

# Georgia Birth Weight EDA

Alejandro Hernandez

2024-04-27

```
birthwt <- read.csv("data/cdc birthwt.csv")

names(birthwt) <- c("mother.id", "birth.order", "birth.weight", "maternal.age", "child.id")

# create maternal age factor
breaks <- c(min(birthwt$maternal.age), 20, 30, 40,
             max(birthwt$maternal.age))
birthwt <- birthwt %>%
  mutate(maternal.age.factor = cut(maternal.age, breaks=breaks,
                                   include.lowest = T))
```

## Validation

```
# validate each mother has 5 children
birthwt %>%
  group_by(mother.id) %>%
  summarize(n = length(birth.order)) %>%
  all(.$n == 5)
```

```
## [1] TRUE
```

```
# summarize numeric variables
birthwt %>%
  select(birth.order, birth.weight, maternal.age) %>%
  mutate_at("birth.order", as.factor) %>%
  summary
```

```
## birth.order birth.weight maternal.age
## 1:878      Min.   : 312   Min.   :12.00
## 2:878      1st Qu.:2850  1st Qu.:18.00
## 3:878      Median :3175  Median :21.00
## 4:878      Mean   :3156  Mean   :21.65
## 5:878      3rd Qu.:3515  3rd Qu.:24.00
##          Max.   :5528  Max.   :42.00
```

# Exploratory data analysis

## Birth weight

**Birth weight** is an infant's weight that is optimally measured in the hours following birth.

The World Health Organization (WHO) defines low birth weight as below 2500 g (5.5 lbs / 5 lbs 8 oz) and very low birth weight as below 1500 g (3.3 lbs / 3 lbs 5 oz). An infant with a low birth weight is often also premature, and at greater risk of health complications (e.g., underdeveloped lungs, inability to maintain body temperature, difficulty gaining weight, intestinal disease, bleeding in the brain, and sudden death).

Factors that are believed to cause low birth weight include maternal health and age, and multiple-baby pregnancies. Effects from race are also believed to impact birth weight: among Americans, Black mothers are twice as likely as white mothers to have low birth weights.

*In some cases, the validity of birth weight data is of concern, as measures may be taken days after birth, after which significant weight loss may have occurred. We assume the source of our data to be valid.*

1. What are the quartiles and average birth weight among all mothers? What quantile defines the threshold of low birth weights? How many births fall into each quartile and each category of (low / not low)?
2. Describe the distribution of all birth weights.
3. Describe the distribution of each mother's average birth weight. How does this distribution compare to that of all birth weights?

```
#####  
# Birth weight  
#####  
  
# numeric summary of birth weight  
round(c(mean = mean(birthwt$birth.weight),  
        sd = sd(birthwt$birth.weight),  
        IQR = IQR(birthwt$birth.weight),  
        range = diff(range(birthwt$birth.weight)),  
        quantile(birthwt$birth.weight)), 2)
```

```
##      mean      sd      IQR    range      0%      25%      50%      75%      100%  
## 3156.30  570.44  665.00 5216.00  312.00 2850.00 3175.00 3515.00 5528.00
```

```
# numeric summary of individual average birth weight  
birthwt %>%  
  group_by(mother.id) %>%  
  summarize(avg.weight = mean(birth.weight)) %>%  
  reframe(c(mean = mean(avg.weight),  
            sd = sd(avg.weight),  
            IQR = IQR(avg.weight),  
            quantile(avg.weight))) %>%  
  round(2) %>% as.list
```

```
## $'c(...)'  
##      mean      sd      IQR      0%      25%      50%      75%      100%  
## 3156.30  416.94  532.65 1690.40 2891.80 3166.80 3424.45 4745.80
```

```

# size of each quartile
birthwt %>%
  mutate(quartile = cut(birth.weight,
                        breaks = quantile(birthwt$birth.weight),
                        include.lowest = T)) %>%
  group_by(quartile) %>%
  summarize(n = length(quartile))

## # A tibble: 4 x 2
##   quartile      n
##   <fct>      <int>
## 1 [312,2.85e+03] 1102
## 2 (2.85e+03,3.18e+03] 1099
## 3 (3.18e+03,3.52e+03] 1113
## 4 (3.52e+03,5.53e+03] 1076

# number of low births
birthwt %>%
  mutate(low.weight = ifelse(birth.weight < 2500,
                             ifelse(birth.weight < 1500, "Very low", "Low"),
                             "Not low"),
         low.weight = factor(low.weight,
                             levels = c("Not low", "Low", "Very low"))) %>%
  group_by(low.weight) %>%
  summarize(n = length(low.weight),
            prop = round(n / nrow(.), 2))

## # A tibble: 3 x 3
##   low.weight      n prop
##   <fct>      <int> <dbl>
## 1 Not low      3941  0.9
## 2 Low          391  0.09
## 3 Very low      58  0.01

# distribution of all birth weights
gg_bw_all <- birthwt %>%
  ggplot(aes(x = birth.weight, y = after_stat(density))) +
  geom_histogram(bins = 30) +
  xlab("Birth weight") + ylab("") +
  geom_vline(xintercept = c(1500, 2500), color = c("red", "blue")) +
  geom_vline(xintercept = quantile(birthwt$birth.weight), color = "grey") +
  theme_bw()

# distribution of mothers' average birth weight
# must compute quartiles separately
quartiles <- birthwt %>%
  group_by(mother.id) %>%
  summarize(avg.weight = mean(birth.weight)) %>%
  reframe(quantile(avg.weight)) %>% as.list

gg_bw_indv <- birthwt %>%
  group_by(mother.id) %>%

