

Exploratory Data Analysis

EDAwR-Slides2.pdf

Read in the data

Read full file *once* then sample, save, comment out.

```
1 # df <- read.csv("~/Downloads/nat2021us.csv")
2 # set.seed(504)
3 # df_samp <- df[sample(nrow(df), 1000),]
4 # write.csv(df_samp, "sample1000.csv", row.names = FALSE)
5 # dat <- df[, c("mager", "fagecomb", "meduc", "dmar")]
6 # write.csv(dat, "births2021.csv", row.names = FALSE)
7 dat <- read.csv("births2021.csv")
```

Load libraries, etc.

```
1 library(tidyverse)
2 options(scipen = 999) # no scientific notation
3 theme_set(theme_bw(16)) # set ggplot2 theme
```

The data

```
1 head(dat)
```

	mager	fagecomb	meduc	dmar
1	22	23	3	1
2	31	27	4	2
3	29	29	4	2
4	39	40	6	2
5	20	28	3	2
6	29	39	3	1

```
1 dim(dat)
```

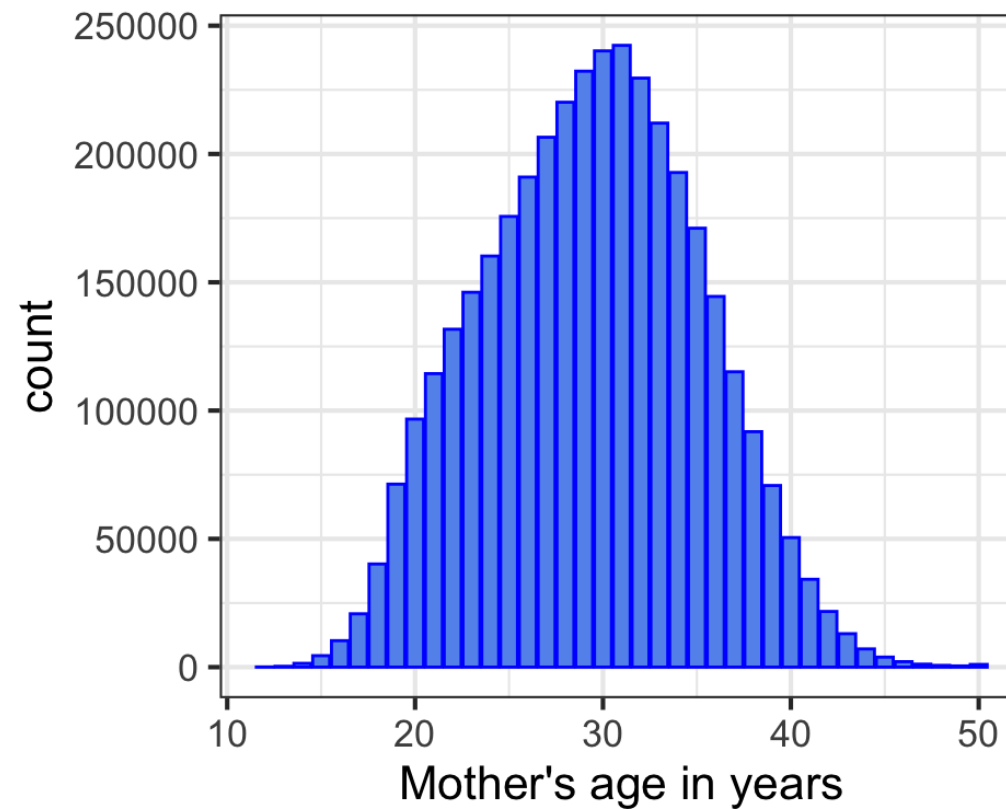
```
[1] 3669928      4
```

format: *cases* or *unbinned*

(no count or frequency column)

