

# Stat 405 Final

Meera Borle, Isis Burgos, Naomi Consiglio, Carson Foster, Aidan Gerber

December 10th, 2022

## Introduction

Spoiler Warning Masaki Kobayashi's *Harakiri* (1962) is a Japanese film about a samurai who asks to commit ritual suicide at a lord's palace. Throughout the film, the audience learns the story of what brings the samurai to the palace. At the palace, the samurai argues with and disrespects the lord's samurais, in revenge for past wrongs. These layers of disrespect lead to conflict and the main samurai kills many of the lord's in combat. The film ends with the lord's history where the events of the film are manipulated and recorded incorrectly to preserve honor. Contrastingly, when a person dies in the modern United States, their causes of death and characteristics are meticulously recorded with substantial effort put in to accuracy — not for honor — but for statistics. We are here to do those statistics. Cue [epic](#) music.

## Primary Dataset

The National Bureau of Economic Research creates and distributes a dataset of US mortality for every year since 1959. This dataset is unique for both its breadth and depth. Each row in the dataset represents a single death, and each column represents a different demographic characteristic of the deceased. The information is derived from death certificates filed by medical professionals in the 50 states plus Washington DC. We made the decision to use

the 2019 edition of the dataset since we did not want to focus on COVID-19. Notable information the dataset contains is education, sex, age classification, day of month, place of death, weekday, manner of death, cause of death, and different risk factors that the deceased had. In 2019, there were 2,861,523 deaths total. The following are the 20 most common groups split up by race, age, education, and sex

count	race	age	education	sex
91215	White	85-89 years	High school graduate or GED	F
88045	White	90-94 years	High school graduate or GED	F
76221	White	80-84 years	High school graduate or GED	F
60991	White	75-79 years	High school graduate or GED	F
60949	White	80-84 years	High school graduate or GED	M
60003	White	75-79 years	High school graduate or GED	M
57183	White	85-89 years	High school graduate or GED	M
53685	White	70-74 years	High school graduate or GED	M
50128	White	65-69 years	High school graduate or GED	M
48356	White	70-74 years	High school graduate or GED	F
47683	White	60-64 years	High school graduate or GED	M
45466	White	95-99 years	High school graduate or GED	F
39068	White	90-94 years	High school graduate or GED	M
35575	White	55-59 years	High school graduate or GED	M
34633	White	65-69 years	High school graduate or GED	F
28678	White	60-64 years	High school graduate or GED	F
23033	White	85-89 years	Bachelor's degree	M
22210	White	70-74 years	Some college credit, but no degree	M
21045	White	75-79 years	Some college credit, but no degree	M
20951	White	80-84 years	Bachelor's degree	M

## Secondary Dataset

For our secondary dataset, we are using the Behavioral Risk Factor Surveillance System Survey. This survey includes different free text survey questions from across the United States and territories with responses broken out by subgroup. There is also information on sample size, percent affirmative response, and confidence interval bounds. We combine the secondary dataset by matching up subgroups between the death dataset and the risk factor dataset and trying to use aggregate statistics to analyze how risk factors can be matched with causes of death.

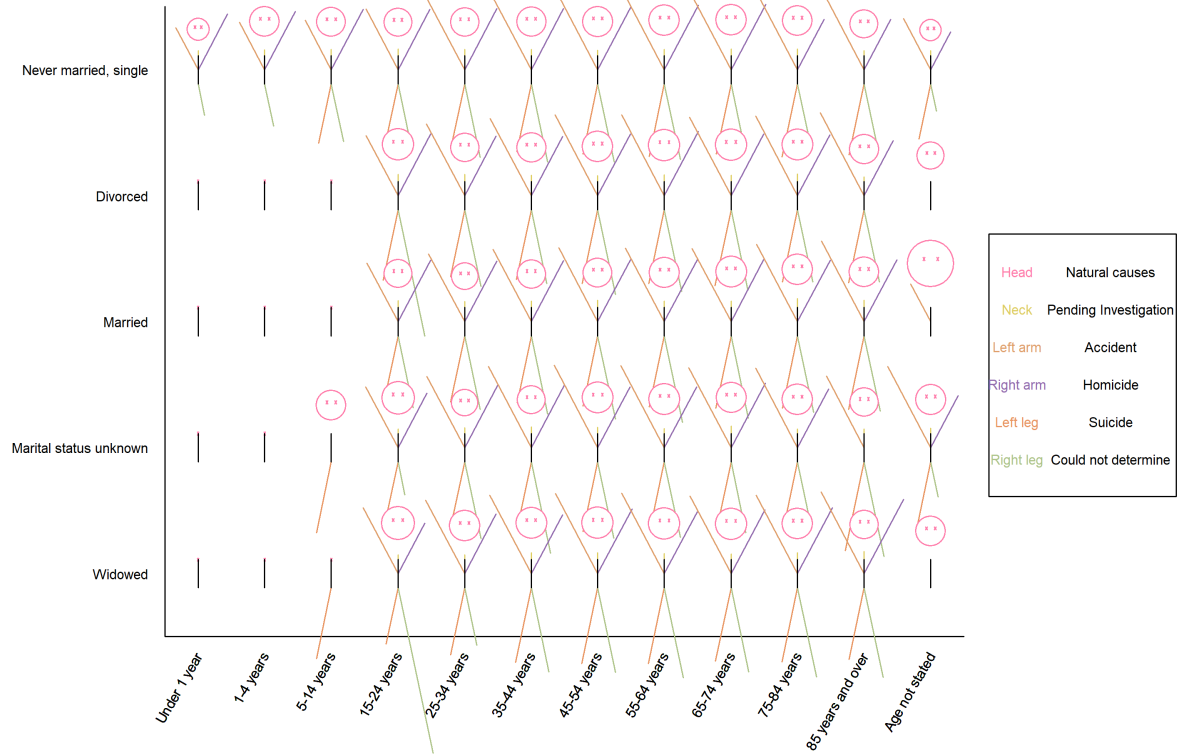
## Questions

We are interested in what different factors are correlated with higher death rates. Depending on the different causes, different policies can be recommended. Moreover, we can break down the different causes of death by demographics such as race, age, and sex to determine where resources should specifically be directed.

