

Chapter 3 Univariate Graphs

Univariate graphs plot the distribution of data from a single variable. The variable can be categorical (e.g., race, sex) or quantitative (e.g., age, weight).

3.1 Categorical

The distribution of a single categorical variable is typically plotted with a bar chart, a pie chart, or (less commonly) a tree map.

3.1.1 Bar chart

The [Marriage](#) dataset contains the marriage records of 98 individuals in Mobile County, Alabama. Below, a bar chart is used to display the distribution of wedding participants by race.

```
library(ggplot2)
data(Marriage, package = "mosaicData")

# plot the distribution of race
ggplot(Marriage, aes(x = race)) +
  geom_bar()
```

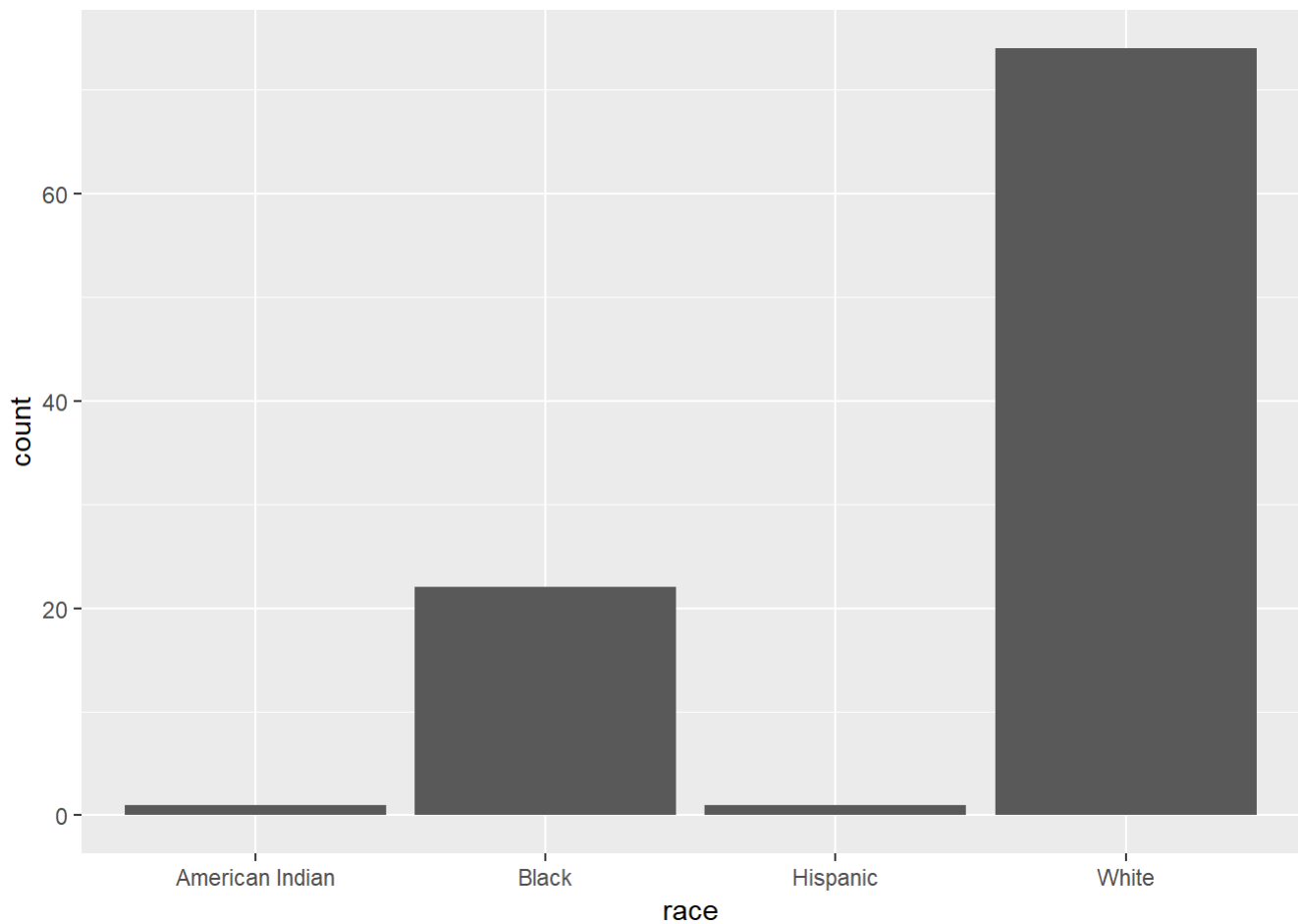


Figure 3.1: Simple barchart

The majority of participants are white, followed by black, with very few Hispanics or American Indians.

You can modify the bar fill and border colors, plot labels, and title by adding options to the `geom_bar` function.

```
# plot the distribution of race with modified colors and labels
ggplot(Marriage, aes(x = race)) +
  geom_bar(fill = "cornflowerblue",
           color="black") +
  labs(x = "Race",
       y = "Frequency",
       title = "Participants by race")
```

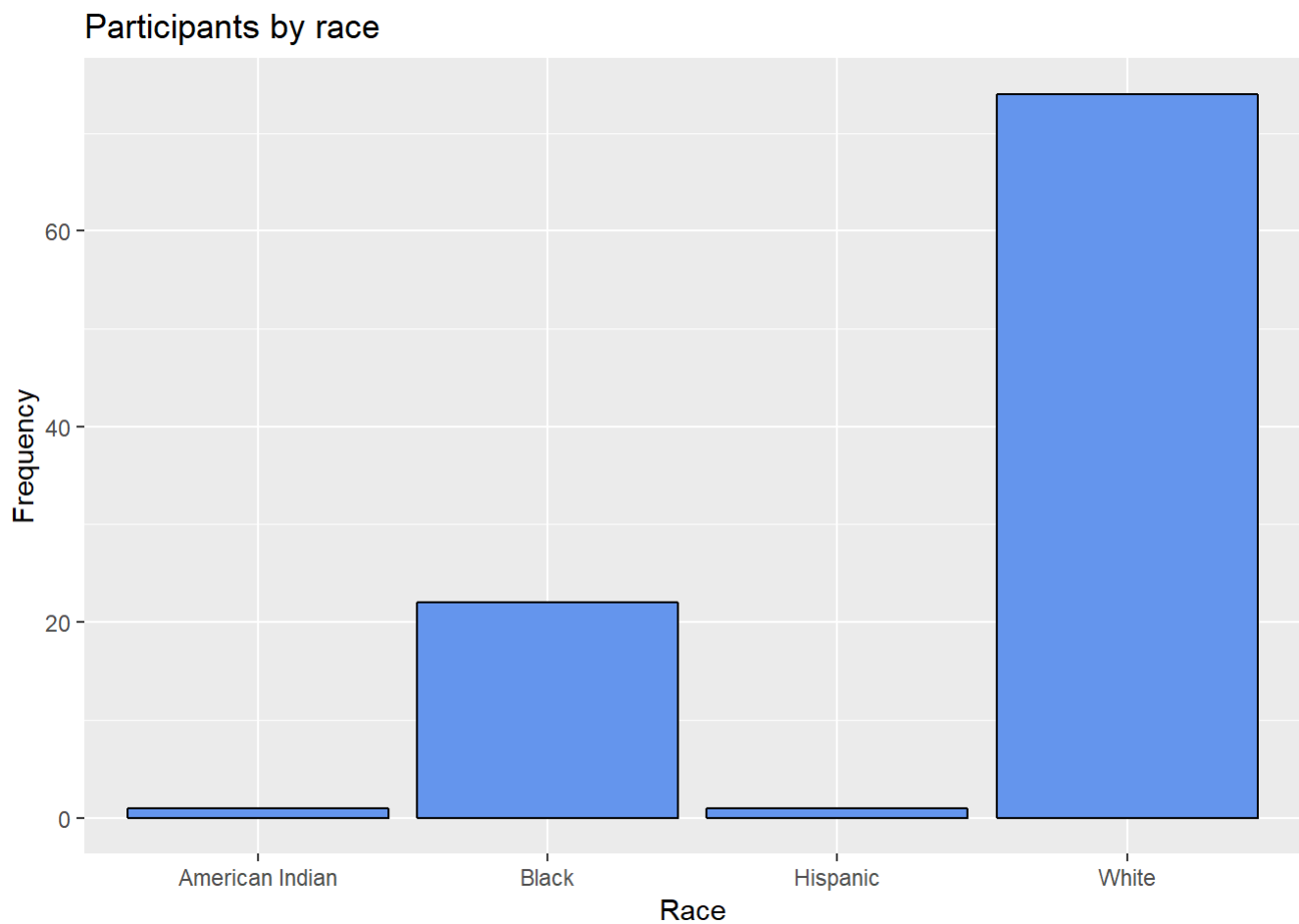


Figure 3.2: Barchart with modified colors, labels, and title

3.1.1.1 Percents

Bars can represent percents rather than counts. For bar charts, the code `aes(x=race)` is actually a shortcut for `aes(x = race, y = ..count..)`, where `..count..` is a special variable representing the frequency within each category. You can use this to calculate percentages, by specifying the `y` variable explicitly.

```
# plot the distribution as percentages
ggplot(Marriage,
       aes(x = race,
           y = ..count.. / sum(..count..))) +
geom_bar() +
labs(x = "Race",
     y = "Percent",
     title = "Participants by race") +
scale_y_continuous(labels = scales::percent)
```

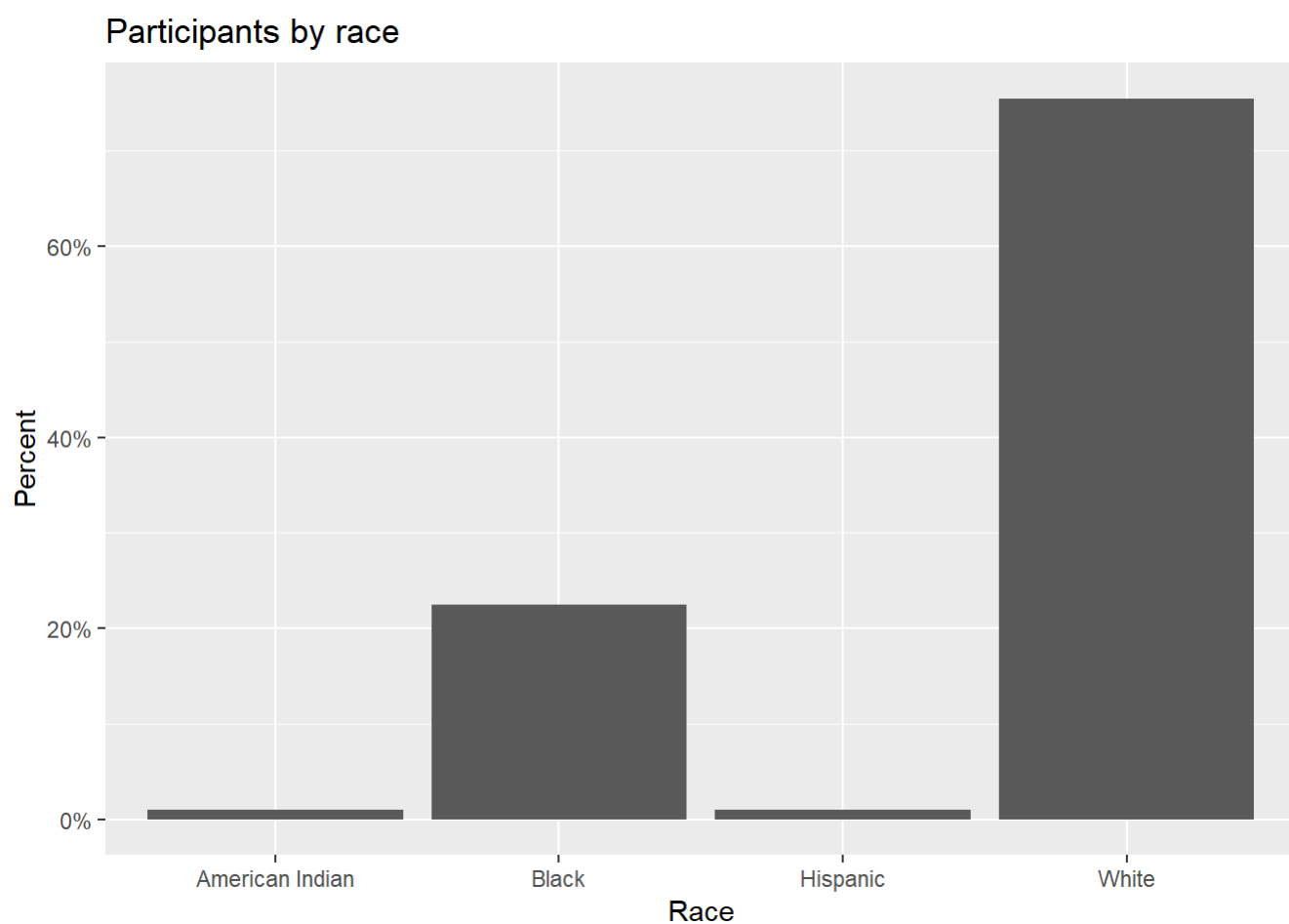


Figure 3.3: Barchart with percentages

In the code above, the `scales` package is used to add % symbols to the y-axis labels.

3.1.1.2 Sorting categories

It is often helpful to sort the bars by frequency. In the code below, the frequencies are calculated explicitly. Then the `reorder` function is used to sort the categories by the frequency. The option `stat="identity"` tells the plotting function not to calculate counts, because they are supplied directly.

```
# calculate number of participants in
# each race category
library(dplyr)
plotdata <- Marriage %>%
  count(race)
```

The resulting dataset is give below.

Table 3.1: plotdata

race	n
American Indian	1
Black	22
Hispanic	1
White	74

This new dataset is then used to create the graph.

```
# plot the bars in ascending order
ggplot(plotdata,
  aes(x = reorder(race, n),
    y = n)) +
  geom_bar(stat = "identity") +
  labs(x = "Race",
    y = "Frequency",
    title = "Participants by race")
```

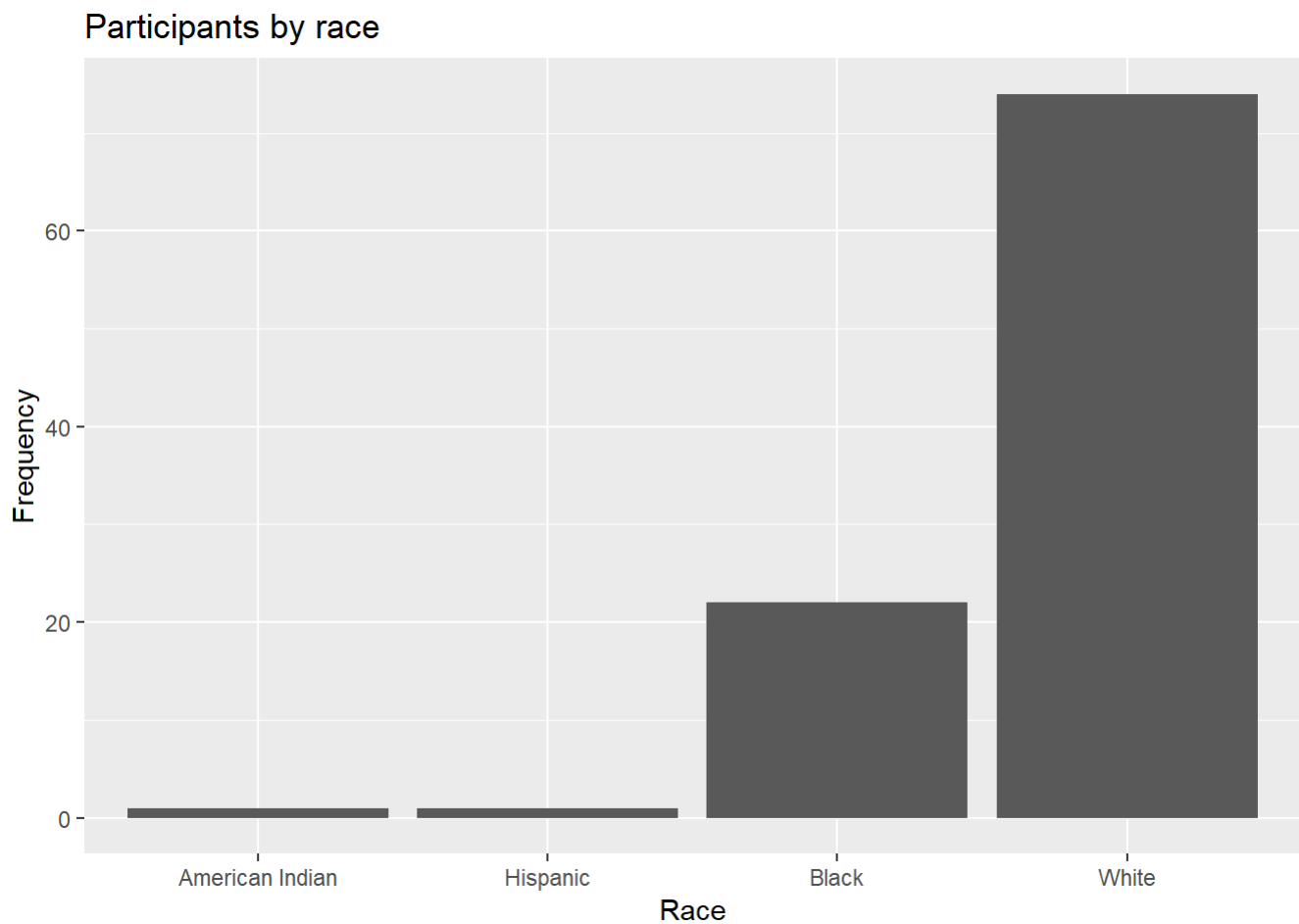


Figure 3.4: Sorted bar chart

The graph bars are sorted in ascending order. Use `reorder(race, -n)` to sort in descending order.

3.1.1.3 Labeling bars

Finally, you may want to label each bar with its numerical value.

```
# plot the bars with numeric labels
ggplot(plotdata,
       aes(x = race,
           y = n)) +
  geom_bar(stat = "identity") +
  geom_text(aes(label = n),
           vjust=-0.5) +
  labs(x = "Race",
       y = "Frequency",
       title = "Participants by race")
```

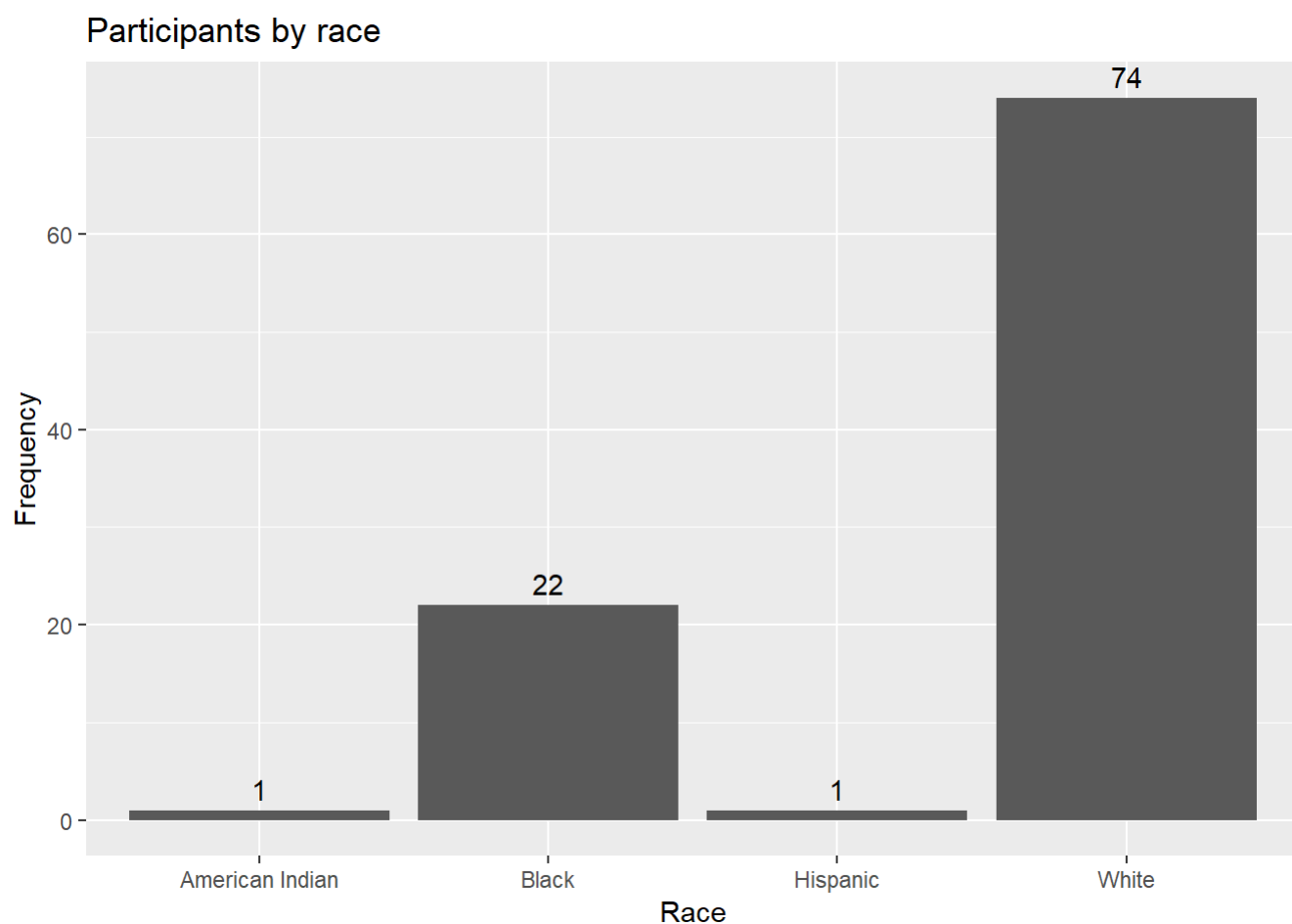


Figure 3.5: Bar chart with numeric labels

Here `geom_text` adds the labels, and `vjust` controls vertical justification. See [Annotations](#) for more details.

Putting these ideas together, you can create a graph like the one below. The minus sign in `reorder(race, -pct)` is used to order the bars in descending order.

```

library(dplyr)
library(scales)
plotdata <- Marriage %>%
  count(race) %>%
  mutate(pct = n / sum(n),
         pctlabel = paste0(round(pct*100), "%"))

# plot the bars as percentages,
# in decending order with bar labels
ggplot(plotdata,
       aes(x = reorder(race, -pct),
          y = pct)) +
  geom_bar(stat = "identity",
          fill = "indianred3",
          color = "black") +
  geom_text(aes(label = pctlabel),
           vjust = -0.25) +
  scale_y_continuous(labels = percent) +
  labs(x = "Race",
       y = "Percent",
       title = "Participants by race")

```

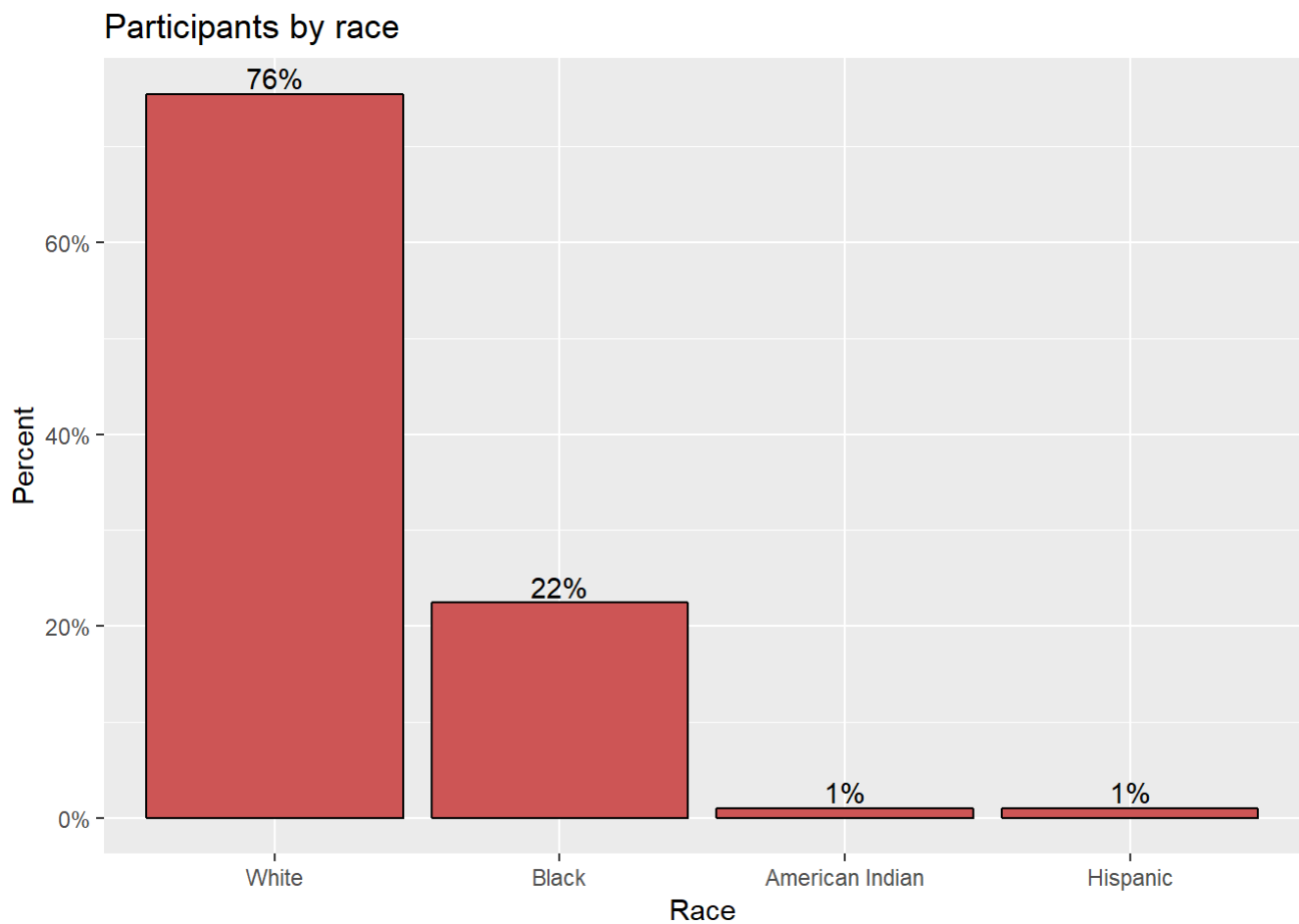



Figure 3.6: Sorted bar chart with percent labels

3.1.1.4 Overlapping labels

Category labels may overlap if (1) there are many categories or (2) the labels are long. Consider the distribution of marriage officials.

```
# basic bar chart with overlapping labels
ggplot(Marriage, aes(x = officialTitle)) +
  geom_bar() +
  labs(x = "Officiate",
       y = "Frequency",
       title = "Marriages by officiate")
```

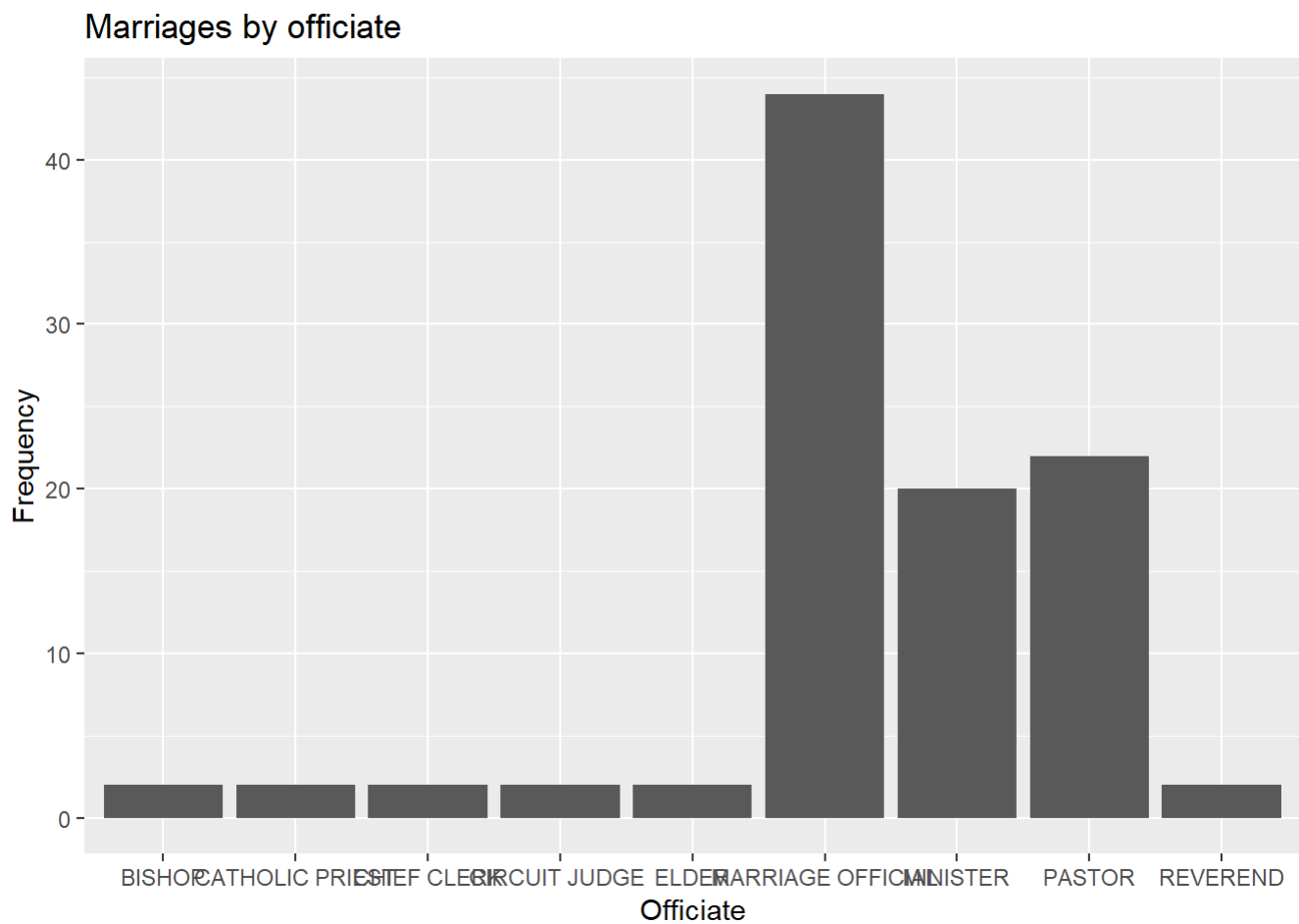


Figure 3.7: Barchart with problematic labels

In this case, you can flip the x and y axes.

```
# horizontal bar chart
ggplot(Marriage, aes(x = officialTitle)) +
  geom_bar() +
  labs(x = "",
       y = "Frequency",
       title = "Marriages by officiate") +
  coord_flip()
```