```

```

summarize(avg.weight = mean(birth.weight)) %>%
ggplot(aes(x = avg.weight, y = after_stat(density))) +
  geom_histogram(bins = 30) +
  xlab("Average birth weight") + ylab("") +
  labs(caption = "Quartiles marked by vertical grey lines") +
  geom_vline(xintercept = c(1500, 2500), color = c("red", "blue")) +
  geom_vline(xintercept = quartiles[[1]], color = "grey") +
  theme_bw()

gridExtra::grid.arrange(gg_bw_all, gg_bw_indv)

```

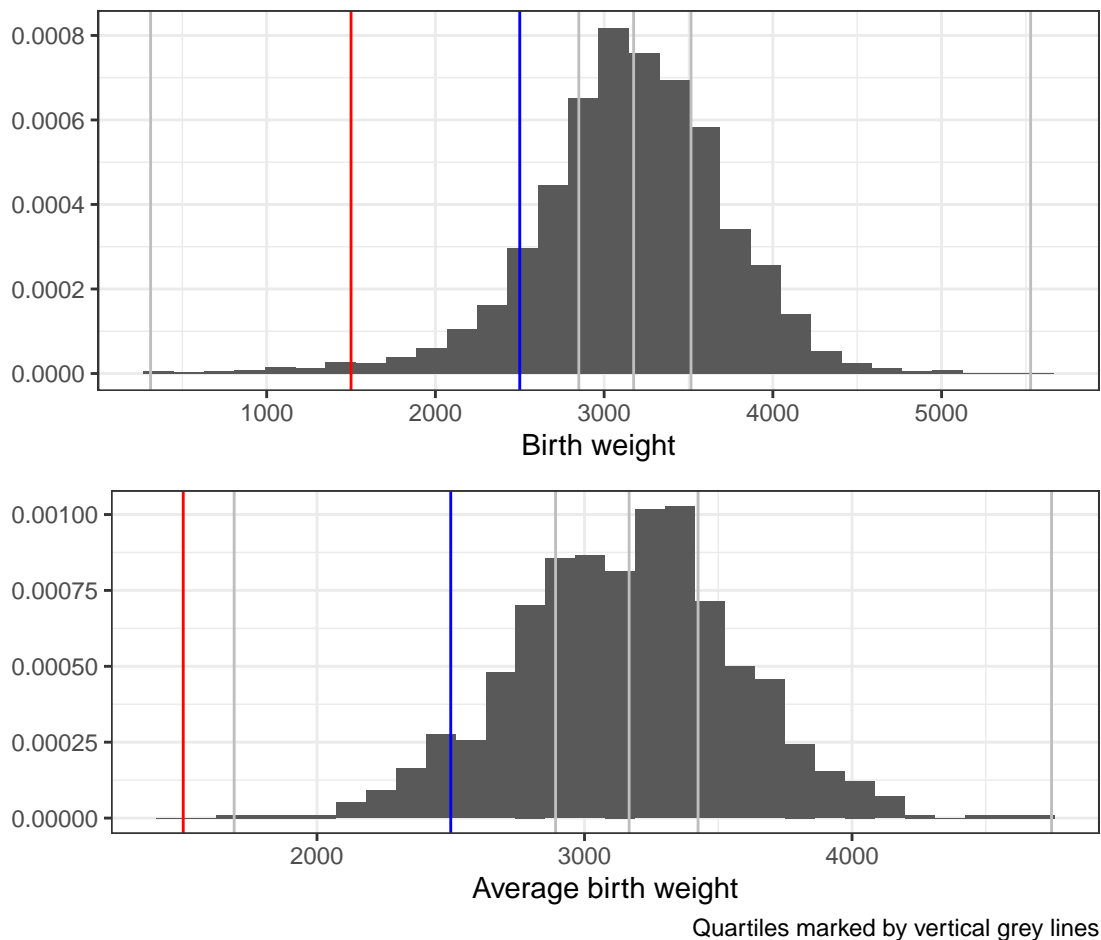


Figure 1: Distribution of birth weight with thresholds of low (blue) and very low (red) birth weight

## Maternal age

**Maternal age** is the age at which a mother gives birth. Previous studies suggest that maternal age is associated with birth weight (younger and older mothers have greater rates of preterm births).