## Killer

This plot demonstrates the most common manners of death among people in different cross sections of age and marriage. Head scale is determined by natural causes. Neck scale is determined by pending investigation. Left arm scale is determined by accident. Right arm scale is determined by homicide. Left leg scale is determined by suicide. Right leg scale is determined by could not determine.

## Average Ailments Attributed to Death by Marital Status, Age, Matter



Age	N	Average Record Count
Under 1 year	21012	2.2
1-4 years	3701	2.8
5-14 years	5541	2.8
15-24 years	29979	3.0
25-34 years	59543	3.3
35-44 years	83472	3.3
45-54 years	161212	3.2
55-64 years	376411	3.2
65-74 years	557075	3.2
75-84 years	689088	3.1
85 years and over	874198	3.0

Age	N	Average Record Count
Age not stated	291	2.6

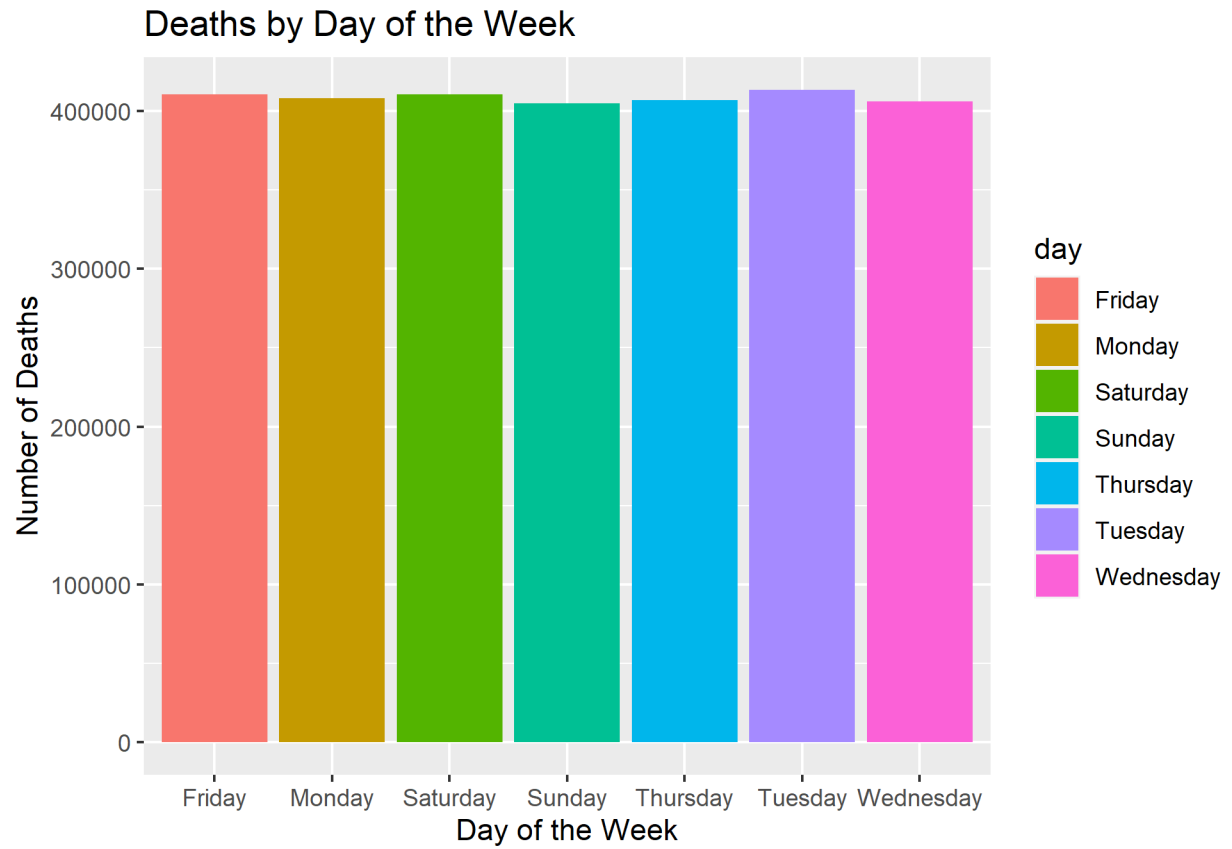
Marital Status	N	Average Record Count
Divorced	478548	3.2
Married	1038238	3.1
Never married, single	403235	3.1
Marital status unknown	23155	3.3
Widowed	918347	3.0

