

# Final Project

Ingrid J. Lu 260773949;

28/10/2021

```
# Set-up
df_fh <- readr::read_csv(here::here("fetal_health.csv")) %>%
  select(!starts_with("histogram"))
```

## Background

Your sister and her partner are expecting a child soon, and she just went to her obstetrician for her routine check. Because she is in her second trimester, her obstetrician asks her to do a fetal cardiogram. The results will not get back to her until a week later. Your sister is a bit of a hypochondriac, so she is afraid that there is something wrong with her child.

So, her partner turns to you, a Masters student who is studying inferential statistics for health, hoping to find some results to calm your sister down.

Luckily, you have identified a publicly available dataset that includes thousands of fetal cardiogram results, and the classification of these babies' health status. You need to convince your sister that she and her child will be safe.

Note: the objective of this exercise is to consolidate all the important concepts covered in EPIB607. When answer each question, be sure to include any units and assumptions and define all parameters, when appropriate. The following questions are based on the publicly available dataset “Fetal Classification”, please find all attribute information of the data from the link.

## Question 1 Data Visualization and Summary Statistics

a)

Is this data tidy? If no, transform it into an untidy data. If yes, provide an explanation.

*Solution:*

b)

Looking at the 3 different classification of fetal health status and each fetus' baseline heart rate provide an appropriate graphic summarizing the distribution of each of baseline heart rates for each class. Be sure to provide the correct title and label for the plot.

