

# MECH476: Engineering Data Analysis in R

## Chapter 7 Homework: Multivariate Exploratory Data Analysis

Brad Portouw

21 October, 2022

### **Load packages**

### **Chapter 7 Homework**

In Chapter 5, we briefly explored data on the salaries of engineering graduates from the National Science Foundation 2017 National Survey of College Graduates from a univariate perspective. Now, let's explore the relationships between multiple variables.

When a question asks you to make a plot, remember to set a theme, title, subtitle, labels, colors, etc. It is up to you how to personalize your plots, but put in some effort, and make the plotting approach consistent throughout the document. For example, you could use the same theme for all plots. I also like to use the subtitle as a place for the main summary for the viewer.

## Question 1: Data wrangling

Within a pipeline, import the data from the .csv file, convert all column names to lowercase text (either “manually” with `dplyr::rename()`, or use `clean_names()` from the `janitor` package), convert `gender` from “numeric” to “factor”, and drop any and all observations with `salary` recorded as 0. Assign this to a dataframe object with a meaningful name.

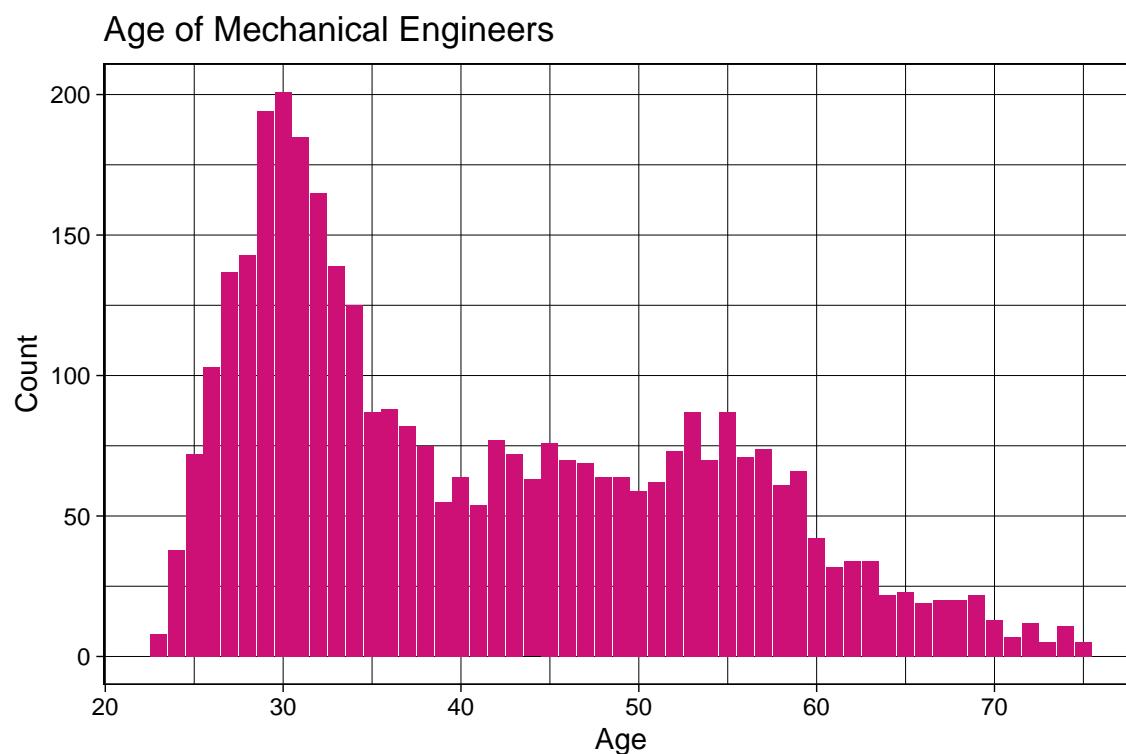
How many observations have a 0 (zero) value for salary? Note: The last question asked you to remove these observations from the resultant data frame.

What are the levels in `gender`? (Ignore the fact that the values for this variable refer to “biological sex”, not “gender”, although the categories often overlap among people).