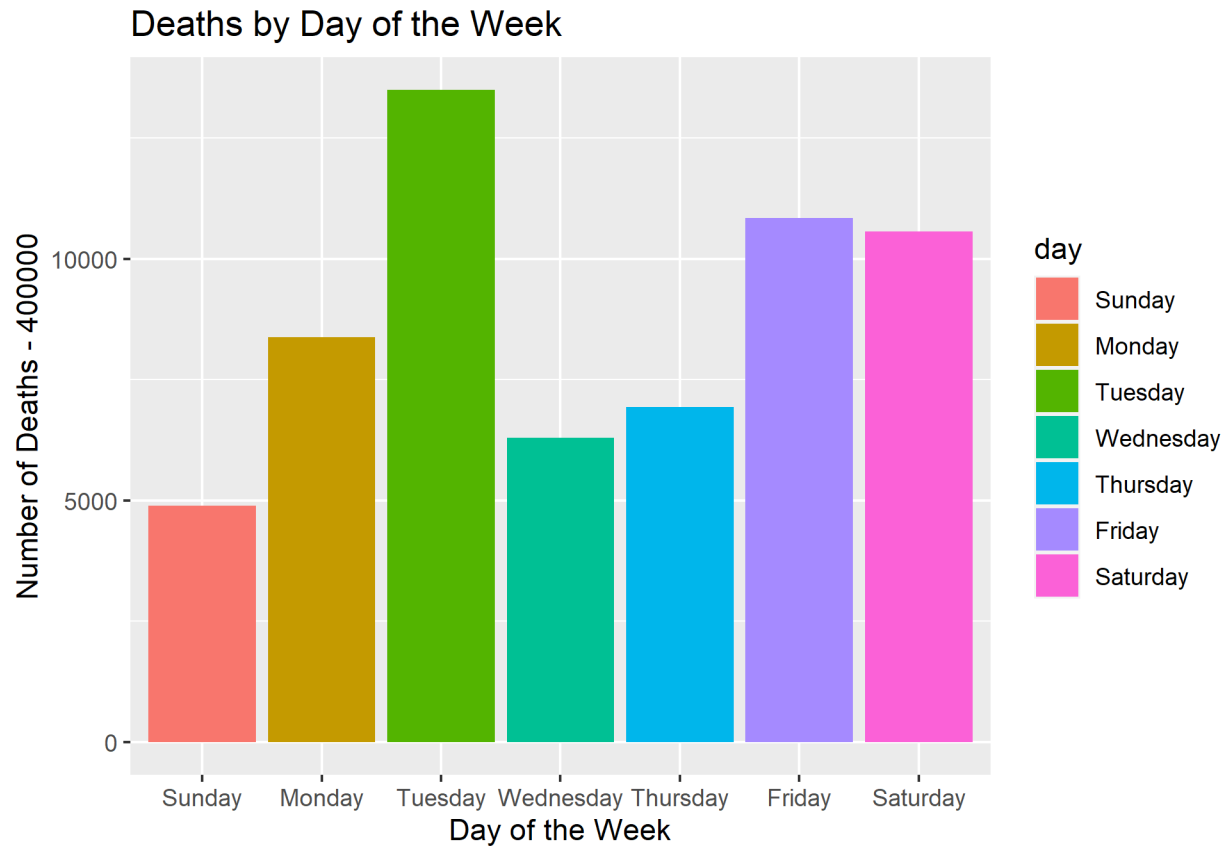
Manner	N	Average Record Count
Accident	173608	4.0
Suicide	47764	2.9
Homicide	20310	3.2
Pending Investigation	4484	1.3
Could Not Determine	11800	2.9
Natural	2327811	3.0
Not Specified	275746	3.4

## Exploration

### Deaths by Weekday

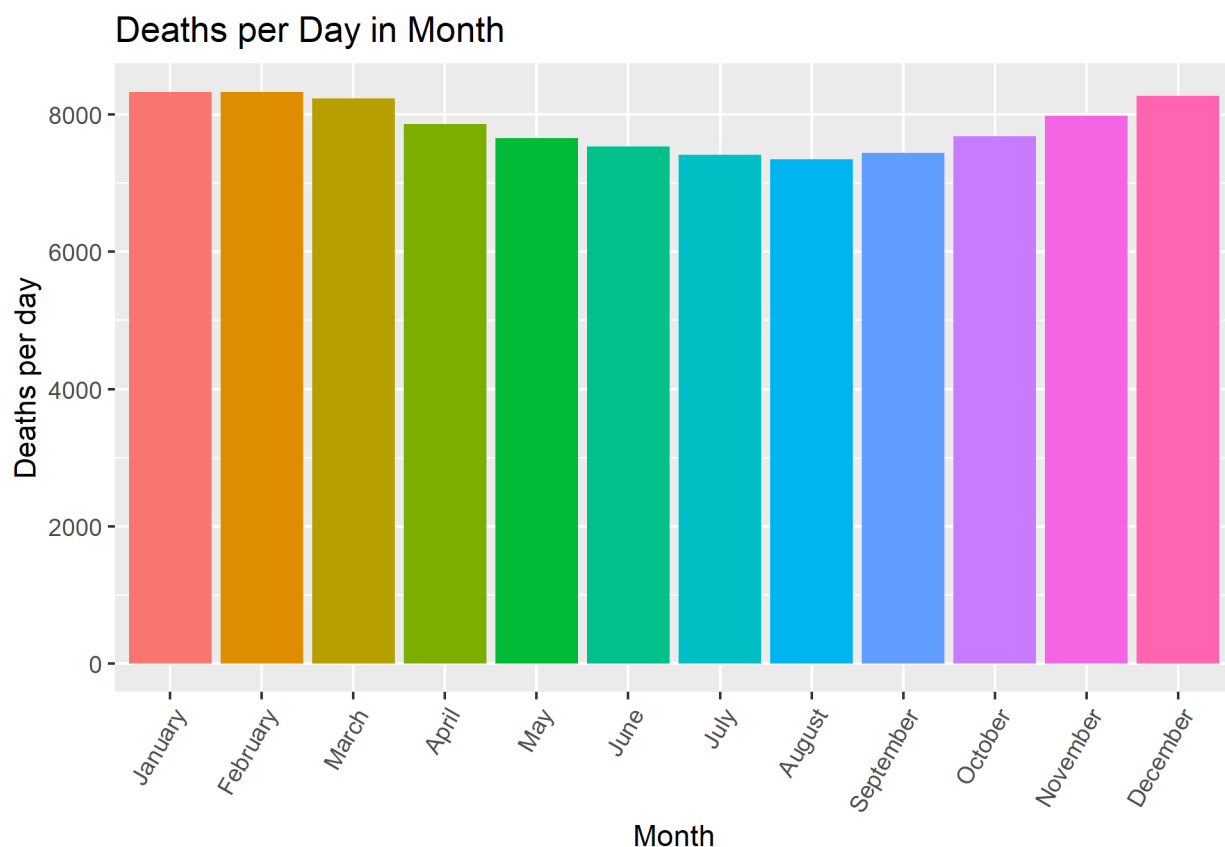
First, we plotted weekday of death versus death counts. There were the most deaths on Tuesday. However, days have an average of 7839.79 deaths and 2019 had an extra Tuesday so adjusting for that, the most deaths were on Fridays.





