

Around Simpson's Paradox

Benjamin Cathelineau

21/11/2021

Introduction

This work is not finished.

Question 1

TODO confidence interval on the chart.

```
df = read.csv("Subject6_smoking.csv")
# Compute the ration of dead for
compute_mortality_ratio <- function(smoker_arg,df_arg, title_arg){
  nb_alive= df_arg %>% filter(Status == "Alive" & Smoker== smoker_arg) %>% nrow()
  nb_dead= df_arg %>% filter(Status == "Dead" & Smoker== smoker_arg) %>% nrow()

  df <- data.frame(
    group = c("Alive", "Dead"),
    value = c(nb_alive,nb_dead)
  )

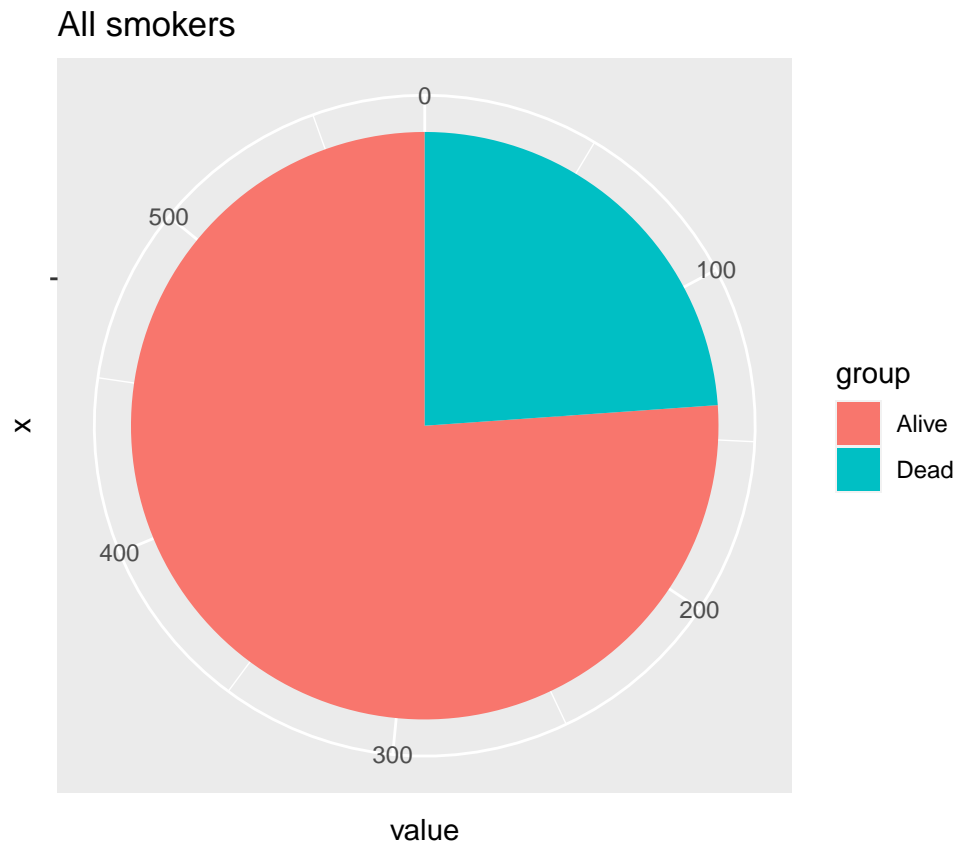
  bp<- ggplot(df, aes(x="", y=value, fill=group))+
    geom_bar(width = 1, stat = "identity") + ggtitle(title_arg) + coord_polar("y", start=0)
  print(bp)

  nb_dead / (nb_alive+nb_dead) # divide the number of dead by the total, the total being the addition of
}

compute_confidence_interval <- function(smoker_arg,df_arg){
  alives= df_arg %>% filter(Status == "Alive"& Smoker== smoker_arg) %>% mutate(death_variable=0)
  dead= df_arg %>% filter(Status == "Dead"& Smoker== smoker_arg) %>% mutate(death_variable=1)
  all_member = rbind(alives,dead)
  # QUESTION for the professor Sample variance * 4 ??
  CI(x=all_member$death_variable,ci=0.95)
}
```

We declare a function so that we can compute both rates (for smokers and non smokers), without repeating our code

```
compute_mortality_ratio(smoker_arg = "Yes", df_arg =df, title_arg = "All smokers" )
```



```
## [1] 0.2388316
```

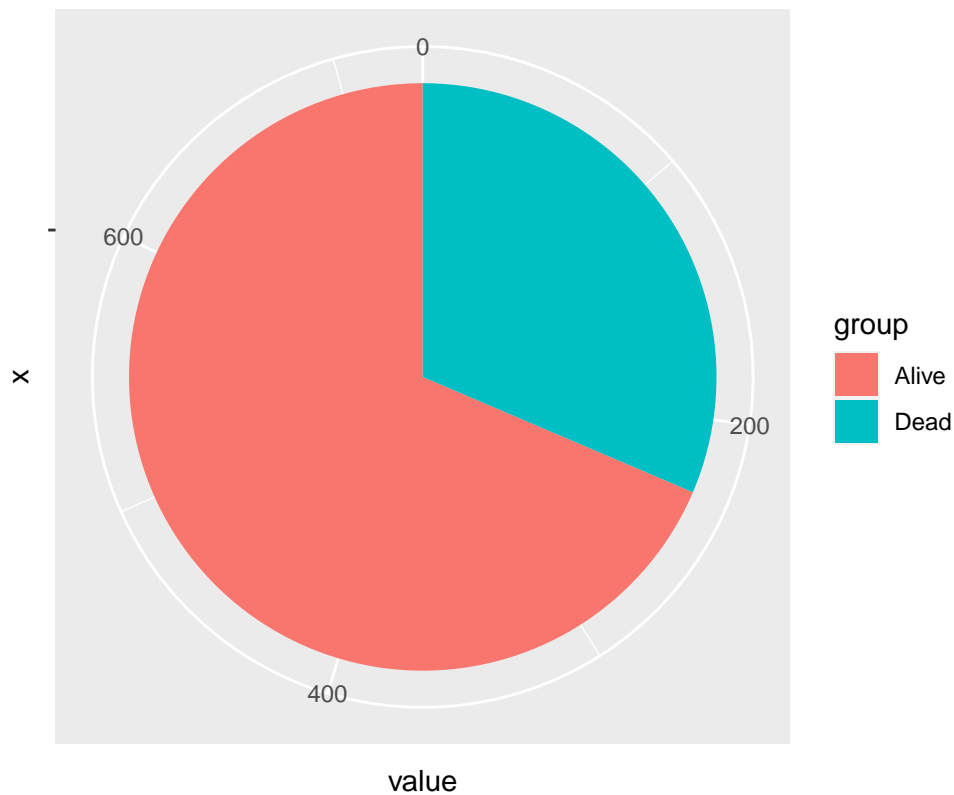
```
compute_confidence_interval(smoker_arg="Yes",df_arg = df)
```

```
##      upper      mean      lower
## 0.2735734 0.2388316 0.2040898
```