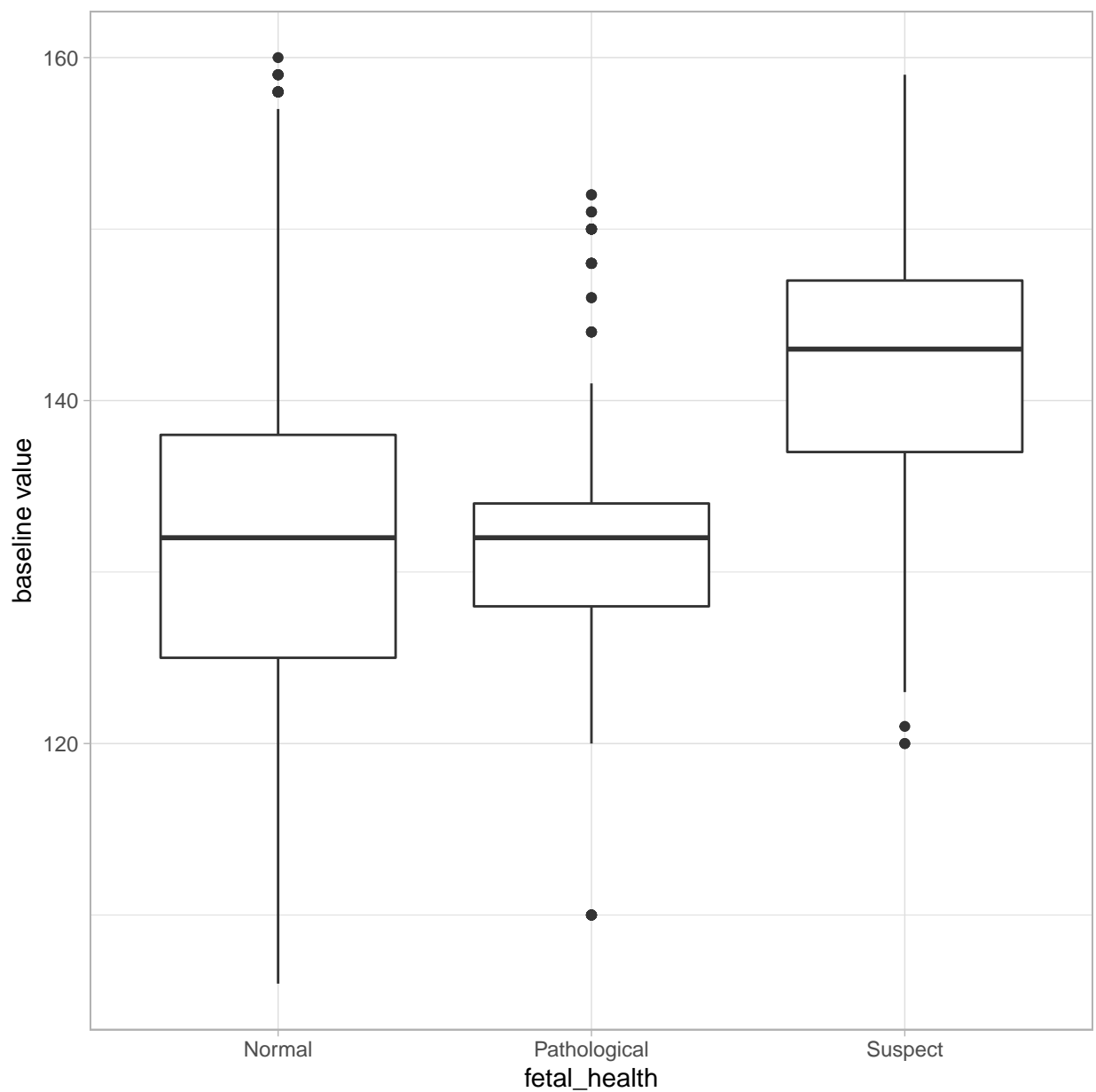
*Solution:*

*#----Question 1-----*

*# 1.b*

```
df_fh$fetal_health <- ifelse(df_fh$fetal_health == 1,  
                             "Normal",  
                             ifelse(df_fh$fetal_health == 2,  
                                     "Suspect",  
                                     "Pathological"))
```

```
df_fh %>%  
  group_by(fetal_health) %>%  
  ggplot(aes(x = fetal_health, y = `baseline value`)) +  
  geom_boxplot()
```



c)

Comment on the boxplot, what are the characteristics of each category?