## Deaths by Month

The most deaths occur in the coldest and darkest months of the year which are February, January, December, and March. Summer months have lower deaths by around 10-11%.

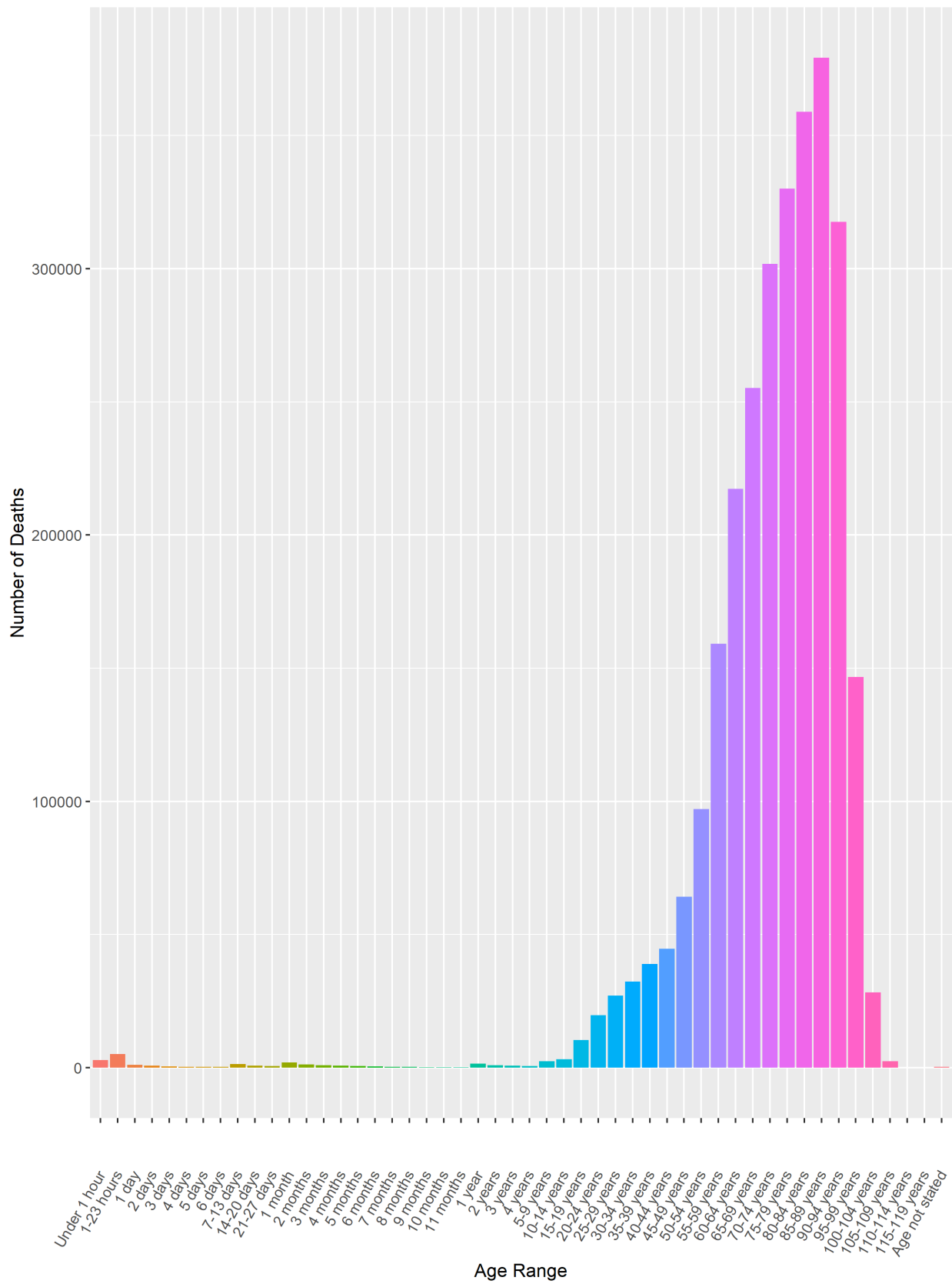


## Deaths by Age

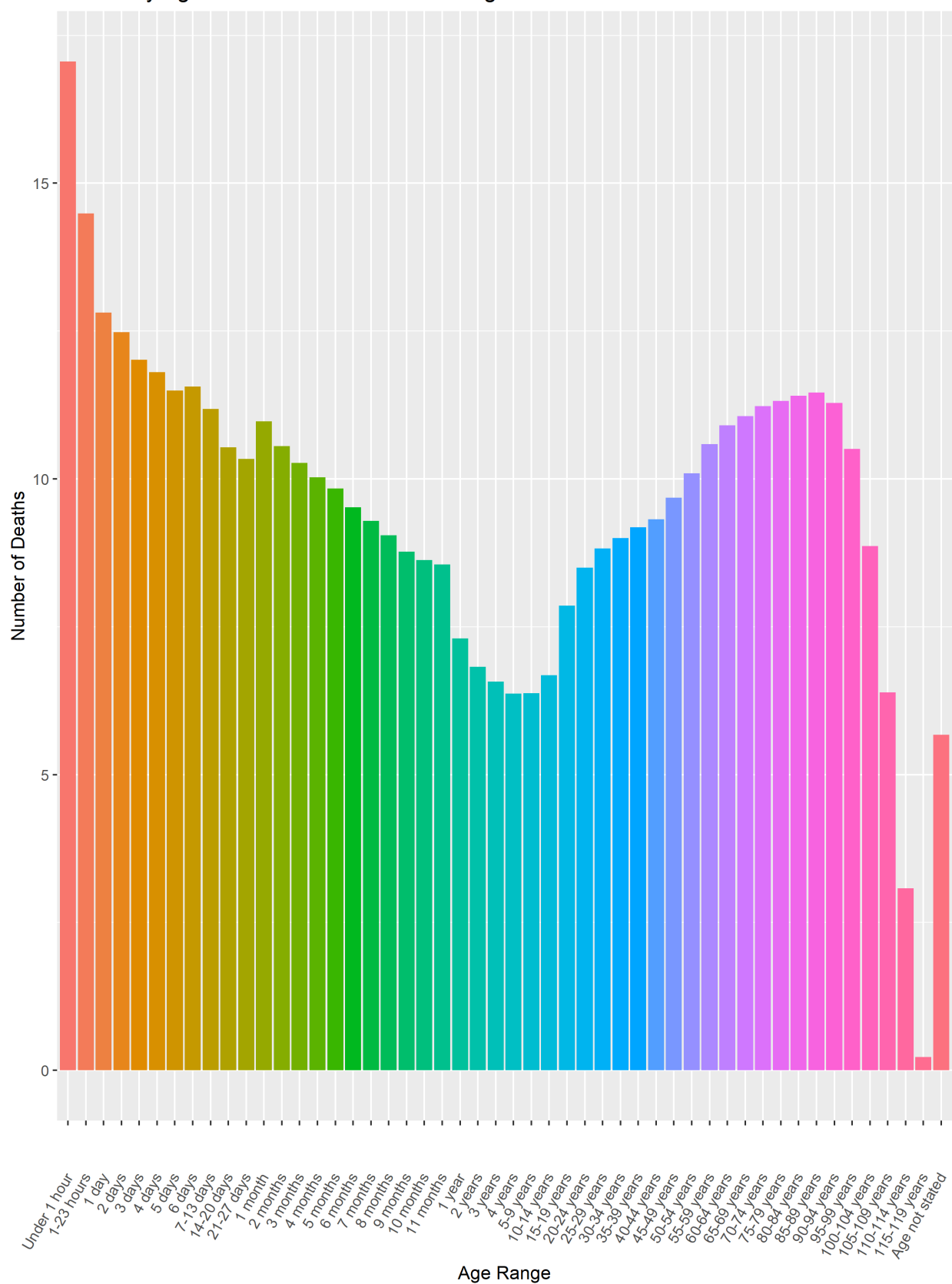
Next, we plotted age versus death counts. Deaths were most prevalent among older age groups such as those between 70 and 84, although deaths start increasing more quickly at age 60. There is also a spike in those less than 1 day old. However, those greater than 1 day old do not frequently die. We also created a version of the plot scaled to bucket size. For privacy reasons, the NBER does not release ages of deaths but rather different buckets that the ages fall into. These bucket are of different lengths of time so we created a rescaled version. This plot was then put on a log scale to better showcase the data.