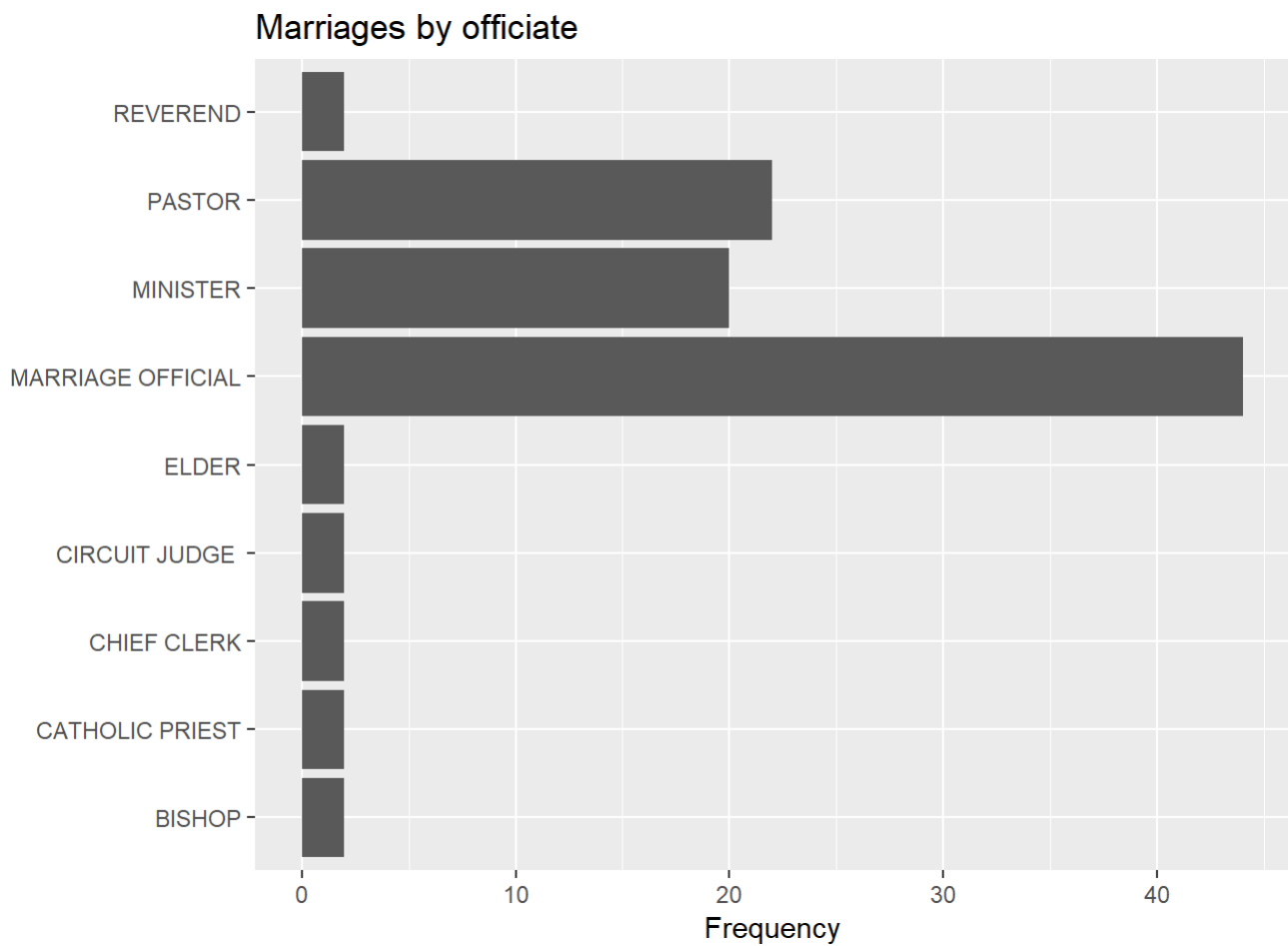


Figure 3.8: Horizontal barchart

Alternatively, you can rotate the axis labels.

```
# bar chart with rotated labels
ggplot(Marriage, aes(x = officialTitle)) +
  geom_bar() +
  labs(x = "",
       y = "Frequency",
       title = "Marriages by officiate") +
  theme(axis.text.x = element_text(angle = 45,
                                    hjust = 1))
```

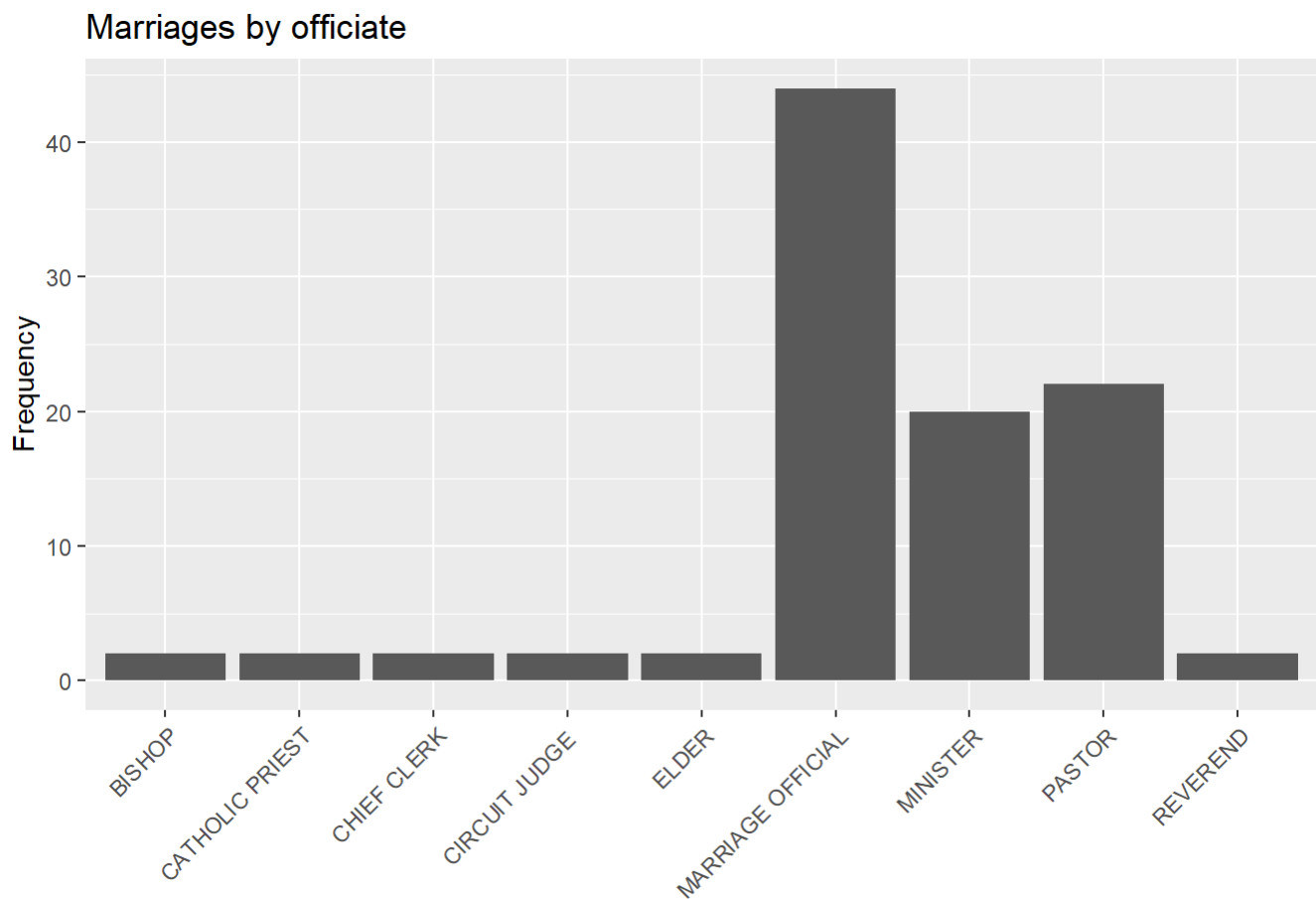


Figure 3.9: Barchart with rotated labels

Finally, you can try staggering the labels. The trick is to add a newline `\n` to every other label.

```
# bar chart with staggered labels
lbls <- paste0(c("", "\n"),
               levels(Marriage$officialTitle))
ggplot(Marriage,
       aes(x=factor(officialTitle,
                    labels = lbls))) +
  geom_bar() +
  labs(x = "",
       y = "Frequency",
       title = "Marriages by officiate")
```

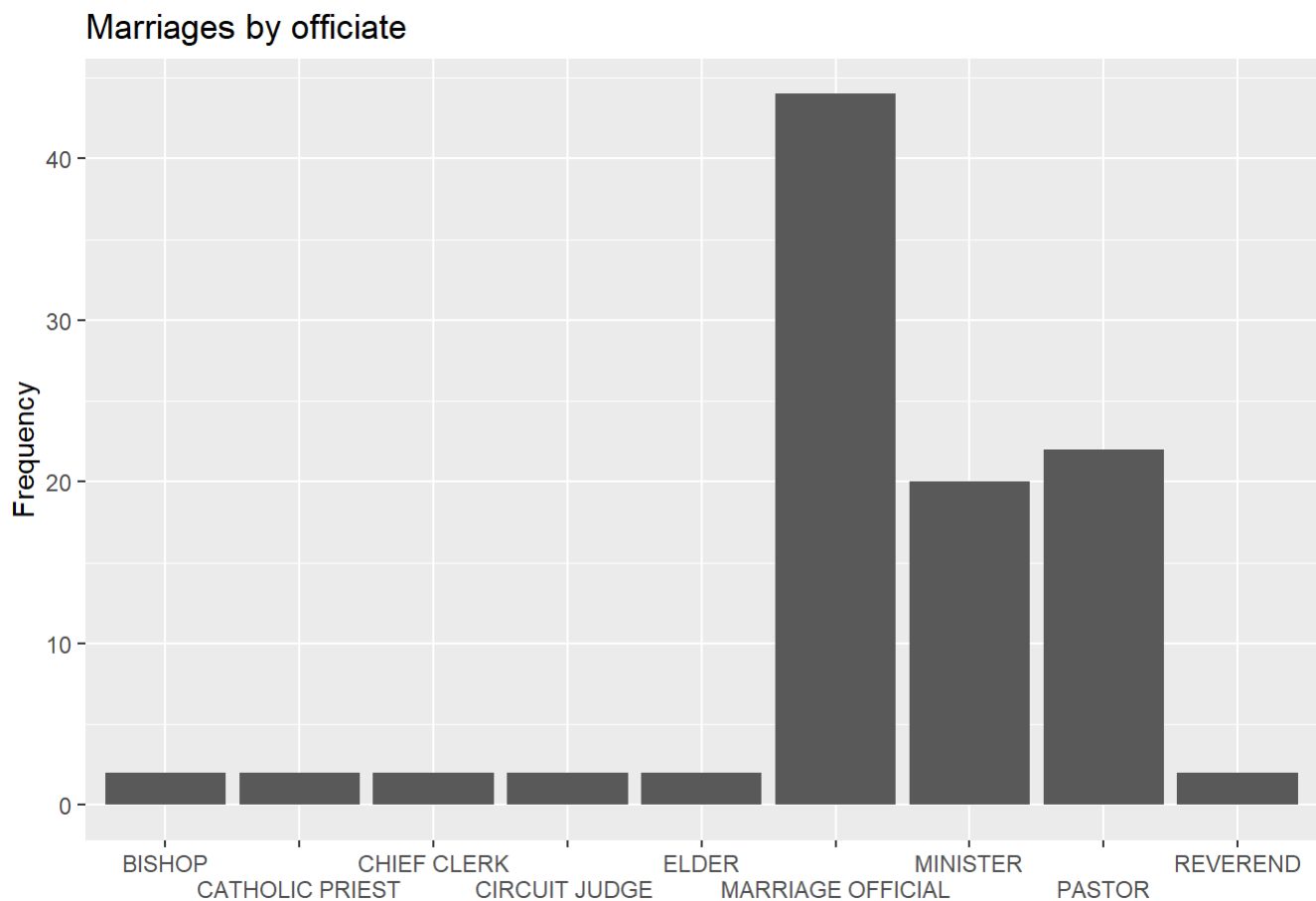


Figure 3.10: Barchart with staggered labels

3.1.2 Pie chart

Pie charts are controversial in statistics. If your goal is to compare the frequency of categories, you are better off with bar charts (humans are better at judging the length of bars than the volume of pie slices). If your goal is compare each category with the the whole (e.g., what portion of participants are Hispanic compared to all participants), and the number of categories is small, then pie charts may work for you. It takes a bit more code to make an attractive pie chart in R.

create a basic ggplot2 pie chart

```
plotdata <- Marriage %>%  
  count(race) %>%  
  arrange(desc(race)) %>%  
  mutate(prop = round(n * 100 / sum(n), 1),  
         lab.ypos = cumsum(prop) - 0.5 * prop)  
  
ggplot(plotdata,  
       aes(x = "",  
          y = prop,  
          fill = race)) +  
  geom_bar(width = 1,  
          stat = "identity",  
          color = "black") +  
  coord_polar("y",  
             start = 0,  
             direction = -1) +  
  theme_void()
```

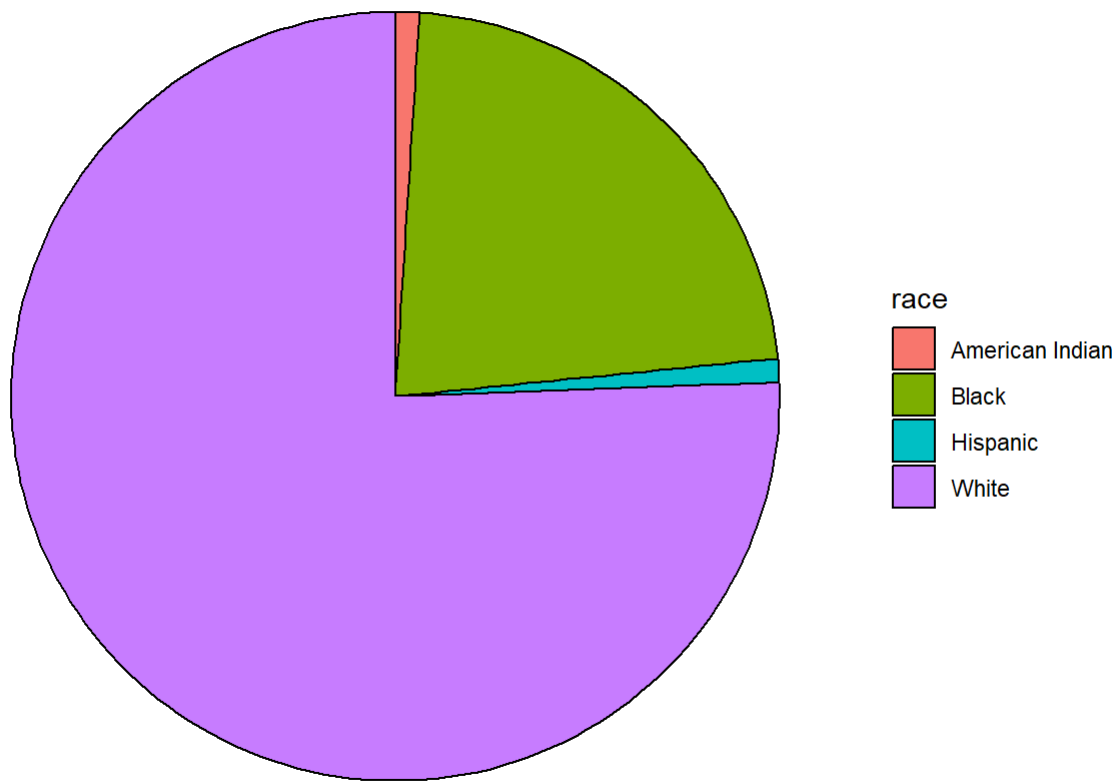


Figure 3.11: Basic pie chart

Now let's get fancy and add labels, while removing the legend.

```

# create a pie chart with slice labels
plotdata <- Marriage %>%
  count(race) %>%
  arrange(desc(race)) %>%
  mutate(prop = round(n*100/sum(n), 1),
         lab.ypos = cumsum(prop) - 0.5*prop)

plotdata$label <- paste0(plotdata$race, "\n",
                        round(plotdata$prop), "%")

ggplot(plotdata,
       aes(x = "",
          y = prop,
          fill = race)) +
  geom_bar(width = 1,
          stat = "identity",
          color = "black") +
  geom_text(aes(y = lab.ypos, label = label),
          color = "black") +
  coord_polar("y",
             start = 0,
             direction = -1) +
  theme_void() +
  theme(legend.position = "FALSE") +
  labs(title = "Participants by race")

```


Participants by race

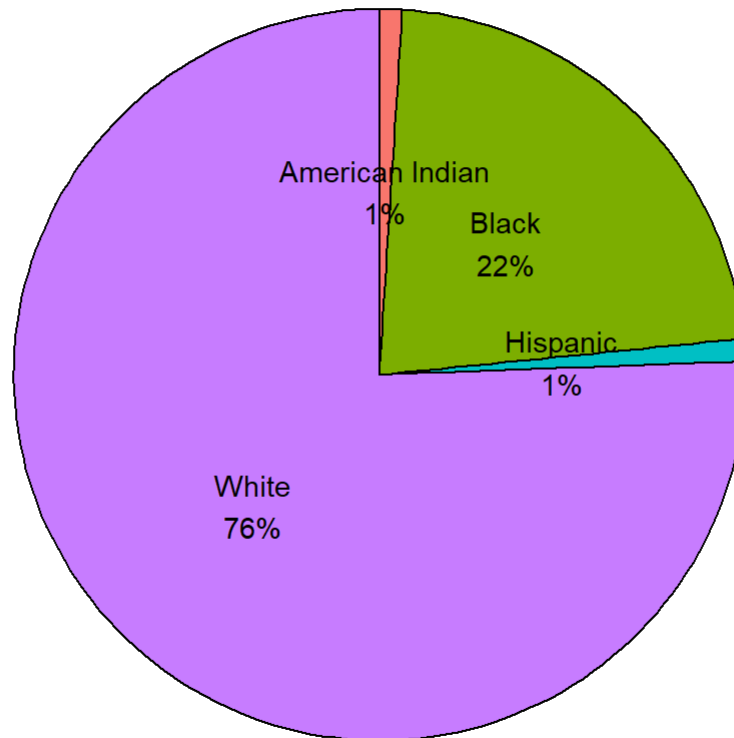


Figure 3.12: Pie chart with percent labels

The pie chart makes it easy to compare each slice with the whole. For example, Black is seen to roughly a quarter of the total participants.

3.1.3 Tree map

An alternative to a pie chart is a tree map. Unlike pie charts, it can handle categorical variables that have *many* levels.

```
library(treemapify)

# create a treemap of marriage officials
plotdata <- Marriage %>%
  count(officialTitle)

ggplot(plotdata,
       aes(fill = officialTitle,
           area = n)) +
  geom_treemap() +
  labs(title = "Marriages by officiate")
```

Marriages by officiate

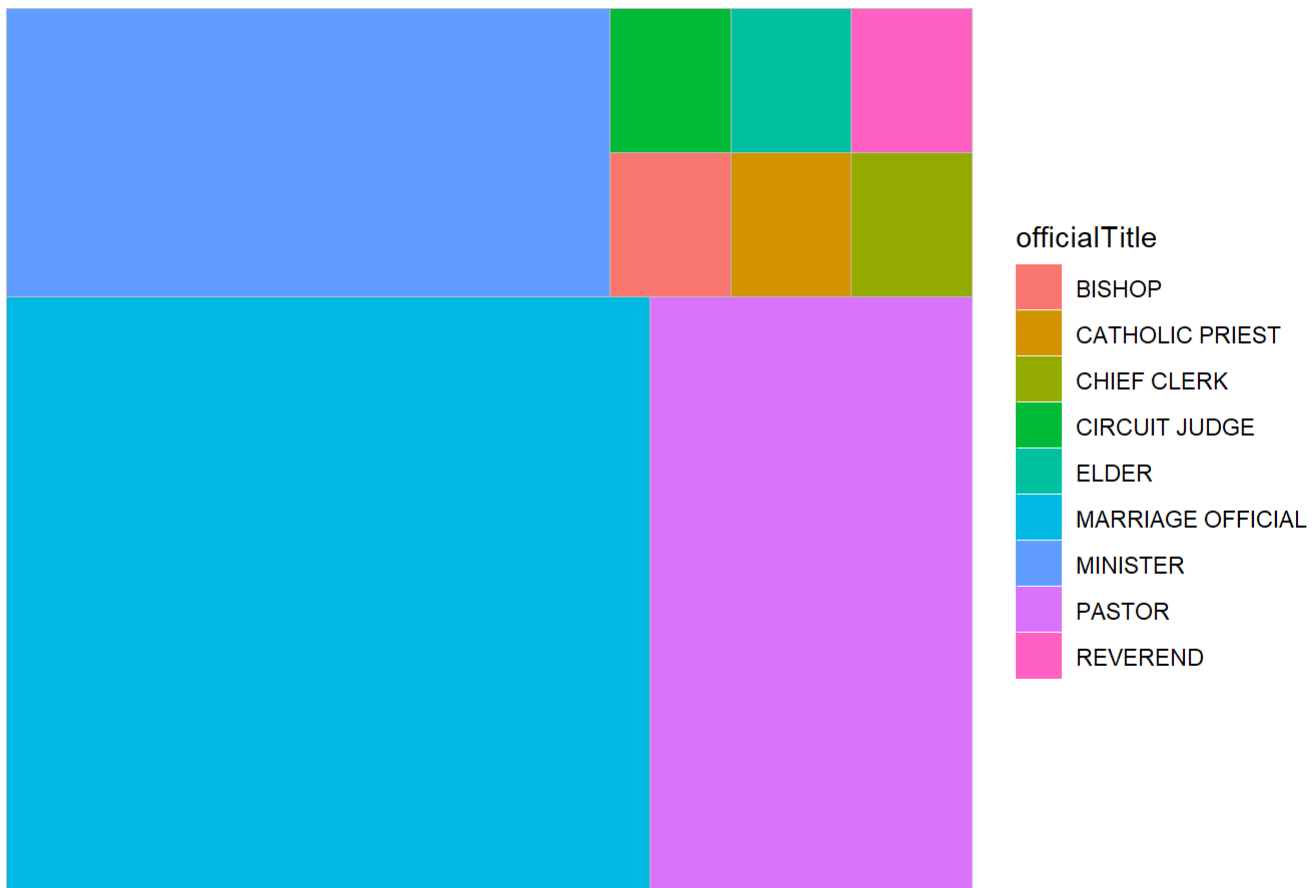


Figure 3.13: Basic treemap

Here is a more useful version with labels.

```

# create a treemap with tile labels
ggplot(plotdata,
      aes(fill = officialTitle,
          area = n,
          label = officialTitle)) +
  geom_treemap() +
  geom_treemap_text(colour = "white",
                   place = "centre") +
  labs(title = "Marriages by officiate") +
  theme(legend.position = "none")

```

Marriages by officiate

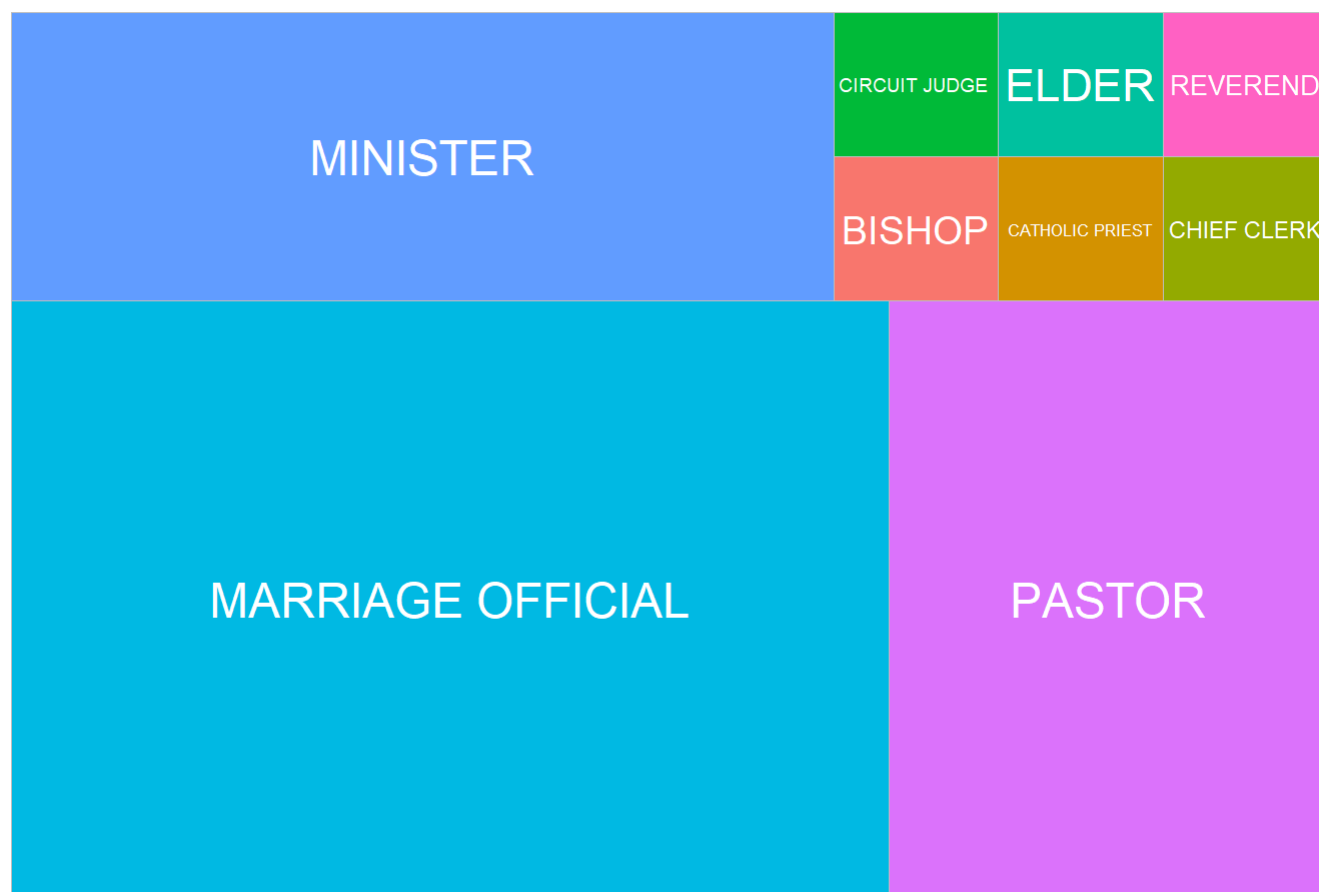


Figure 3.14: Treemap with labels

3.2 Quantitative

The distribution of a single quantitative variable is typically plotted with a histogram, kernel density plot, or dot plot.