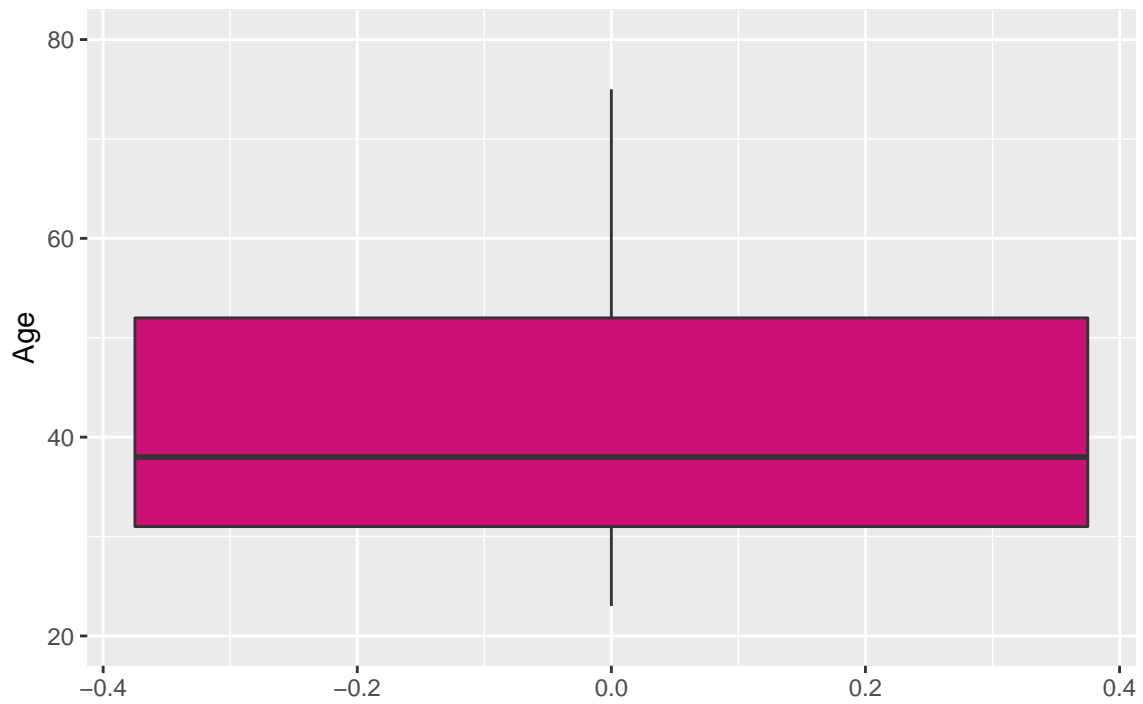
```
## [1] "F" "M"
```

## Question 2: Univariate EDA

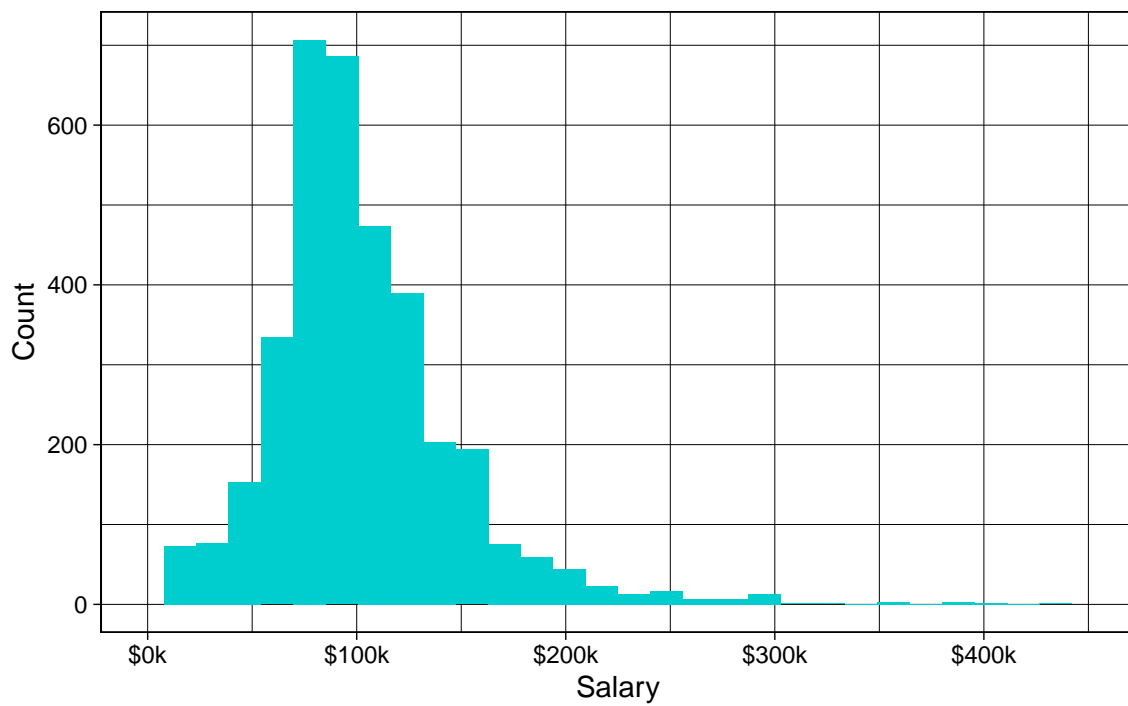
Using what you learned in Chapter 5, generate basic plots and/or descriptive statistics to explore `age`, `gender`, and `salary`. List whether each variable is continuous or categorical, and explain how and why you adjusted your EDA approach accordingly.