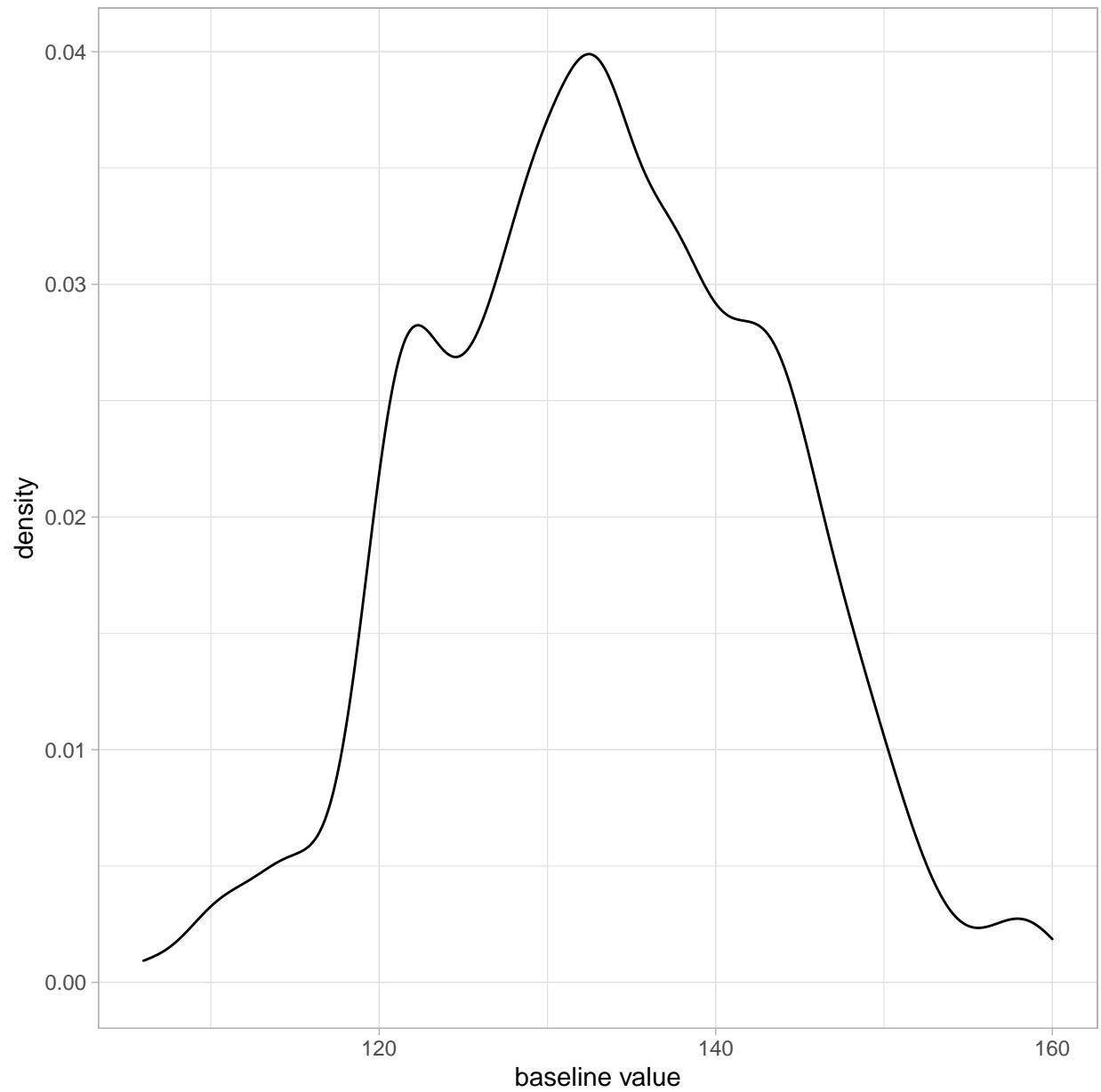
d)

Describe the distribution of the baseline heart rate for all participants in this sample, is the baseline heart rate normally distributed, comment on any skewness. What about each class? Use an appropriate graph to answer this question.