3.2.1 Histogram

Using the [Marriage](#) dataset, let's plot the ages of the wedding participants.

```
library(ggplot2)

# plot the age distribution using a histogram
ggplot(Marriage, aes(x = age)) +
  geom_histogram() +
  labs(title = "Participants by age",
       x = "Age")
```

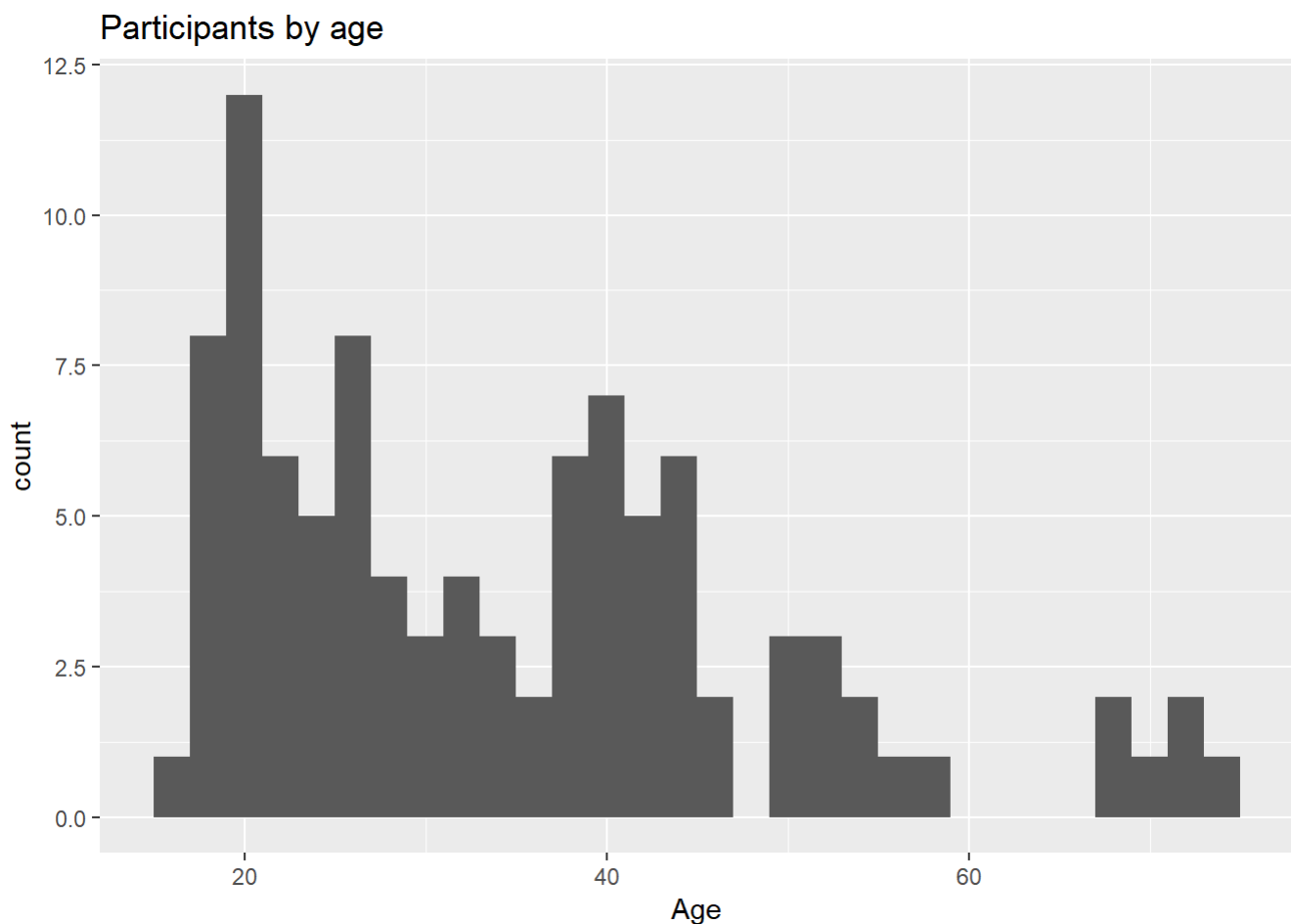


Figure 3.15: Basic histogram

Most participants appear to be in their early 20's with another group in their 40's, and a much smaller group in their later sixties and early seventies. This would be a *multimodal* distribution.

Histogram colors can be modified using two options

- `fill` - fill color for the bars
- `color` - border color around the bars

```
# plot the histogram with blue bars and white borders
```

```
ggplot(Marriage, aes(x = age)) +  
  geom_histogram(fill = "cornflowerblue",  
                 color = "white") +  
  labs(title="Participants by age",  
       x = "Age")
```

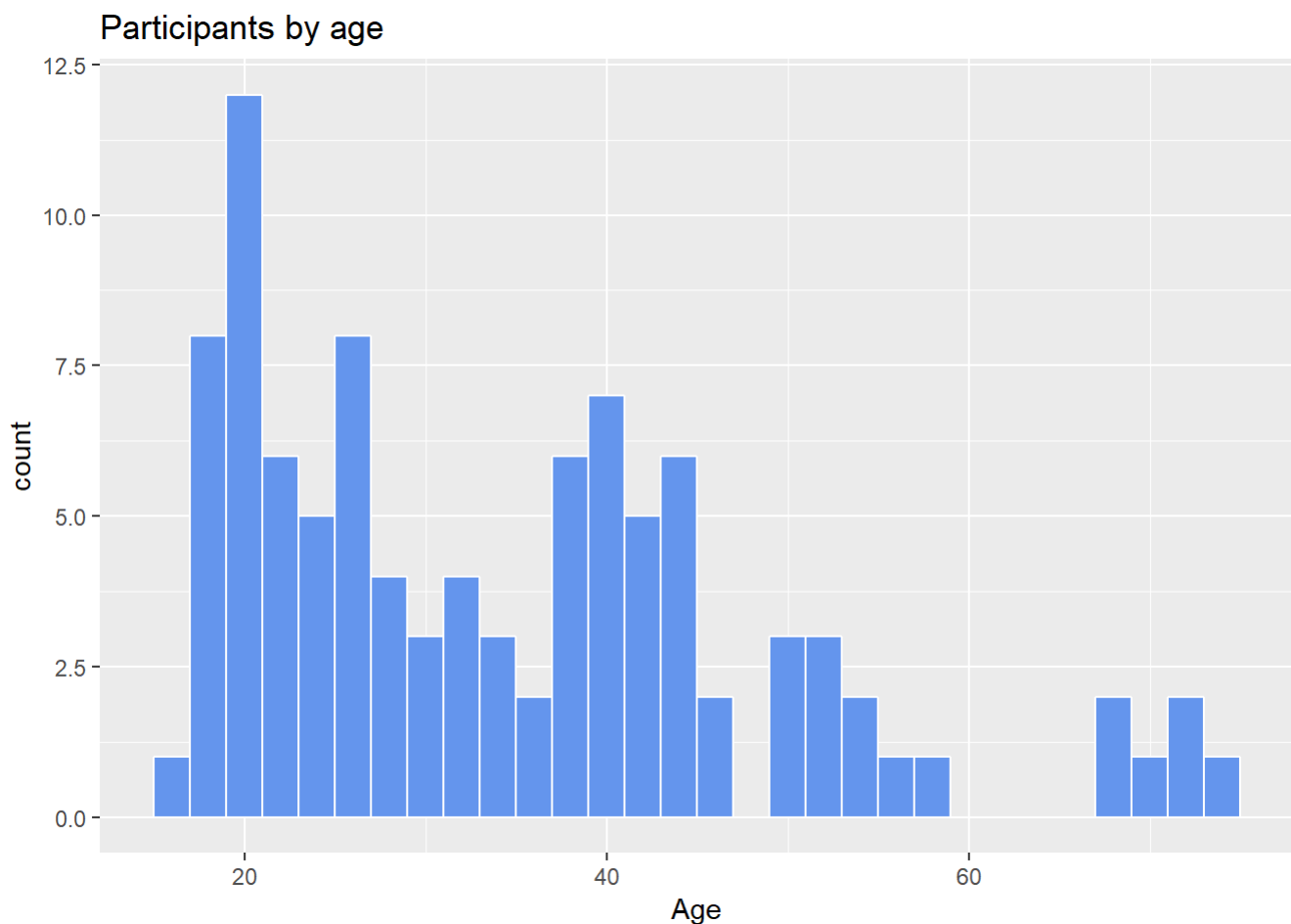


Figure 3.16: Histogram with specified fill and border colors

3.2.1.1 Bins and bandwidths

One of the most important histogram options is `bins`, which controls the number of bins into which the numeric variable is divided (i.e., the number of bars in the plot). The default is 30, but it is helpful to try smaller and larger numbers to get a better impression of the shape of the distribution.

```
# plot the histogram with 20 bins
ggplot(Marriage, aes(x = age)) +
  geom_histogram(fill = "cornflowerblue",
                 color = "white",
                 bins = 20) +
  labs(title="Participants by age",
       subtitle = "number of bins = 20",
       x = "Age")
```

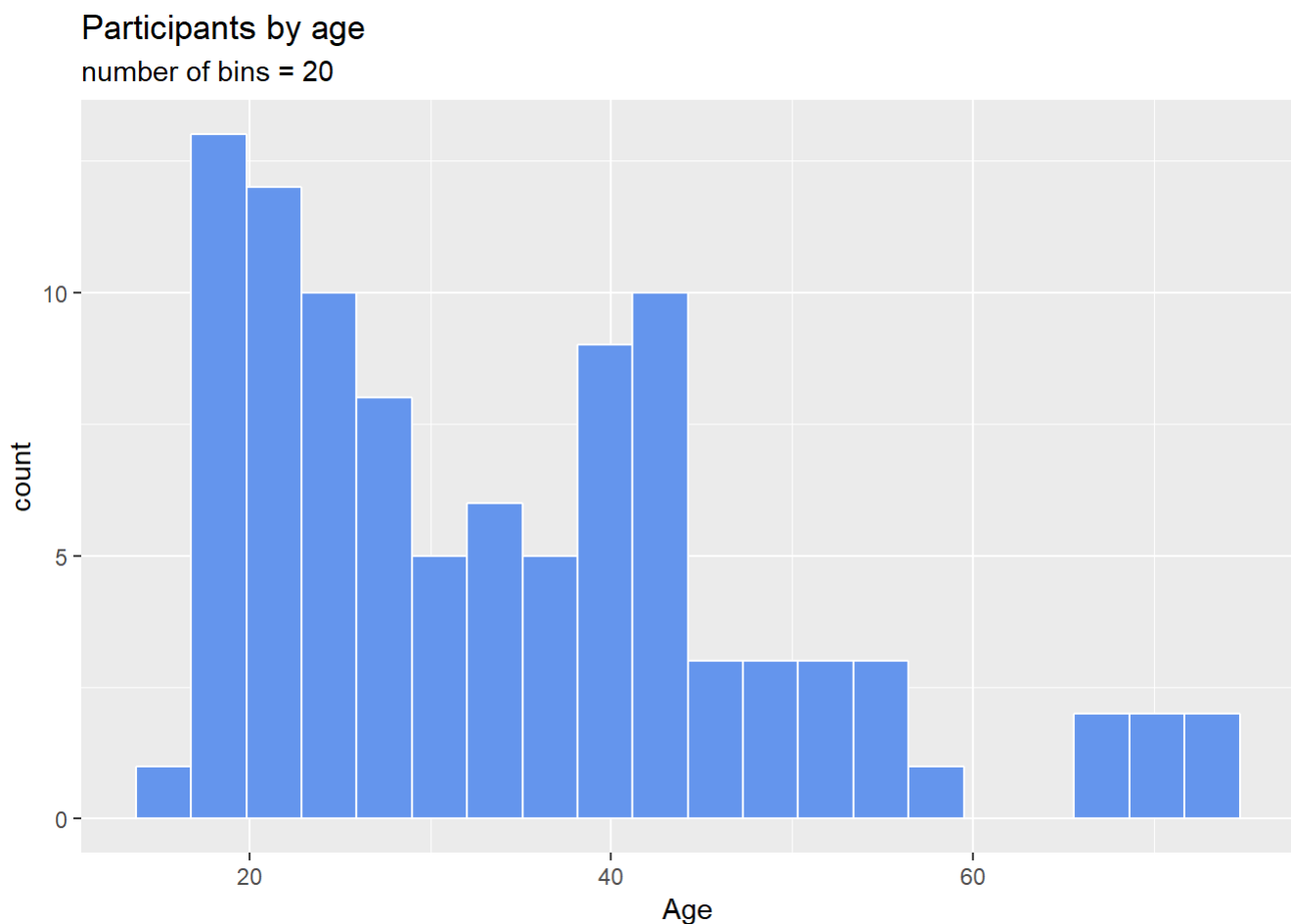


Figure 3.17: Histogram with a specified number of bins

Alternatively, you can specify the `binwidth`, the width of the bins represented by the bars.

```
# plot the histogram with a binwidth of 5
ggplot(Marriage, aes(x = age)) +
  geom_histogram(fill = "cornflowerblue",
                 color = "white",
                 binwidth = 5) +
  labs(title="Participants by age",
        subtitle = "binwidth = 5 years",
        x = "Age")
```

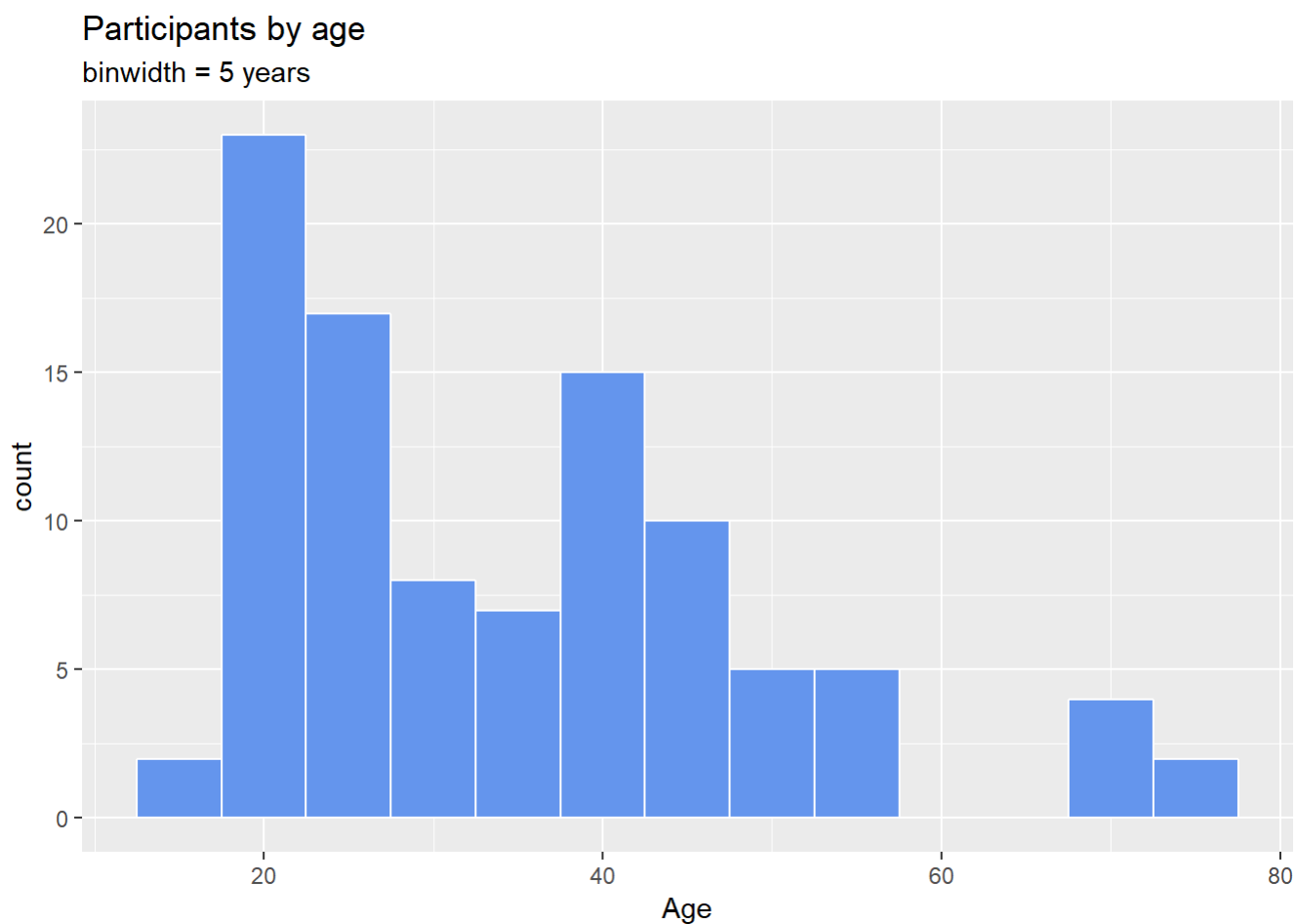


Figure 3.18: Histogram with specified a bin width

As with bar charts, the y-axis can represent counts or percent of the total.

```
# plot the histogram with percentages on the y-axis
library(scales)
ggplot(Marriage,
       aes(x = age,
           y= ..count.. / sum(..count..))) +
  geom_histogram(fill = "cornflowerblue",
                color = "white",
                binwidth = 5) +
  labs(title="Participants by age",
       y = "Percent",
       x = "Age") +
  scale_y_continuous(labels = percent)
```

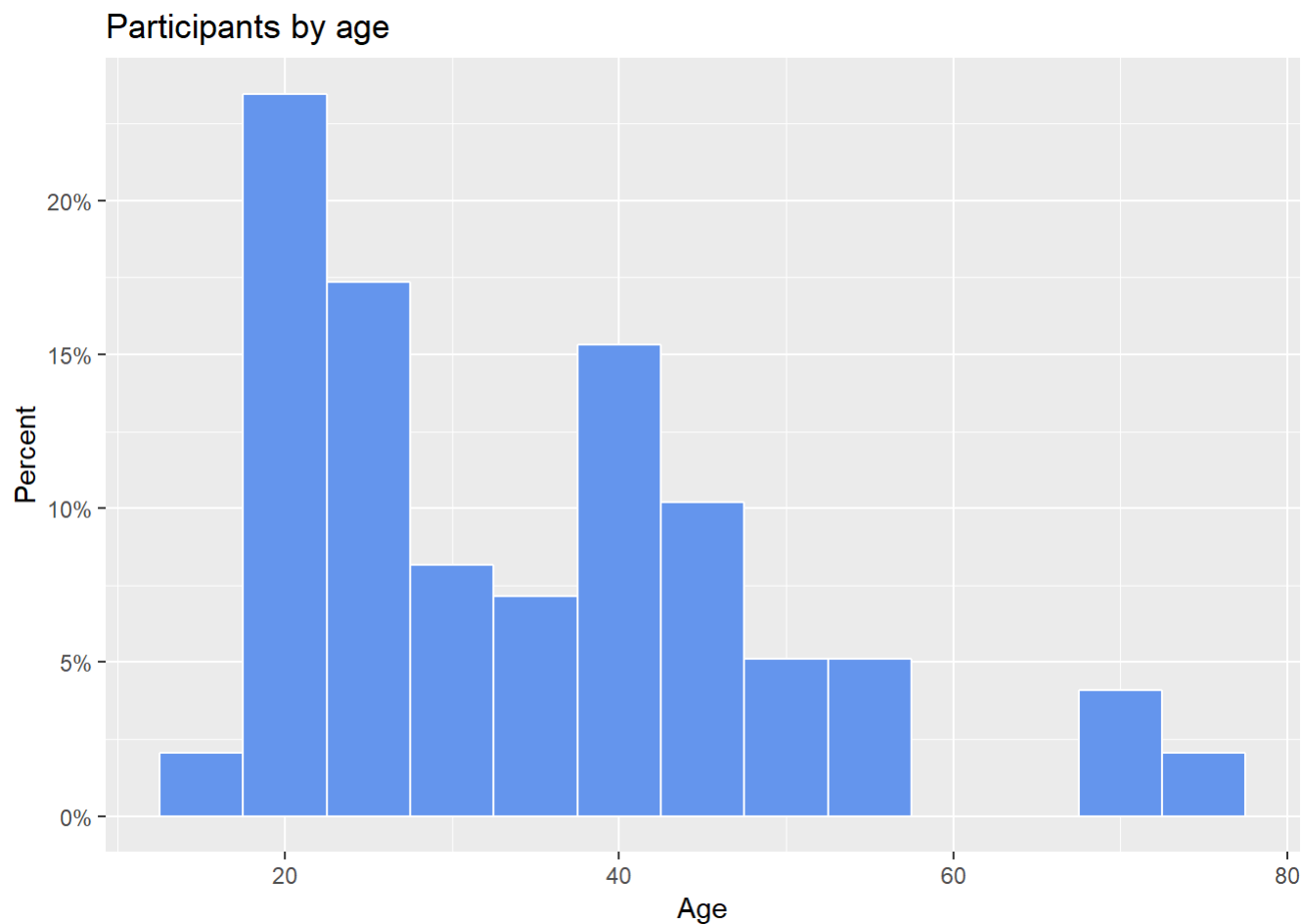


Figure 3.19: Histogram with percentages on the y-axis

3.2.2 Kernel Density plot

An alternative to a histogram is the kernel density plot. Technically, kernel density estimation is a nonparametric method for estimating the probability density function of a continuous random variable. (What??) Basically, we are trying to draw a smoothed histogram, where the area under the curve equals one.

```
# Create a kernel density plot of age
ggplot(Marriage, aes(x = age)) +
  geom_density() +
  labs(title = "Participants by age")
```

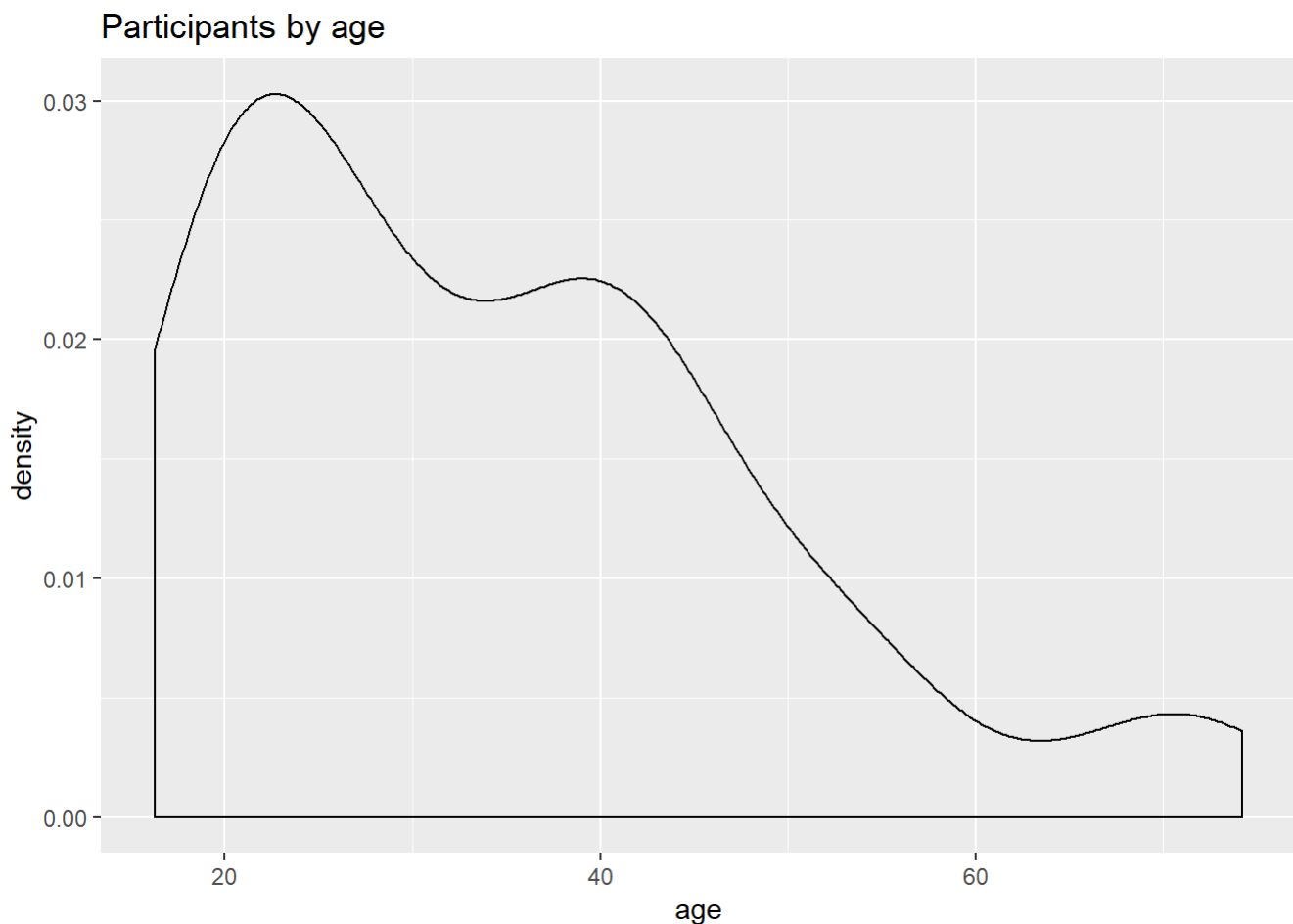


Figure 3.20: Basic kernel density plot

The graph shows the distribution of scores. For example, the proportion of cases between 20 and 40 years old would be represented by the area under the curve between 20 and 40 on the x-axis.

As with previous charts, we can use `fill` and `color` to specify the fill and border colors.

```
# Create a kernel density plot of age
ggplot(Marriage, aes(x = age)) +
  geom_density(fill = "indianred3") +
  labs(title = "Participants by age")
```

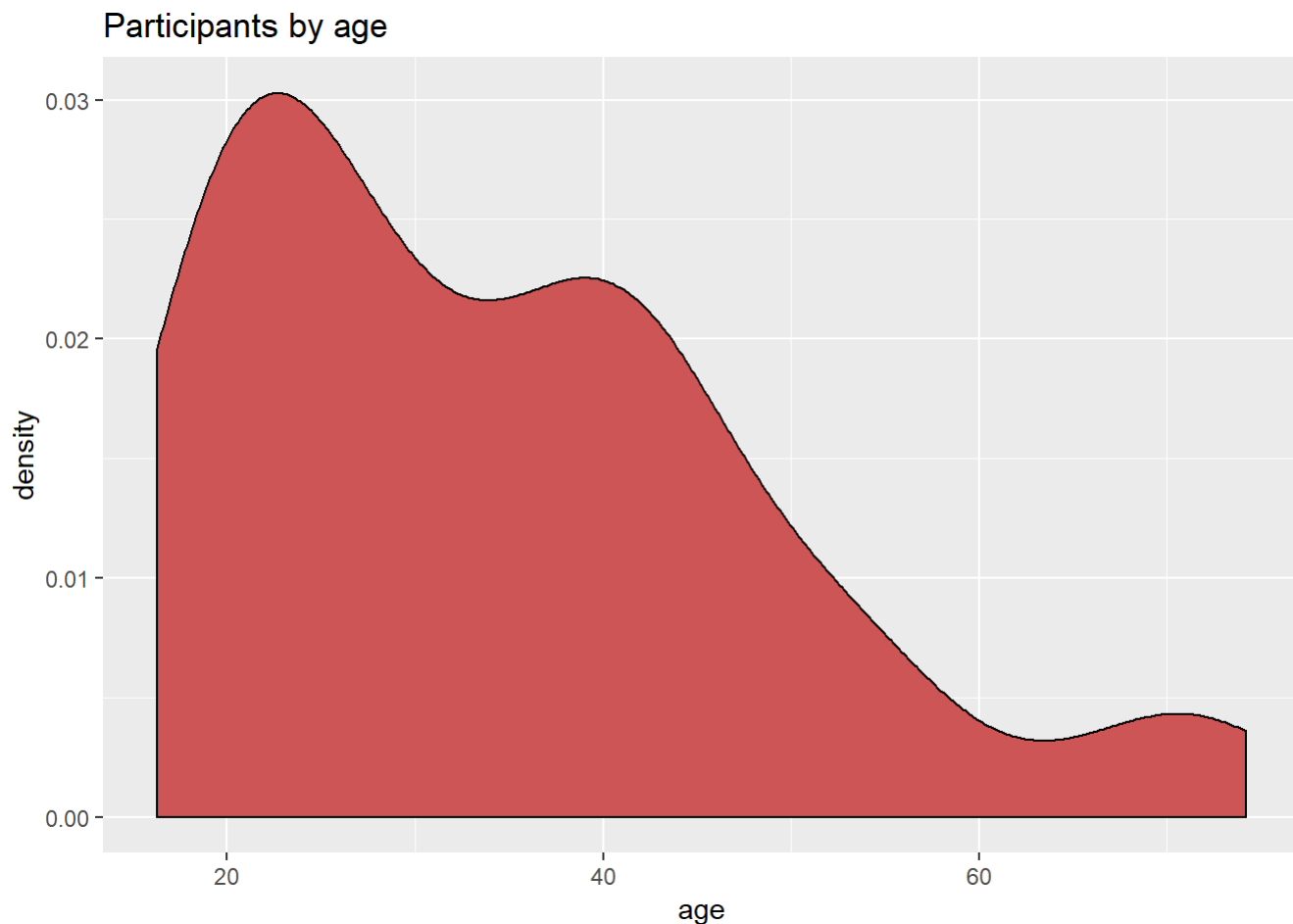


Figure 3.21: Kernel density plot with fill

3.2.2.1 Smoothing parameter

The degree of smoothness is controlled by the bandwidth parameter `bw`. To find the default value for a particular variable, use the `bw.nrd0` function. Values that are larger will result in more smoothing, while values that are smaller will produce less smoothing.

```
# default bandwidth for the age variable
bw.nrd0(Marriage$age)
```

```
## [1] 5.181946
```

```
# Create a kernel density plot of age  
ggplot(Marriage, aes(x = age)) +  
  geom_density(fill = "deepskyblue",  
               bw = 1) +  
  labs(title = "Participants by age",  
        subtitle = "bandwidth = 1")
```

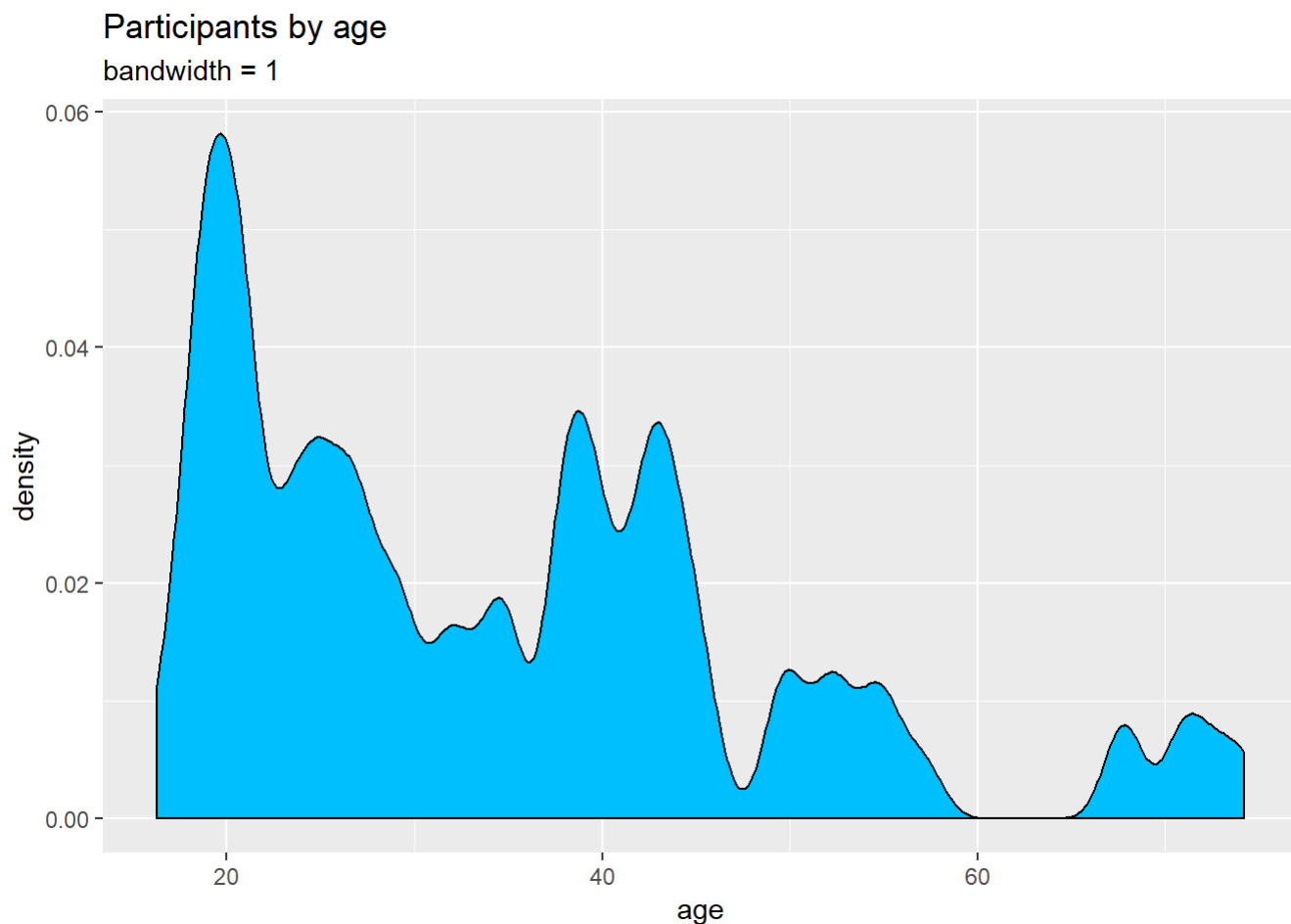


Figure 3.22: Kernel density plot with a specified bandwidth

In this example, the default bandwidth for age is 5.18. Choosing a value of 1 resulted in less smoothing and more detail.