Age of Mechanical Engineers

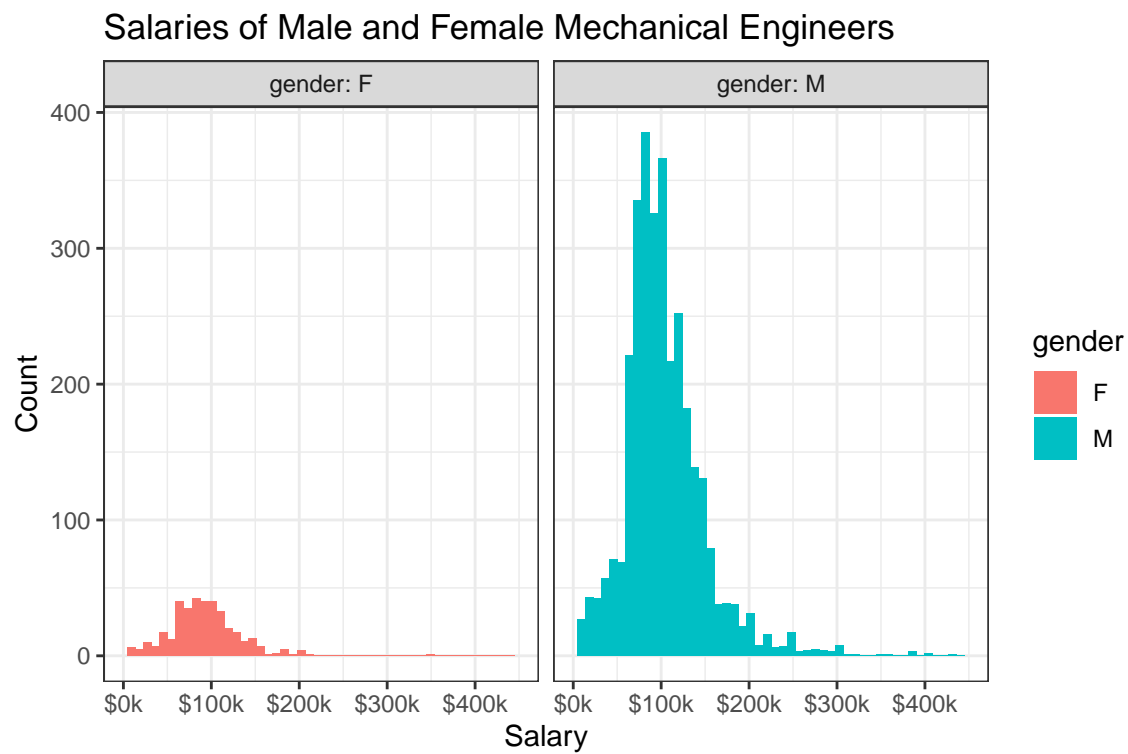


Salary Distribution for Mechanical Engineers

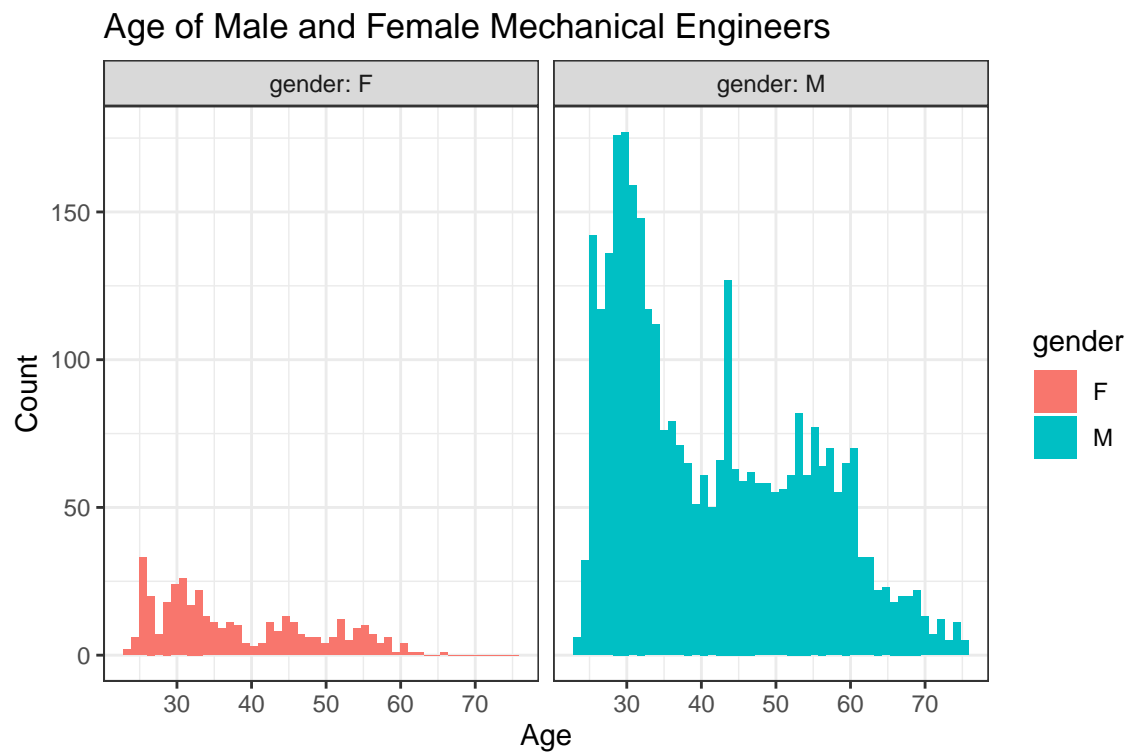


### Question 3: Multivariate