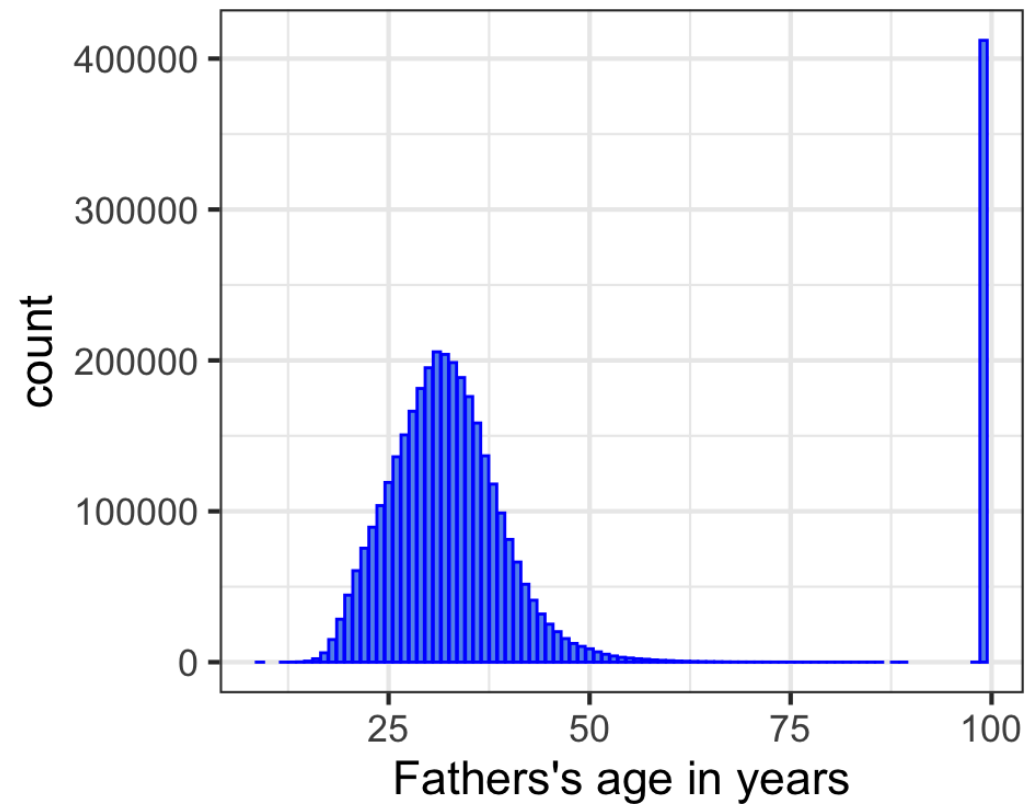
Mother's age

```
1 ggplot(dat, aes(mager)) +  
2   geom_bar(color = "blue", fill = "cornflowerblue") +  
3   xlab("Mother's age in years")
```



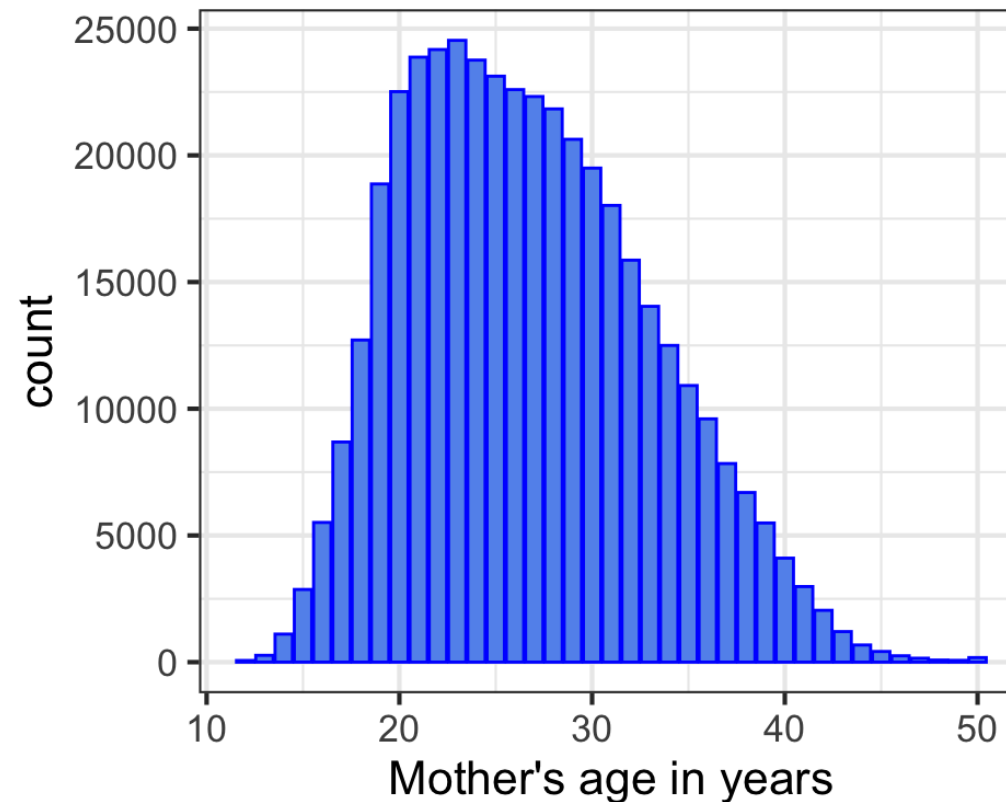
Father's age

```
1 ggplot(dat, aes(fagecomb)) +  
2   geom_bar(color = "blue", fill = "cornflowerblue") +  
3   xlab("Fathers's age in years")
```



Mother's age (father's age is missing)