The rate for the smoker group

```
compute_mortality_ratio(smoker_arg = "No", df_arg =df, title_arg = "All non smokers")
```

All non smokers



```
## [1] 0.3142077
```

```
compute_confidence_interval(smoker_arg="No",df_arg = df)
```

```
##      upper      mean      lower
## 0.3479142 0.3142077 0.2805011
```

The rate for the non smoking group

The mortality rate is significantly higher for the group that is not smoking. In other words, in with this data, a woman who smoked in 1977 is less likely to have died in 1995 than a woman who did not smoke in 1977.

Of course, this is very surprising because it is now known that smoking cigarette increases the risk of death, trough various mechanisms, such as increased risk of cancer and cardiovascular disease. For more details, consult the relevant [wikipedia](#) article.

Question 2

We will use the recommended age grouping.

```
class1834 = df %>% filter(Age >= 18 & Age < 34)
class3454 = df %>% filter(Age >= 34 & Age < 54)
class5464 = df %>% filter(Age >= 54 & Age < 64)
```

```

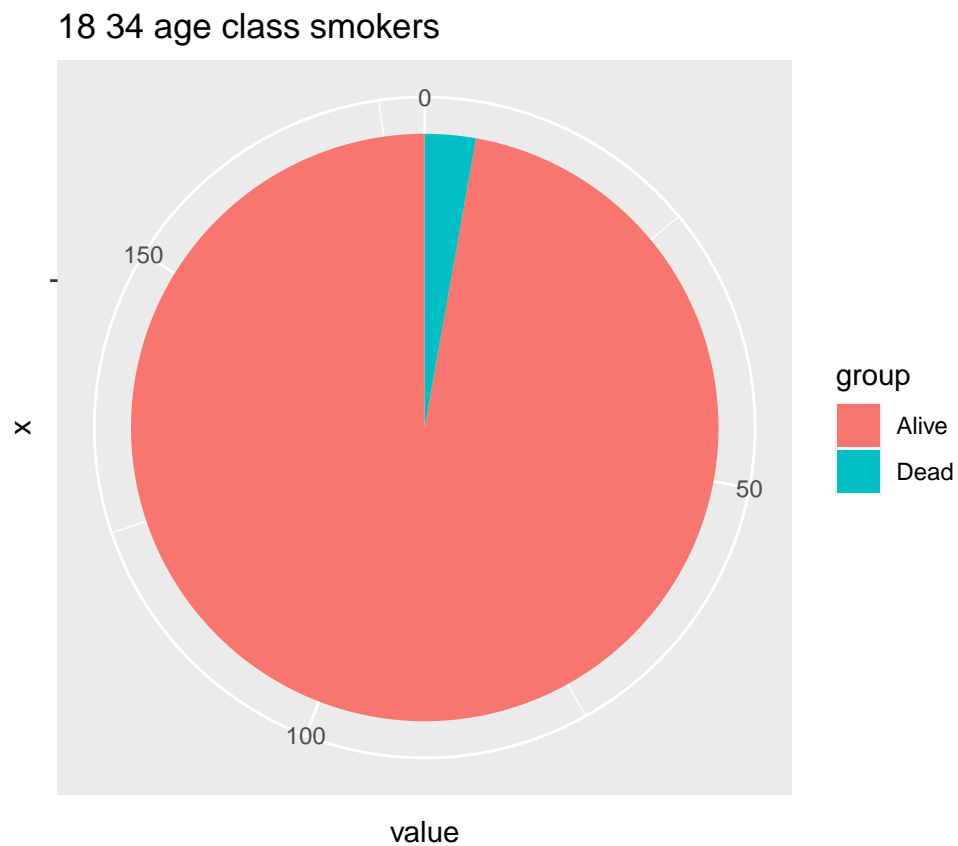
class64 = df %>% filter(Age >= 64)

mortality_ratio_and_ci <- function(df_arg, title_arg){
  compute_mortality_ratio(smoker_arg = "Yes", df_arg = df_arg, title_arg = paste(title_arg, "smokers"))
  compute_confidence_interval(smoker_arg="Yes", df_arg = df_arg)

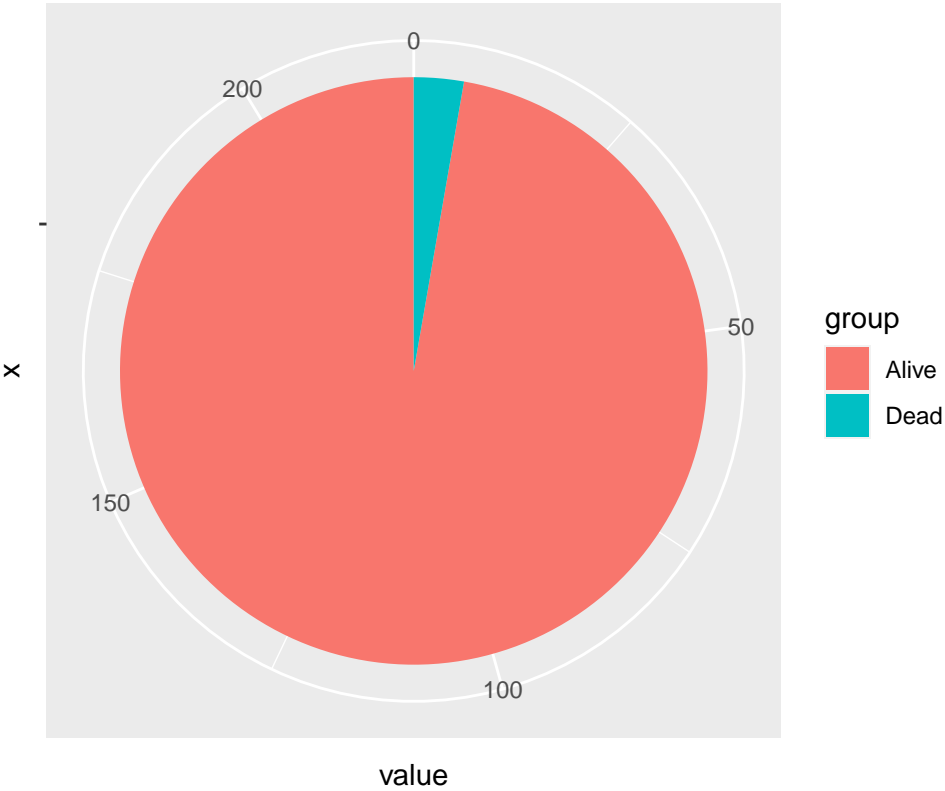
  compute_mortality_ratio(smoker_arg = "No", df_arg = df_arg, title_arg = paste(title_arg, "non smokers"))
  compute_confidence_interval(smoker_arg="No", df_arg = df_arg)
}

mortality_ratio_and_ci(class1834, "18 34 age class")