histograms

Create a histogram of `salary`, faceted by `gender`. Add `bins = 50`.

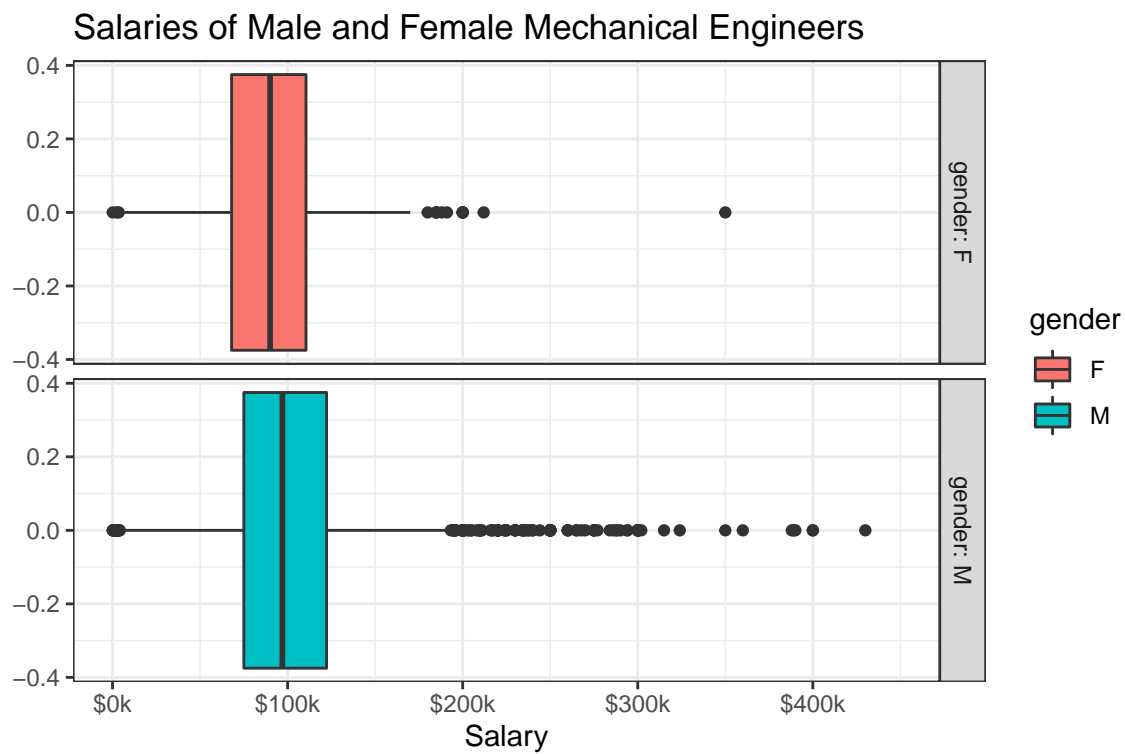


Create a histogram of `age`, faceted by `gender`. Add `bins = 50`.