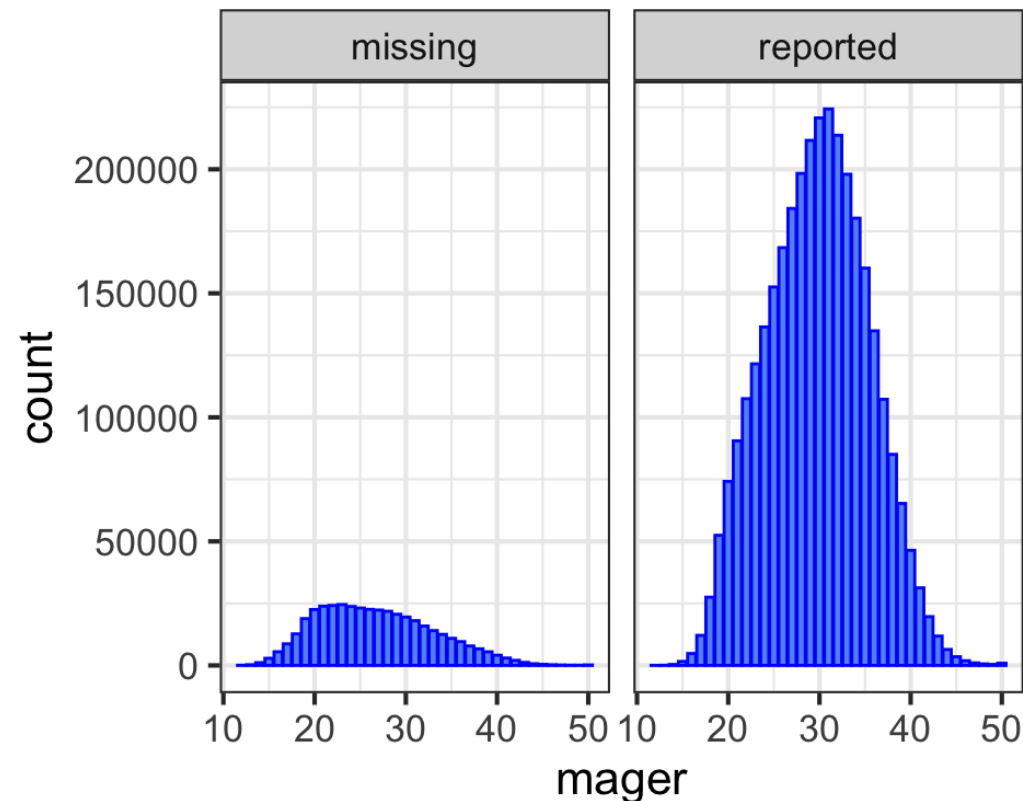
```
1 dat |> filter(fagecomb == 99) |>  
2 ggplot(aes(mager)) +  
3   geom_bar(color = "blue", fill = "cornflowerblue") +  
4   xlab("Mother's age in years")
```



Mother's age by father's age reporting

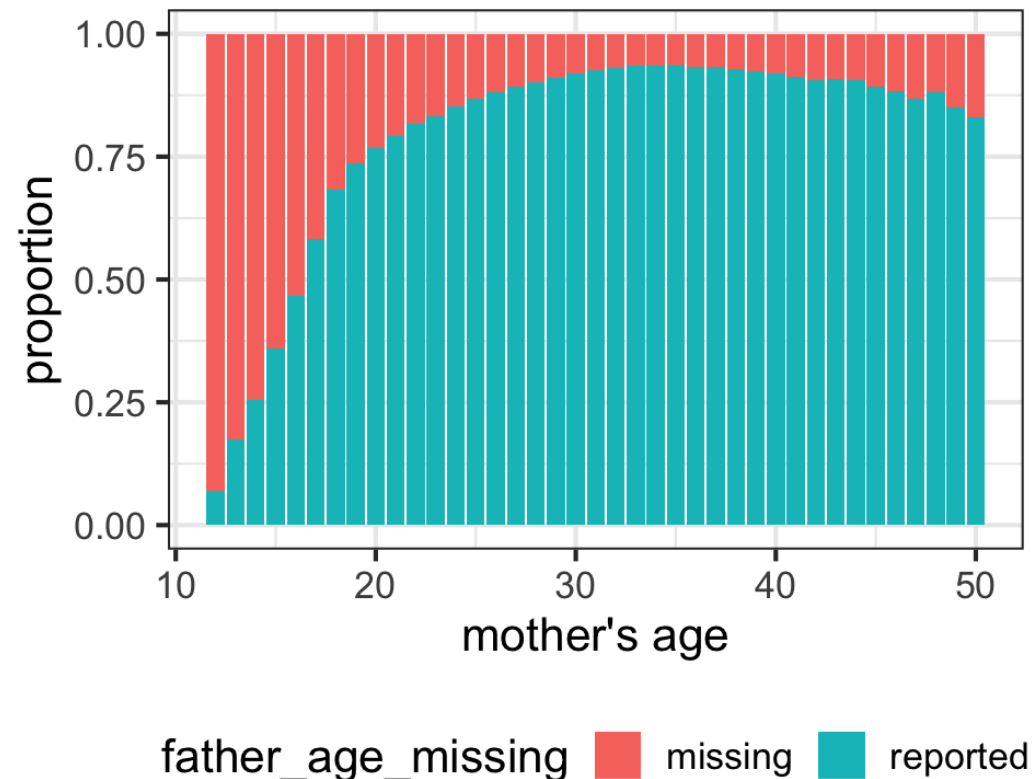
Create a new variable to indicate missing status

```
1 dat$father_age_missing <- ifelse(dat$fagecomb == 99, "missing", "reported")
2 ggplot(dat, aes(mager)) + geom_bar(color = "blue", fill = "cornflowerblue")
3   facet_wrap(~father_age_missing)
```