```



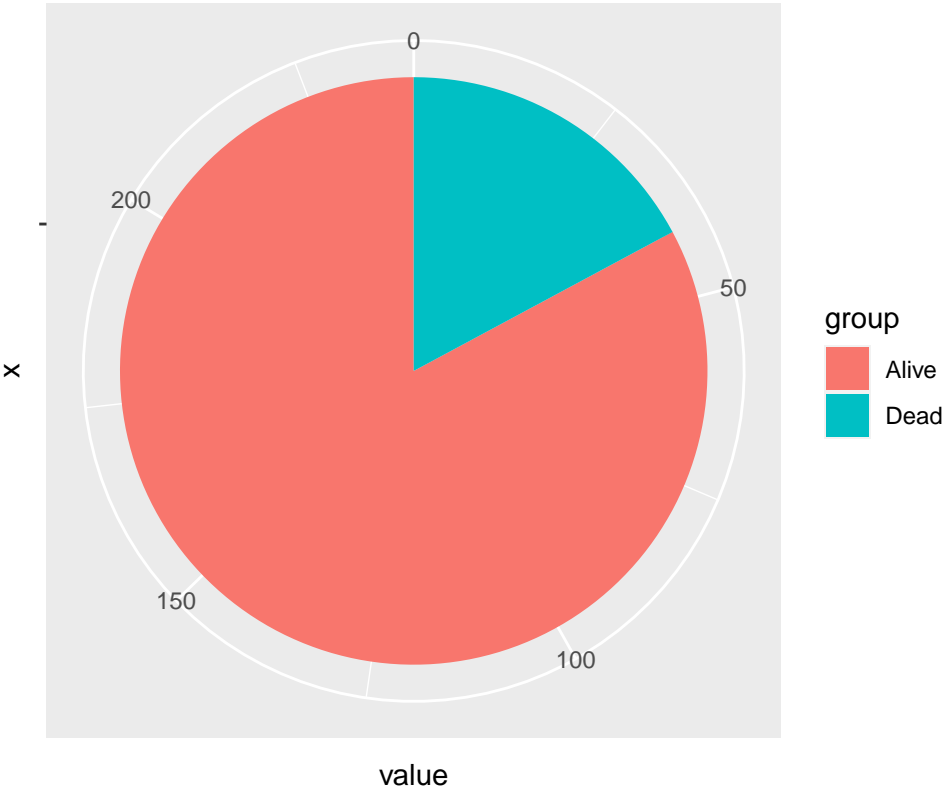
18 34 age class non smokers



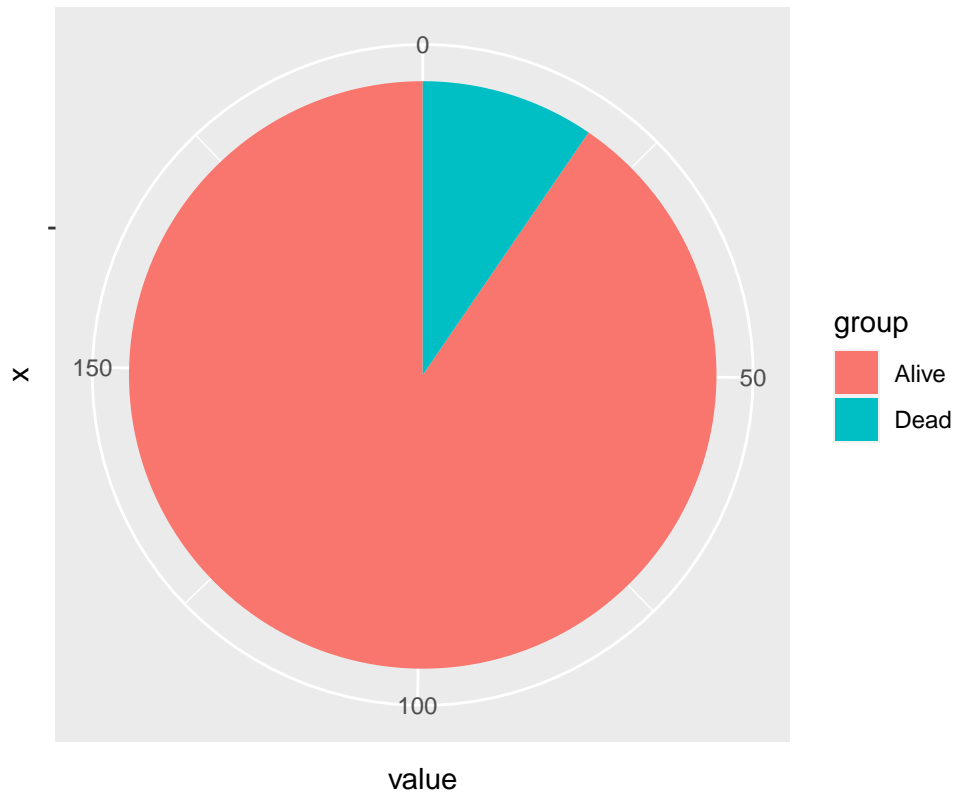
```
##      upper      mean      lower
## 0.049187343 0.027397260 0.005607177

mortality_ratio_and_ci(class3454,"34 54 age class")
```