Deaths by Age



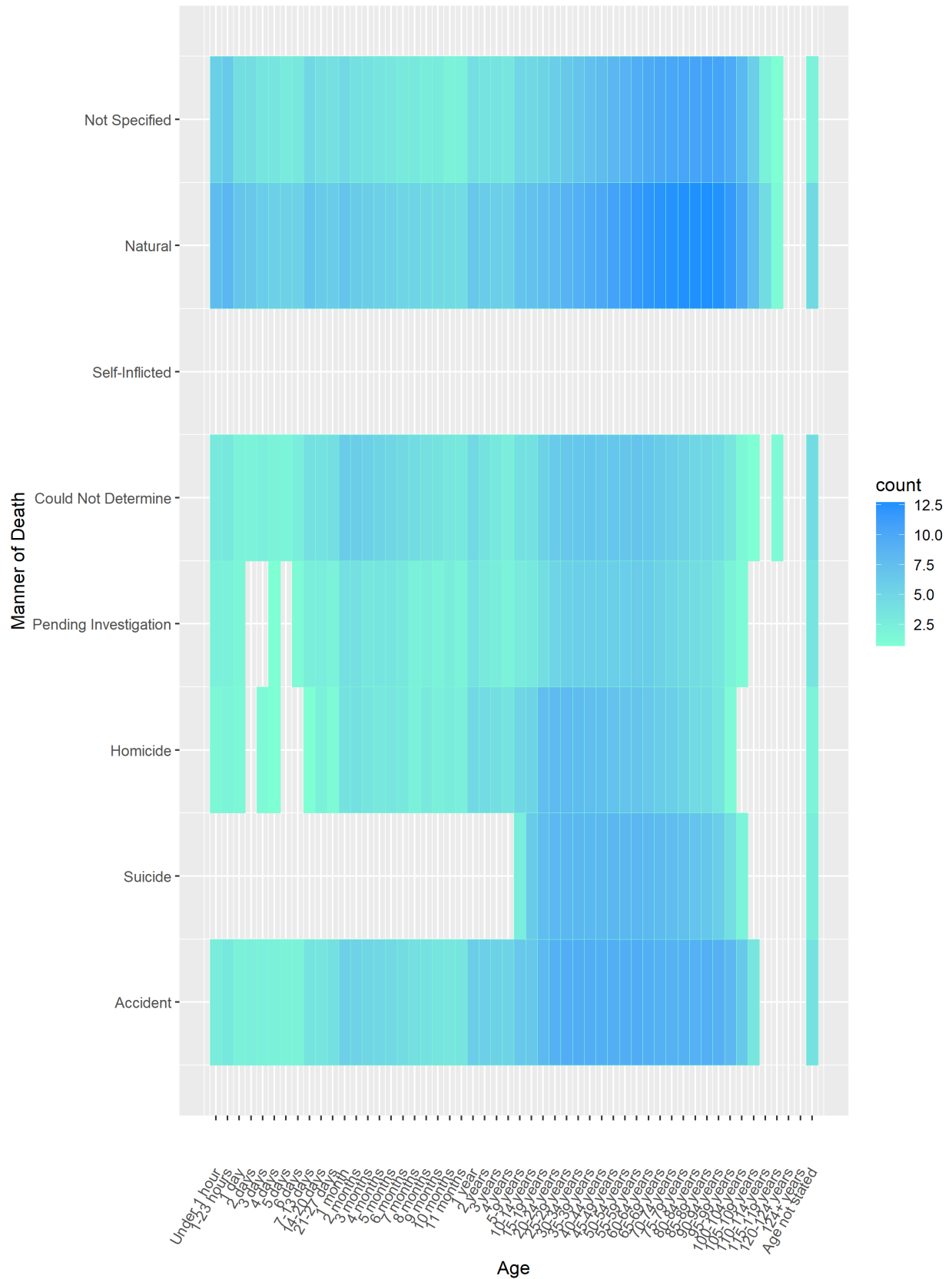
Deaths by Age Scaled to Bucket Size on Log Scale



## Deaths by Manner

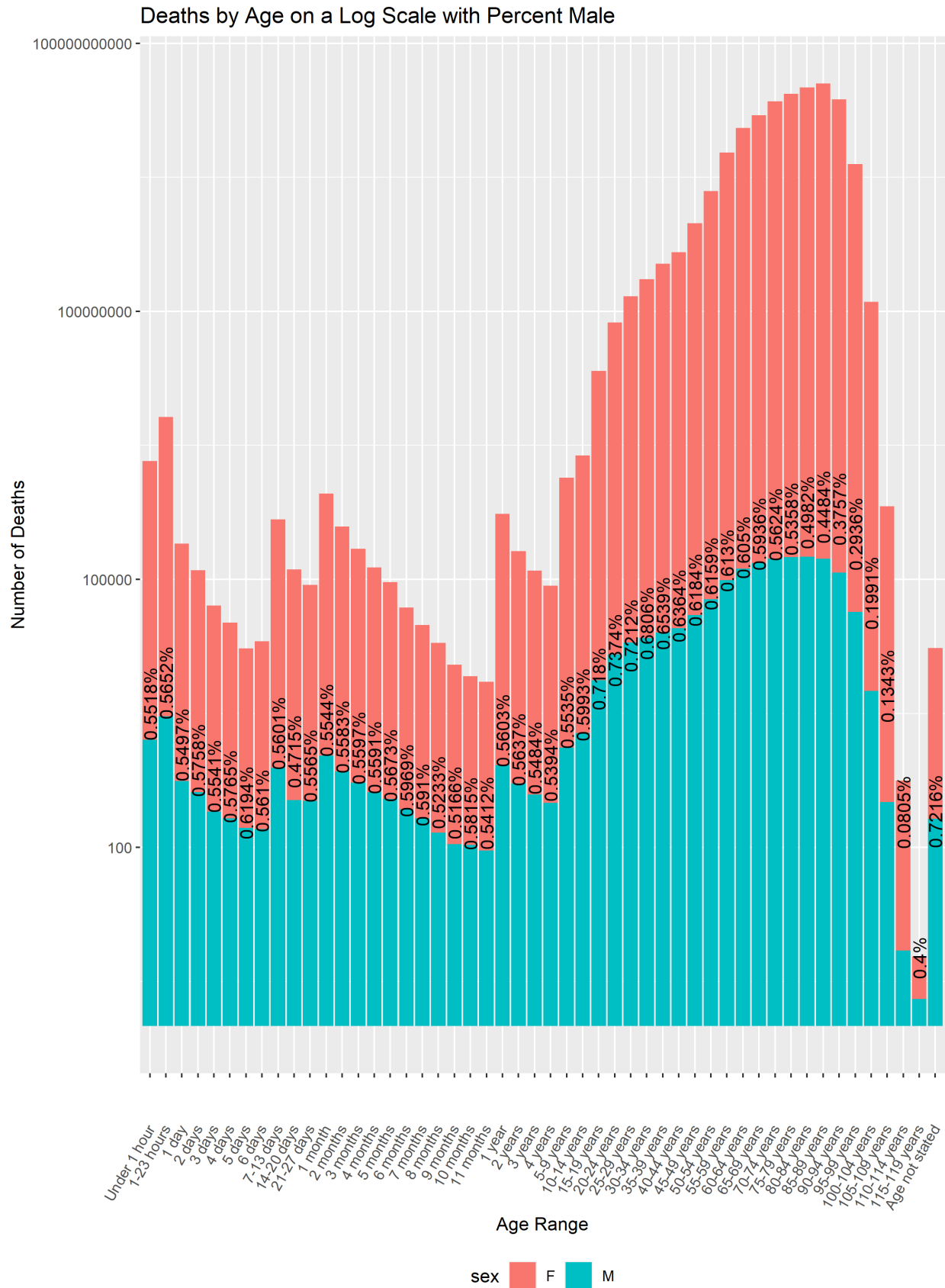
Here, we plotted the manner of death versus age and counted how many people of a certain aged died based on a certain manner of death. A few key finding of this analysis shows that the majority of people die from natural causes, especially those aged 60+ and less than 1 day old, and accidental causes, spanning across all age groups. What this plot may help to inform us about is the behavior and activities that people in a general age group may commonly engage in that may have lead to their manner of passing. By being observant of the manners of death based on age group, preventative methods can be used to decrease the number of accidental related deaths if we are able to determine commonly engaged activities for age groups. Using this plot will help us answer the cause of death among the different age groups, and further promote research in what actual activities people are participating in that lead to their manner of death.