Mother's age by father's age reporting

```
1 dat |>
2   ggplot(aes(mager, fill = father_age_missing)) +
3   geom_bar(position = "fill") +
4   labs(x = "mother's age", y = "proportion") +
5   theme(legend.position = "bottom")
```



Father's age vs. mother's age

Bin data first. (Not practical to make a scatterplot with almost 4 million data points.)

```
1 df_count <- dat |> group_by(mager, fagecomb) |> count()  
2 head(df_count)
```

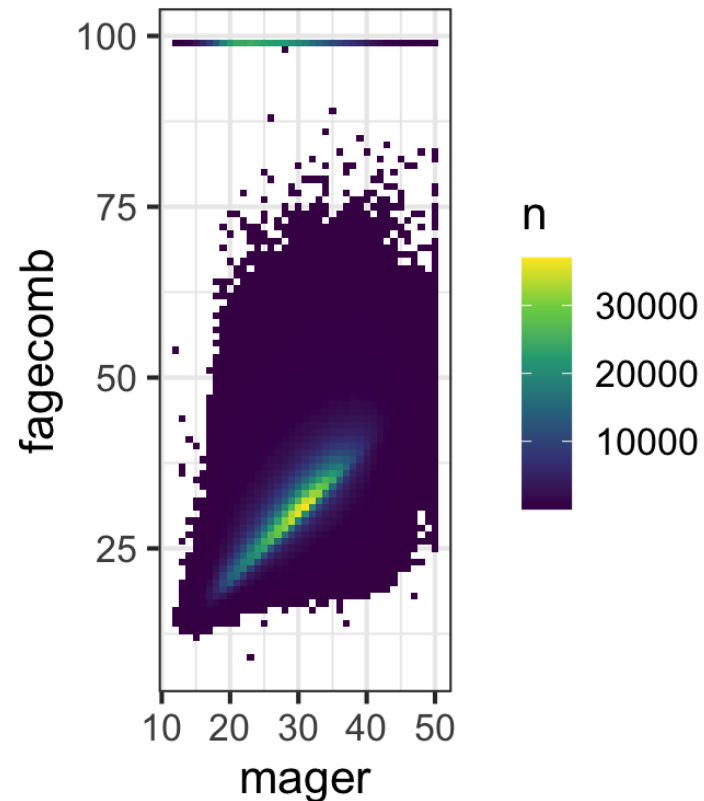
```
# A tibble: 6 × 3  
# Groups:   mager, fagecomb [6]  
  mager fagecomb      n  
  <int>   <int> <int>  
1     12      14      1  
2     12      15      1  
3     12      16      2  
4     12      54      1  
5     12      99     66  
6     13      13      8
```

```
1 dim(df_count)
```

```
[1] 1999      3
```

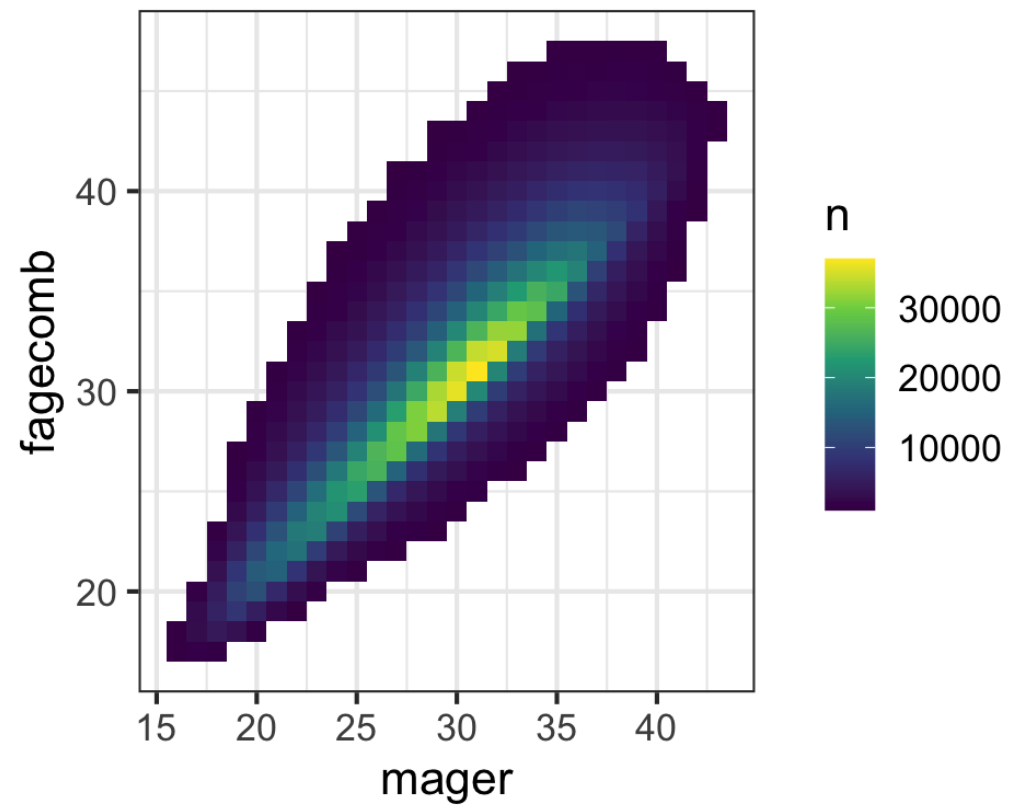
Father's age vs. mother's age

```
1 ggplot(df_count, aes(mager, fagecomb, fill = n)) +  
2   geom_tile() +  
3   scale_fill_viridis_c() +  
4   coord_fixed()
```