Kernel density plots allow you to easily see which scores are most frequent and which are relatively rare. However it can be difficult to explain the meaning of the y-axis to a non-statistician. (But it will make you look really smart at parties!)

3.2.3 Dot Chart

Another alternative to the histogram is the dot chart. Again, the quantitative variable is divided into bins, but rather than summary bars, each observation is represented by a dot. By default, the width of a dot corresponds to the bin width, and dots are stacked, with each dot representing one observation. This works best when the number of observations is small (say, less than 150).

```
# plot the age distribution using a dotplot  
ggplot(Marriage, aes(x = age)) +  
  geom_dotplot() +  
  labs(title = "Participants by age",  
        y = "Proportion",  
        x = "Age")
```

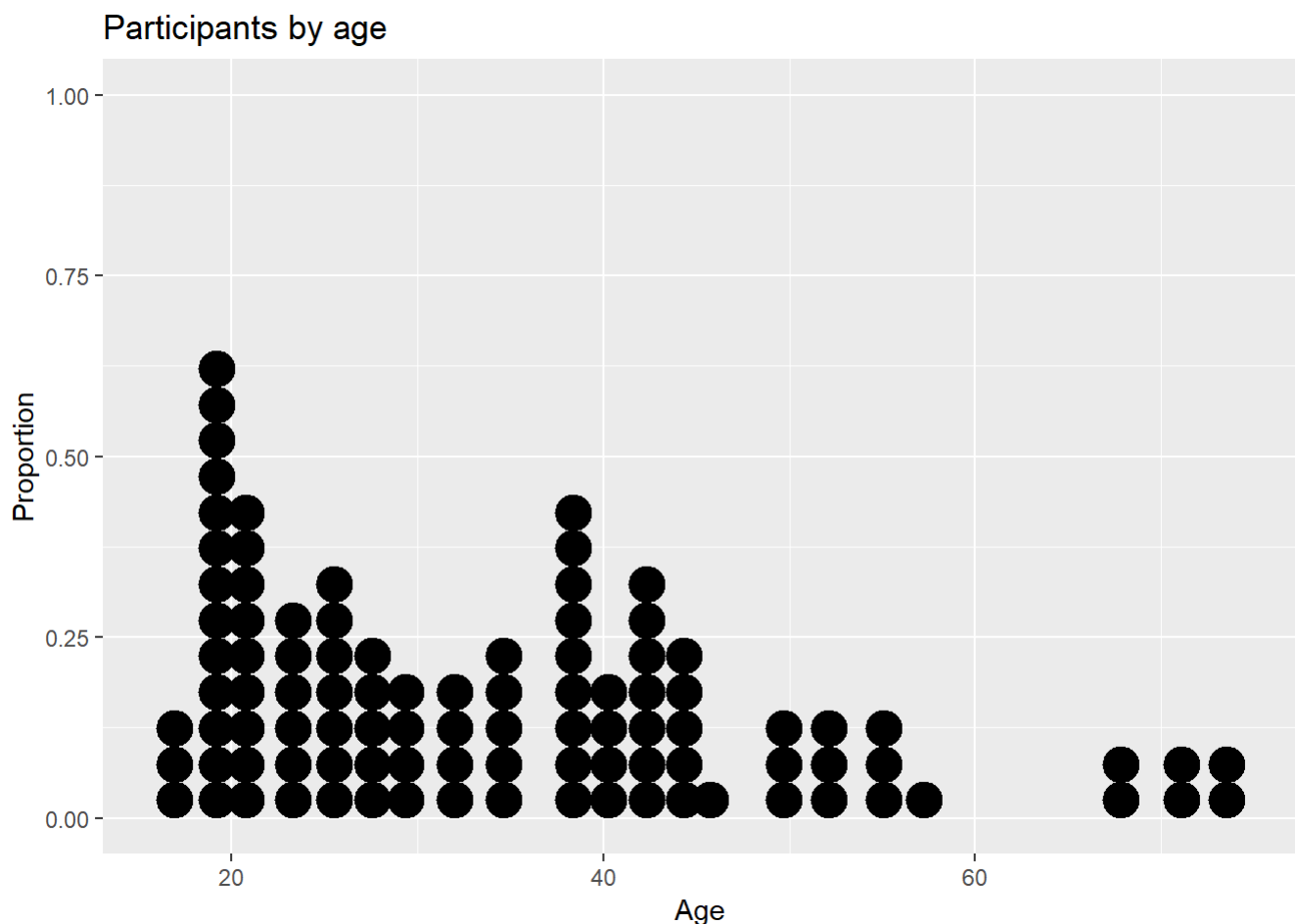


Figure 3.23: Basic dotplot

The `fill` and `color` options can be used to specify the fill and border color of each dot respectively.

```
# Plot ages as a dot plot using
# gold dots with black borders
ggplot(Marriage, aes(x = age)) +
  geom_dotplot(fill = "gold",
               color = "black") +
  labs(title = "Participants by age",
       y = "Proportion",
       x = "Age")
```

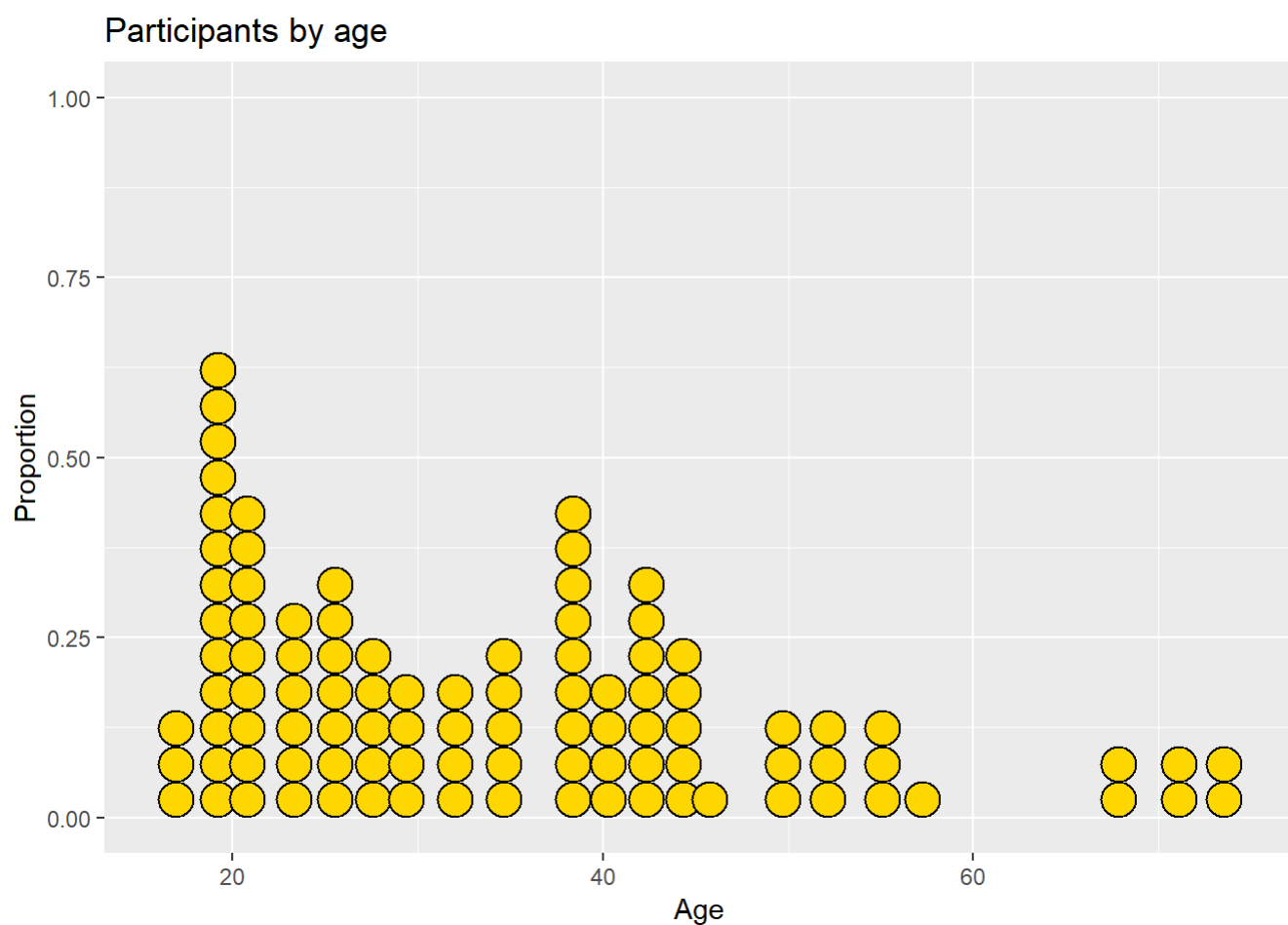


Figure 3.24: Dotplot with a specified color scheme

There are many more options available. See the [help](#) for details and examples.

Chapter 4 Bivariate Graphs

Bivariate graphs display the relationship between two variables. The type of graph will depend on the measurement level of the variables (categorical or quantitative).

4.1 Categorical vs. Categorical

When plotting the relationship between two categorical variables, stacked, grouped, or segmented bar charts are typically used. A less common approach is the [mosaic](#) chart.

4.1.1 Stacked bar chart

Let's plot the relationship between automobile class and drive type (front-wheel, rear-wheel, or 4-wheel drive) for the automobiles in the [Fuel economy](#) dataset.

```
library(ggplot2)

# stacked bar chart
ggplot(mpg,
       aes(x = class,
           fill = drv)) +
  geom_bar(position = "stack")
```

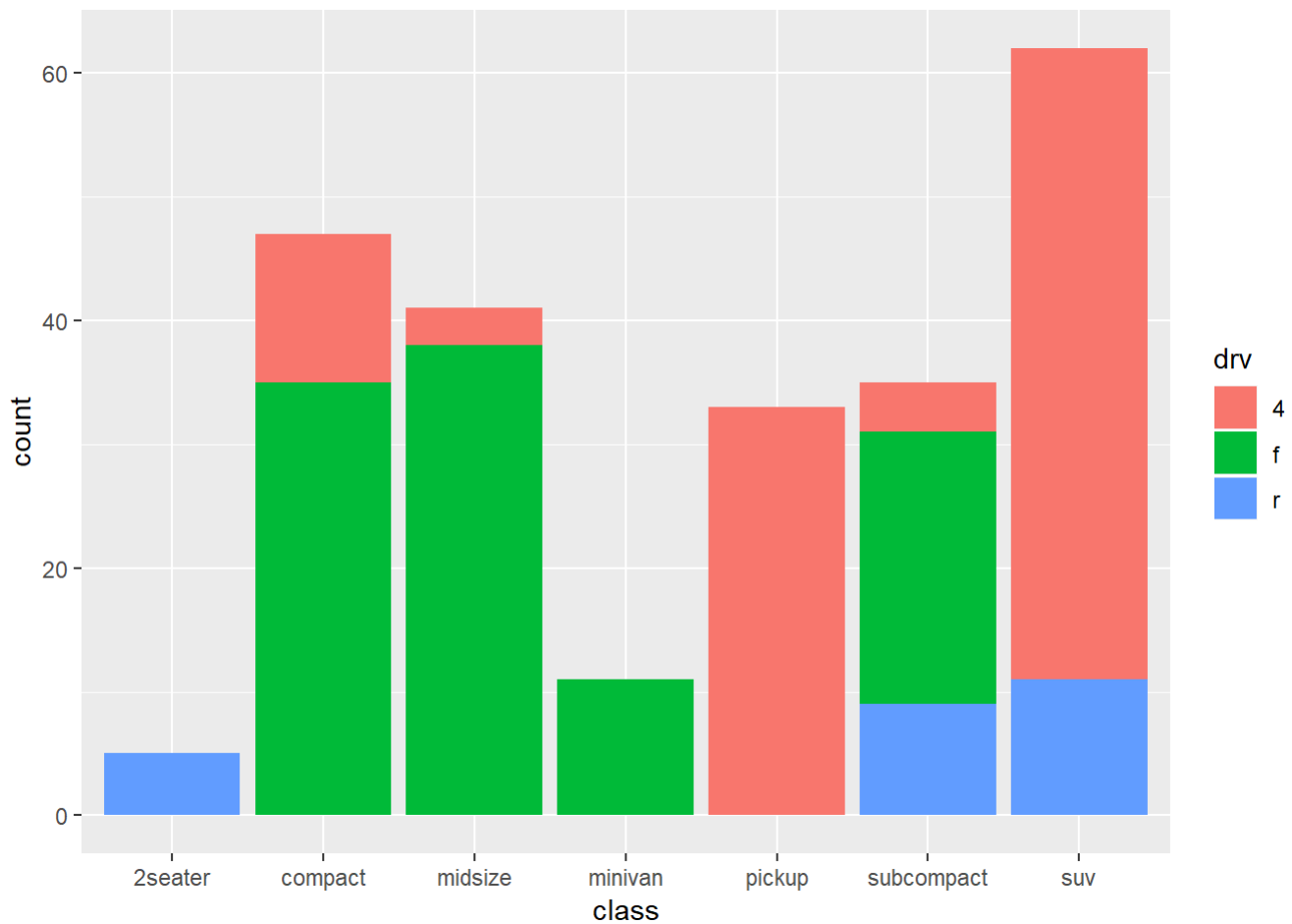


Figure 4.1: Stacked bar chart

From the chart, we can see for example, that the most common vehicle is the SUV. All 2seater cars are rear wheel drive, while most, but not all SUVs are 4-wheel drive.

Stacked is the default, so the last line could have also been written as `geom_bar()` .

4.1.2 Grouped bar chart

Grouped bar charts place bars for the second categorical variable side-by-side. To create a grouped bar plot use the `position = "dodge"` option.

```
library(ggplot2)

# grouped bar plot
ggplot(mpg,
       aes(x = class,
           fill = drv)) +
  geom_bar(position = "dodge")
```

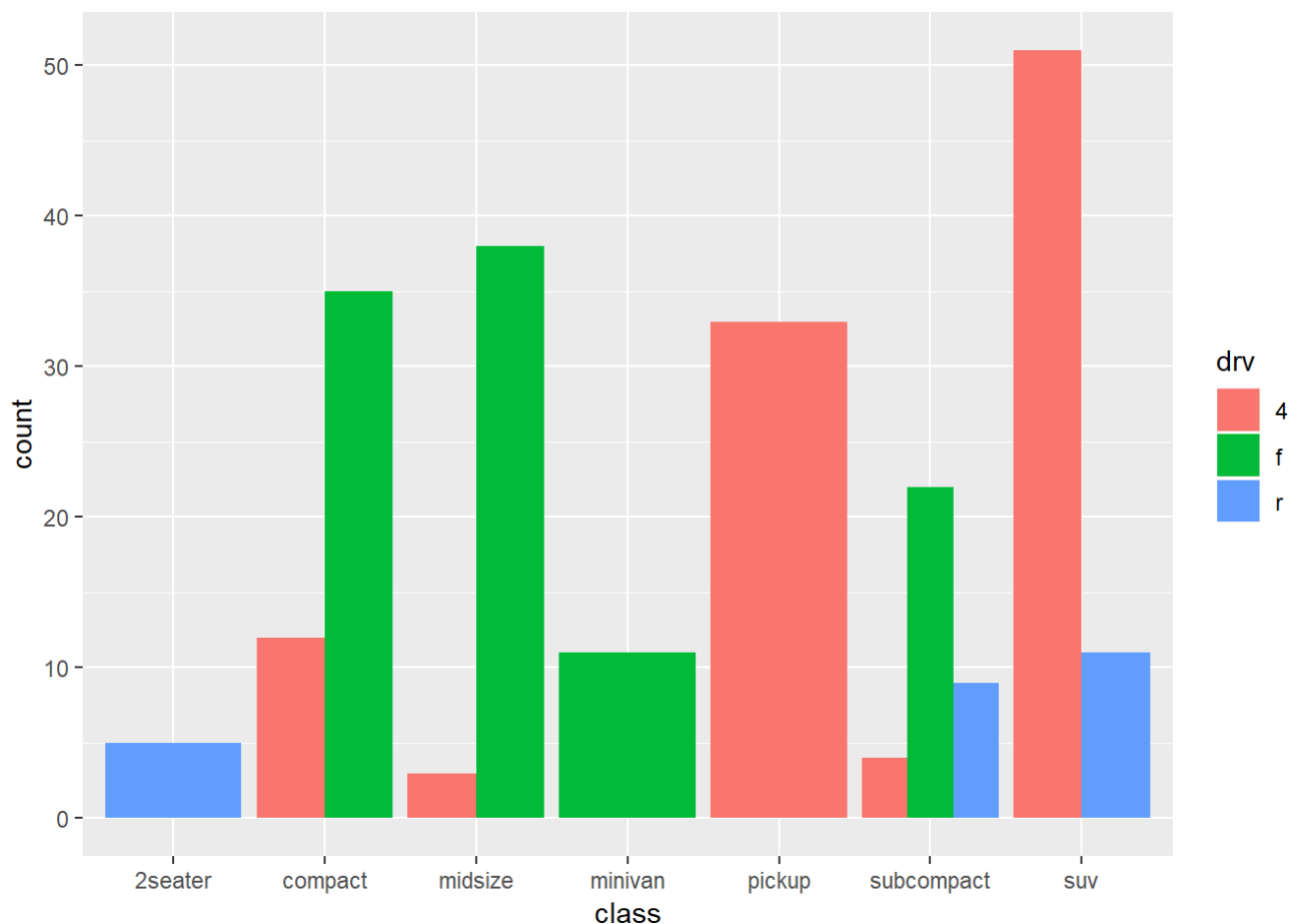


Figure 4.2: Side-by-side bar chart

Notice that all Minivans are front-wheel drive. By default, zero count bars are dropped and the remaining bars are made wider. This may not be the behavior you want. You can modify this using the `position = position_dodge(preserve = "single")` option.


```
library(ggplot2)

# grouped bar plot preserving zero count bars
ggplot(mpg,
       aes(x = class,
           fill = drv)) +
  geom_bar(position = position_dodge(preserve = "single"))
```

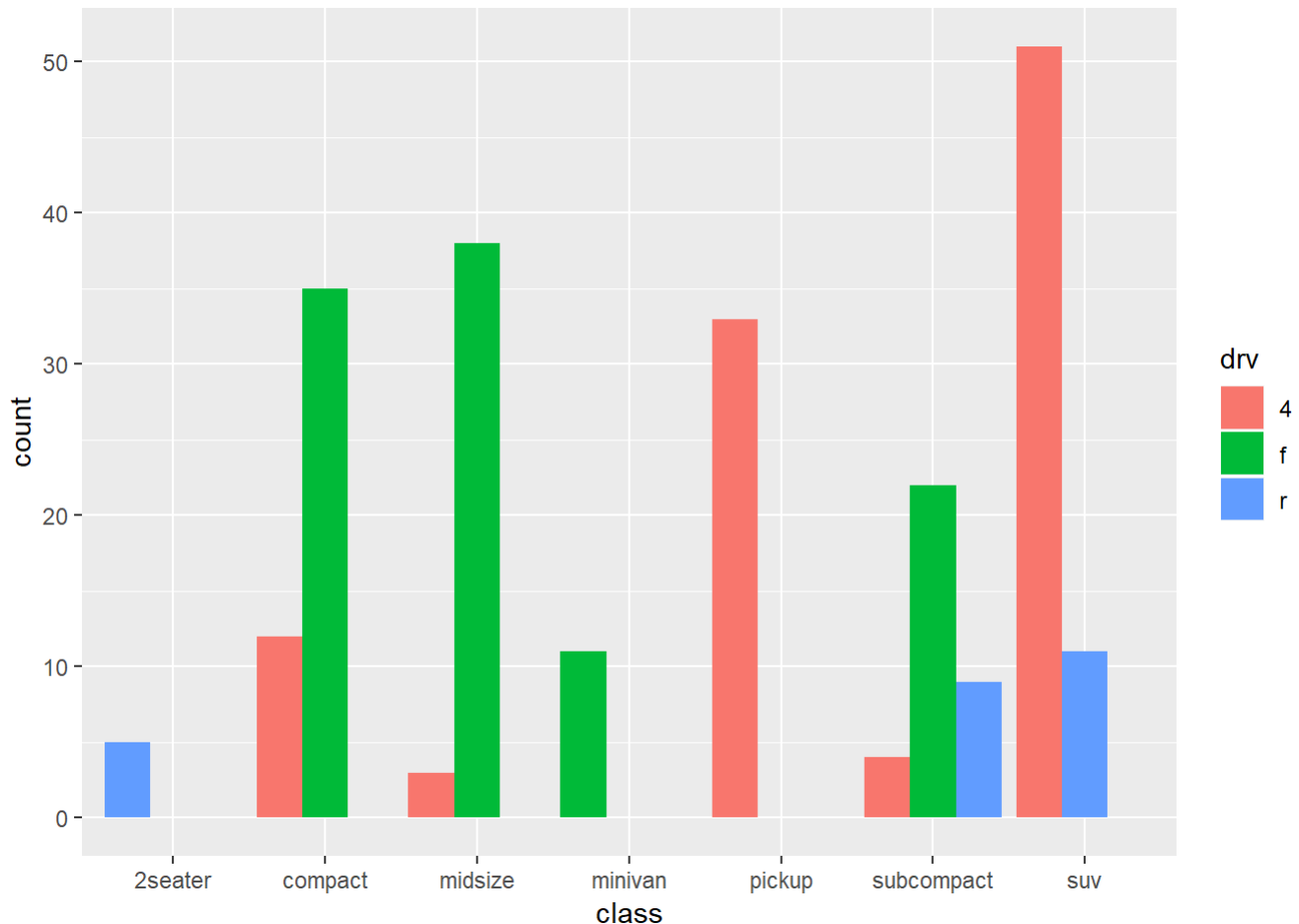


Figure 4.3: Side-by-side bar chart with zero count bars retained

Note that this option is only available in the latest development version of `ggplot2`, but should be generally available shortly.

4.1.3 Segmented bar chart

A segmented bar plot is a stacked bar plot where each bar represents 100 percent. You can create a segmented bar chart using the `position = "filled"` option.

```
library(ggplot2)

# bar plot, with each bar representing 100%
ggplot(mpg,
       aes(x = class,
           fill = drv)) +
  geom_bar(position = "fill") +
  labs(y = "Proportion")
```

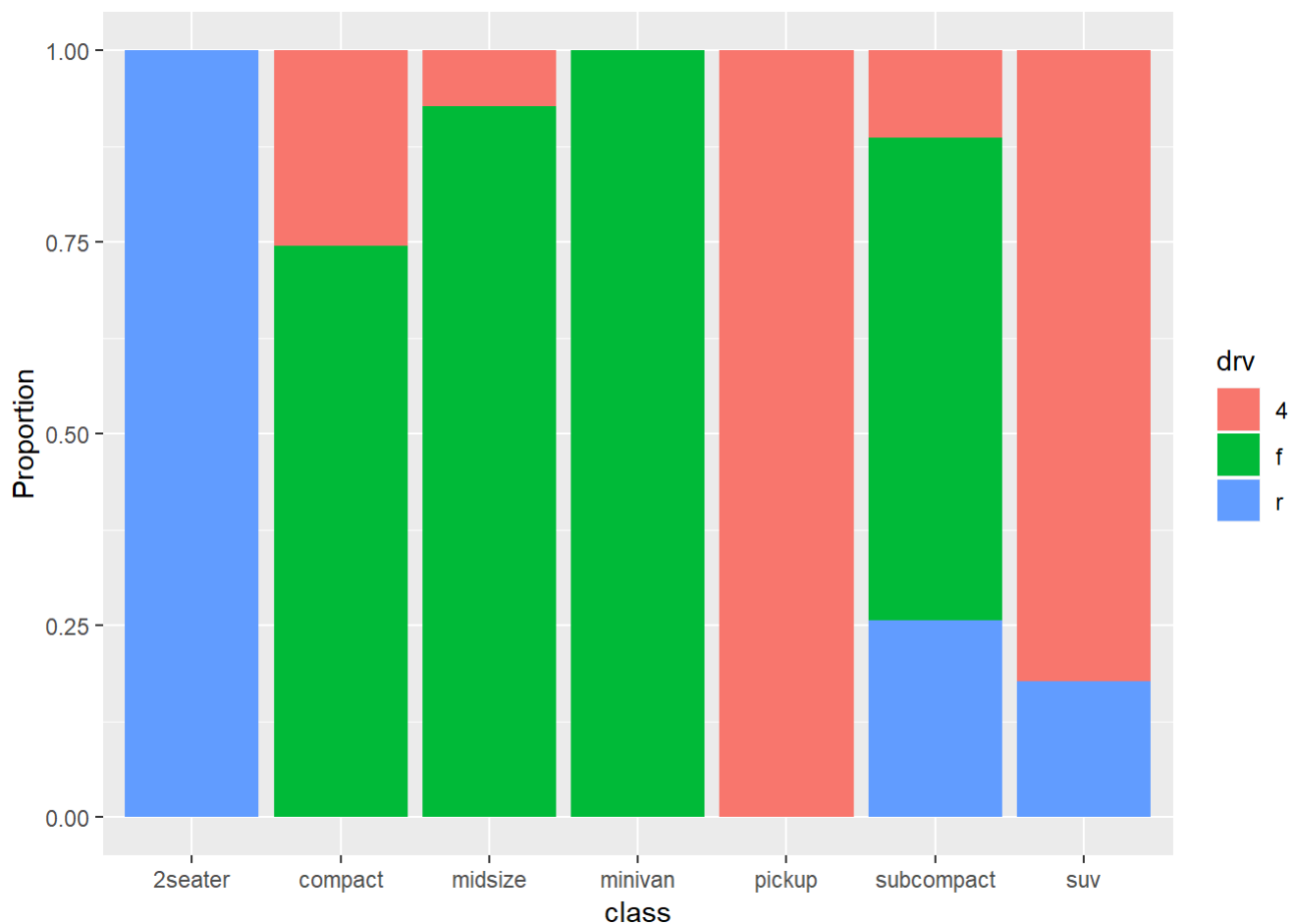


Figure 4.4: Segmented bar chart

This type of plot is particularly useful if the goal is to compare the percentage of a category in one variable across each level of another variable. For example, the proportion of front-wheel drive cars go up as you move from compact, to midsize, to minivan.