1. What are the quartiles and average maternal age among all mothers? What quantile defines the thresholds of middle aged? How many births fall into each quantile and each category of (young / middle age)?

2. Describe the distribution of all maternal ages.
3. Describe the distribution of each mother's average maternal age. How does this distribution compare to that of all maternal ages?

```
#####
# Maternal age
#####

# numeric summary of maternal age
round(c(mean = mean(birthwt$maternal.age),
  sd = sd(birthwt$maternal.age),
  IQR = IQR(birthwt$maternal.age),
  range = diff(range(birthwt$maternal.age)),
  quantile(birthwt$maternal.age)), 2)
```

```
## mean    sd    IQR range    0%    25%    50%    75%    100%
## 21.65   4.63   6.00 30.00 12.00 18.00 21.00 24.00 42.00
```

```
# numeric summary of individual average maternal age
birthwt %>%
  group_by(mother.id) %>%
  summarize(avg.age = mean(maternal.age)) %>%
  reframe(c(mean = mean(avg.age),
    sd = sd(avg.age),
    IQR = IQR(avg.age),
    quantile(avg.age))) %>%
  round(2) %>% as.list
```

```
## $'c(...)'
## mean    sd    IQR    0%    25%    50%    75%    100%
## 21.65   3.69   4.20 15.40 19.00 20.80 23.20 38.20
```

```
# size of each quartile
birthwt %>%
  mutate(quartile = cut(maternal.age,
    breaks = quantile(birthwt$maternal.age),
    include.lowest = T)) %>%
  group_by(quartile) %>%
  summarize(n = length(quartile))
```

```
## # A tibble: 4 x 2
##   quartile      n
##   <fct>    <int>
## 1 [12,18]   1212
## 2 (18,21]  1193
## 3 (21,24]   958
## 4 (24,42]  1027
```

```
# condensed version of above code (this has worse readability):
# table(cut(birthwt$maternal.age,
#           breaks = quantile(birthwt$maternal.age),
```

```
#           include.lowest = T))

# size of each age group
birthwt %>%
  mutate(age.group = cut(maternal.age,
                        breaks = c(0, 18, 25, 35, 45),
                        right = F)) %>%
  group_by(age.group) %>%
  summarize(n = length(age.group),
            prop = round(n / nrow(), 2))
```

```
## # A tibble: 4 x 3
##   age.group      n prop
##   <fct>      <int> <dbl>
## 1 [0,18)       790  0.18
## 2 [18,25)     2573  0.59
## 3 [25,35)      960  0.22
## 4 [35,45)       67  0.02
```

```
# distribution of maternal age
birthwt %>%
  ggplot(aes(x = maternal.age, y = ..density..)) +
  geom_histogram()
```

```
## Warning: The dot-dot notation ('..density..') was deprecated in ggplot2 3.4.0.
## i Please use 'after_stat(density)' instead.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```

```
# distribution of all maternal ages
gg_ma_all <- birthwt %>%
  ggplot(aes(x = maternal.age, y = after_stat(density))) +
  geom_histogram(bins = 30) +
  xlab("Maternal age") + ylab("") +
  scale_x_continuous(breaks = seq(15, 40, by = 5)) +
  geom_vline(xintercept = quantile(birthwt$maternal.age), color = "grey") +
  theme_bw()
```

```
# distribution of mothers' average maternal age
# must compute quartiles separately
quartiles <- birthwt %>%
  group_by(mother.id) %>%
  summarize(avg.age = mean(maternal.age)) %>%
  reframe(quantile(avg.age)) %>% as.list
```