34 54 age class smokers



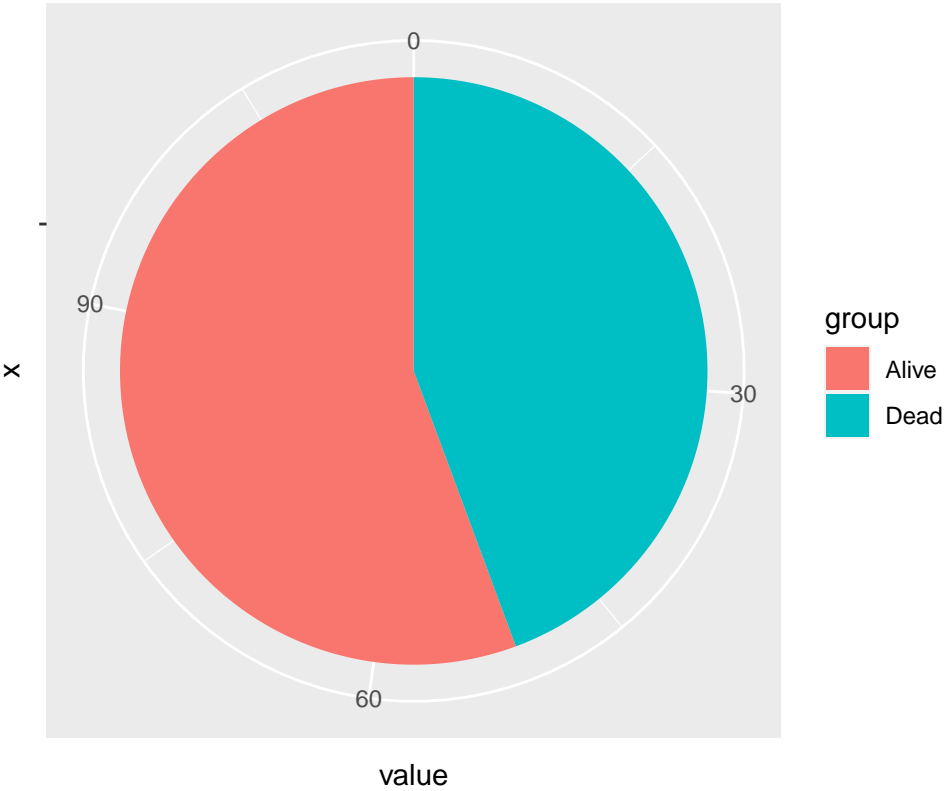
34 54 age class non smokers



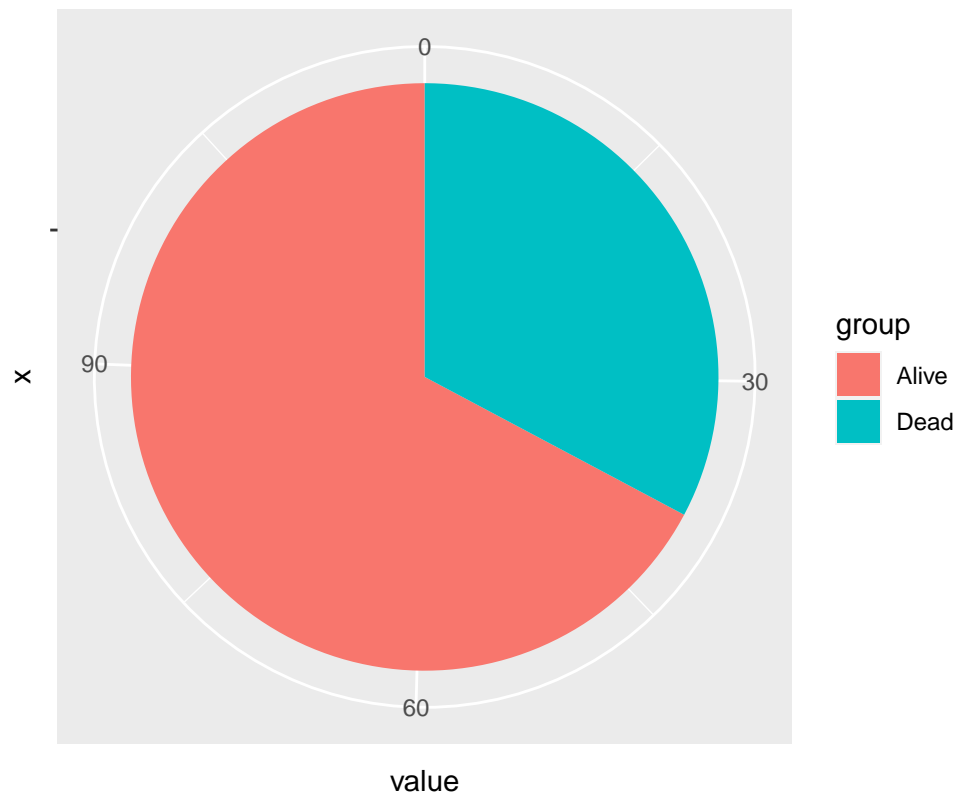
```
##      upper      mean      lower
## 0.13666230 0.09547739 0.05429248
```

```
mortality_ratio_and_ci(class5464, "54 64 age class")
```