Manner of Death by Age on a log scale



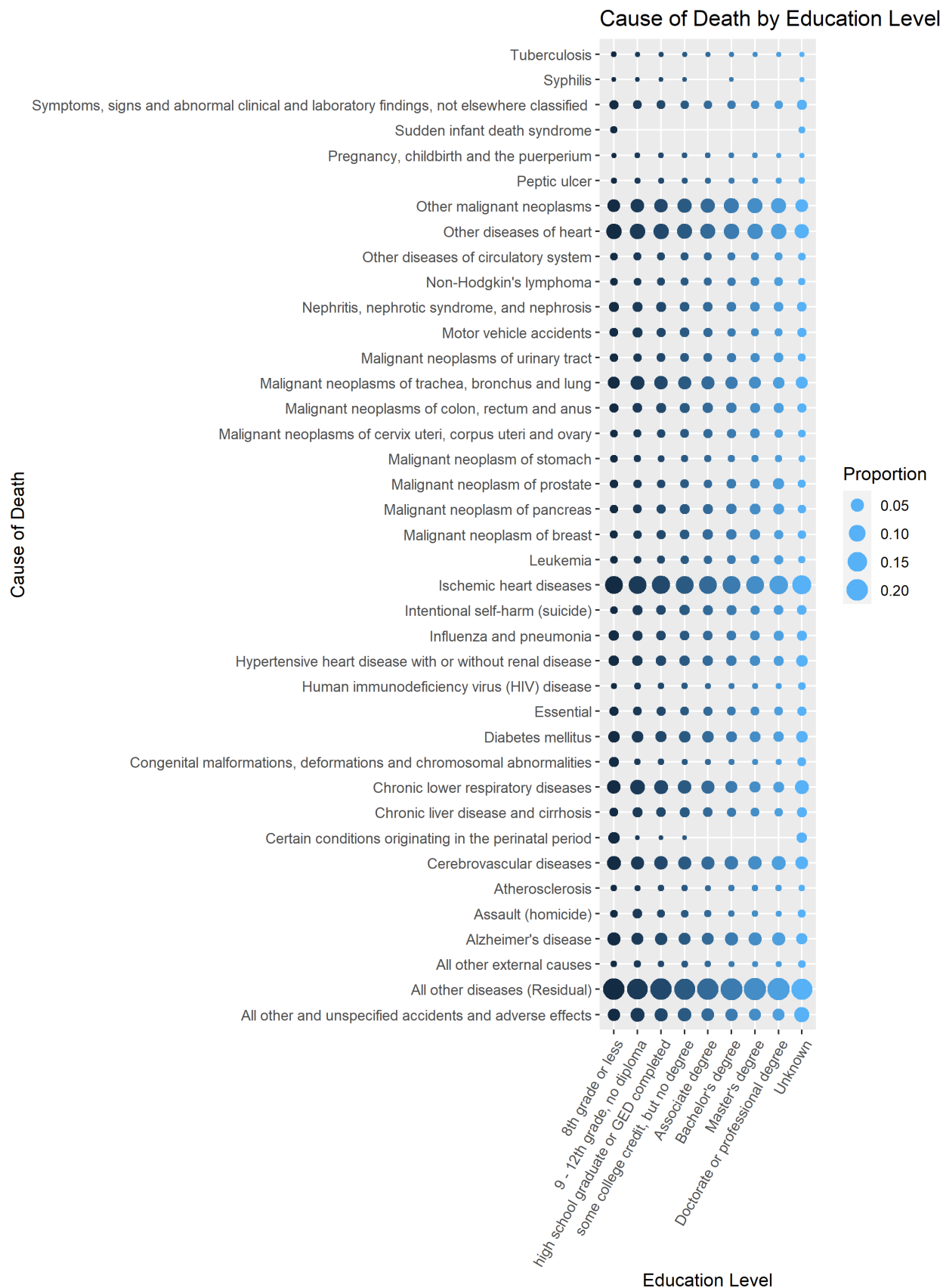
## Deaths by Age and Gender

In this plot, we plotted the number of deaths versus age ranges while demonstrating how many men compared to women passed away in each age category. In each of the bars, the red fill represents the amount of women who passed away in that particular age range while the blue accounts for the amount of men. The percentage seen in each bar represents the proportion of men in a given age range that passed away. This analysis shows that the majority of people under the age of 80 who pass away tend to be men, as nearly every bar from ages 0-80 shows the proportion of male deaths to be above 50%. This proportion of male deaths goes down after 80 years of age, and is likely because women who are of an older age tend to live to a complete life expectancy. What this plot may help to inform us about is the differences in male and females lives and life expectancies. Further research into differences in lifestyle choices for men versus women as a whole may help better explain why women tend to live longer than men. Furthermore, this plot accompanied with a plot on cause of death by gender, may assist in determining what kind of, potentially more risky, behaviors men may partake in during their lifetimes that lead to an earlier death than women.



