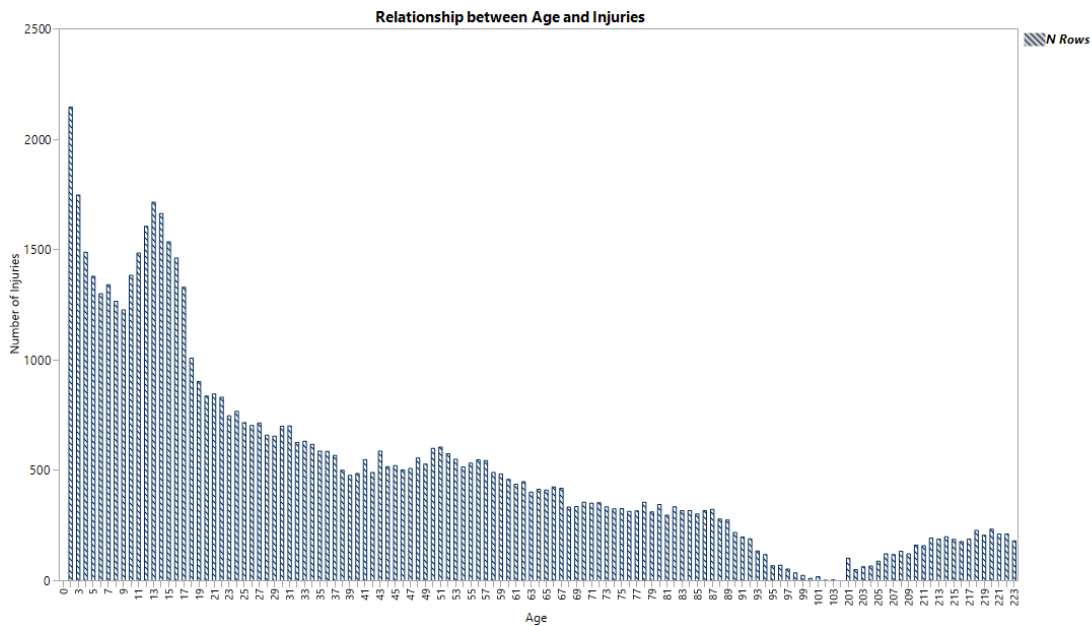
**Explanation:**
The below graphs are visualizations showing the number of injuries with respect to the age (Same using R and JMP). It can be seen that the maximum injuries recorded were at the age of 2. Then the number of injuries decreased and again peaked between the age of 10 – 17 years. From the age of 18, the number of injuries had a slow decrease, peaking a little at 30-31, 41, 43, 49- 53, 56-57. Also, it can be seen that comparatively there were less incidents below the age of 2.

The graph makes sense, as kids below 2 years are generally taken utmost care which implies they are prone to have less injuries. At 2 years, a possible explanation could be

that, since kids start walking/running on their own and are generally admitted into day-cares, they are prone to experience a huge amount of incidents involving injuries. Later, we can say that teens are more prone to have accidents, as they now start using a lot of products and are less cautious at that age.  As age increases, people become more cautious which decreases the risk of having injuries.  Probably, at middle-age due to health problems we can see the little peaks in the graph.



Using R



Using JMP

**5)** **Investigate the data however you like and discuss any interesting insights you can find in the data**

I have imported and joined all the four datasets in JMP.
Also, I have changed the "Data Type" and "Modelling Type" of few variables in the dataset for the below visualizations.
This is because, Ordinal/Nominal gives the frequency counts and Continuous give the average calculations.
Below is the current dataset variable list:

# Investigating individual variables with respect to the number of Injuries:

## Weight



## Sex



| Sex | Count |
|---|---|
| Female | 29996 |
| Male | 35503 |

## Race



| Race | Count |
|---|---|
| American Indian/Alaska Native | 249 |
| Asian | 621 |
| Black/African American | 9935 |
| Native Hawaiian/Pacific Islander | 36 |
| None listed | 19593 |
| Other / Mixed Race | 3389 |
| White | 31676 |