54 64 age class smokers



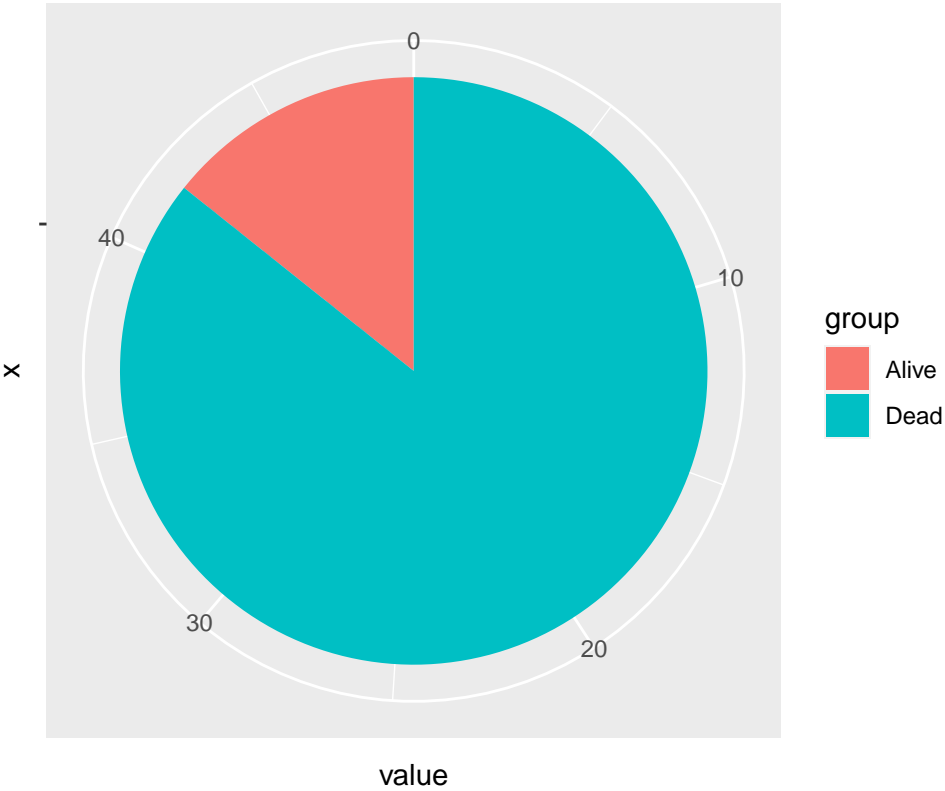
54 64 age class non smokers



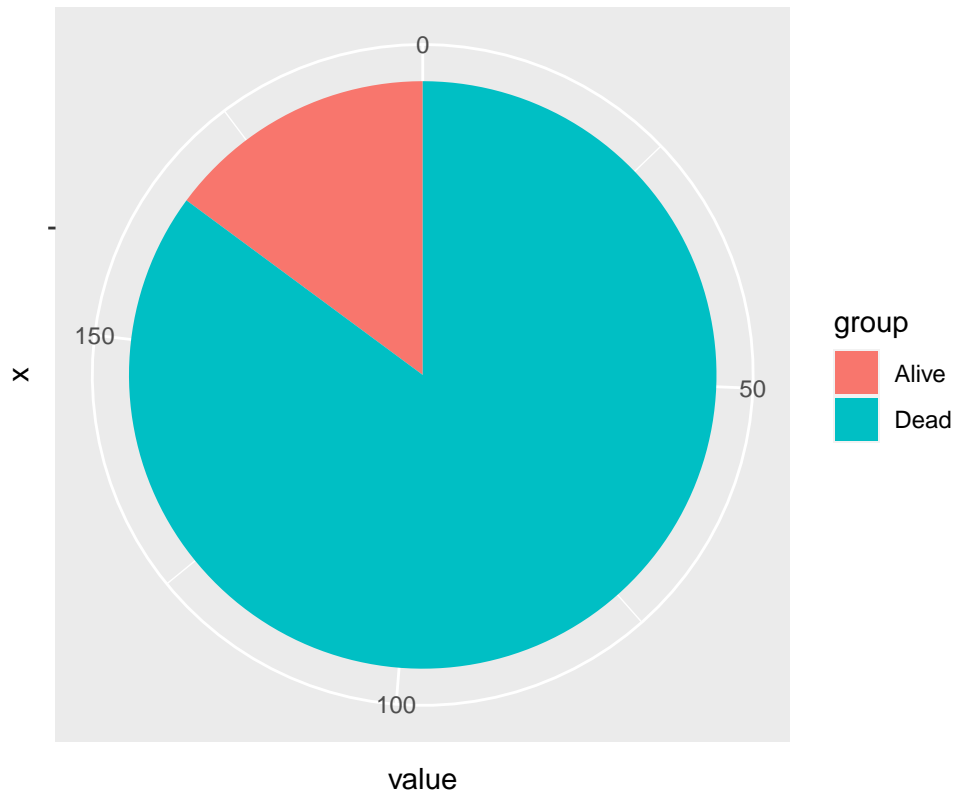
```
##      upper      mean      lower
## 0.4132997 0.3277311 0.2421625
```

```
mortality_ratio_and_ci(class64,"64 + age class")
```

64 + age class smokers



64 + age class non smokers



```
##      upper      mean      lower
## 0.9016650 0.8512821 0.8008991
```

This is very surprising, because, as we saw in question 1, the mortality rate was higher for *non smoker*. But now, after organizing the data in age classes, for every single class, the mortality is higher for the *smoker* group. So there is seemingly direct contradiction.