```
gg_ma_indv <- birthwt %>%
  group_by(mother.id) %>%
  summarize(avg.age = mean(maternal.age)) %>%
```

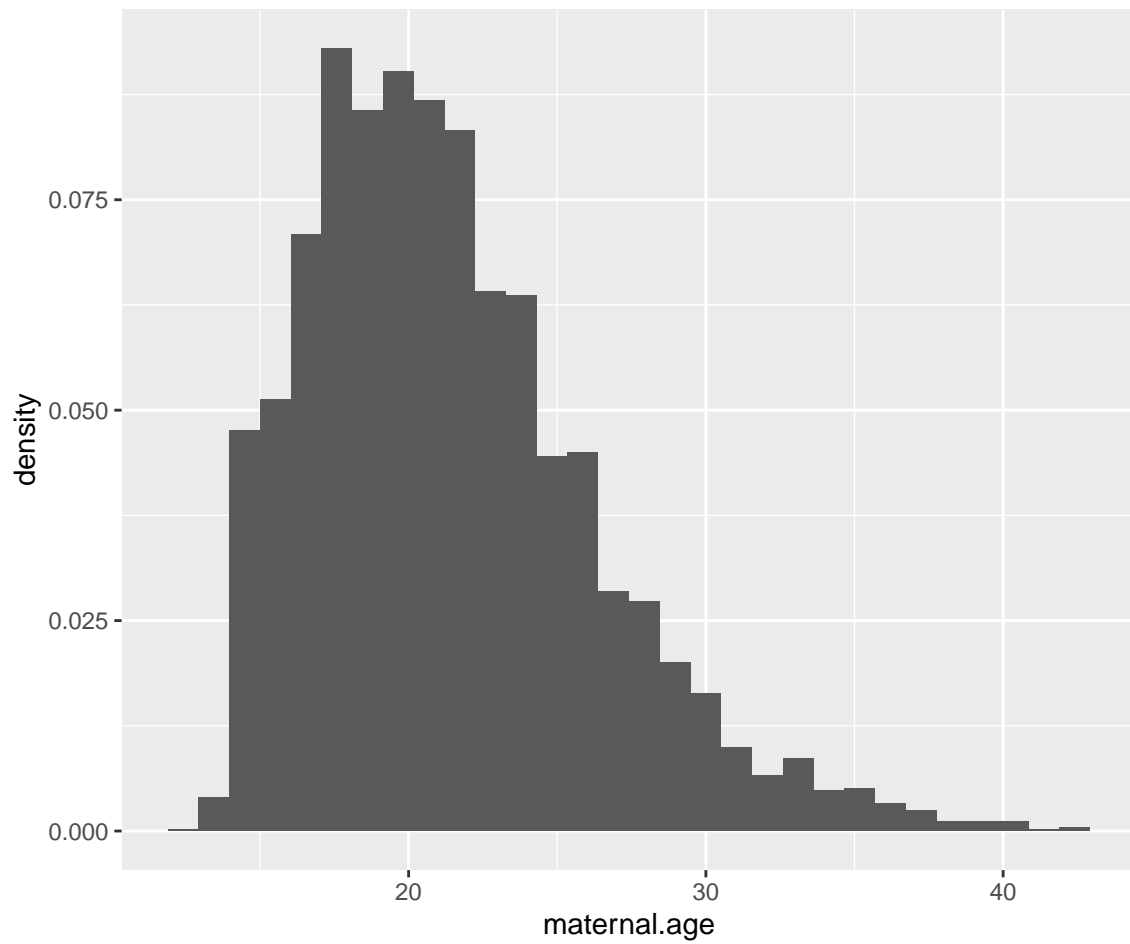


Figure 2: Distribution of maternal age

```
ggplot(aes(x = avg.age, y = after_stat(density))) +
  geom_histogram(bins = 30) +
  xlab("Average maternal age") + ylab("") +
  labs(caption = "Quartiles marked by vertical grey lines") +
  geom_vline(xintercept = quartiles[[1]], color = "grey") +
  theme_bw()

gridExtra::grid.arrange(gg_ma_all, gg_ma_indv)
```

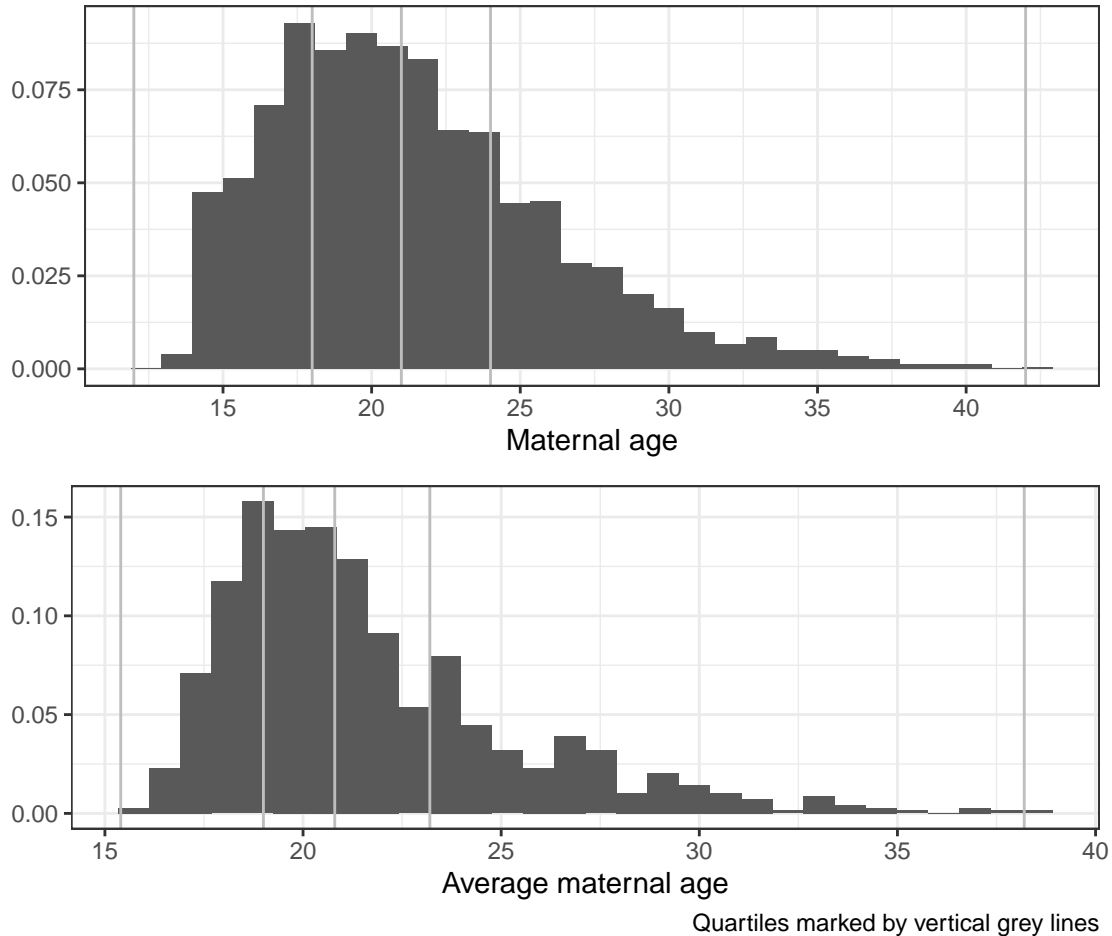


Figure 3: Distribution of maternal age

## Interpregnancy interval

**Interpregnancy interval** is the time elapsed between one child's birth and the subsequent child's conception, often measured in months. This interval can be brief: mothers may become pregnant as early as four weeks after delivery.<sup>4</sup> Similar to maternal age, extremely short and long intervals are associated with health risks for the mother and second-born child. Interpregnancy intervals less than 18 months introduce moderate risk to children and significant risk is associated with intervals shorter than 6 months; intervals greater than 5-10 years is associated with increased risk of adverse health outcomes for both mother and child.<sup>4</sup>



With available data, our best estimate of this interval is the difference in maternal age between two subsequent births, measured in years, which in certain cases may be a slight underestimate.

1. What are the group sizes and mode for all interval values?
2. Describe the distribution of all intervals.
3. Describe the distribution of each mother's interval mode. How does this distribution compare to that of all intervals?

*Interdelivery interval is the time elapsed between two subsequent births. May we more accurately estimate this interval using our data? Not necessarily- the accuracy of using the difference in maternal ages to estimate an interpregnancy versus an interdelivery interval varies by situation and relies on information unavailable to us. We have no workaround for this issue, other than to define short and long intervals according to respective thresholds (which differ by ?? 9 months?) and examine the extent to which results agree.*