**Race_Other**

Count

American Indian 3, _ 1, Arabic 1, Armenian 1, Biracial 8, Brazilian 1, Declined 2, Guatemalan 1, Hindi 1, Hipanic 1, Hisa panic 1, Hispanic 3, Hispanic 161, Hisp px 1, Hispan jc 2549, Hispanic/Latino 1, Hispanic 1, Hispanic 1, Hispanic 11, Mixed 6, Multi 14, Multi racial 1, Multi-racial 10, Multiracial 252, Native Hawaiian 1, Nepali 1, Not stated 52, NS 1, Pa 15, Somali 1, Tigrinya 1, Uknown 27, UNK 30, UNKN 410, Unknown 1, Vietnamese

**Diagnosis**

Count

| Diagnosis | Count |
|---|---|
| Amputation | 125 |
| Anoxia | 170 |
| Aspirated foreign object | 66 |
| Avulsion | 337 |
| Burns, chemical (caustics, etc.) | 86 |
| Burns, electrical | 25 |
| Burns, not specified | 19 |
| Burns, radiation (includes all cell damage by ultraviolet, … | 35 |
| Burns, scald (from hot liquids or steam) | 392 |
| Burns, thermal (from flames or hot surface) | 512 |
| Concussions | 1495 |
| Contusions, Abrasions | 10646 |
| Crushing | 149 |
| Dental injury | 251 |
| Dermatitis, Conjunctivitis | 414 |
| Dislocation | 988 |
| Electric shock | 32 |
| Foreign body | 1270 |
| Fracture | 9735 |
| Hematoma | 534 |
| Hemorrhage | 64 |
| Ingested foreign object | 483 |
| Internal organ injury | 5306 |
| Laceration | 12307 |
| Nerve damage | 198 |
| Other/Not Stated | 8120 |
| Poisoning | 724 |
| Puncture | 636 |
| Strain or Sprain | 10326 |
| Submersion (including Drowning) | 54 |

**Diagnosis_Other**

Count

**BodyPart**

Count

>50% of body: 1422; 25-50% of body: 4; Ankle: 3781; Arm, lower: 2561; Arm, upper: 745; Ear: 782; Elbow: 1612; Eyeball: 847; Face: 5786; Finger: 5783; Foot: 3090; Hand: 3369; Head: 9891; Internal: 549; Knee: 3616; Leg, lower: 2239; Leg, upper: 756; Mouth: 1254; Neck: 1080; Not Recorded: 390; Pubic region: 286; Shoulder: 2675; Toe: 1280; Trunk, lower: 5717; Trunk, upper: 3868; Wrist: 2116

**Disposition**

Count

Fatality, including DOA, died in the ED: 28; Held for observation (includes admitted for observation): 415; Left without being seen/Left against medical advice: 619; Treated and admitted for hospitalization (within same fa...): 3979; Treated and released, or examined and released without treatment: 59807; Treated and transferred to another hospital: 651

**Location**

Count

0: 19245; 1: 28953; 2: 18; 4: 1483; 5: 4072; 6: 15; 7: 5; 8: 3348; 9: 8360

## Results:

- People weighing $0 - 25$ Units have faced more injuries, followed by people weighing $75 - 100$, $50 - 75$, $25 - 50$ and $100 - 125$ Units. This matches with the age relationship.
- Males have a little more number of injuries as compared to Females. This is obvious, as their population is higher.
- White people have more injuries, followed by Black/African American people, followed by Hispanic people. This makes sense, as the data is from the hospitals in the US and its territories. Also, a major portion of injuries (19,593) have the race unlisted. So, this can be taken as a partial result.
- The top five list of diagnosis are:
  - Laceration
  - Contusions, Abrasion
  - Strain or Sprain
  - Fracture
  - Internal Organ Injury
- Head is the most often injured part of the body. Face, finger and lower trunk parts of the body are also injured often, after head and almost in an equal proportion.
- Maximum proportion of the cases were treated and released or examined and released without treatment.
- It is surprising to see that most of the injuries happened in a farm/ranch followed by home; industrial place being the least.
- Majority of the incidents did not involve fire at all. The one's which involved fire had an approximately equal share of the fire department involvement (or not).
- Maximum injuries reported (7282) have been because of the Floors or Flooring materials, followed by Stairs or Steps (excluding pull-down and folding stairs).
  [**Note**: This was achieved by grouping by product 1 and product 2 and merging data to find out the total number of injuries caused by the product.
  Please find the attached dataset.]

## Future Work:

Relationship between multiple variables and the number of injuries can be found out and visualized for more interesting facts.