4.1.4 Improving the color and labeling

You can use additional options to improve color and labeling. In the graph below

- `factor` modifies the order of the categories for the class variable and both the order and the labels for the drive variable
- `scale_y_continuous` modifies the y-axis tick mark labels
- `labs` provides a title and changed the labels for the x and y axes and the legend
- `scale_fill_brewer` changes the fill color scheme
- `theme_minimal` removes the grey background and changed the grid color

```
library(ggplot2)

# bar plot, with each bar representing 100%,
# reordered bars, and better labels and colors
library(scales)
ggplot(mpg,
       aes(x = factor(class,
                      levels = c("2seater", "subcompact",
                                "compact", "midsize",
                                "minivan", "suv", "pickup")),
          fill = factor(drv,
                      levels = c("f", "r", "4"),
                      labels = c("front-wheel",
                                "rear-wheel",
                                "4-wheel")))) +
  geom_bar(position = "fill") +
  scale_y_continuous(breaks = seq(0, 1, .2),
                    label = percent) +
  scale_fill_brewer(palette = "Set2") +
  labs(y = "Percent",
       fill = "Drive Train",
       x = "Class",
       title = "Automobile Drive by Class") +
  theme_minimal()
```

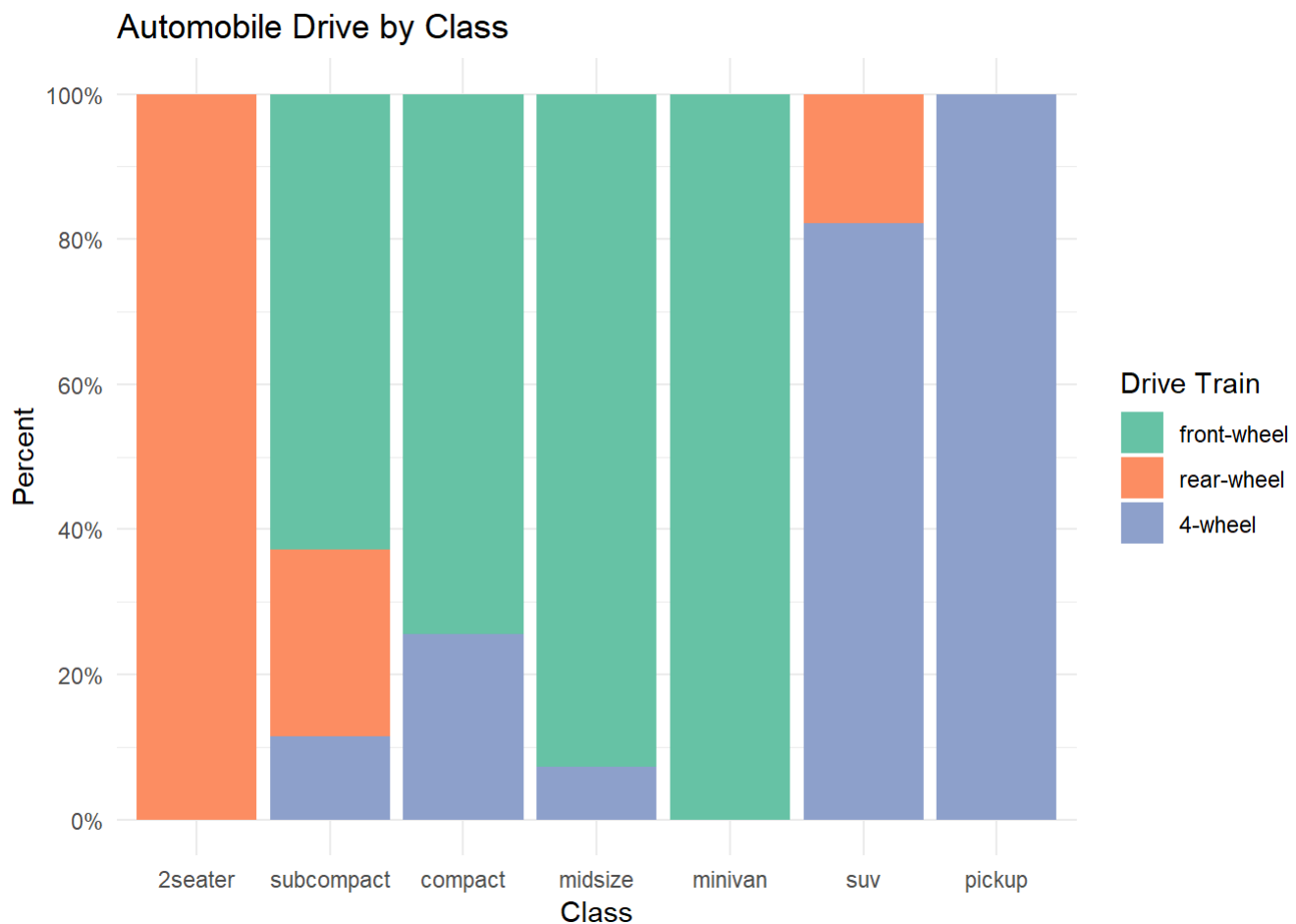


Figure 4.5: Segmented bar chart with improved labeling and color

In the graph above, the `factor` function was used to reorder and/or rename the levels of a categorical variable. You could also apply this to the original dataset, making these changes permanent. It would then apply to all future graphs using that dataset. For example:

```
# change the order the levels for the categorical variable "class"
mpg$class = factor(mpg$class,
  levels = c("2seater", "subcompact",
             "compact", "midsize",
             "minivan", "suv", "pickup"))
```

I placed the `factor` function within the `ggplot` function to demonstrate that, if desired, you can change the order of the categories and labels for the categories for a single graph.

The other functions are discussed more fully in the section on [Customizing](#) graphs.

Next, let's add percent labels to each segment. First, we'll create a summary dataset that has the necessary labels.

```

# create a summary dataset
library(dplyr)
plotdata <- mpg %>%
  group_by(class, drv) %>%
  summarize(n = n()) %>%
  mutate(pct = n/sum(n),
         lbl = scales::percent(pct))

plotdata

## # A tibble: 12 x 5
## # Groups:   class [7]
##   class      drv      n    pct lbl
##   <chr>    <chr> <int> <dbl> <chr>
## 1 2seater    r         5 1.00  100%
## 2 compact   4        12 0.255  25.5%
## 3 compact   f        35 0.745  74.5%
## 4 midsize    4         3 0.0732 7.3%
## 5 midsize    f        38 0.927  92.7%
## 6 minivan    f        11 1.00  100%
## 7 pickup     4        33 1.00  100%
## 8 subcompact 4         4 0.114  11.4%
## 9 subcompact f        22 0.629  62.9%
## 10 subcompact r         9 0.257  25.7%
## 11 suv       4        51 0.823  82.3%
## 12 suv       r        11 0.177  17.7%

```

Next, we'll use this dataset and the `geom_text` function to add labels to each bar segment.

```

# create segmented bar chart
# adding labels to each segment

ggplot(plotdata,
      aes(x = factor(class,
                     levels = c("2seater", "subcompact",
                                "compact", "midsize",
                                "minivan", "suv", "pickup")),
          y = pct,
          fill = factor(drv,
                       levels = c("f", "r", "4"),
                       labels = c("front-wheel",
                                  "rear-wheel",
                                  "4-wheel")))) +
  geom_bar(stat = "identity",
          position = "fill") +
  scale_y_continuous(breaks = seq(0, 1, .2),
                    label = percent) +
  geom_text(aes(label = lbl),
            size = 3,
            position = position_stack(vjust = 0.5)) +
  scale_fill_brewer(palette = "Set2") +
  labs(y = "Percent",
       fill = "Drive Train",
       x = "Class",
       title = "Automobile Drive by Class") +
  theme_minimal()

```

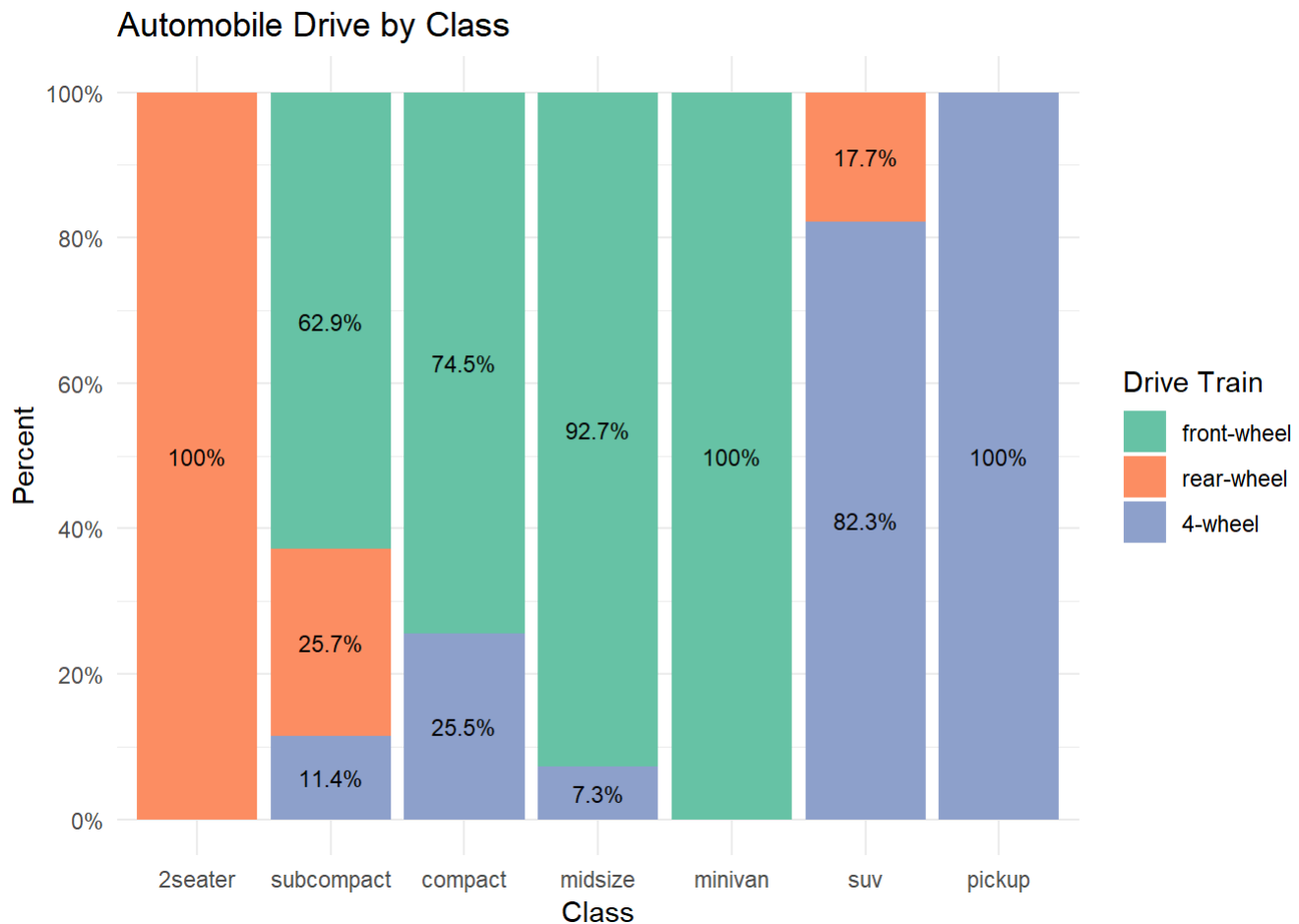


Figure 4.6: Segmented bar chart with value labeling

Now we have a graph that is easy to read and interpret.

4.1.5 Other plots

[Mosaic plots](#) provide an alternative to stacked bar charts for displaying the relationship between categorical variables. They can also provide more sophisticated statistical information.

4.2 Quantitative vs. Quantitative

The relationship between two quantitative variables is typically displayed using scatterplots and line graphs.

4.2.1 Scatterplot

The simplest display of two quantitative variables is a scatterplot, with each variable represented on an axis. For example, using the [Salaries](#) dataset, we can plot experience (*yrs.since.phd*) vs. academic salary (*salary*) for college professors.

```
library(ggplot2)
data(Salaries, package="carData")

# simple scatterplot
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point()
```

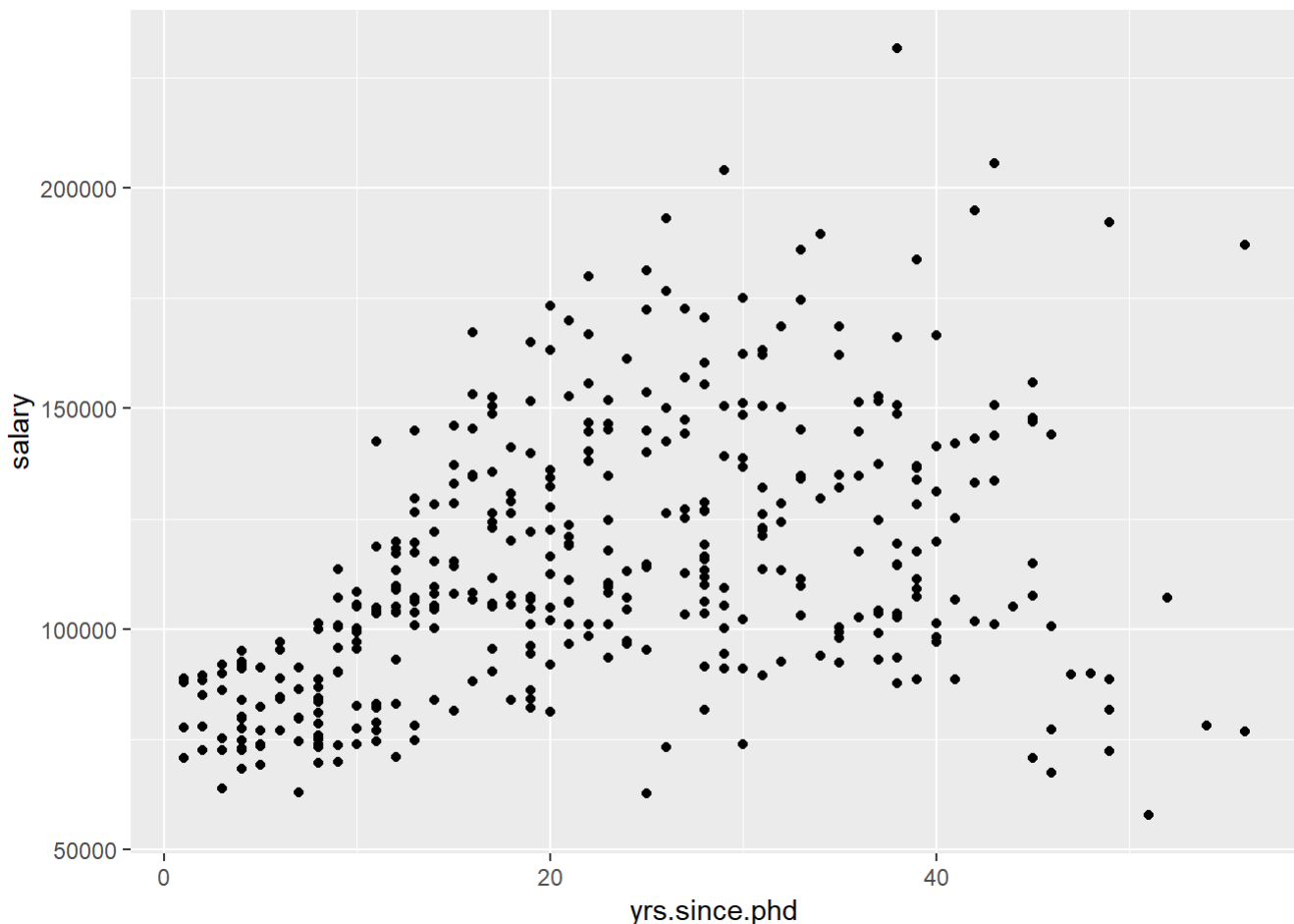


Figure 4.7: Simple scatterplot

`geom_point` options can be used to change the

- `color` - point color

- `size` - point size
- `shape` - point shape
- `alpha` - point transparency. Transparency ranges from 0 (transparent) to 1 (opaque), and is a useful parameter when points overlap.

The functions `scale_x_continuous` and `scale_y_continuous` control the scaling on x and y axes respectively.

See [Customizing](#) graphs for details.

We can use these options and functions to create a more attractive scatterplot.

```
# enhanced scatter plot
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point(color="cornflowerblue",
             size = 2,
             alpha=.8) +
  scale_y_continuous(label = scales::dollar,
                    limits = c(50000, 250000)) +
  scale_x_continuous(breaks = seq(0, 60, 10),
                    limits=c(0, 60)) +
  labs(x = "Years Since PhD",
       y = "",
       title = "Experience vs. Salary",
       subtitle = "9-month salary for 2008-2009")
```

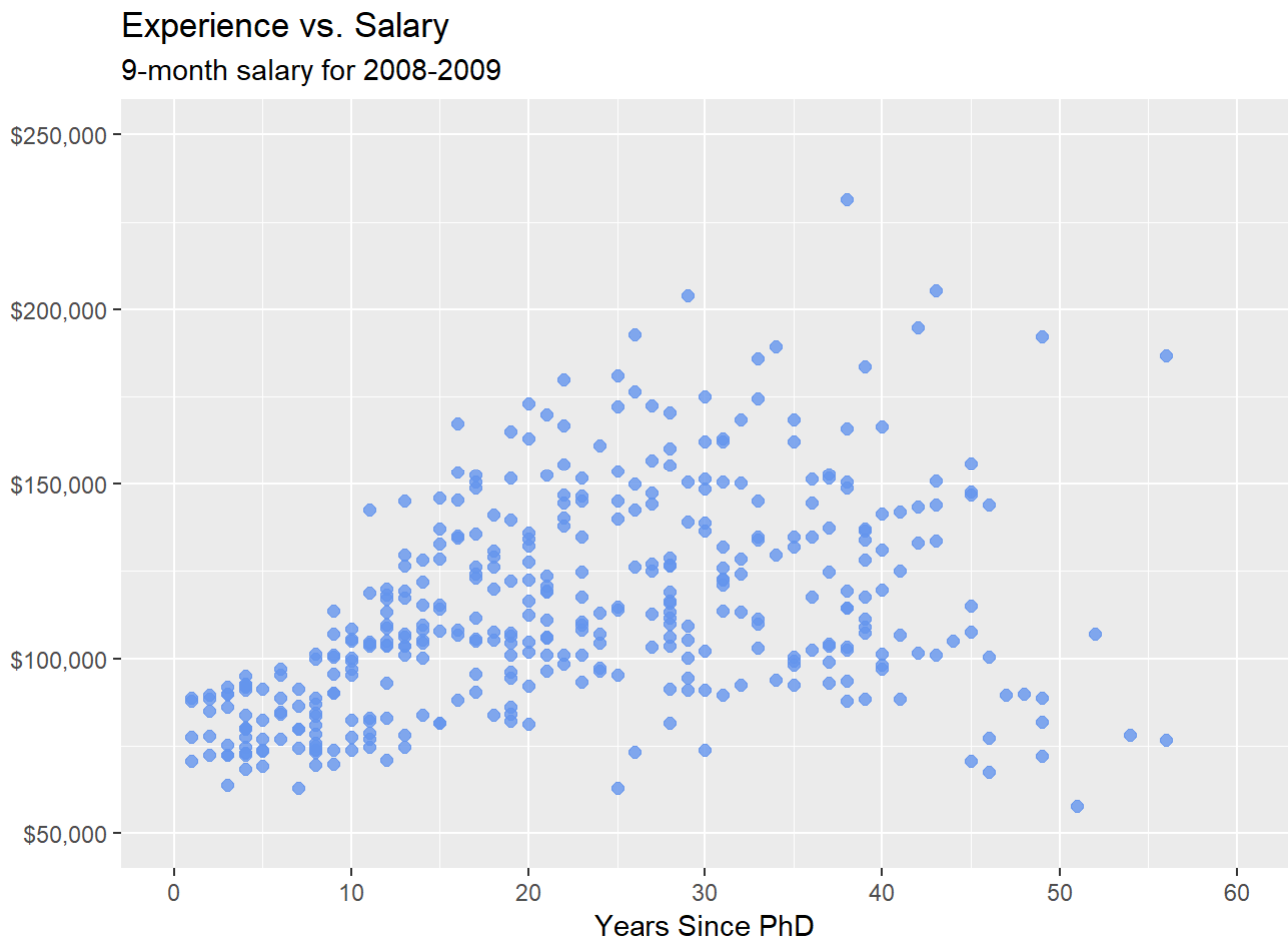


Figure 4.8: Scatterplot with color, transparency, and axis scaling

4.2.1.1 Adding best fit lines

It is often useful to summarize the relationship displayed in the scatterplot, using a best fit line. Many types of lines are supported, including linear, polynomial, and nonparametric (loess). By default, 95% confidence limits for these lines are displayed.

```
# scatterplot with linear fit line
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point(color= "steelblue") +
  geom_smooth(method = "lm")
```

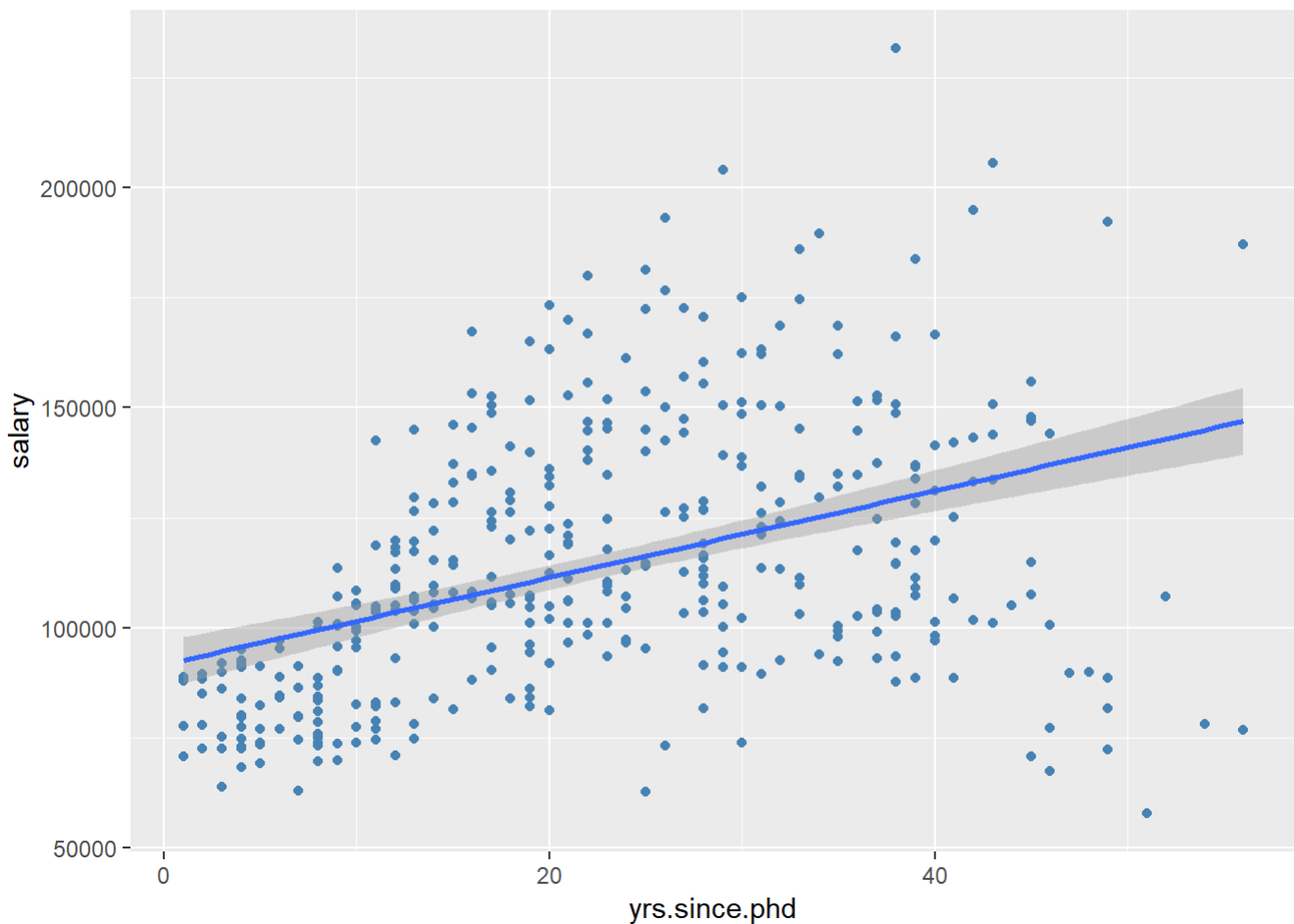


Figure 4.9: Scatterplot with linear fit line

Clearly, salary increases with experience. However, there seems to be a dip at the right end - professors with significant experience, earning lower salaries. A straight line does not capture this non-linear effect. A line with a bend will fit better here.