```
#####
# Interpregnancy interval
#####

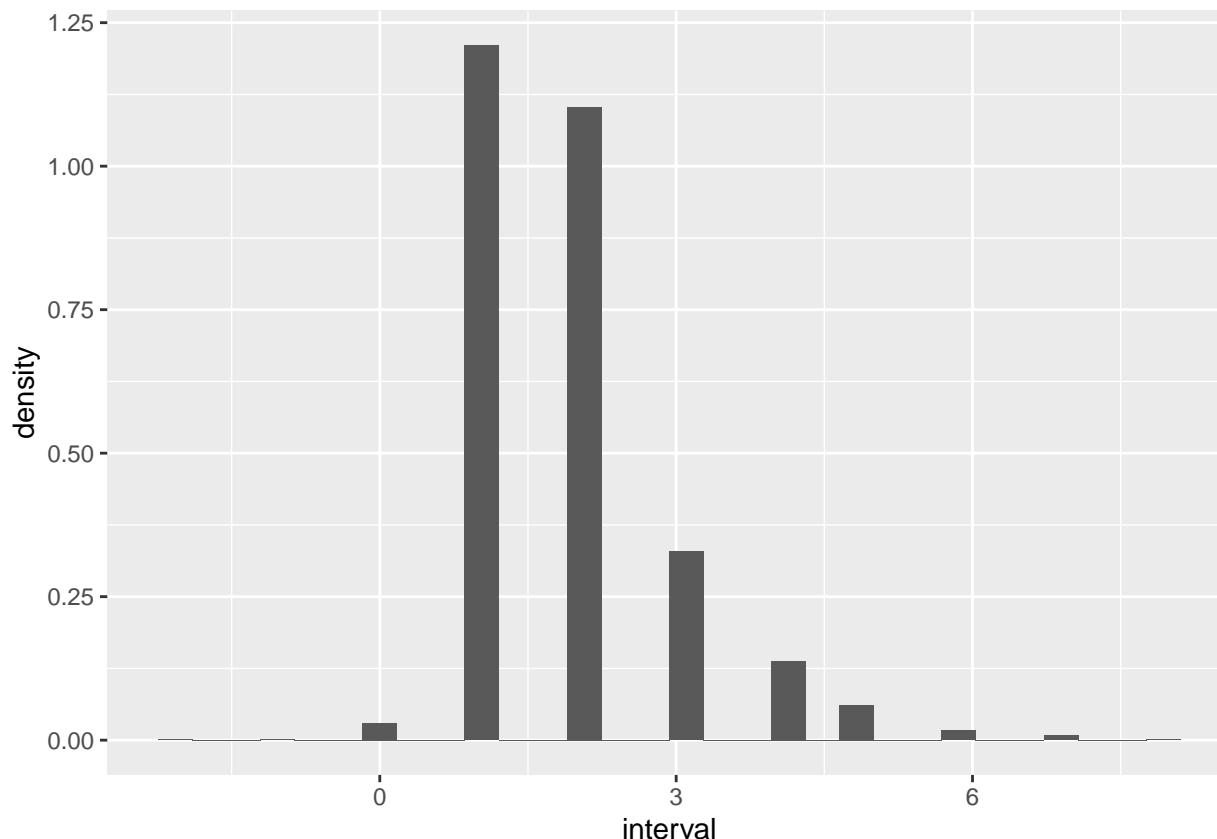
# create an interval variable
birthwt <- birthwt %>%
  arrange(mother.id, birth.order) %>%
  group_by(mother.id) %>%
  mutate(interval = maternal.age - lag(maternal.age))