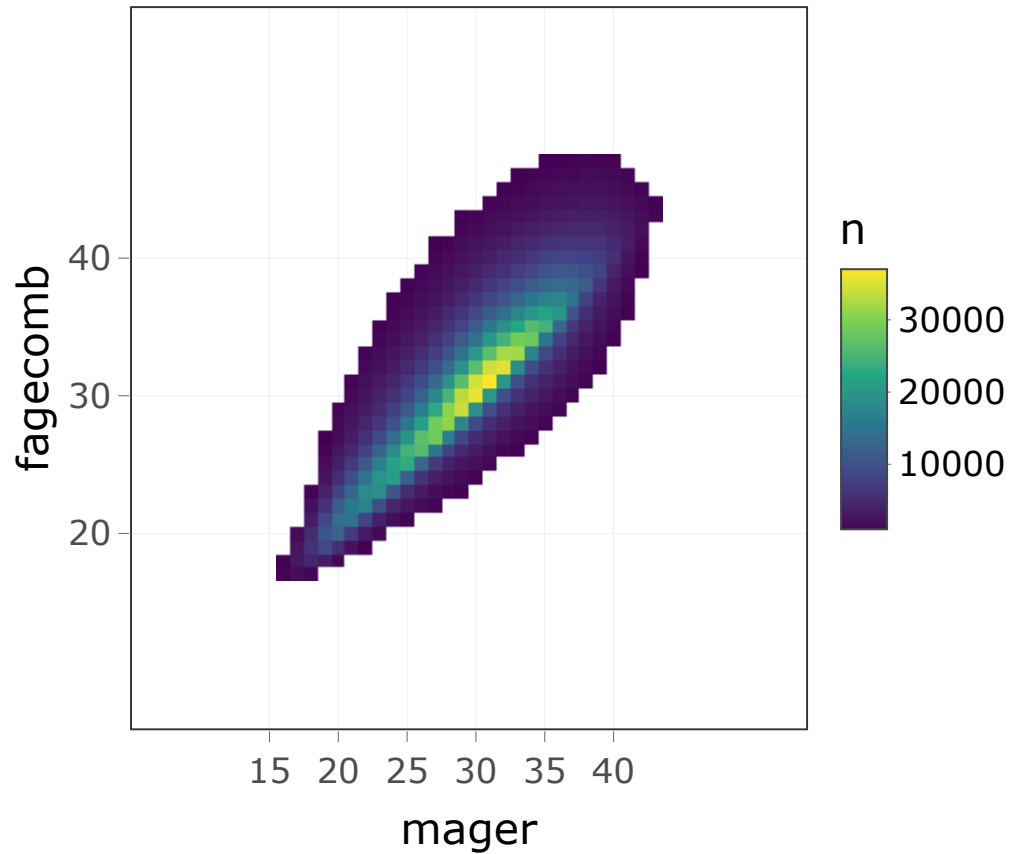
Father's age vs. mother's age

```
1 g <- df_count |> filter(n > 1000, fagecomb != 99) |>
2   ggplot(aes(mager, fagecomb, fill = n)) +
3   geom_tile() +
4   scale_fill_viridis_c() +
5   coord_fixed()
6 g
```



Father's age vs. mother's age

```
1 plotly::ggplotly(g)
```



Most popular ages

```
1 df_count |> group_by(mager) |> summarize(count = sum(n)) |>
2   slice_max(count, n = 4)
```