A polynomial regression line provides a fit line of the form

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4 + \dots$$

Typically either a quadratic (one bend), or cubic (two bends) line is used. It is rarely necessary to use a higher order(>3) polynomials. Applying a quadratic fit to the salary dataset produces the following result.

```
# scatterplot with quadratic line of best fit
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point(color= "steelblue") +
  geom_smooth(method = "lm",
             formula = y ~ poly(x, 2),
             color = "indianred3")
```

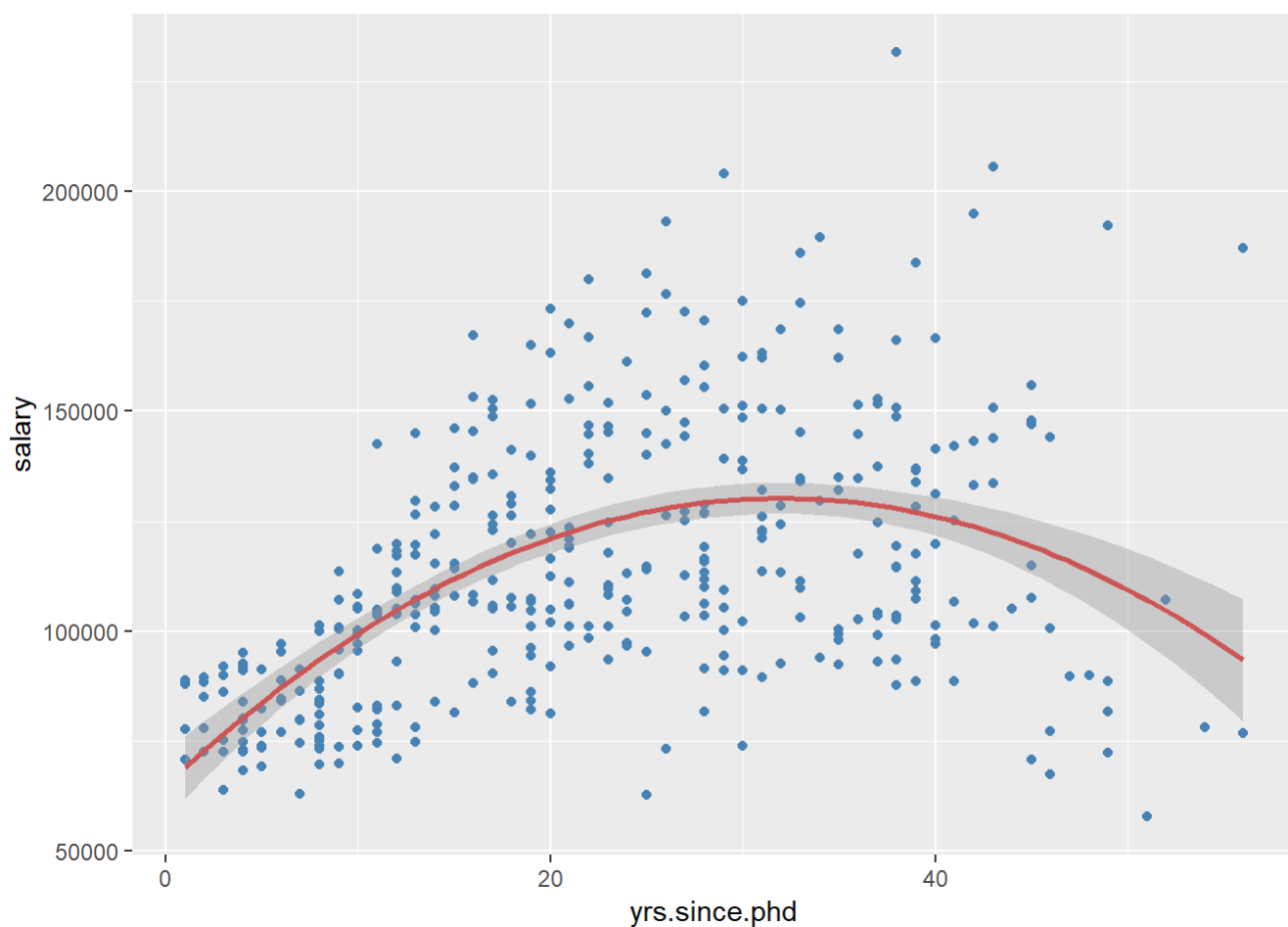


Figure 4.10: Scatterplot with quadratic fit line

Finally, a smoothed nonparametric fit line can often provide a good picture of the relationship. The default in `ggplot2` is a `loess` line which stands for for locally weighted scatterplot smoothing.

```
# scatterplot with loess smoothed line
```

```
ggplot(Salaries,  
       aes(x = yrs.since.phd,  
           y = salary)) +  
  geom_point(color= "steelblue") +  
  geom_smooth(color = "tomato")
```

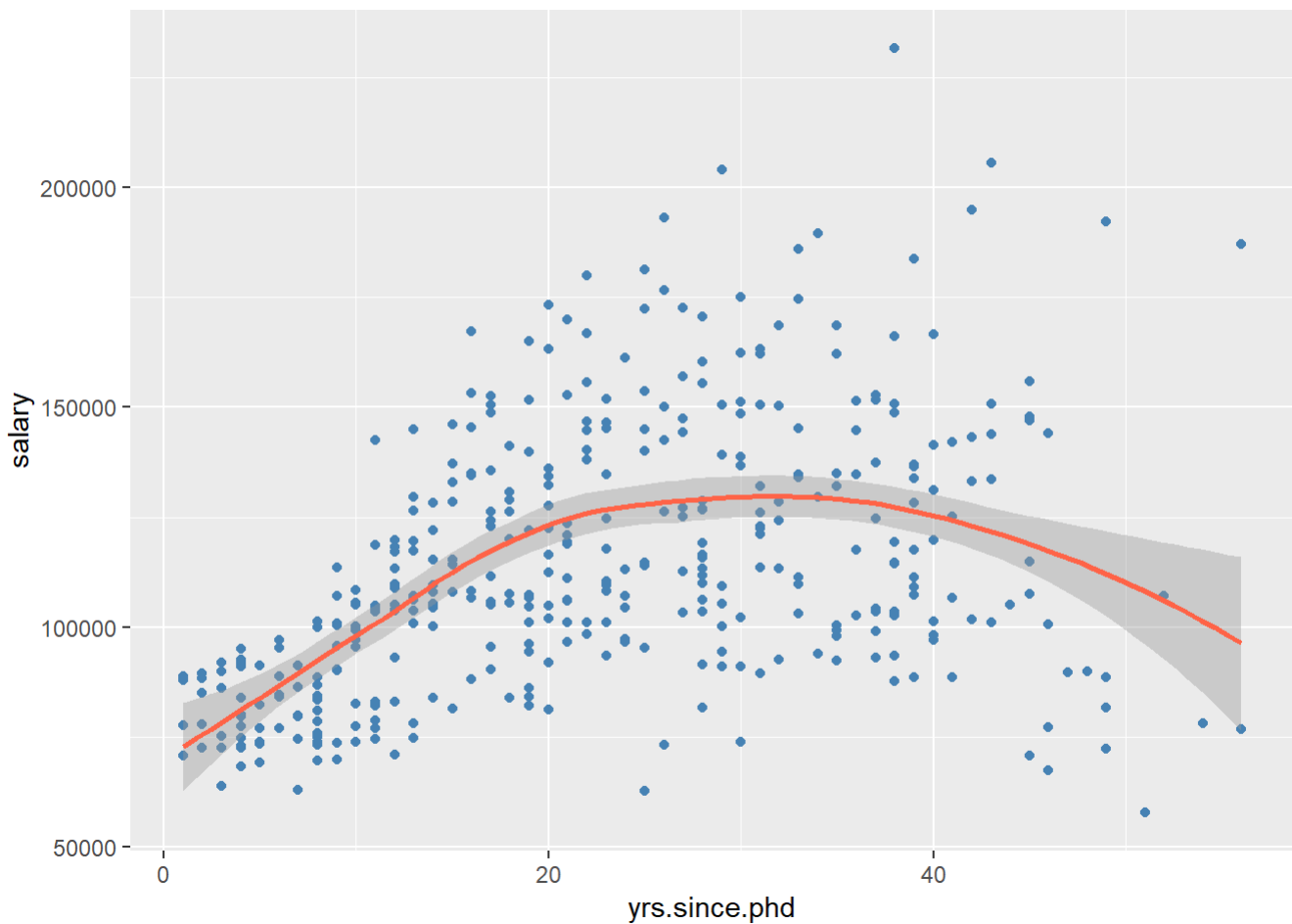


Figure 4.11: Scatterplot with nonparametric fit line

You can suppress the confidence bands by including the option `se = FALSE`.

Here is a complete (and more attractive) plot.

```

# scatterplot with loess smoothed line
# and better labeling and color
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point(color="cornflowerblue",
            size = 2,
            alpha = .6) +
  geom_smooth(size = 1.5,
            color = "darkgrey") +
  scale_y_continuous(label = scales::dollar,
                    limits = c(50000, 250000)) +
  scale_x_continuous(breaks = seq(0, 60, 10),
                    limits = c(0, 60)) +
  labs(x = "Years Since PhD",
       y = "",
       title = "Experience vs. Salary",
       subtitle = "9-month salary for 2008-2009") +
  theme_minimal()

```

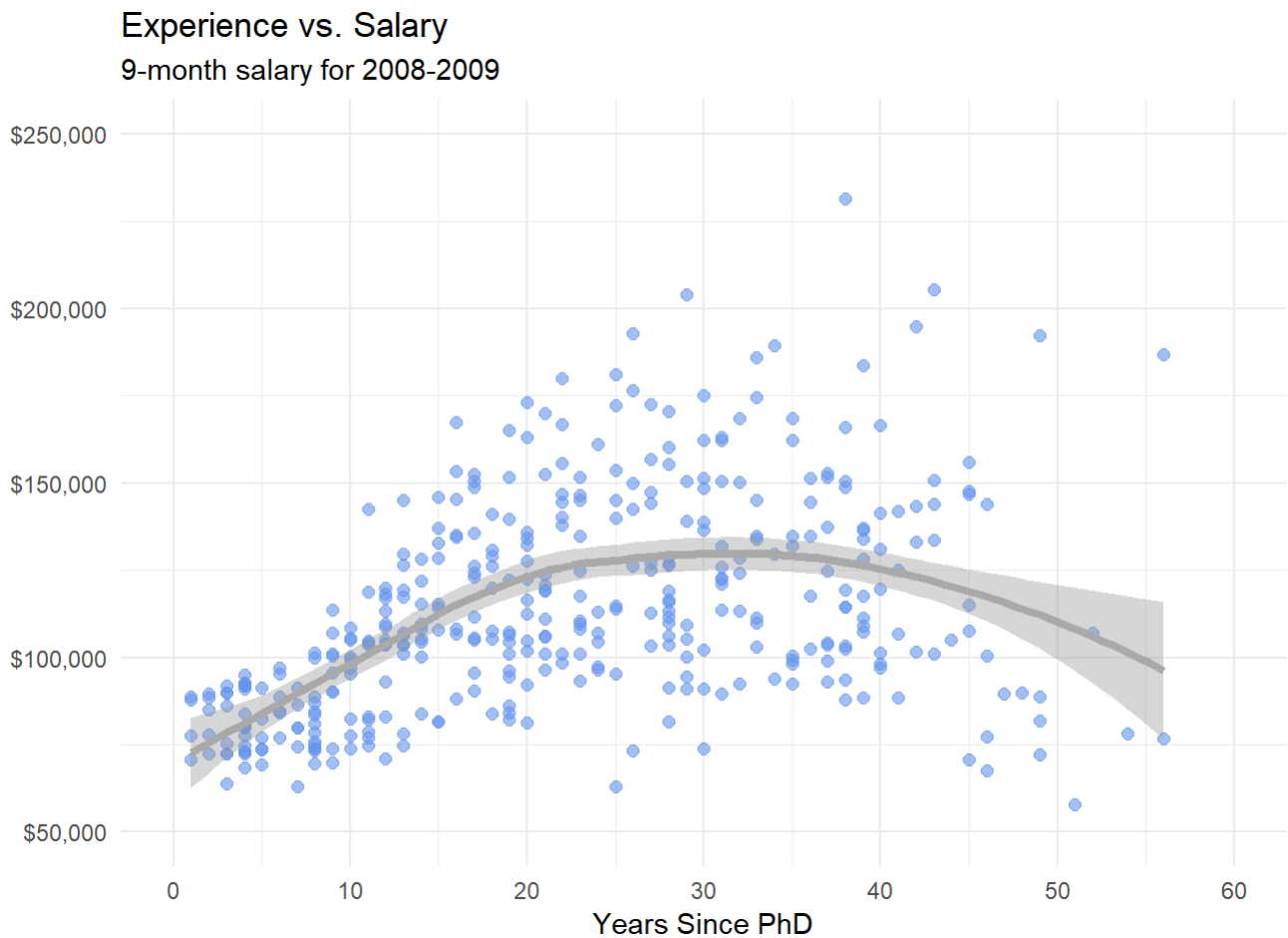


Figure 4.12: Scatterplot with nonparametric fit line

4.2.2 Line plot

When one of the two variables represents time, a line plot can be an effective method of displaying relationship. For example, the code below displays the relationship between time (*year*) and life expectancy (*lifeExp*) in the United States between 1952 and 2007. The data comes from the [gapminder](#) dataset.

```

data(gapminder, package="gapminder")

# Select US cases
library(dplyr)
plotdata <- filter(gapminder,
                    country == "United States")

# simple line plot
ggplot(plotdata,
       aes(x = year,
           y = lifeExp)) +
  geom_line()

```

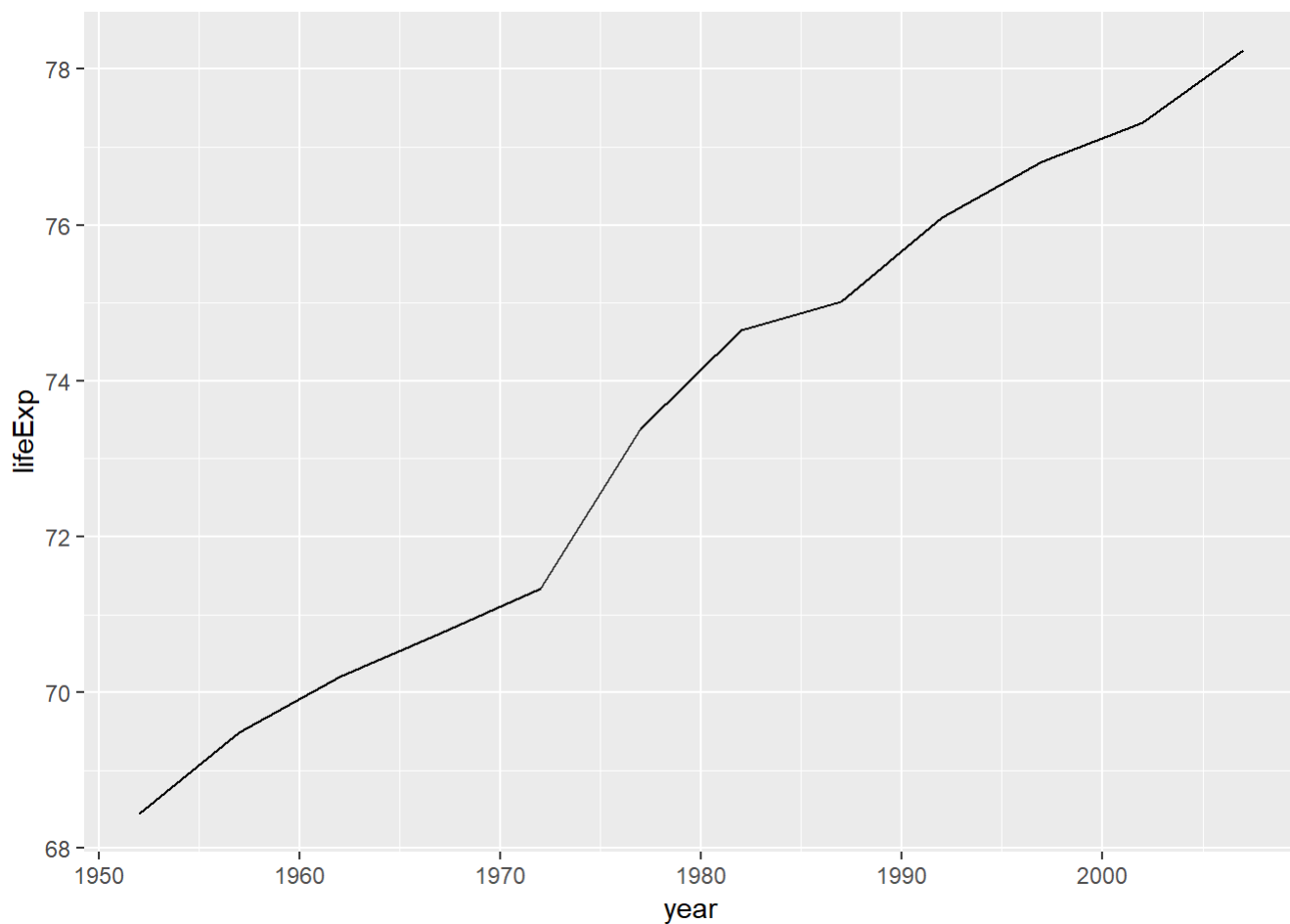


Figure 4.13: Simple line plot

It is hard to read individual values in the graph above. In the next plot, we'll add points as well.


```

# Line plot with points
# and improved labeling
ggplot(plotdata,
       aes(x = year,
           y = lifeExp)) +
  geom_line(size = 1.5,
            color = "lightgrey") +
  geom_point(size = 3,
             color = "steelblue") +
  labs(y = "Life Expectancy (years)",
       x = "Year",
       title = "Life expectancy changes over time",
       subtitle = "United States (1952-2007)",
       caption = "Source: http://www.gapminder.org/data/")

```

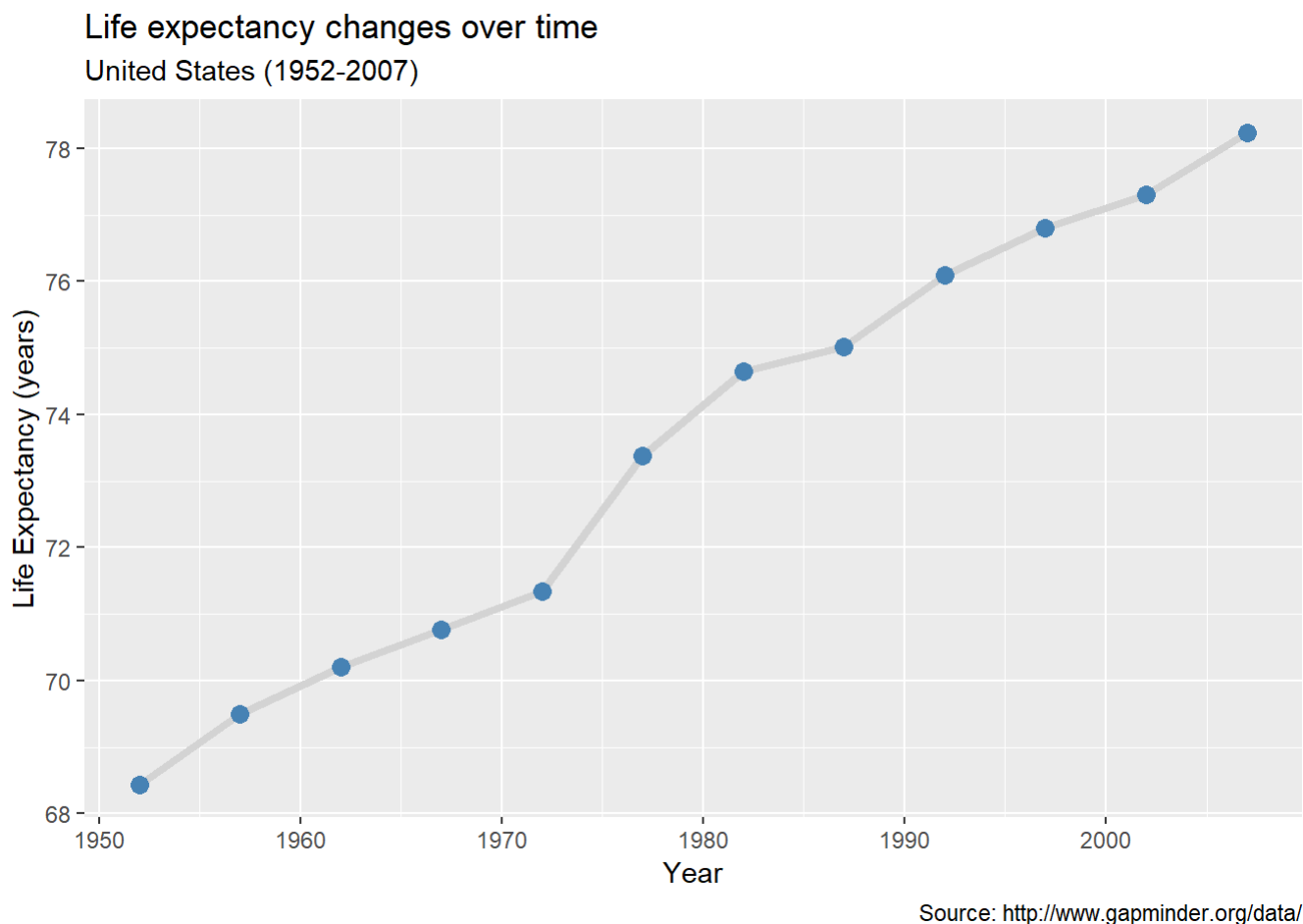


Figure 4.14: Line plot with points and labels

Time dependent data is covered in more detail under [Time series](#). Customizing line graphs is covered in the [Customizing graphs](#) section.

4.3 Categorical vs. Quantitative

When plotting the relationship between a categorical variable and a quantitative variable, a large number of graph types are available. These include bar charts using summary statistics, grouped kernel density plots, side-by-side box plots, side-by-side violin plots, mean/sem plots, ridgeline plots, and Cleveland plots.

4.3.1 Bar chart (on summary statistics)

In previous sections, bar charts were used to display the number of cases by category for a [single variable](#) or for [two variables](#). You can also use bar charts to display other summary statistics (e.g., means or medians) on a quantitative variable for each level of a categorical variable.

For example, the following graph displays the mean salary for a sample of university professors by their academic rank.

```
data(Salaries, package="carData")

# calculate mean salary for each rank
library(dplyr)
plotdata <- Salaries %>%
  group_by(rank) %>%
  summarize(mean_salary = mean(salary))

# plot mean salaries
ggplot(plotdata,
       aes(x = rank,
           y = mean_salary)) +
  geom_bar(stat = "identity")
```

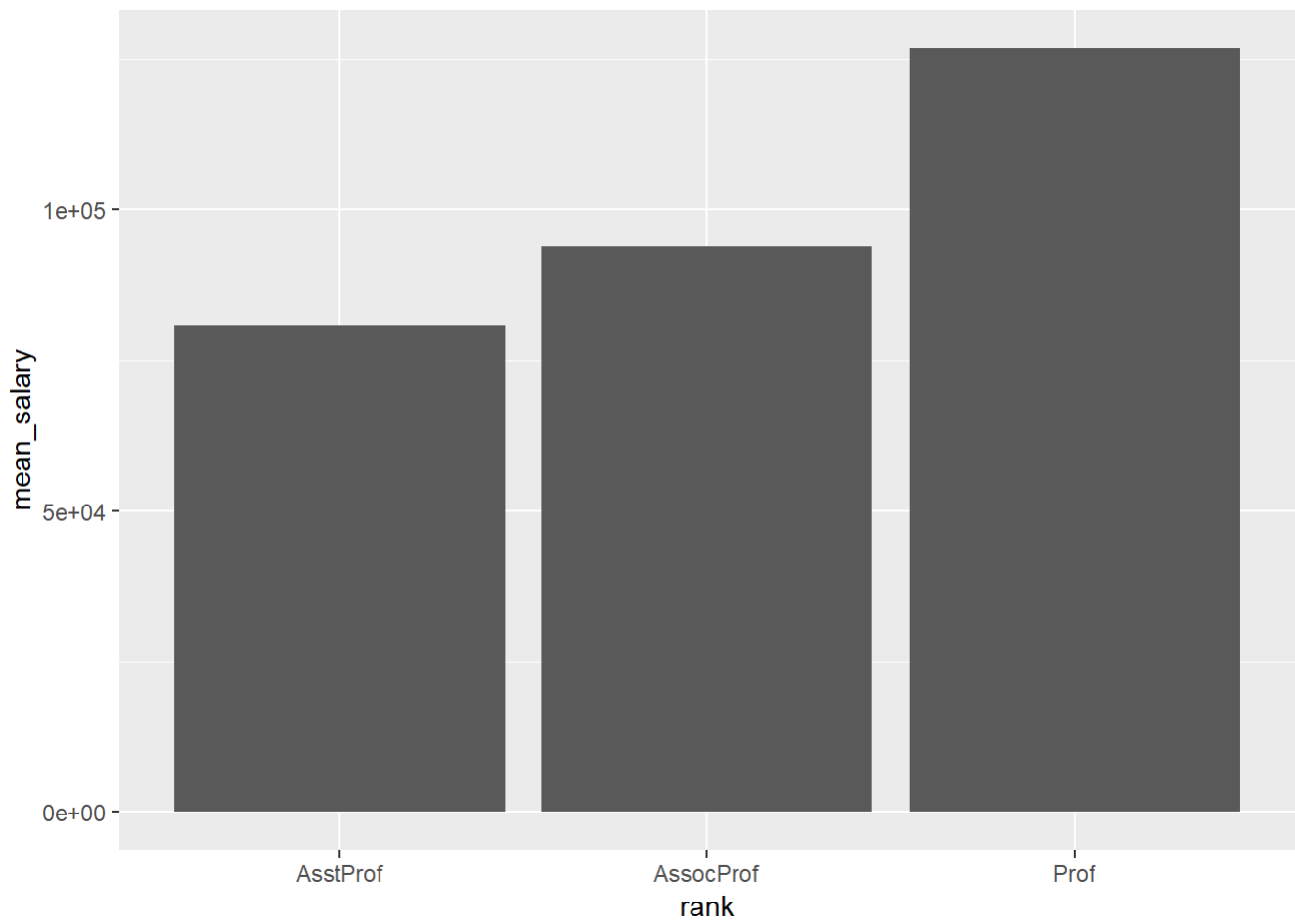


Figure 4.15: Bar chart displaying means

We can make it more attractive with some options.

```

# plot mean salaries in a more attractive fashion
library(scales)
ggplot(plotdata,
       aes(x = factor(rank,
                       labels = c("Assistant\nProfessor",
                                   "Associate\nProfessor",
                                   "Full\nProfessor")),
          y = mean_salary)) +
  geom_bar(stat = "identity",
          fill = "cornflowerblue") +
  geom_text(aes(label = dollar(mean_salary)),
          vjust = -0.25) +
  scale_y_continuous(breaks = seq(0, 130000, 20000),
                    label = dollar) +
  labs(title = "Mean Salary by Rank",
       subtitle = "9-month academic salary for 2008-2009",
       x = "",
       y = "")

```

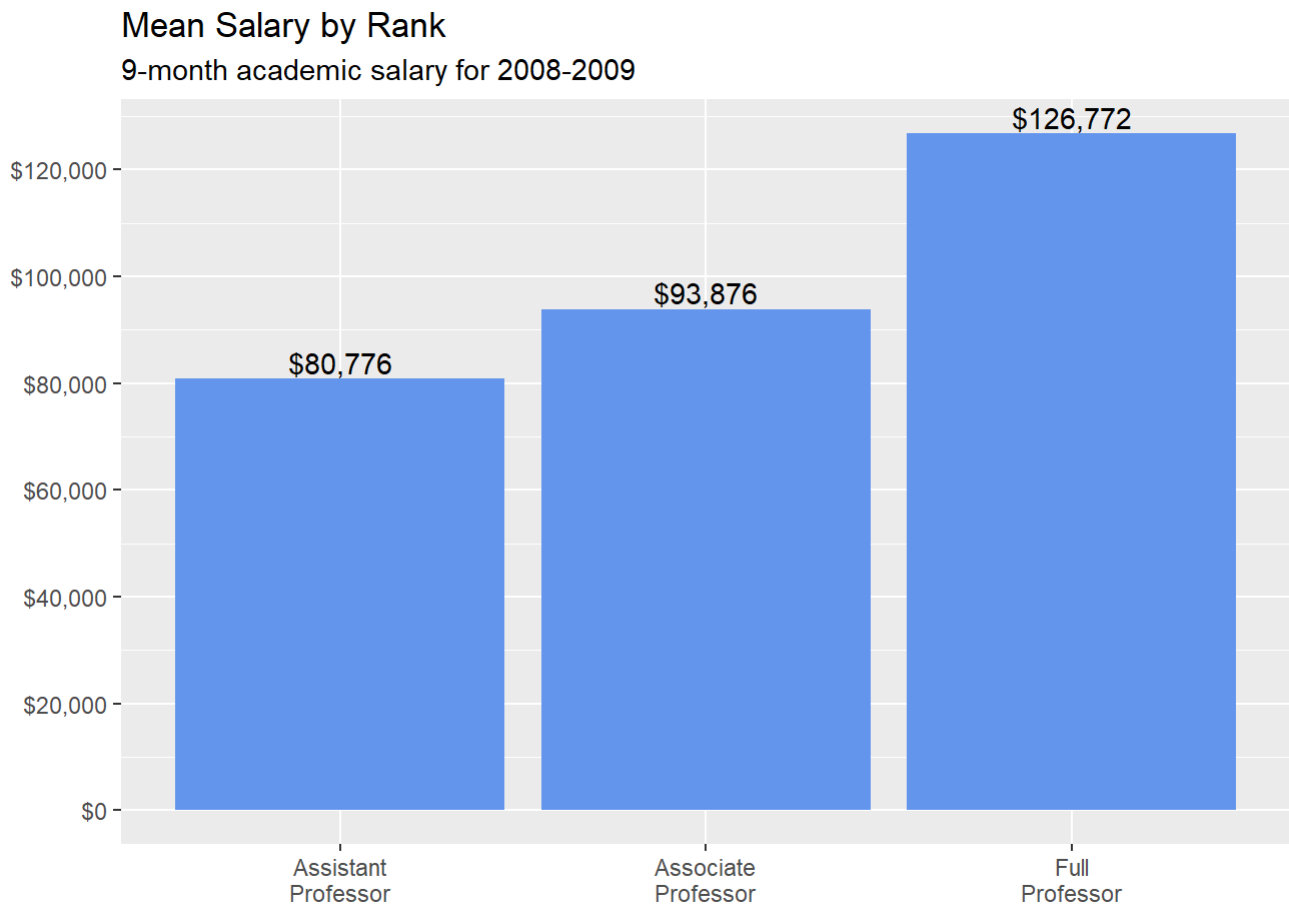


Figure 4.16: Bar chart displaying means

One limitation of such plots is that they do not display the distribution of the data - only the summary statistic for each group. The plots below correct this limitation to some extent.