What is happening

To figure out what is happening, we need more information. Currently, we have the global death ratio for smoker and non smoker. We also have the death ratio for separate age class, and that's where we find a contradiction. An interesting information to have would be the ratio of smoker, especially in separate age class because we might find that people of different age have different smoking habit. I developed a function to compute to ratio of smoker.

```
# TODO confidence interval here as well ?
compute_smoker_ratio <- function(df_arg, title_arg){
  nb_smoker= df_arg %>% filter(Smoker== "Yes") %>% nrow()
  nb_non_smoker= df_arg %>% filter(Smoker== "No") %>% nrow()

  df <- data.frame(
    group = c("Smoker", "Non Smoker"),
    value = c(nb_smoker, nb_non_smoker)
```

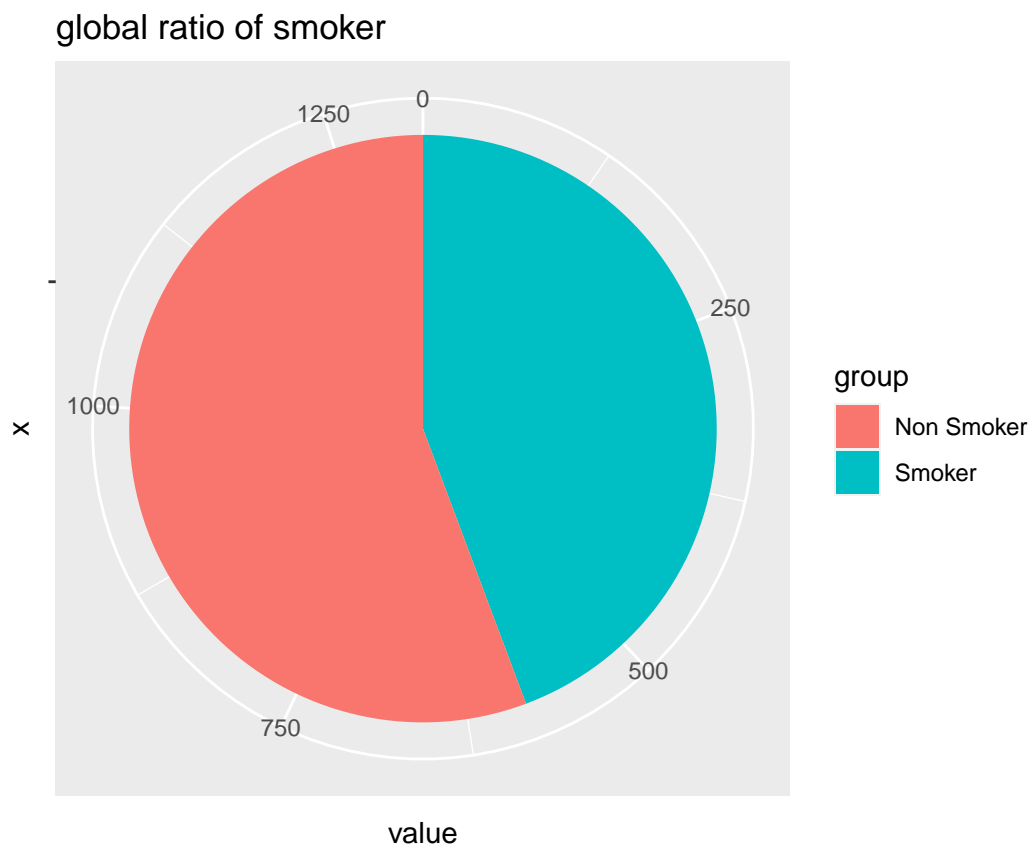
```

)

bp<- ggplot(df, aes(x="", y=value, fill=group))+
  geom_bar(width = 1, stat = "identity") + ggtitle(title_arg) +coord_polar("y", start=0)
print(bp)

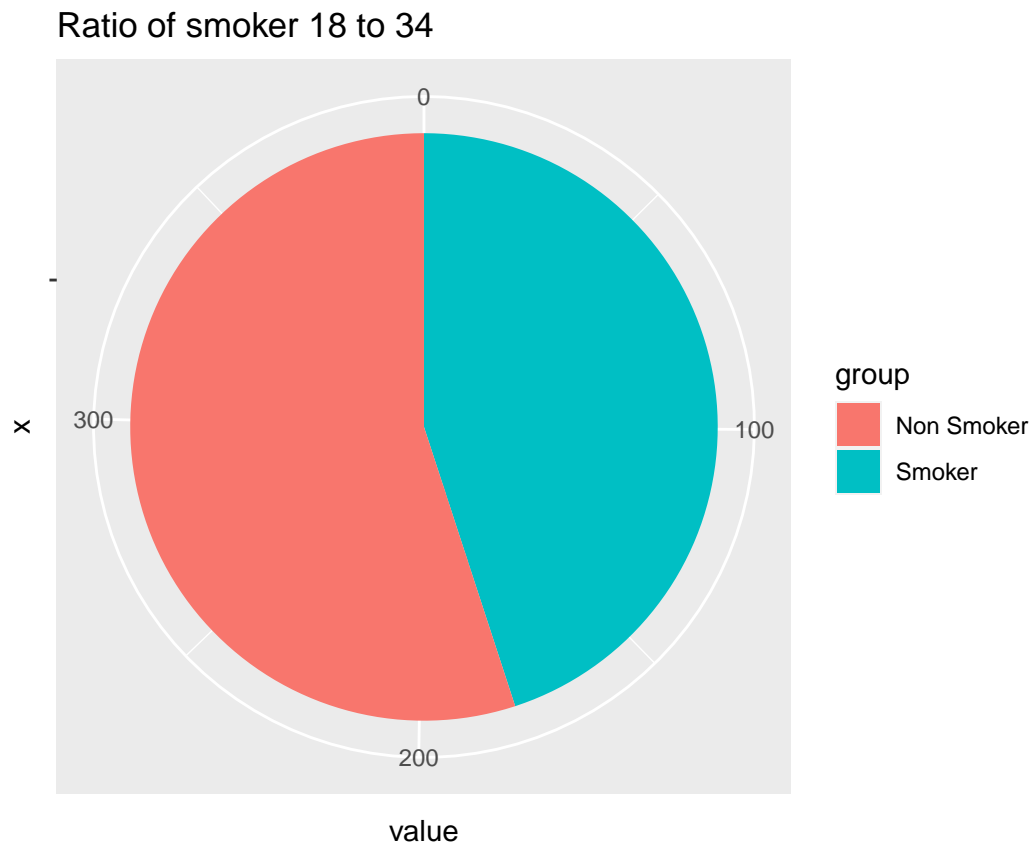
nb_smoker / (nb_smoker+nb_non_smoker) # divide the number of dead by the total, the total being the a
}
compute_smoker_ratio(df_arg = df, title_arg = "global ratio of smoker")