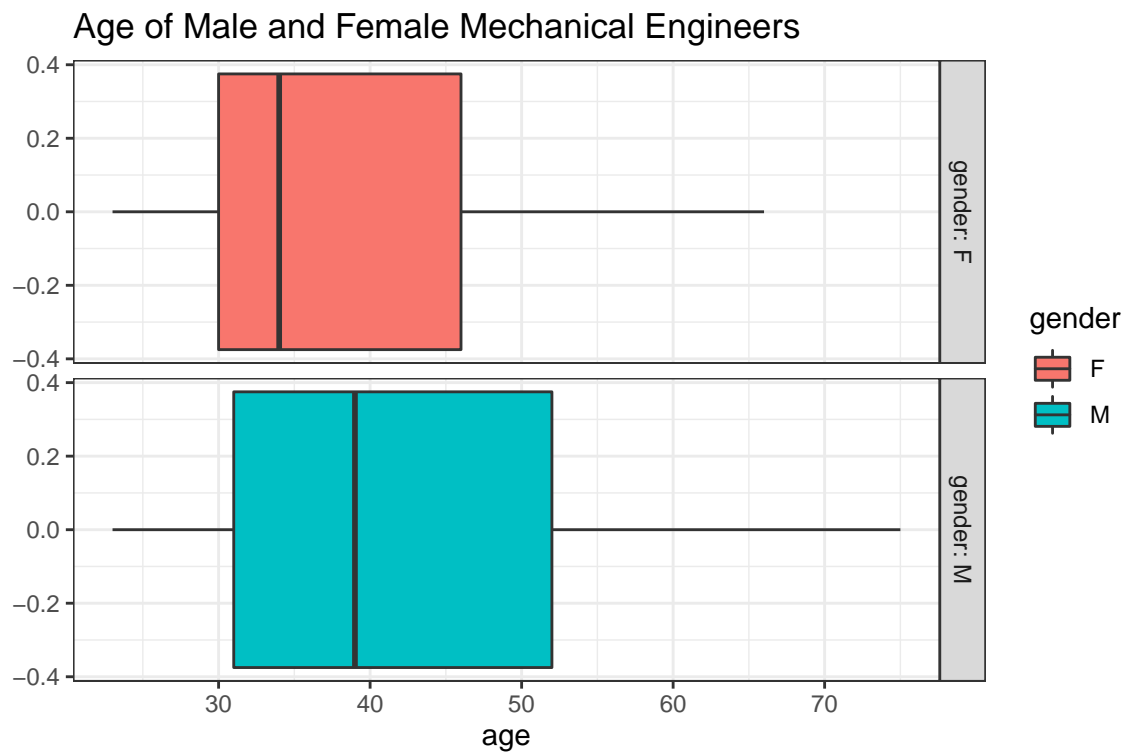


#### Question 4: Multivariate boxplots

Create a boxplot of `salary`, faceted by `gender`.