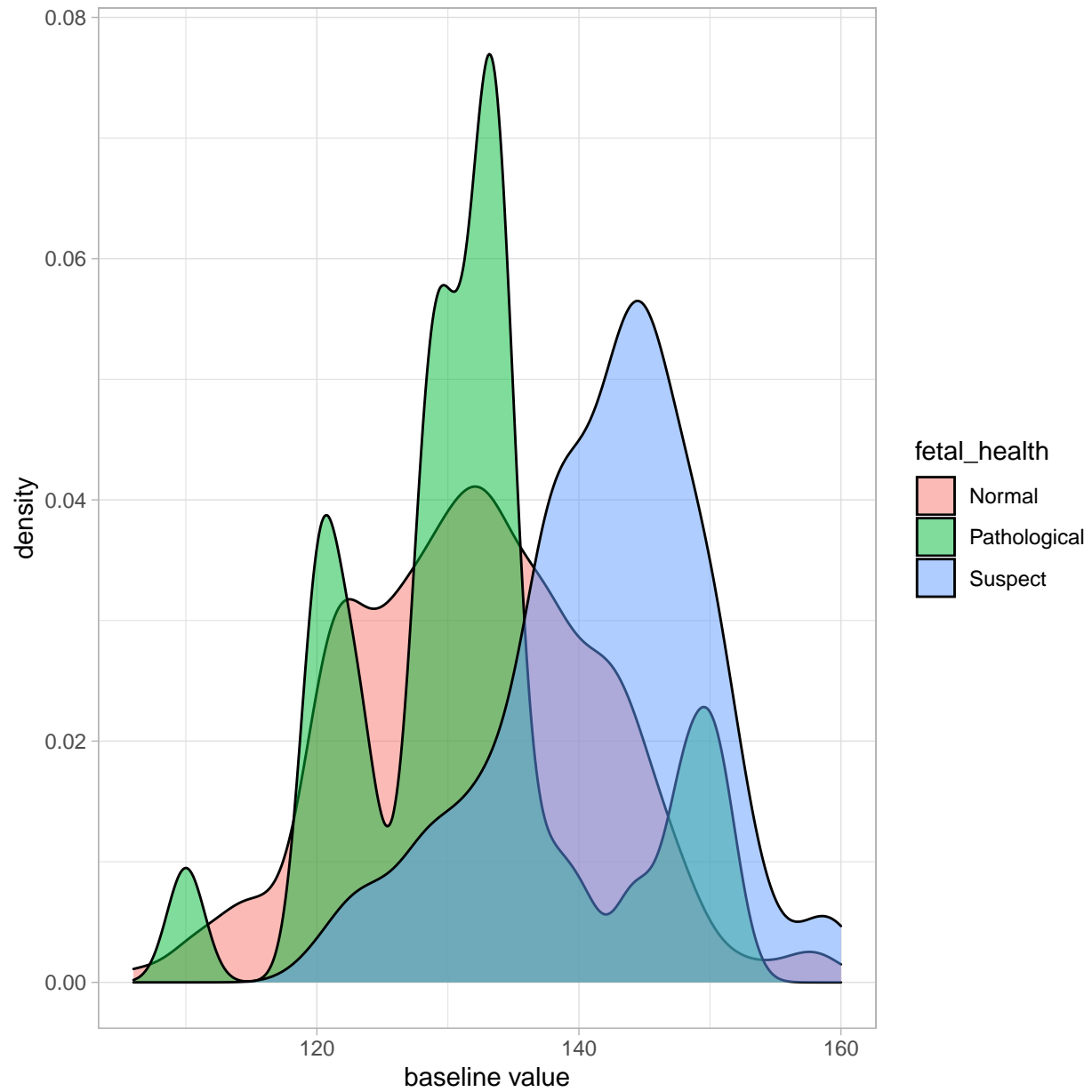
*Solution:*

```
## 1.d

df_fh %>%
  ggplot(aes(x = `baseline value`))+
  geom_density()
```



```
df_fh %>%  
  ggplot(aes(x = `baseline value`, fill = fetal_health))+  
  geom_density(alpha = 0.5)
```



## Question 2

a)

Your sister thinks that more uterine contraction means that her fetus is unhealthy and she is more likely to give a pre-term birth.

Is this true? What is the median uterine contraction for each class?

*Solution:*

```
#----Question 2-----
df_fh_n <- df_fh %>% filter(fetal_health == "Normal")
median(df_fh_n$uterine_contractions)
```

```
## [1] 0.005
```

```
df_fh_s <- df_fh %>% filter(fetal_health == "Suspect")
median(df_fh_s$uterine_contractions)
```

```
## [1] 0.001
```

```
df_fh_p <- df_fh %>% filter(fetal_health == "Pathological")
median(df_fh_p$uterine_contractions)
```

```
## [1] 0.003
```

b)

Since we have a small sample size for those who are suspected to be pathological and those who are determined to be pathological, what is one method that we can use to artificially create a pseudo-population and calculate the median of these two groups? State your assumptions.

*Solution:* Bootstrapping

c)

What is one weakness of using the bootstrapping method?