```



```
## [1] 0.4429224
```

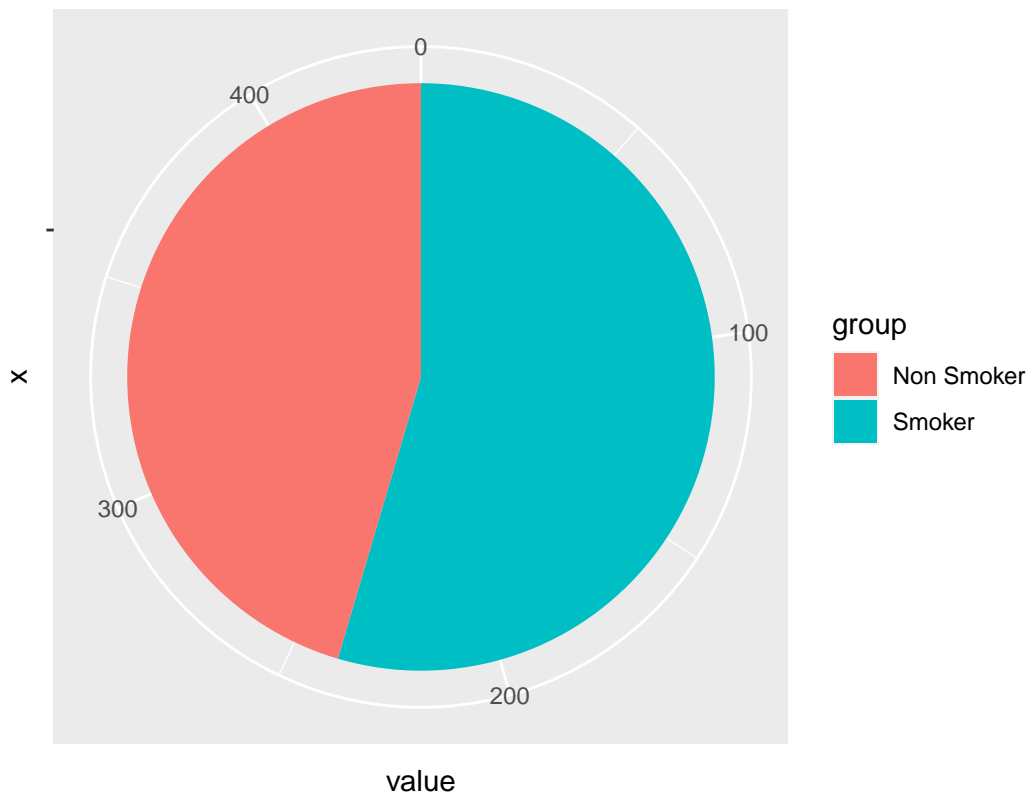
```
compute_smoker_ratio(df_arg = class1834, title_arg = "Ratio of smoker 18 to 34")
```



```
## [1] 0.4497487
```

```
compute_smoker_ratio(df_arg = class3454, title_arg = "Ratio of smoker 34 to 54")
```