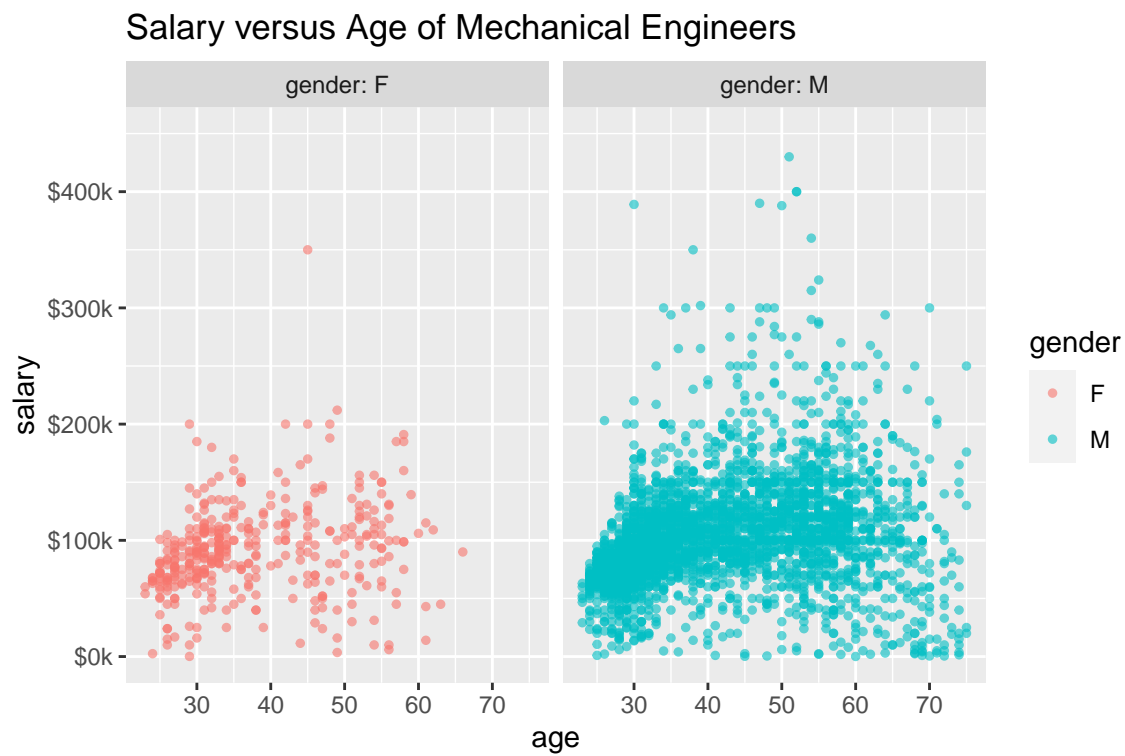


Create a boxplot of `age`, faceted by `gender`.



## Question 5: Scatterplot and correlation

Create a scatterplot of `age` (x-axis) and `salary`, differentiating by `gender`.



*Bonus point:* Is there a correlation between an engineer's salary and age? What is the estimated Pearson correlation coefficient  $r$ ? Run a formal test.

```
## [1] 0.2077451
```

## Question 6: Cumulative distribution function

Plot the cumulative distribution function of `salary` by `gender`. Adjust the x-axis with `scale_x_log10(limits = c(5e4, 5e5))` to zoom in a bit. What do you notice about the salaries for men and women? Hint: Remember there are greater differences the farther up you go on a log scale axis.

## Question 7: Quantiles

Calculate the quantiles of `salary` by `gender`. You can either subset the data with `dplyr::filter()` and dataframe assignment, or you can group by, summarize by quantile, and ungroup.

*Bonus point:* Assign the output to a dataframe, and use inline code to call individual values when answering the following questions. Do not let R use scientific notation in the text output; check the knitted document.

What is the difference in salary between men and women at the median?

- Median salary for women is  $\$9 \times 10^4$
- Median salary for men is  $\$9.7 \times 10^4$
- The difference at the median is \$7000

At the top percentile (maximum)?

- Maximum salary for women is  $\$3.5 \times 10^5$
- Maximum salary for men is  $\$1.027653 \times 10^6$
- The difference at the maximum is \$677653

Do you think there is a salary difference by gender across the pay scale? What other information would you need to test your hypothesis?

## Question 8: Hypothetical analysis

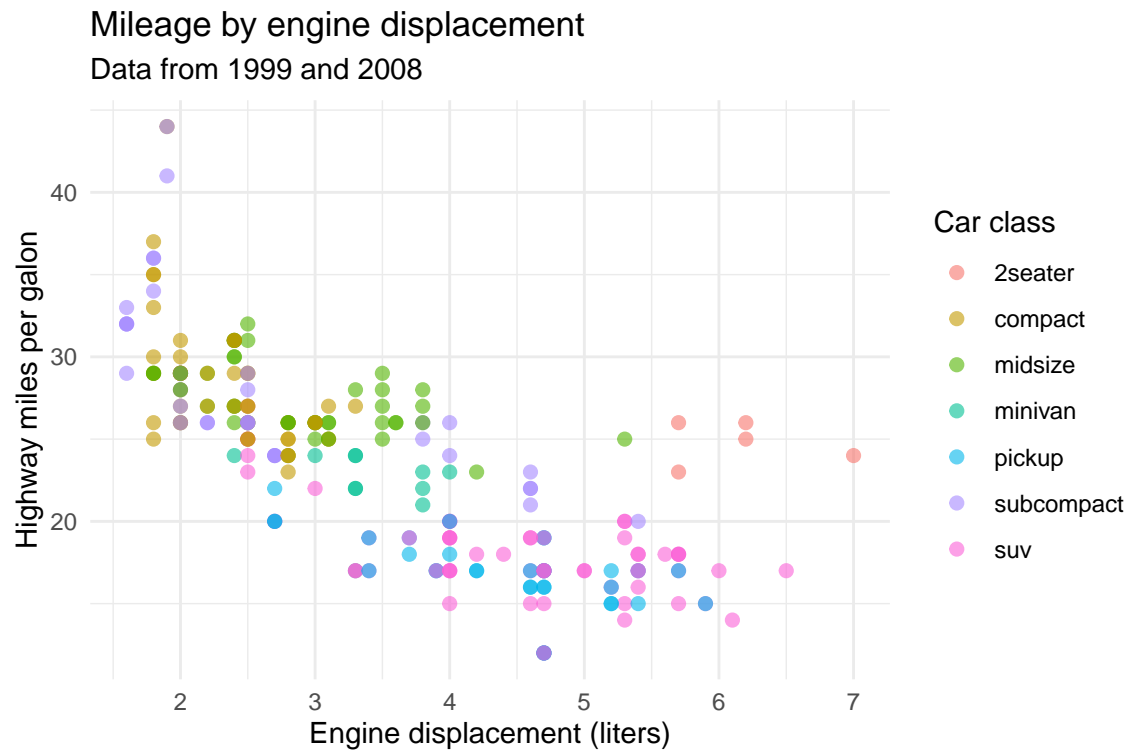
Think about what other variables you would like to include in an hypothetical analysis. From your perspective, what are the most important individual, family, and workforce factors related to salary—beyond gender and age?