birthwt

## # A tibble: 4,390 x 7
## # Groups:   mother.id [878]
##   mother.id birth.order birth.weight maternal.age child.id maternal.age.factor
##   <int>      <int>      <int>      <int>      <int>      <fct>
## 1         80         1        3175         18         1 [12,20]
## 2         80         2        3572         21         2 (20,30]
## 3         80         3        3317         24         3 (20,30]
## 4         80         4        4281         26         4 (20,30]
## 5         80         5        3827         28         5 (20,30]
## 6         84         1        2892         14         6 [12,20]
## 7         84         2        3204         16         7 [12,20]
## 8         84         3        4253         20         8 [12,20]
## 9         84         4        2948         22         9 (20,30]
## 10        84         5        3402         23        10 (20,30]
## # i 4,380 more rows
## # i 1 more variable: interval <int>

# distribution of maternal age interval
birthwt %>%
  filter(!is.na(interval)) %>%
  ggplot(aes(x = interval, y = ..density..)) +
  geom_histogram()

## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



## Visualizing relationships with birth weight

1. Produce a spaghetti plot of birth weights versus birth order, marking the threshold for low birth weight; include individual and overall averages. Does baseline birth weight seem to inform subsequent values?
  - (a) Color the plot by (i) age group, (ii) interpregnancy interval group?
2. Produce boxplots of birth weights across birth order for all mothers, marking the threshold for low birth weight. Is there a visual trend?
3. Produce boxplots of maternal age across birth order for all mothers, marking the threshold for low birth weight. Is there a visual trend?
  - (a) Color the plot by (i) low / not low birth weight
4. Produce a spaghetti plot of birth weights versus maternal age, marking the threshold for low birth weight; include individual and overall averages. Does baseline birth weight seem to inform subsequent values?
  - (a) Color the plot by (i) age group, (ii) interpregnancy interval group?
5. Produce a spaghetti plot of percent in birth weight group versus maternal age. Repeat for number in birth weight group versus maternal age. Compare the two. Does the most common birth weight group change as mothers age?

# Modeling

## Modeling birth weight from interpregnancy interval (and more?)

Consider two response variables: birth weight, which is a continuous variable measured in grams, and low-weight birth, a binary indicator variable.