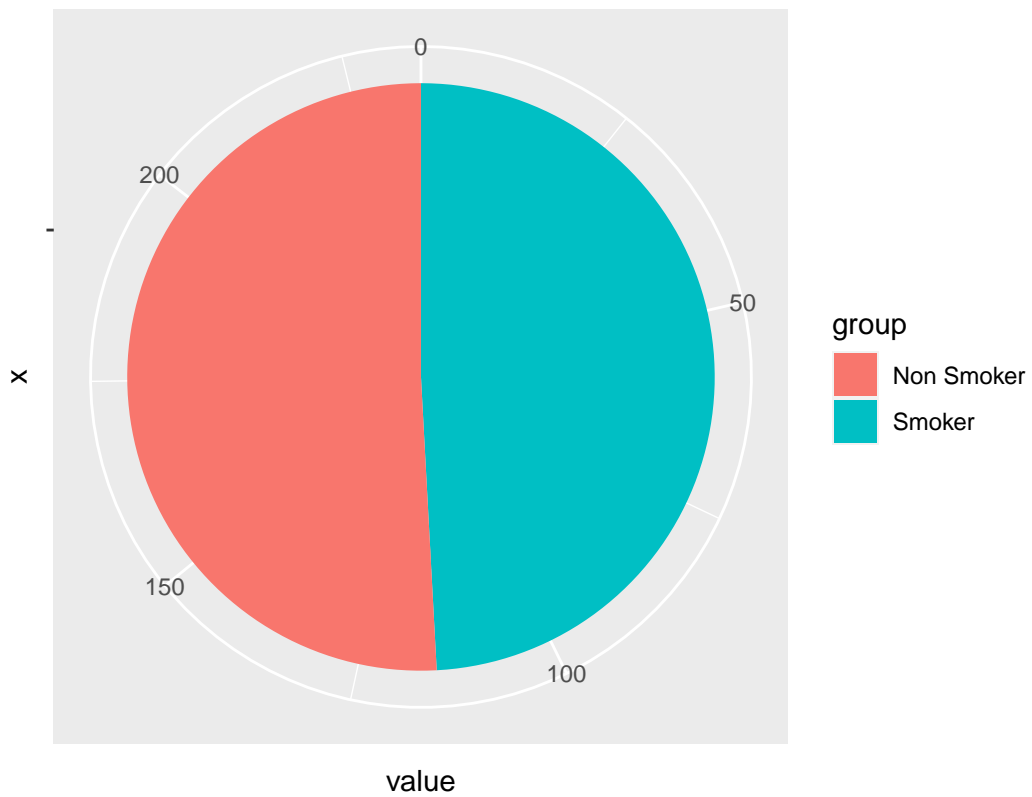
Ratio of smoker 34 to 54



```
## [1] 0.5456621
```

```
compute_smoker_ratio(df_arg = class5464, title_arg = "Ratio of smoker 54 to 64")
```

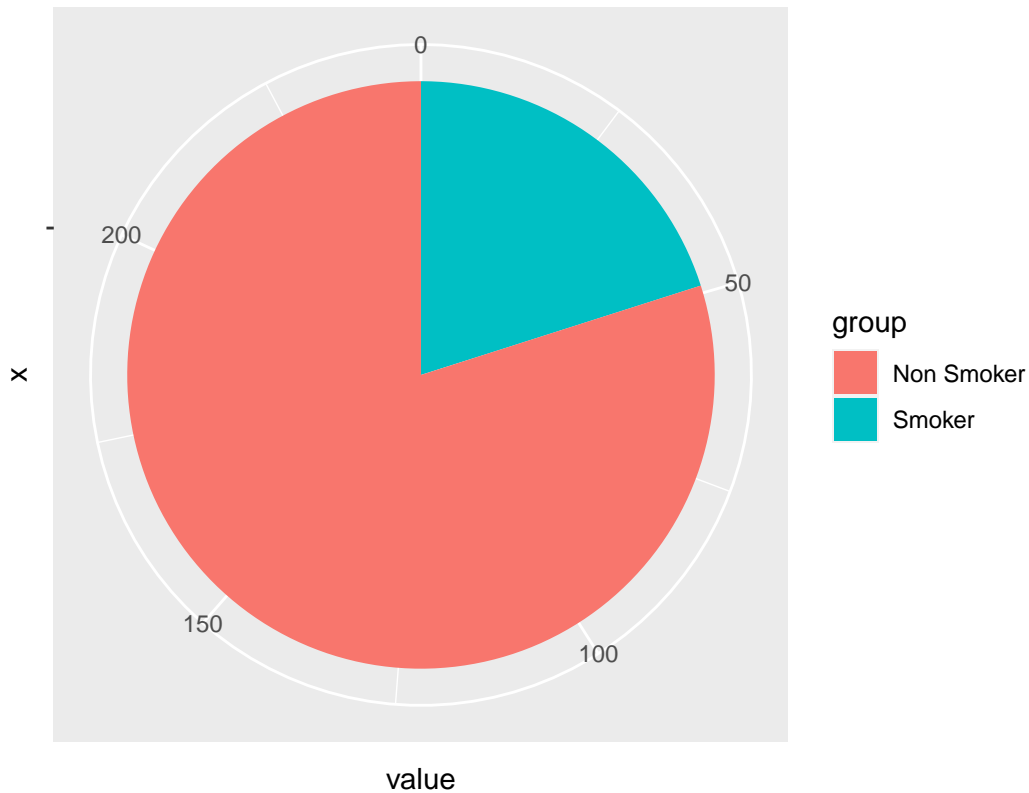
Ratio of smoker 54 to 64



```
## [1] 0.491453
```

```
compute_smoker_ratio(df_arg = class64,title_arg = "Ratio of smoker 64 +" )
```

Ratio of smoker 64 +



```
## [1] 0.2008197
```

It is immediately apparent that the oldest age class smokes much less than all the other age classes.

Given this, I can give a tentative explanation of what is going on: - Age is, as far as I know, the most important factor in increasing the probability of death, in the general population. - People of older age are therefore, everything else being equal, much more likely to die.

- In this study it is revealed that older people (in the 64 + age class), are much less likely to be smokers than the younger age classes.
- However, due to their age, they are much more likely to die than the younger age classes.

The age classes are not exactly the same size but they are all in the same order of magnitude.

And that's why we see this effect. Old people are much more likely to die in the first place.

The conclusion is that Age is more important to your health than smoking, even if smoking is still very important and should be avoided for a good health.

According to this [youtube video](#) that I watched recently, the same phenomena will most likely occur at some point with the covid 19 vaccine. The vaccine lowers the risk of death from covid 19 (lower risk of infection and lower risk of death when infected). But since older people are more vaccinated and much more likely to die in the first place, we might see a negative association between the vaccine and covid 19 protection. In other words, unvaccinated people will have less probability of dying of covid 19 than vaccinated people, even if the vaccine is effective in protecting people. That is, again, because older people are much more likely to die in the first place, and they are more vaccinated than younger people.

Question 3

Conclusion