There are a large number of potential factors in play for overall pay of employees, but the largest one might be the company and occupation of each engineer. If a company and type of engineering could be found, it would be helpful. Also it seems that the pay for many engineers plateaus at around 100k. It could be helpful to know which engineers have management positions. The school each engineer graduated at may play a small role in salary, as new hires out of prestigious schools may be offered more to start. Also where the engineer lives would play a role in pay considering cost of living, and potential benefits.

### Question 9: Recreate plot

Recreate this plot with the `mpg` dataset. Remember to use `?mpg` for information on the dataset and the variables. How would you describe the correlation between the independent variable and dependent variable? Do you see any patterns when considering the third variable?

(View R Markdown PDF for image)





## Appendix

```
# set global options for figures, code, warnings, and messages
knitr::opts_chunk$set(fig.width=6, fig.height=4, fig.path="../figs/",
                      echo=FALSE, warning=FALSE, message=FALSE)

# load packages for current session
library("tidyverse")

# import and tidy salary data
ME_salaries <- read.csv("../Data/ME_salaries.csv")%>%
  dplyr::rename(salary = SALARY, age=AGE, gender = GENDER) %>%
  dplyr::filter(salary != 0 & !is.na(salary))

ME_salaries$gender <- as.factor(ME_salaries$gender)
#check to see how many zero salaries were removed:

# 15 zero salary rows were removed.

# number of observations with salary as 0
ME_sal_zeroes <- read.csv("../Data/ME_salaries.csv")%>%
  dplyr::filter(SALARY == 0)%>%
  nrow()

# there were 15 rows removed for having 0 salaries.
# number of factor levels
levels(ME_salaries$gender)
#the levels associated with gender here are either Female (F) or Male(M), not considering other possible
# univariate eda

# age of Mechanical Engineers in the data frame. Excluding data with no salary listed (0).
ggplot(data = ME_salaries, aes(x = age))+
  geom_bar( fill= "deeppink3")+
  ggtitle("Age of Mechanical Engineers")+
  xlab("Age")+
  ylab("Count")+
  theme_linedraw()

# including a bar chart for age for more information:
ggplot(data = ME_salaries, aes(y = age))+
  geom_boxplot(fill = "deeppink3")+
  ylab("Age")+
  ylim(20,80)+
  ggtitle("Age of Mechanical Engineers")

#Basic plot for salary distribution. Outlier exists at 1 million dollars, so range is reduced to make plot
ggplot(data = ME_salaries, aes(x= salary))+
  geom_histogram( fill = "cyan3")+
  labs(x= "Salary",y="Count")+
  scale_x_continuous(limits = c(0, 450000),
                    labels = scales::label_dollar(scale = 0.001,
                                                  prefix = '$',
                                                  suffix = 'k'))+
  ggtitle("Salary Distribution for Mechanical Engineers")+
  theme_linedraw()

# histogram of salaries split by gender, again removing most extreme outliers, being capped at a salary
```

```

sal_gen_hist <- ggplot(ME_salaries, aes(x = salary, fill = gender))+
  geom_histogram(bins = 50)+
  labs(x= "Salary",y="Count")+
  scale_x_continuous(limits = c(0, 450000),
    labels = scales::label_dollar(scale = 0.001,
                                prefix = '$',
                                suffix = 'k'))+

  theme_bw()+
  ggtitle("Salaries of Male and Female Mechanical Engineers")+
  facet_grid(cols = vars(gender),
    labeller= label_both)

sal_gen_hist
# histogram of ages split by gender
age_gen_hist <- ggplot(ME_salaries, aes(x= age, fill = gender))+
  geom_histogram(bins = 50)+
  labs(x="Age", y="Count")+
  theme_bw()+
  ggtitle("Age of Male and Female Mechanical Engineers")+
  facet_grid(cols = vars(gender),
    labeller= label_both)

age_gen_hist
# boxplots of salary data by gender
sal_gen_box <- ggplot(ME_salaries, aes(x = salary, fill = gender))+
  geom_boxplot()+
  labs(x= "Salary")+
  scale_x_continuous(limits = c(0, 450000),
    labels = scales::label_dollar(scale = 0.001,
                                prefix = '$',
                                suffix = 'k'))+

  theme_bw()+
  ggtitle("Salaries of Male and Female Mechanical Engineers")+
  facet_grid(rows = vars(gender),
    labeller= label_both)

sal_gen_box
# boxplots of age data by gender
age_gen_box <- ggplot(ME_salaries, aes(x = age, fill = gender))+
  geom_boxplot()+
  labs(x= "age")+
  theme_bw()+
  ggtitle("Age of Male and Female Mechanical Engineers")+
  facet_grid(rows = vars(gender),
    labeller= label_both)

age_gen_box

# scatterplot of salary across age by gender
# due to the large overlap of data, along with a much greater concentration of male entries, the female
# I wasn't entirely sure if the focus should be on the difference in pay for different ages or genders
sal_vs_age_facet <- ggplot(ME_salaries)+
  geom_point(aes(x=age,y = salary, color = gender),
    width = 0.4,
    alpha = 0.6,
    size = 1)+