4.3.2 Grouped kernel density plots

One can compare groups on a numeric variable by superimposing [kernel density](#) plots in a single graph.

```
# plot the distribution of salaries
# by rank using kernel density plots
ggplot(Salaries,
       aes(x = salary,
           fill = rank)) +
  geom_density(alpha = 0.4) +
  labs(title = "Salary distribution by rank")
```

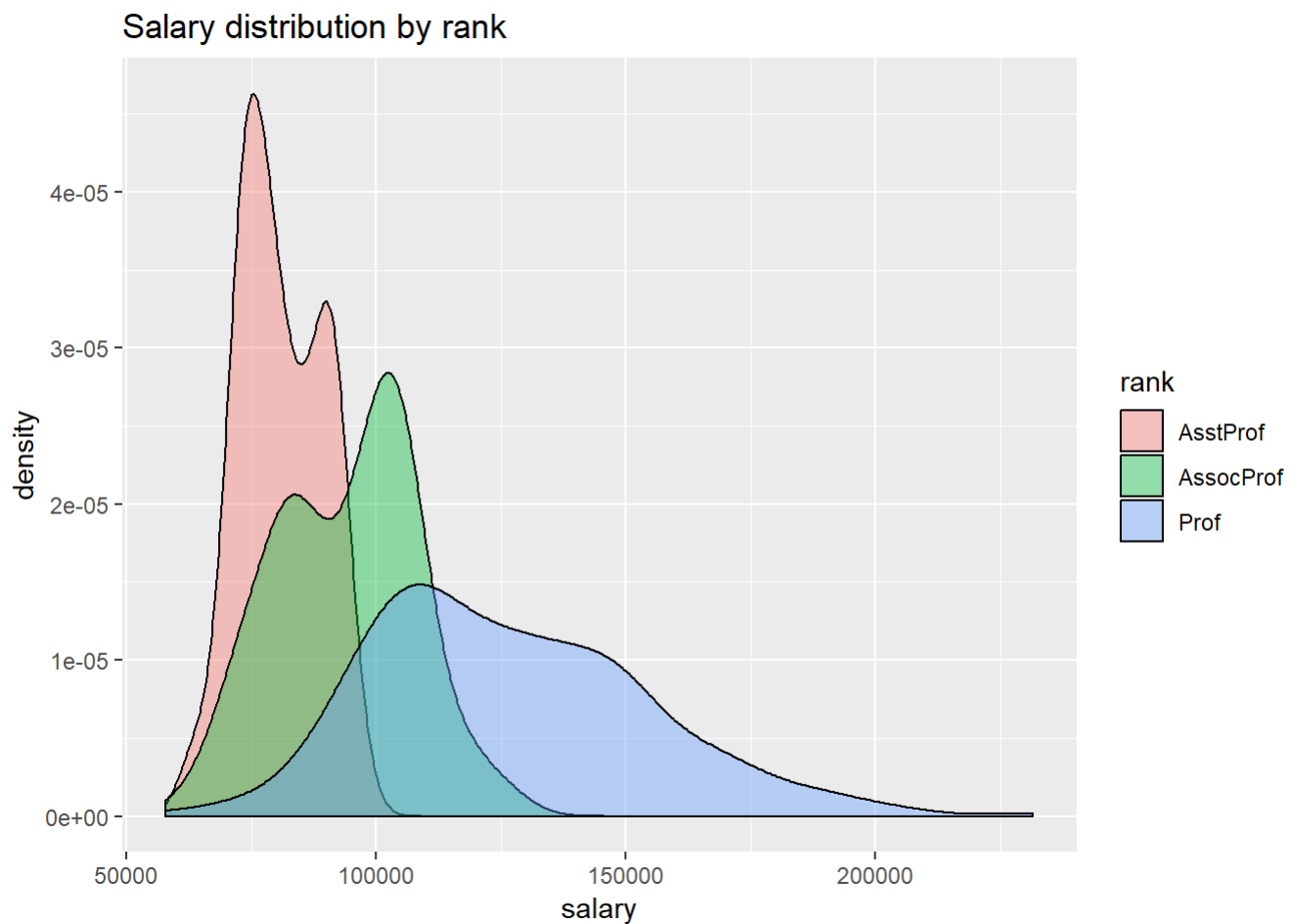
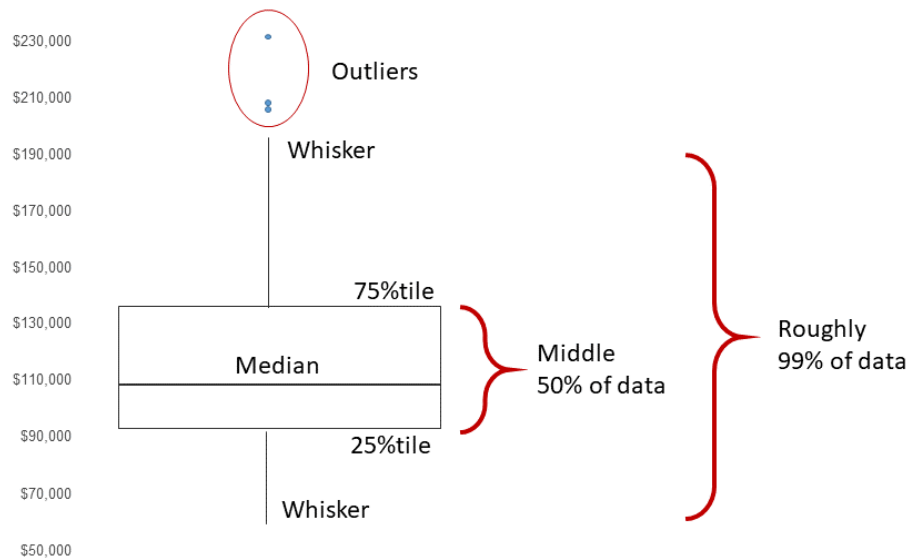


Figure 4.17: Grouped kernel density plots

The `alpha` option makes the density plots partially transparent, so that we can see what is happening under the overlaps. Alpha values range from 0 (transparent) to 1 (opaque). The graph makes clear that, in general, salary goes up with rank. However, the salary range for full professors is *very* wide.

4.3.3 Box plots

A boxplot displays the 25th percentile, median, and 75th percentile of a distribution. The whiskers (vertical lines) capture roughly 99% of a normal distribution, and observations outside this range are plotted as points representing outliers (see the figure below).



Side-by-side box plots are very useful for comparing groups (i.e., the levels of a categorical variable) on a numerical variable.

plot the distribution of salaries by rank using boxplots

```
ggplot(Salaries,
       aes(x = rank,
           y = salary)) +
  geom_boxplot() +
  labs(title = "Salary distribution by rank")
```

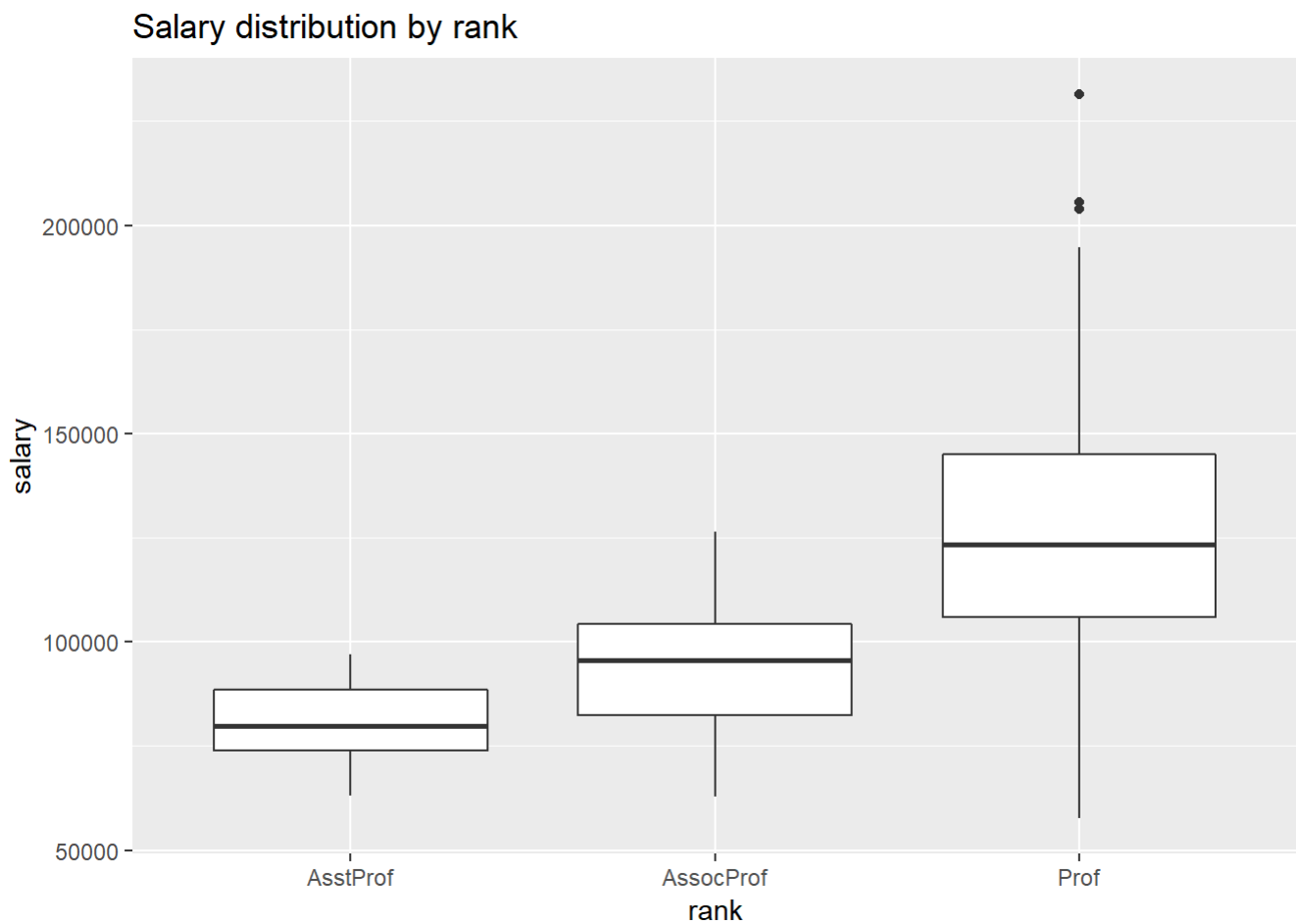


Figure 4.18: Side-by-side boxplots

Notched boxplots provide an approximate method for visualizing whether groups differ. Although not a formal test, if the notches of two boxplots do not overlap, there is strong evidence (95% confidence) that the medians of the two groups differ.

```
# plot the distribution of salaries by rank using boxplots
ggplot(Salaries, aes(x = rank,
                     y = salary)) +
  geom_boxplot(notch = TRUE,
               fill = "cornflowerblue",
               alpha = .7) +
  labs(title = "Salary distribution by rank")
```

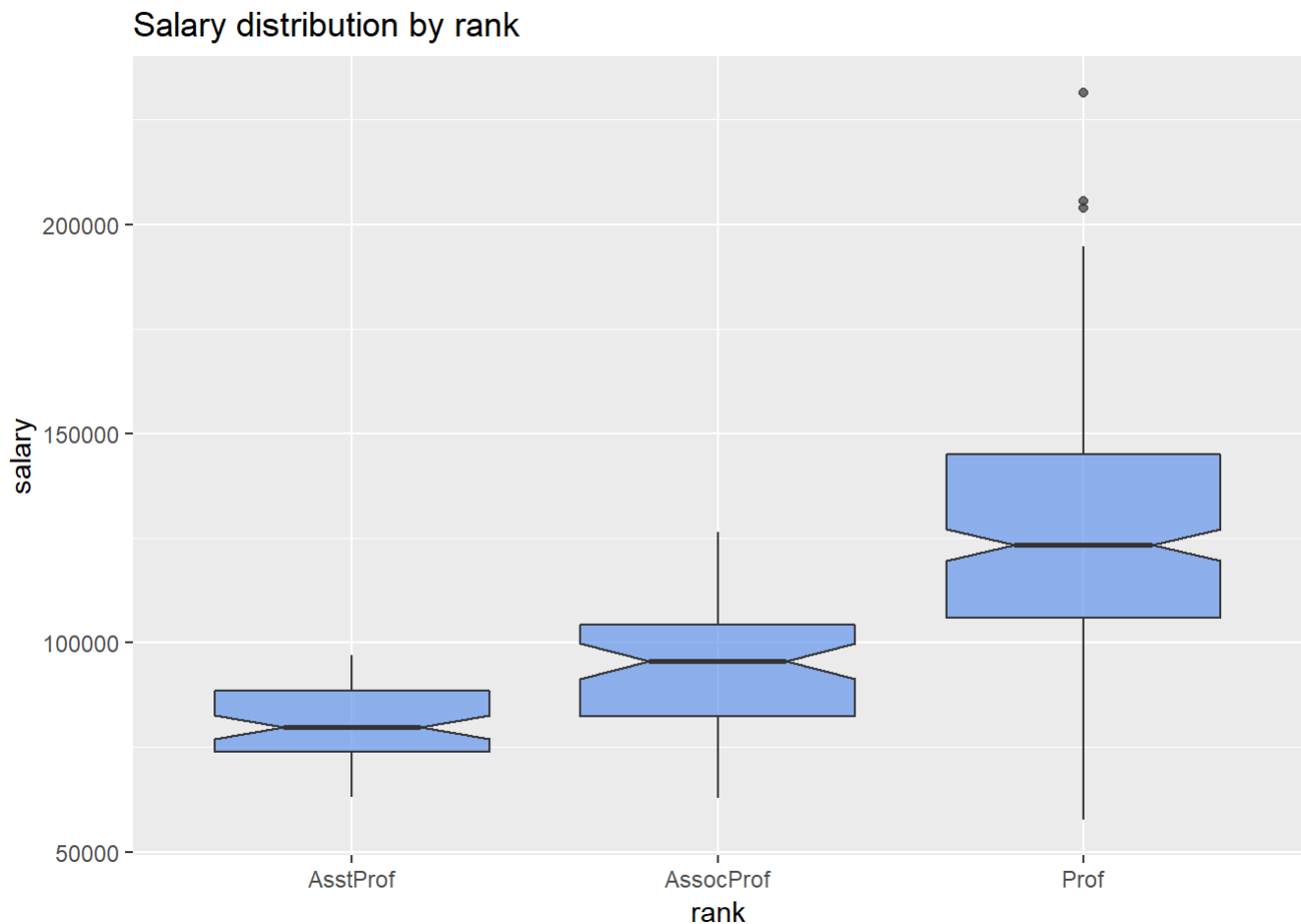



Figure 4.19: Side-by-side notched boxplots

In the example above, all three groups appear to differ.

One of the advantages of boxplots is that their widths are not usually meaningful. This allows you to compare the distribution of many groups in a single graph.

4.3.4 Violin plots

Violin plots are similar to [kernel density plots](#), but are mirrored and rotated 90°.

```
# plot the distribution of salaries
# by rank using violin plots
ggplot(Salaries,
       aes(x = rank,
           y = salary)) +
  geom_violin() +
  labs(title = "Salary distribution by rank")
```

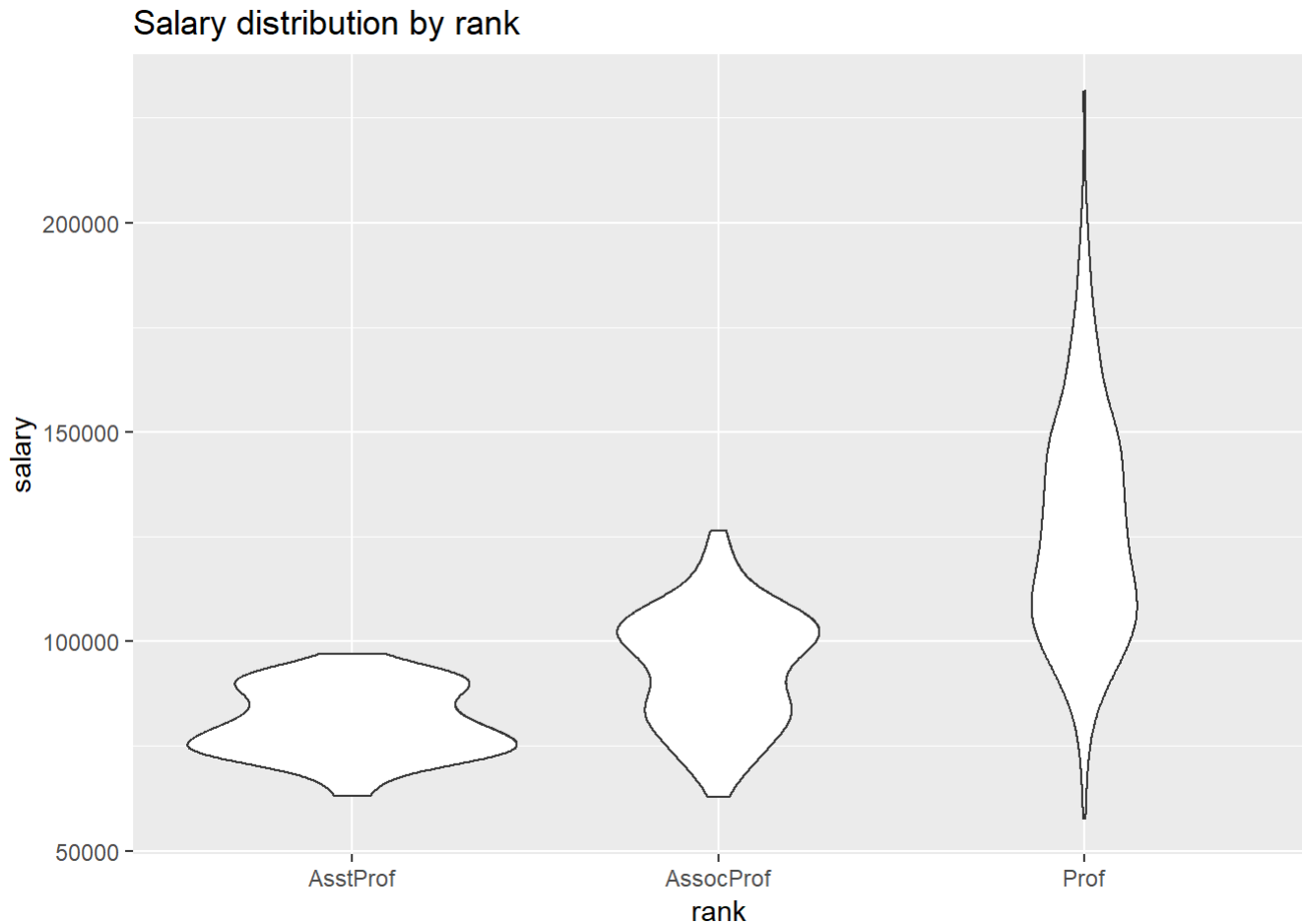


Figure 4.20: Side-by-side violin plots

A useful variation is to superimpose boxplots on violin plots.

```
# plot the distribution using violin and boxplots
ggplot(Salaries,
  aes(x = rank,
    y = salary)) +
  geom_violin(fill = "cornflowerblue") +
  geom_boxplot(width = .2,
    fill = "orange",
    outlier.color = "orange",
    outlier.size = 2) +
  labs(title = "Salary distribution by rank")
```

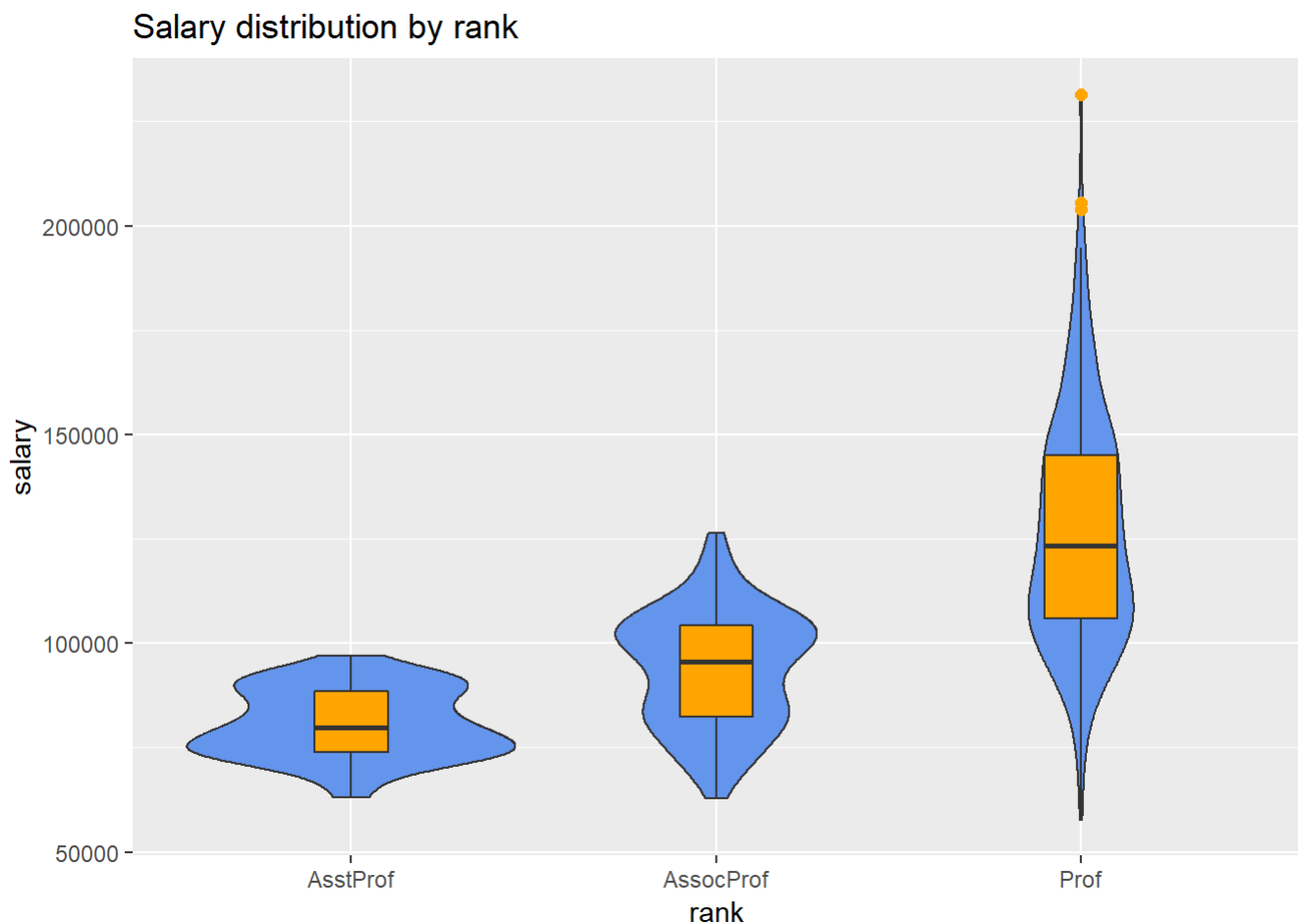


Figure 4.21: Side-by-side violin/box plots

4.3.5 Ridgeline plots

A ridgeline plot (also called a joyplot) displays the distribution of a quantitative variable for several groups. They're similar to [kernel density](#) plots with vertical [faceting](#), but take up less room.

Ridgeline plots are created with the `ggridges` package.

Using the [Fuel economy](#) dataset, let's plot the distribution of city driving miles per gallon by car class.

```
# create ridgeline graph
library(ggplot2)
library(ggribes)

ggplot(mpg,
       aes(x = cty,
           y = class,
           fill = class)) +
  geom_density_ridges() +
  theme_ridges() +
  labs("Highway mileage by auto class") +
  theme(legend.position = "none")
```

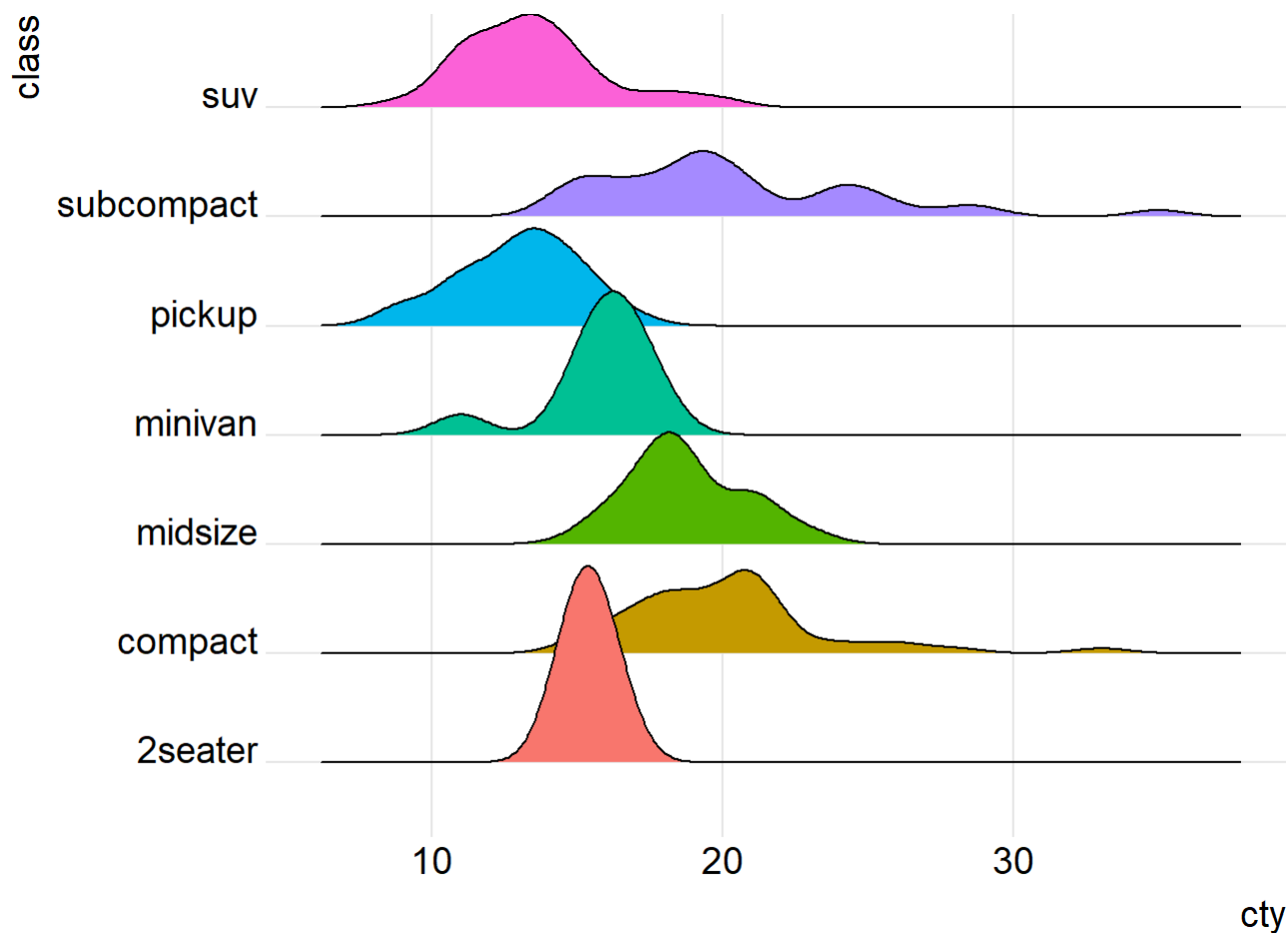


Figure 4.22: Ridgeline graph with color fill

I've suppressed the legend here because it's redundant (the distributions are already labeled on the y-axis). Unsurprisingly, pickup trucks have the poorest mileage, while subcompacts and compact cars tend to achieve ratings. However, there is a very wide range of gas mileage scores for these smaller cars.

Note the the possible overlap of distributions is the trade-off for a more compact graph. You can add transparency if the the overlap is severe using `geom_density_ridges(alpha = n)`, where n ranges from 0 (transparent) to 1 (opaque). See the [package vingnette](#) for more details.

4.3.6 Mean/SEM plots

A popular method for comparing groups on a numeric variable is the mean plot with error bars. Error bars can represent standard deviations, standard error of the mean, or confidence intervals. In this section, we'll plot means and standard errors.

```
# calculate means, standard deviations,  
# standard errors, and 95% confidence  
# intervals by rank  
library(dplyr)  
plotdata <- Salaries %>%  
  group_by(rank) %>%  
  summarize(n = n(),  
            mean = mean(salary),  
            sd = sd(salary),  
            se = sd / sqrt(n),  
            ci = qt(0.975, df = n - 1) * sd / sqrt(n))
```

The resulting dataset is given below.

Table 4.1: Plot data

rank	n	mean	sd	se	ci
AsstProf	67	80775.99	8174.113	998.6268	1993.823
AssocProf	64	93876.44	13831.700	1728.9625	3455.056
Prof	266	126772.11	27718.675	1699.5410	3346.322

```
# plot the means and standard errors
ggplot(plotdata,
       aes(x = rank,
           y = mean,
           group = 1)) +
  geom_point(size = 3) +
  geom_line() +
  geom_errorbar(aes(ymin = mean - se,
                    ymax = mean + se),
                width = .1)
```

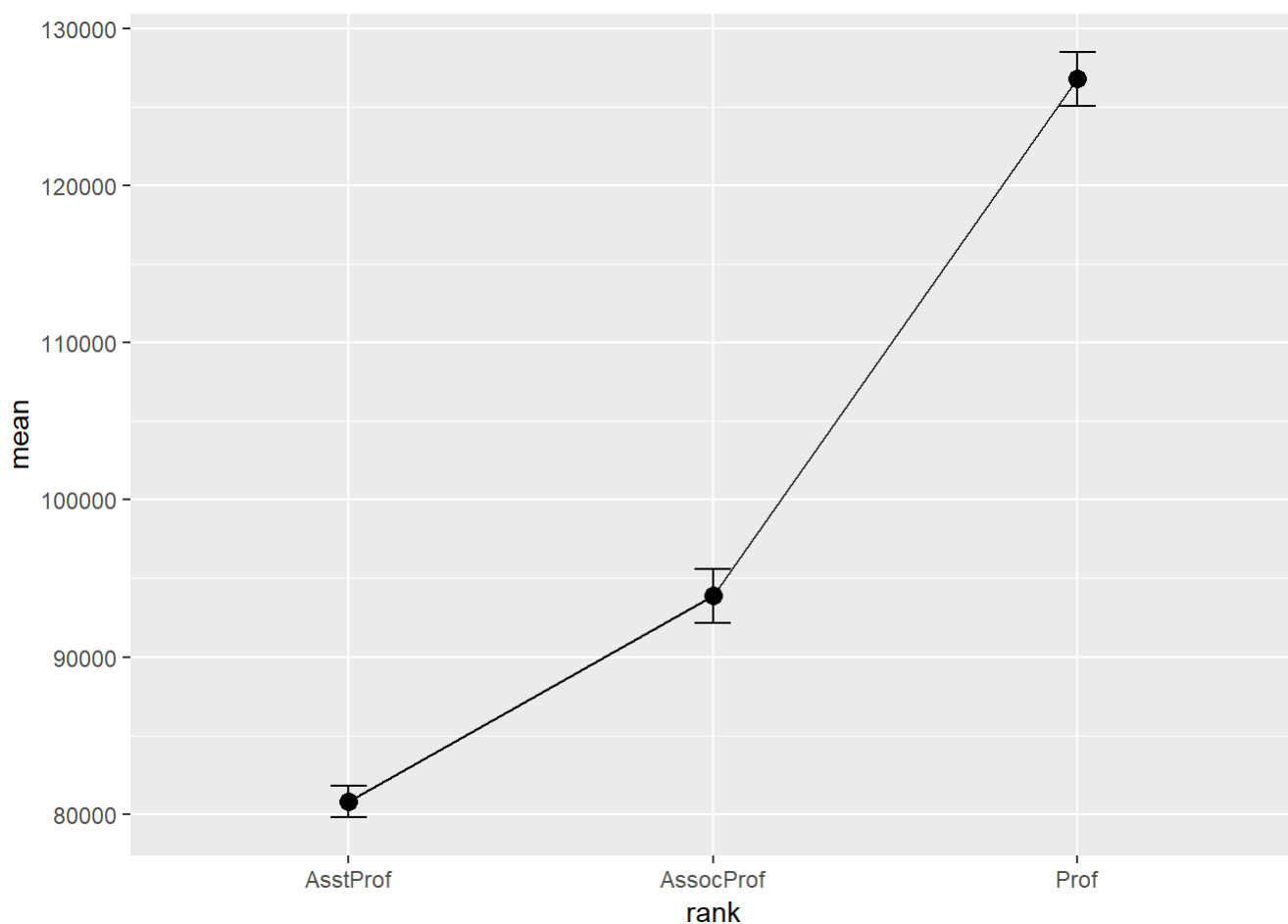


Figure 4.23: Mean plots with standard error bars

Although we plotted error bars representing the standard error, we could have plotted standard deviations or 95% confidence intervals. Simply replace `se` with `sd` or `error` in the `aes` option.