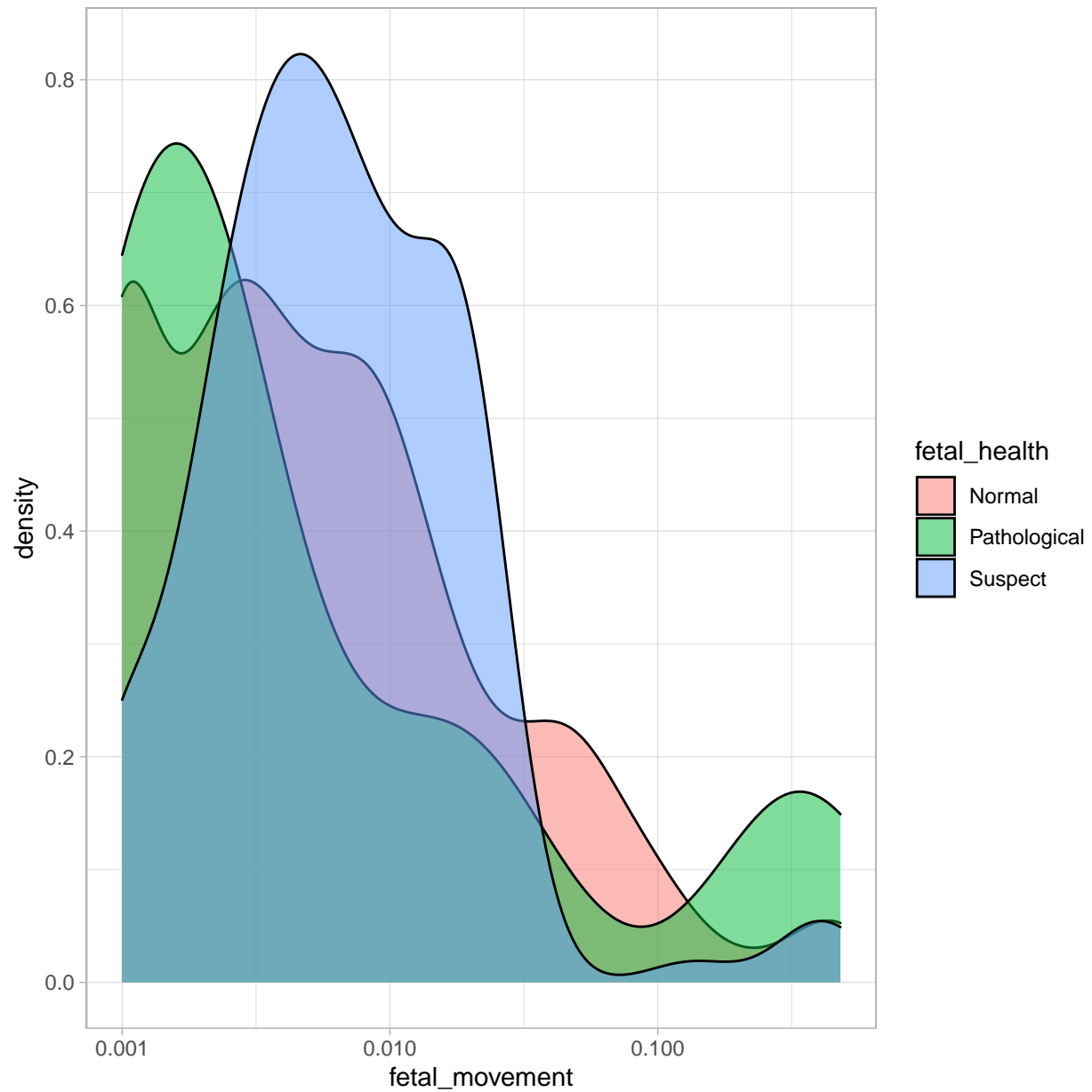
*Solution:* The bootstrapping method assumes that the sample is representative of the whole population. However, when we have a small sample size, that is not always the case. There could be sampling errors which can skew the distribution of the sample. And the bootstrapping method cannot correct this error as it can only sample within this small set of data. As the saying goes, “garbage in, garbage out”. If the sample is not representative of the population, the results obtained from bootstrap will not be accurate and will result in an incorrect inference of the population parameter.

d)

Your sister is also concerned that her baby does not move as much, which could be a sign of unhealthy pregnancy.