```

```

scale_y_continuous(limits = c(0, 450000),
                   labels = scales::label_dollar(scale = 0.001,
                                                prefix = '$',
                                                suffix = 'k'))+

ggtitle("Salary versus Age of Mechanical Engineers")+
facet_grid(cols = vars(gender),
           labeller= label_both)
sal_vs_age_facet

# This plot is throwing an error of

#Warning: Ignoring unknown parameters: widthWarning: Ignoring unknown parameters: width

# I haven't been able to fix this error, and have moved on to the other question. Hopefully the plot ab

#sal_vs_age <- ggplot(ME_salaries)+
# geom_point(aes(x=age,y = salary, color = gender),
#           width = 1)+
#scale_y_continuous(limits = c(0, 450000),
#                   labels = scales::label_dollar(scale = 0.001,
#                                                prefix = '$',
#                                                suffix = 'k'))+
# ggtitle("Salary versus Age of Mechanical Engineers")

# correlation test
# the pearson coefficient is determined by the function, cor(), using age and salary across all enginee

r_age_salary <-cor(ME_salaries$age,ME_salaries$salary)

r_age_salary
# the pearson coefficient for age vs salary is about 0.2077, which isn't very signifigant. Looking at t
# plot cdf of salary by gender
#making cumulative fractions of salaries from ME_salaries
ordered_salaries <- ME_salaries%>%
  dplyr::filter(salary<500000)%>%
  dplyr::arrange(salary)
# dplyr::mutate(cum_pct = seq.int(from = #1/length(ordered_salaries$salary),
#                               to = 1,
#                               by = #1/length(ordered_salaries$salary)))
# calculate quantiles of salary by gender
#Making new dataframe from orginal, converting gender from factor into intergers for filtering.
gen_ints <- ME_salaries
gen_ints$gender <- as.integer(ME_salaries$gender)

#filter into Male and Female datasets
F_sal <- gen_ints %>%
  dplyr::filter(gender == 1)

F_quantile <- as.integer(quantile(F_sal$salary,
                                probs= seq(0,1,0.1)))%>%
  round(0)#%>%
  #format(scientific = FALSE)

```

```

M_sal <- gen_ints %>%
  dplyr::filter(gender == 2)

M_quantile <- as.integer(quantile(M_sal$salary,
                                probs=seq(0,1,0.1)))%>%
  round(0)#%>%

  #format(scientific = FALSE)
med_diff <- as.integer(M_quantile[6]-F_quantile[6])%>%
  format(scientific = FALSE)
max_diff <- as.integer(M_quantile[11]-F_quantile[11])%>%
  format(scientific = FALSE)
#
H <- ggplot(mpg, aes(x= displ,y = hwy, color = class))+
  geom_point(alpha = 0.6,
            size = 2)+
  labs(title = "Mileage by engine displacement",
        subtitle = "Data from 1999 and 2008",
        x="Engine displacement (liters)",
        y="Highway miles per gallon")+
  guides(color=guide_legend(title="Car class"))+
  theme_minimal()
H

```