## Cause of Death by Education

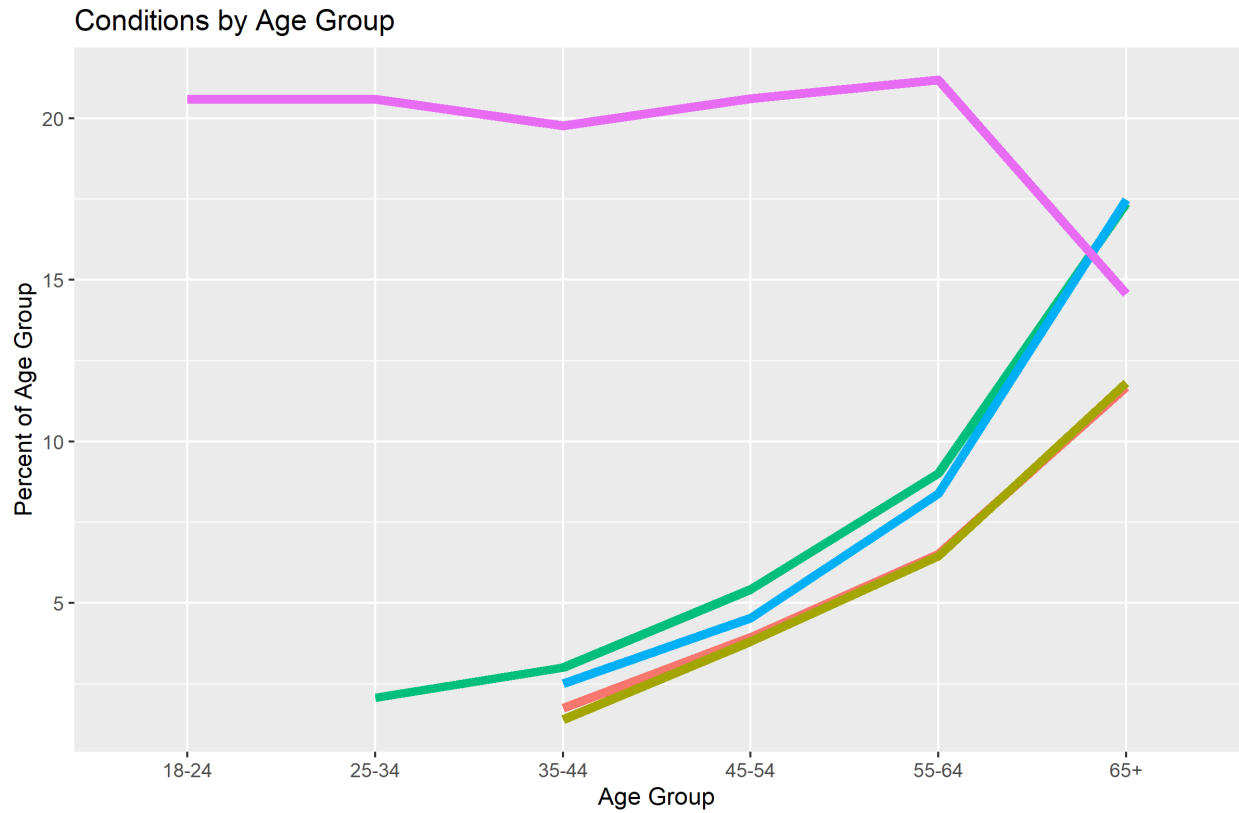
For most causes of death, level of education does not have an impact on what proportion of people have that cause of death. The largest difference belongs to “Certain conditions originating in the perinatal period” with high occurrences in those with 8th grade or less education and those with unknown education and nearly no occurrences in all others. Another large proportion difference is in “Congenital malformations” where 8th grade or less has a much higher mortality proportion than other education levels. For causes of death that are not highly tied to conditions at birth, “Syphilis” and “Assault (homicide)” have the highest differing proportions. “Syphilis” has a quite small sample size but unknown education has the highest mortality proportion. For “Assault (homicide)”, 9 - 12th grade, no diploma has the highest mortality proportion.





## Free Text Analysis: Selected Health Conditions by Age Group

Using regular expressions, we polled the BRFSS data set for questions related to heart conditions, cancer, and depression, while grouping by the age of respondents. Furthermore, we restricted entries to those where participants responded positively, indicating that they did have those conditions. Interestingly, all age groups except 65+ had a high incidence of depression, hovering around the 20% mark. This dips significantly to 15% for the 65+ age group, perhaps because mental health was more stigmatized during their lives and psychological diagnoses were less readily available. In addition, the orange and olive green lines (for heart attack and coronary heart disease, respectively) have a significant degree of overlap, which makes sense, given the conditions. The green line (corresponding to non-skin cancer), is slightly higher than the blue line (corresponding to skin cancer) for all age groups. Furthermore, both cancer have positive slopes, indicating that as your age increases, you become more likely to be diagnosed with cancer.



#### Question

- Ever told you had a heart attack (myocardial infarction)?
- Ever told you had angina or coronary heart disease?
- Ever told you had any other types of cancer?
- Ever told you had skin cancer?
- Ever told you that you have a form of depression?