Transform x-axis to a logarithmic scale.

```
# 2.d
df_fh %>%
  ggplot(aes(x = fetal_movement, fill = fetal_health))+
  geom_density(alpha = 0.5) +
  scale_x_log10()
```



# Question 3

# Question 4

## Question 5

Suppose this data represents the population of newborns in one hospital and you take a simple random sample of 100 babies from the population.

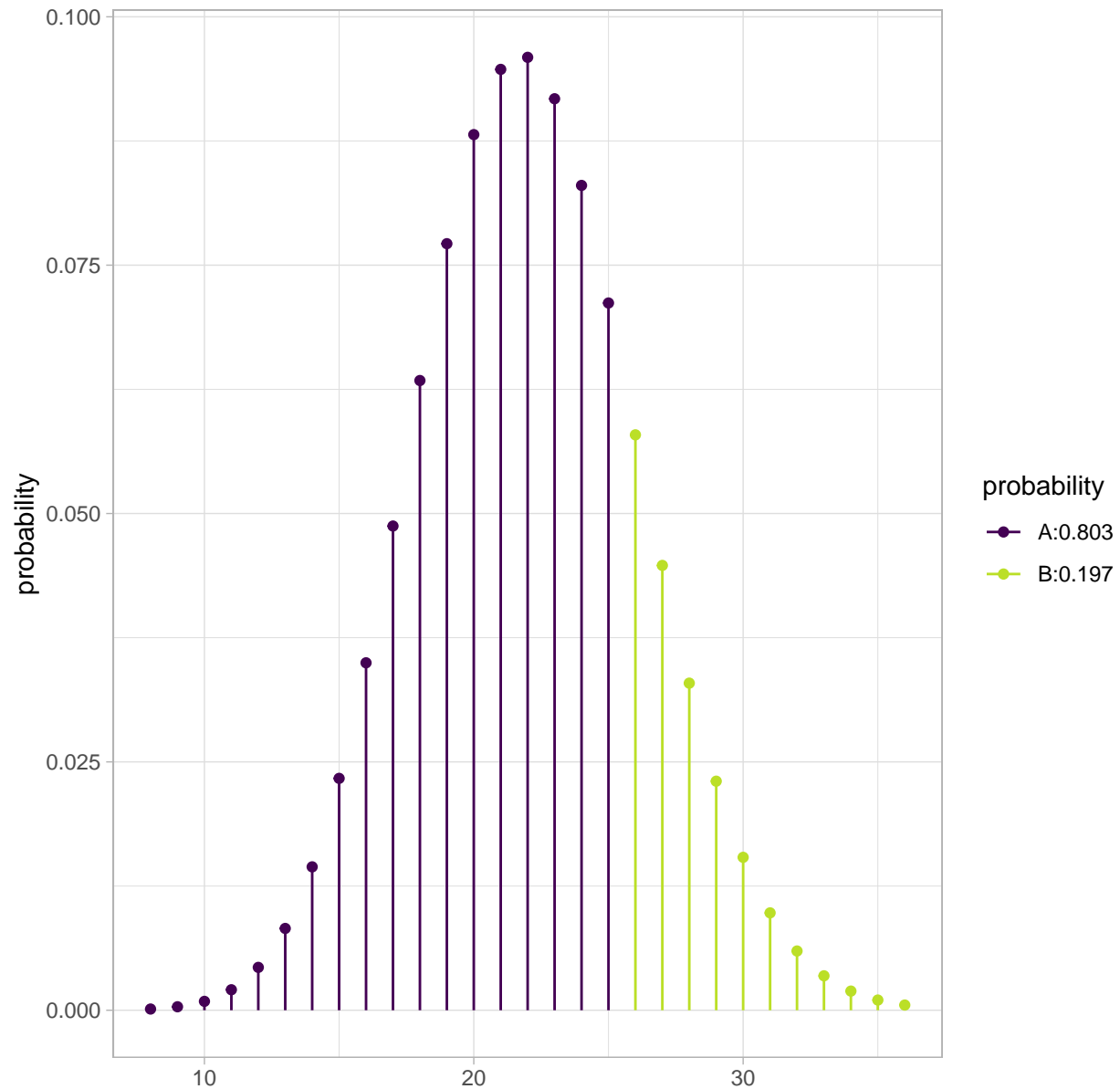
a) What is the probability that your sample contains more than 25 babies has abnormal health status (1 = health, 2 = suspect, 3 = Pathological)?

```
## ---- Question-5a -----  
df_fh %>%  
  select(fetal_health)%>%  
  filter(fetal_health != 1)%>%  
  nrow()
```

```
## [1] 2126
```

```
#there are 471 abnormal in this population  
#the probability of having abnormal is 471/2126 = 0.22  
1 - mosaic::xpbinom(q = 25, size = 100, prob = 0.22)
```





```
## [1] 0.1972269
```