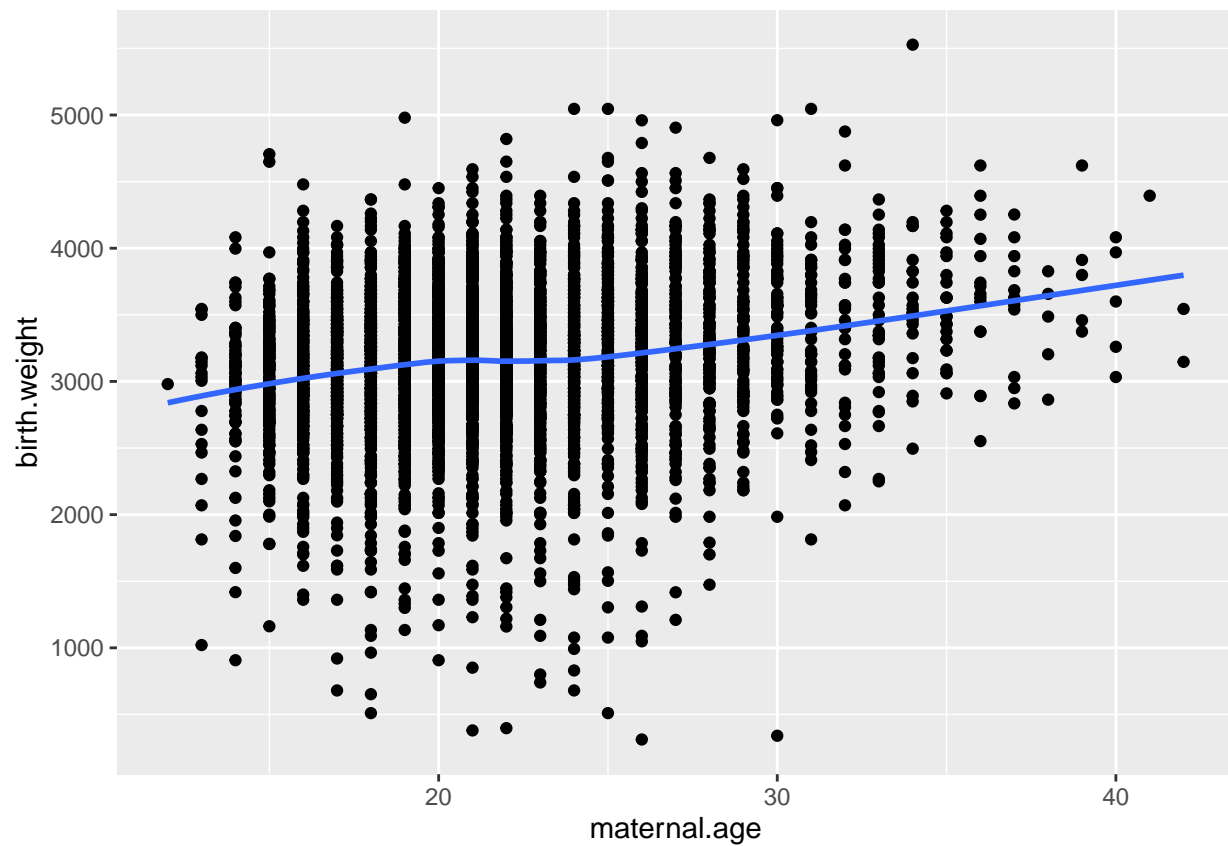
1. Model birth weight using a linear mixed effects model and interpret its coefficients.
2. Model birth weight using a linear fixed effects model and interpret its coefficients.
3. Compare the estimates and standard errors of the two models. Do they suggest different conclusions about effects? Do you believe a nonlinear model is well-motivated?
4. Model low-weight birth using logistic regression and interpret its coefficients.