A tibble: 4 × 2

	mager	count
	<int>	<int>
1	31	242336
2	30	240183
3	29	232272
4	32	229576

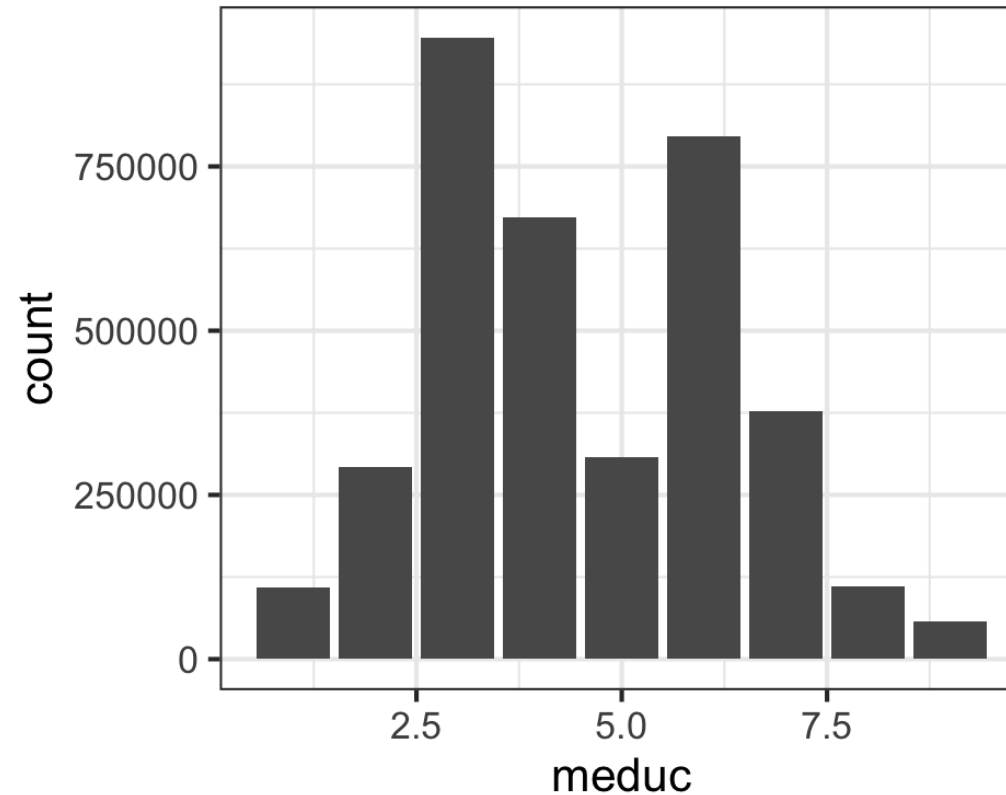
```
1 df_count |> group_by(fagecomb) |> summarize(count = sum(n)) |>
2   slice_max(count, n = 4)
```

A tibble: 4 × 2

	fagecomb	count
	<int>	<int>
1	99	412109
2	31	205647
3	32	203960
4	33	198552

Mother's education

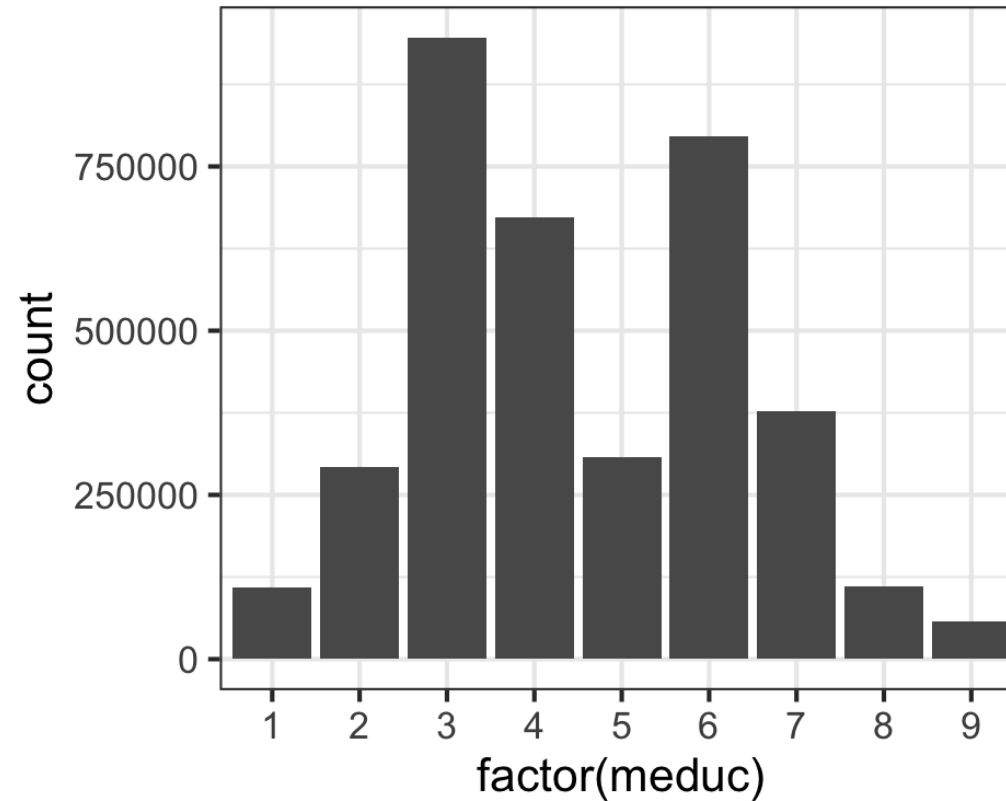
```
1 ggplot(dat, aes(meduc)) + geom_bar()
```