The probability of having more than 30 abnormal is 0.197.

b) Turns out that your sample actually contains 20 babies with abnormal health status. What is the 95% confidence interval of this proportion? Can you use a normal approximation for this sample? Why or why not?

```
## ---- Question-5b -----
mosaic::binom.test(x = 20, n = 100, ci.method = "Clopper-Pearson")
```

```
##
##
##
## data: 200 out of 1000
## number of successes = 200, number of trials = 1000, p-value < 2.2e-16
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
## 0.1756206 0.2261594
## sample estimates:
## probability of success
## 0.2
```

```
mosaic::binom.test(x = 200, n = 1000, ci.method = "Wald")
```

```
##
## Exact binomial test (Wald CI)
##
## data: 200 out of 1000
## number of successes = 200, number of trials = 1000, p-value < 2.2e-16
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
## 0.1752082 0.2247918
## sample estimates:
## probability of success
## 0.2
```

```
mosaic::binom.test(x = 20, n = 100, ci.method = "Clopper-Pearson")
```

```
##
##
##
## data: 20 out of 100
## number of successes = 20, number of trials = 100, p-value = 1.116e-09
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
## 0.1266556 0.2918427
## sample estimates:
## probability of success
## 0.2
```

```
mosaic::binom.test(x = 20, n = 100, ci.method = "Wald")
```

```
##
## Exact binomial test (Wald CI)
##
## data: 20 out of 100
## number of successes = 20, number of trials = 100, p-value = 1.116e-09
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
## 0.1216014 0.2783986
## sample estimates:
## probability of success
## 0.2
```

The 95% CI using Clopper-Pearson method: [0.127,0.292] The 95% CI using exact method(ie: normal approximation): [0.122,0.278] The two methods give similar 95% CIs. Normal approximation can be used here since the sample size is large enough to generate a binomial distribution approximating the normal distribution and for CLT to kick in.

c) Another sample taken have the same proportion of event but the sample size is now only 10 and the count of abnormal is 2. Calculate the 95% CI using this sample and compare it with the one you have in b). Describe their difference and the reason why.

```
## ---- Question-5c -----
mosaic::binom.test(x = 2, n = 10, ci.method = "Clopper-Pearson")

##
##
##
## data:  2 out of 10
## number of successes = 2, number of trials = 10, p-value = 0.1094
## alternative hypothesis: true probability of success is not equal to 0.5
## 95 percent confidence interval:
##  0.02521073 0.55609546
## sample estimates:
## probability of success
##                                0.2
```

The 95% CI for this sample is: [0.0252, 0.556]. 95% CI in c) with sample size = 10 is wider than the 95% CI in b) with sample size = 100. The sample size is different so the standard error is different. The 95% CI is calculated using the formula:

$$\bar{y} - 1.96\left(\frac{\sigma}{\sqrt{n}}\right), \bar{y} + 1.96\left(\frac{\sigma}{\sqrt{n}}\right)$$

If the sample size n is larger, the standard error( $\sigma/\sqrt{n}$ ) is smaller. The value 1.96\*standard error is also smaller, and results in a narrower confident interval.

*# Question 6*