## Birth weight vs covariates

Covariates: Maternal age, interpregnancy interval, birth order

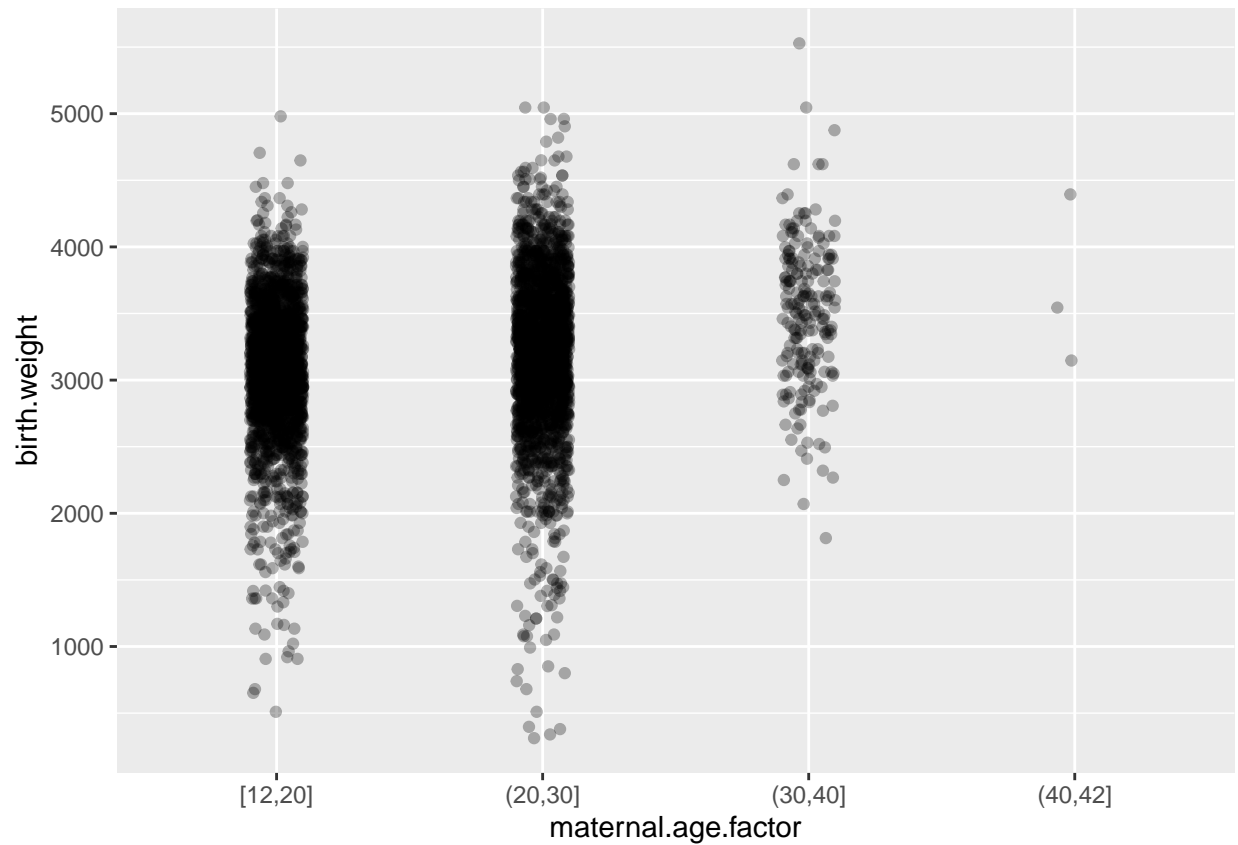
```
# plot birth weight vs maternal age
birthwt %>%
  ggplot(aes(x = maternal.age, y = birth.weight)) +
  geom_point() +
  geom_smooth(method = "loess", se = F)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



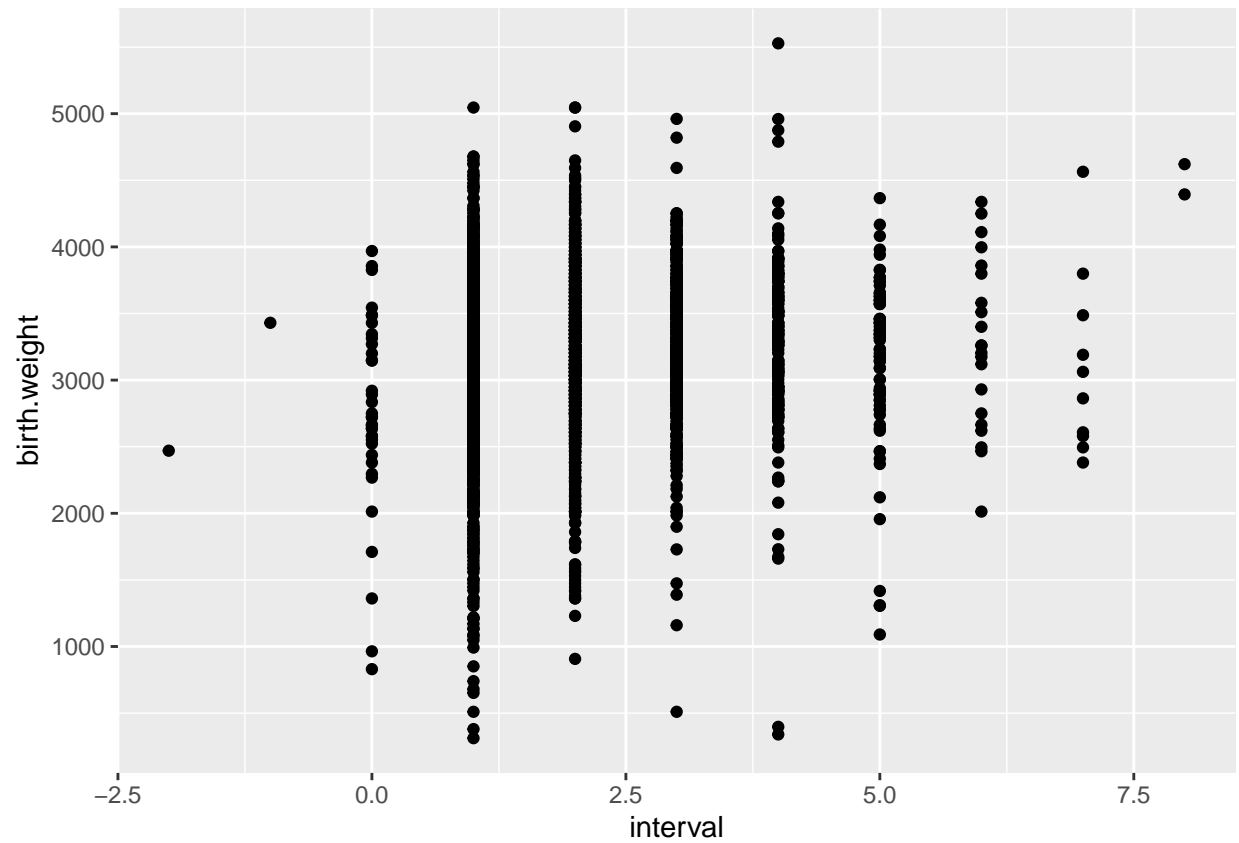
```
birthwt %>%
  ggplot(aes(x = maternal.age.factor, y = birth.weight)) +
  geom_jitter(width = 0.1, height = 0, alpha = 0.3) +
  geom_smooth(method = "loess", se = F)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



```
# plot birth weight vs interpregnancy interval
birthwt %>%
  ggplot(aes(y = birth.weight, x = interval)) +
  geom_point()
```

```
## Warning: Removed 878 rows containing missing values ('geom_point()').
```



```
# plot birth weight vs birth order  
birthwt %>%  
  ggplot(aes(y = birth.weight, x = birth.order)) +  
  geom_point()
```