(x-axis is incorrect)

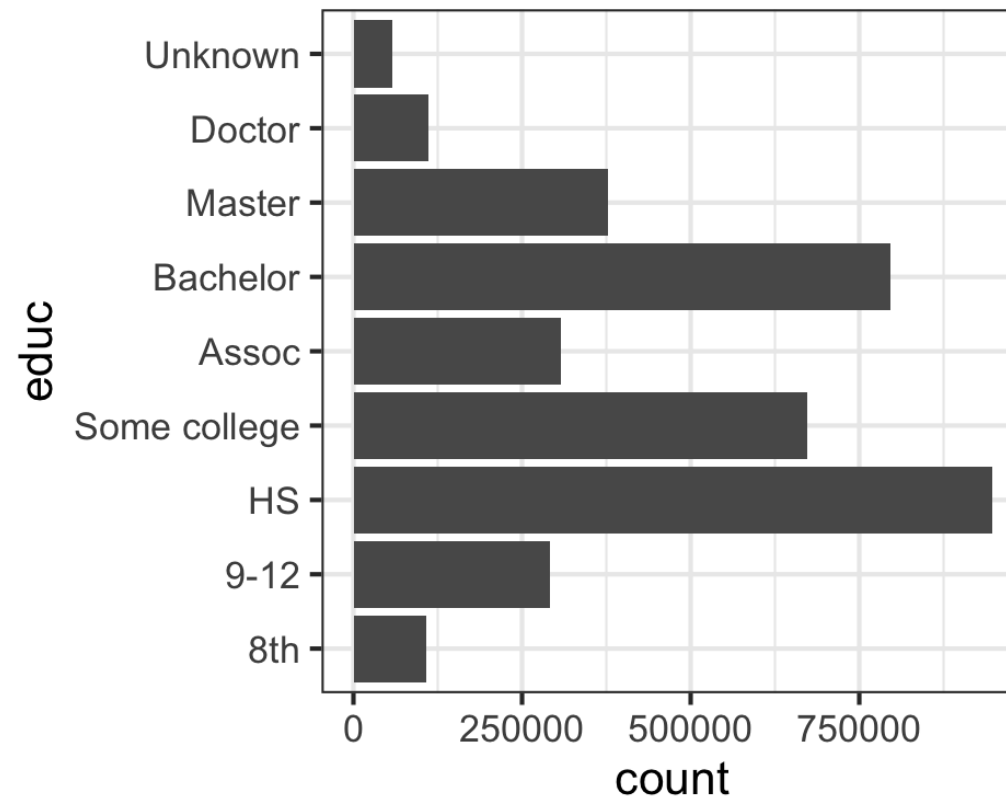
Mother's education

```
1 ggplot(dat, aes(factor(meduc))) + geom_bar()
```



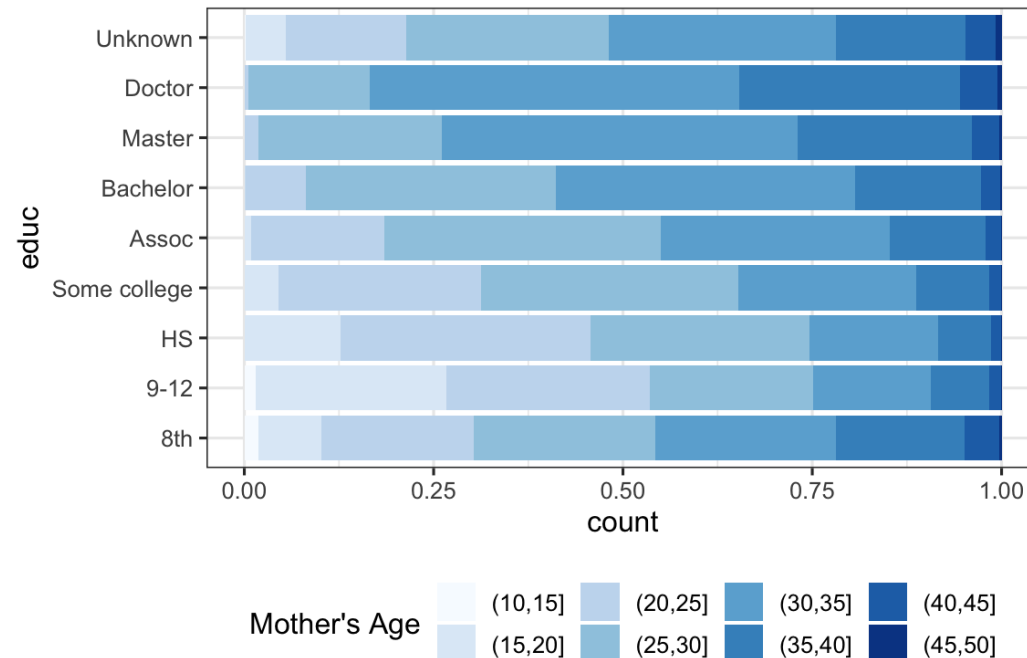
Mother's education

```
1 dat$educ <- fct_recode(factor(dat$meduc), "8th" = "1", "9-12" = "2",  
2                               "HS" = "3", "Some college" = "4", "Assoc" = "5",  
3                               "Bachelor" = "6", "Master" = "7", "Doctor" = "8",  
4                               "Unknown" = "9")  
5 ggplot(dat, aes(educ)) + geom_bar() + coord_flip()
```



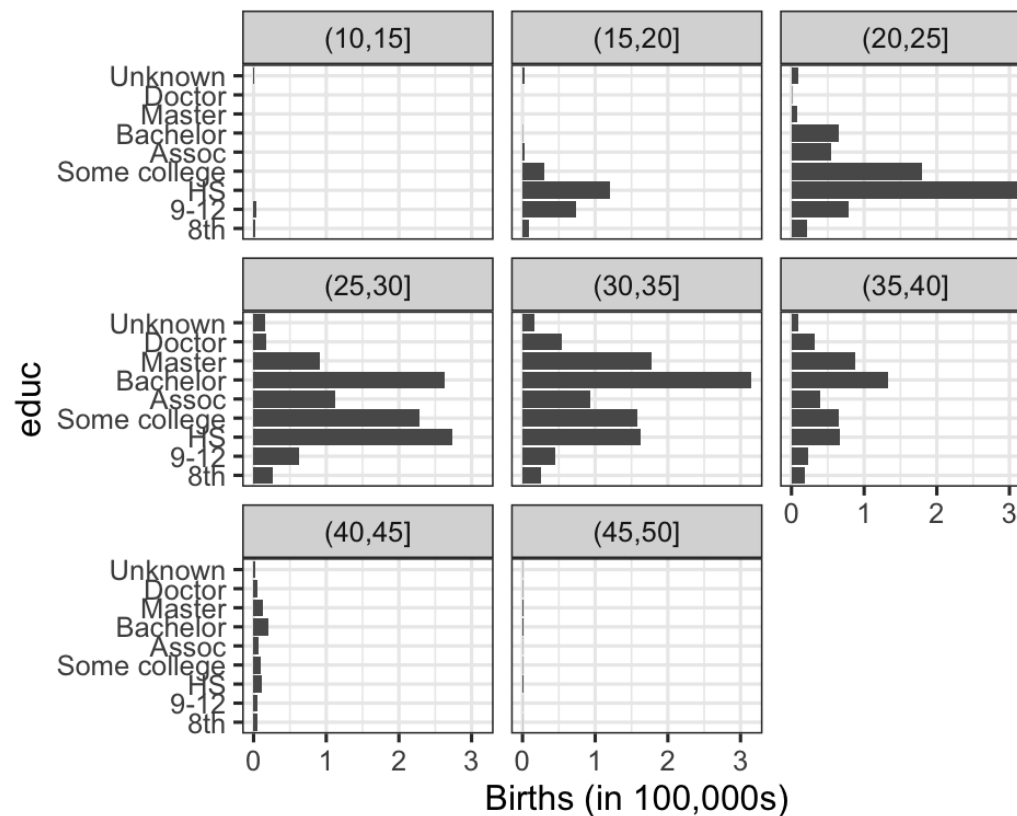
Mother's age and education

```
1 dat$m_age <- cut(dat$mager, breaks = seq(10, 50, 5))
2 ggplot(dat, aes(educ, fill = fct_rev(m_age))) +
3   geom_bar(position = "fill") +
4   coord_flip() +
5   scale_fill_brewer(direction = -1, breaks = levels(dat$m_age)) +
6   guides(fill=guide_legend(title="Mother's Age")) +
7   theme_bw(12) +
8   theme(legend.position = "bottom")
```



Mother's age and education

```
1 ggplot(dat, aes(educ)) +  
2   geom_bar() + facet_wrap(~m_age) + coord_flip() +  
3   scale_y_continuous(labels = ~.x/100000) +  
4   ylab("Births (in 100,000s)") + theme_bw(12)
```



Mother's age and education

```
1 ggplot(dat, aes(m_age)) + geom_bar() + facet_wrap(~educ) +  
2   coord_flip() + scale_y_continuous(labels = ~.x/100000) +  
3   ylab("Births (in 100,000s)") + theme_bw(12)
```