We can use the same technique to compare salary across rank and sex. (Technically, this is not bivariate since we're plotting rank, sex, and salary, but it seems to fit here)

```
# calculate means and standard errors by rank and sex
```

```
plotdata <- Salaries %>%  
  group_by(rank, sex) %>%  
  summarize(n = n(),  
            mean = mean(salary),  
            sd = sd(salary),  
            se = sd/sqrt(n))
```

```
# plot the means and standard errors by sex
```

```
ggplot(plotdata, aes(x = rank,  
                     y = mean,  
                     group=sex,  
                     color=sex)) +  
  geom_point(size = 3) +  
  geom_line(size = 1) +  
  geom_errorbar(aes(ymin = mean - se,  
                   ymax = mean+se),  
               width = .1)
```

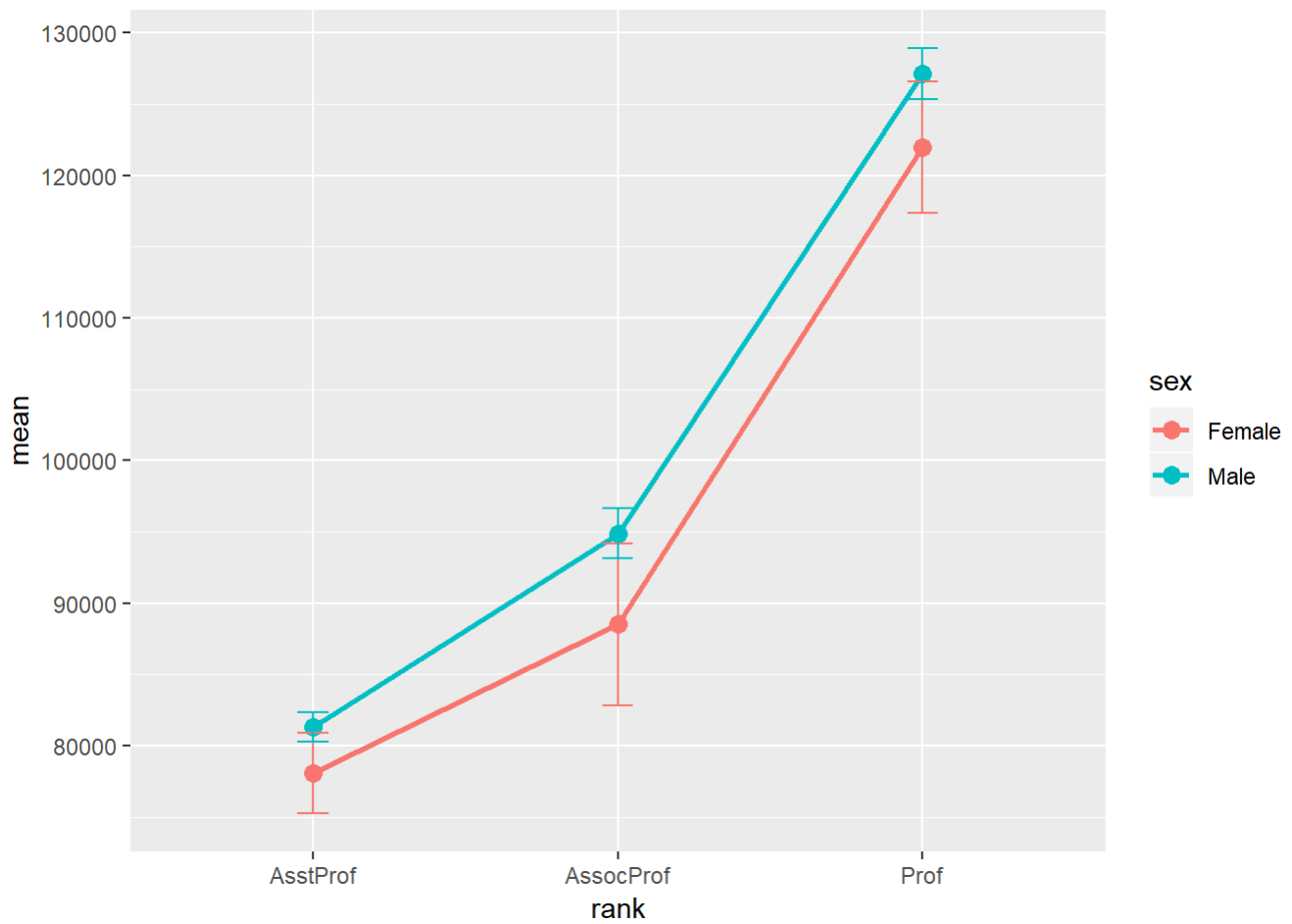


Figure 4.24: Mean plots with standard error bars by sex

Unfortunately, the error bars overlap. We can dodge the horizontal positions a bit to overcome this.


```

# plot the means and standard errors by sex (dodged)
pd <- position_dodge(0.2)
ggplot(plotdata,
       aes(x = rank,
           y = mean,
           group=sex,
           color=sex)) +
  geom_point(position = pd,
             size = 3) +
  geom_line(position = pd,
            size = 1) +
  geom_errorbar(aes(ymin = mean - se,
                    ymax = mean + se),
               width = .1,
               position= pd)

```

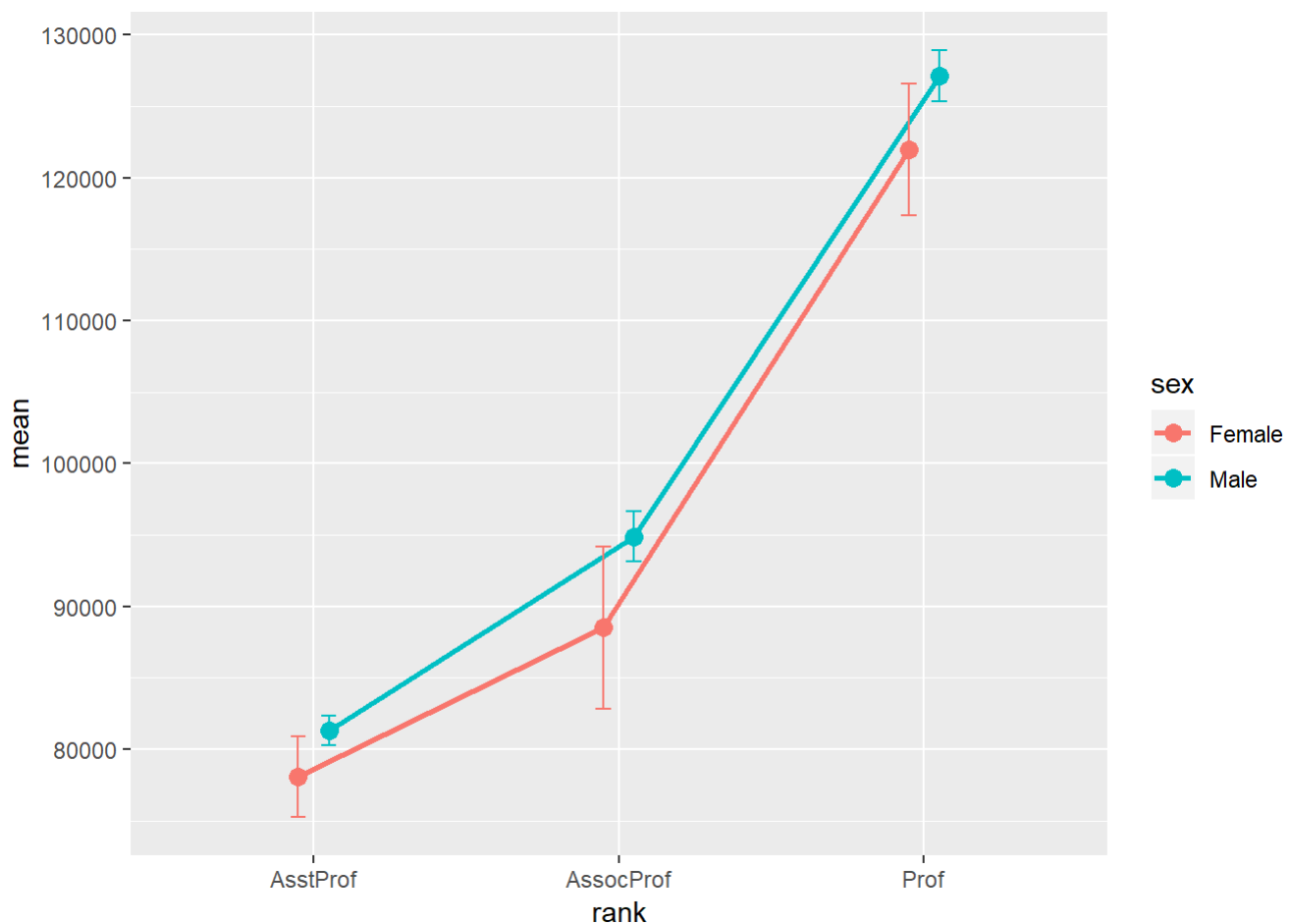


Figure 4.25: Mean plots with standard error bars (dodged)

Finally, lets add some options to make the graph more attractive.

```
# improved means/standard error plot
pd <- position_dodge(0.2)
ggplot(plotdata,
       aes(x = factor(rank,
                      labels = c("Assistant\nProfessor",
                                "Associate\nProfessor",
                                "Full\nProfessor")),
          y = mean,
          group=sex,
          color=sex)) +
  geom_point(position=pd,
            size = 3) +
  geom_line(position = pd,
            size = 1) +
  geom_errorbar(aes(ymin = mean - se,
                   ymax = mean + se),
               width = .1,
               position = pd,
               size = 1) +
  scale_y_continuous(label = scales::dollar) +
  scale_color_brewer(palette="Set1") +
  theme_minimal() +
  labs(title = "Mean salary by rank and sex",
       subtitle = "(mean +/- standard error)",
       x = "",
       y = "",
       color = "Gender")
```



Figure 4.26: Mean/se plot with better labels and colors

4.3.7 Strip plots

The relationship between a grouping variable and a numeric variable can be displayed with a scatter plot. For example

```
# plot the distribution of salaries
# by rank using strip plots
ggplot(Salaries,
       aes(y = rank,
           x = salary)) +
  geom_point() +
  labs(title = "Salary distribution by rank")
```

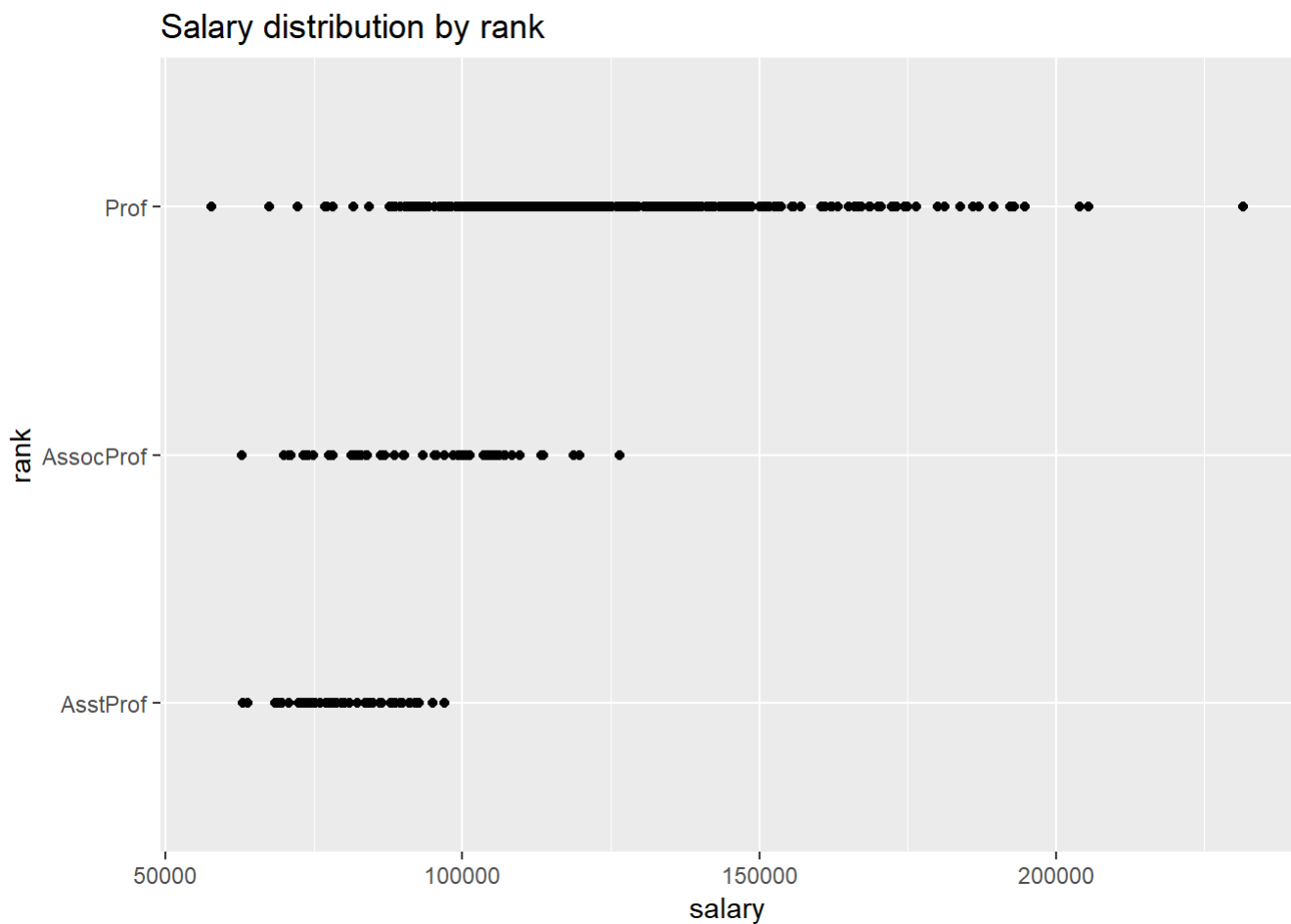


Figure 4.27: Categorical by quantitative scatterplot

These one-dimensional scatterplots are called strip plots. Unfortunately, overprinting of points makes interpretation difficult. The relationship is easier to see if the the points are jittered. Basically a small random number is added to each y-coordinate.

```
# plot the distribution of salaries
# by rank using jittering
ggplot(Salaries,
       aes(y = rank,
           x = salary)) +
  geom_jitter() +
  labs(title = "Salary distribution by rank")
```

Salary distribution by rank

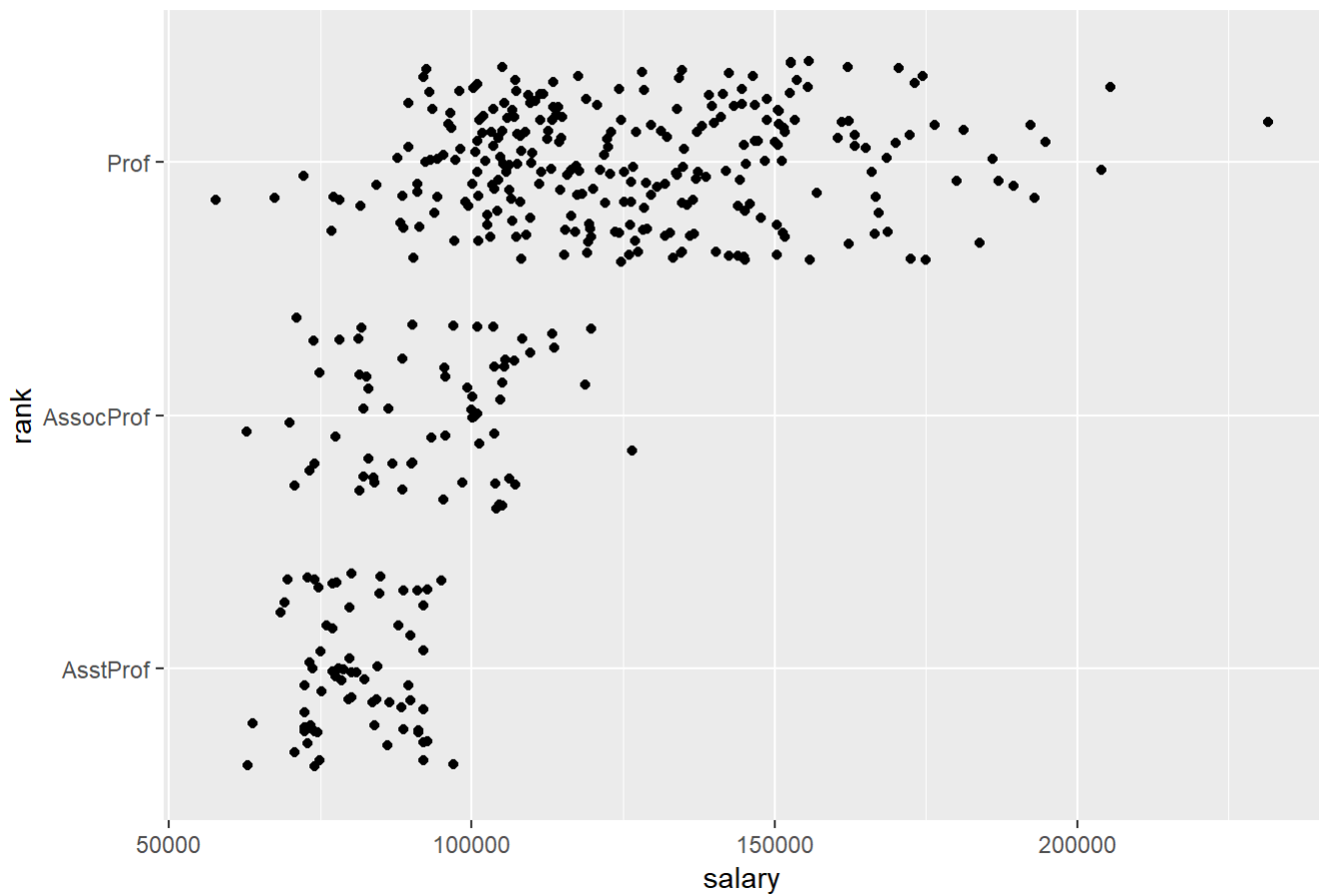


Figure 4.28: Jittered plot

It is easier to compare groups if we use color.

```

# plot the distribution of salaries
# by rank using jittering
library(scales)
ggplot(Salaries,
       aes(y = factor(rank,
                      labels = c("Assistant\nProfessor",
                                "Associate\nProfessor",
                                "Full\nProfessor")),
          x = salary,
          color = rank)) +
  geom_jitter(alpha = 0.7,
             size = 1.5) +
  scale_x_continuous(label = dollar) +
  labs(title = "Academic Salary by Rank",
       subtitle = "9-month salary for 2008-2009",
       x = "",
       y = "") +
  theme_minimal() +
  theme(legend.position = "none")

```

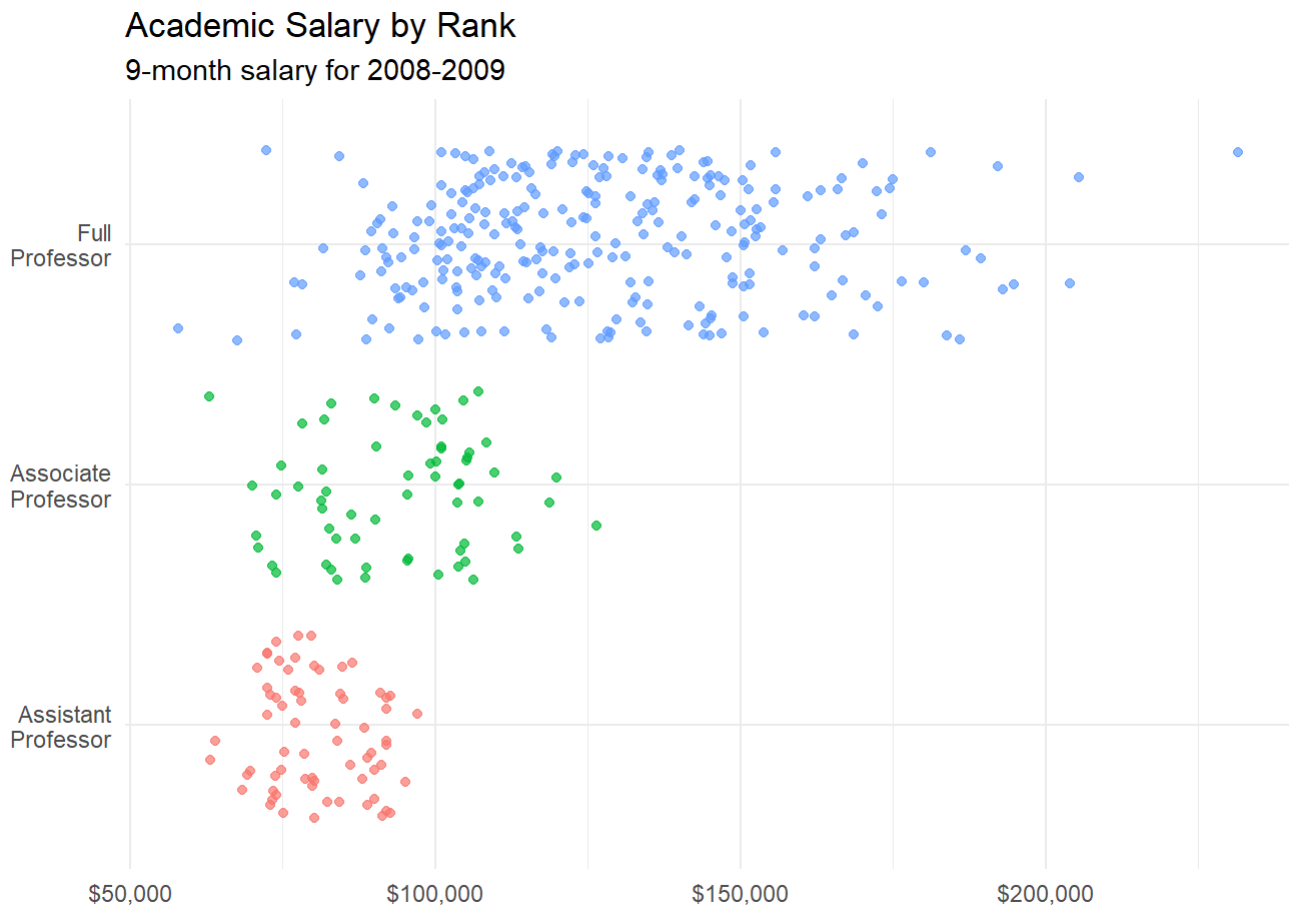


Figure 4.29: Fancy jittered plot

The option `legend.position = "none"` is used to suppress the legend (which is not needed here). Jittered plots work well when the number of points is not overly large.

4.3.7.1 Combining jitter and boxplots

It may be easier to visualize distributions if we add boxplots to the jitter plots.

```

# plot the distribution of salaries
# by rank using jittering
library(scales)
ggplot(Salaries,
       aes(x = factor(rank,
                      labels = c("Assistant\nProfessor",
                                "Associate\nProfessor",
                                "Full\nProfessor")),
          y = salary,
          color = rank)) +
  geom_boxplot(size=1,
              outlier.shape = 1,
              outlier.color = "black",
              outlier.size = 3) +
  geom_jitter(alpha = 0.5,
              width=.2) +
  scale_y_continuous(label = dollar) +
  labs(title = "Academic Salary by Rank",
       subtitle = "9-month salary for 2008-2009",
       x = "",
       y = "") +
  theme_minimal() +
  theme(legend.position = "none") +
  coord_flip()

```


Academic Salary by Rank

9-month salary for 2008-2009

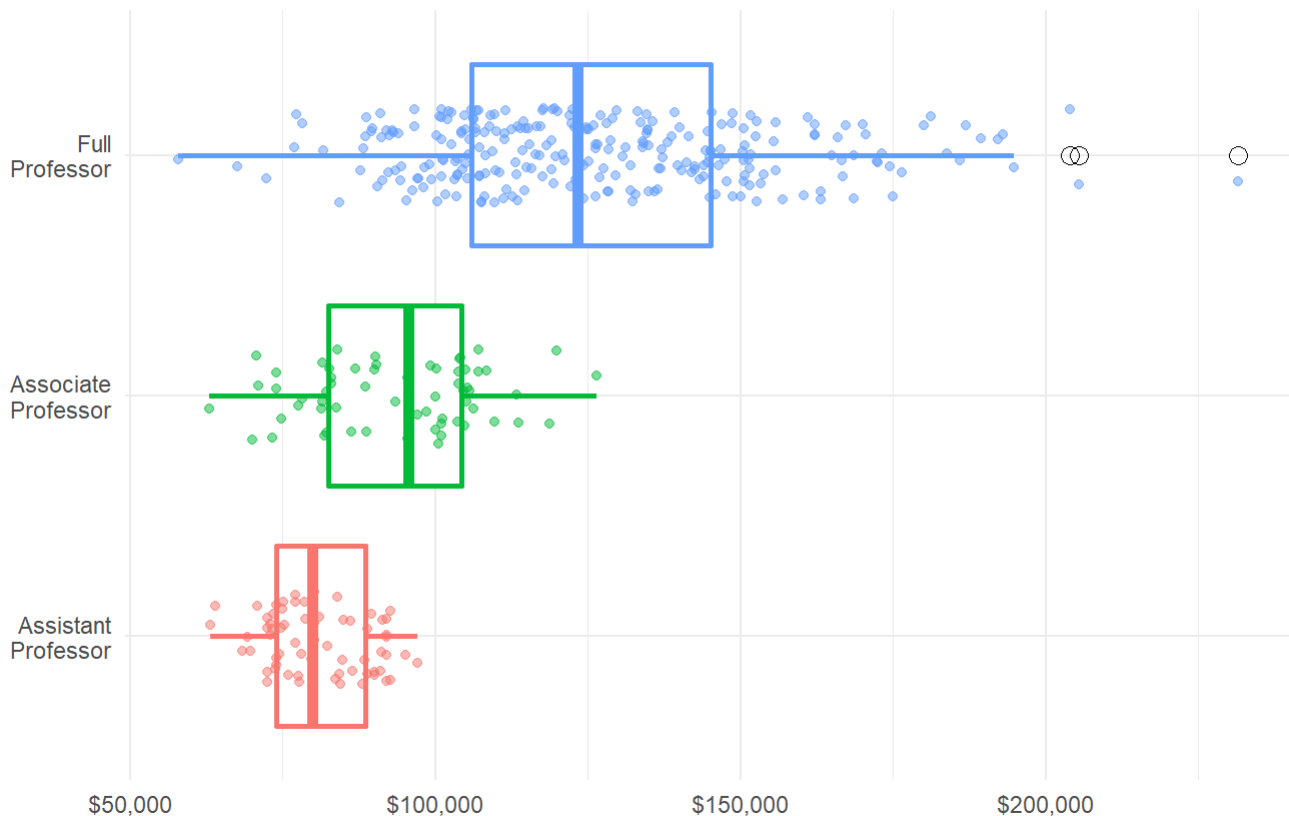


Figure 4.30: Jitter plot with superimposed box plots

Several options were added to create this plot.

For the boxplot

- `size = 1` makes the lines thicker
- `outlier.color = "black"` makes outliers black
- `outlier.shape = 1` specifies circles for outliers
- `outlier.size = 3` increases the size of the outlier symbol

For the jitter

- `alpha = 0.5` makes the points more transparent
- `width = .2` decreases the amount of jitter (.4 is the default)

Finally, the `x` and `y` axes are reversed using the `coord_flip` function (i.e., the graph is turned on its side).

Before moving on, it is worth mentioning the `geom_boxjitter` function provided in the `ggplot` package. It creates a hybrid boxplot - half boxplot, half scatterplot.

```

# plot the distribution of salaries
# by rank using jittering
library(ggplot)
library(scales)
ggplot(Salaries,
       aes(x = factor(rank,
                      labels = c("Assistant\nProfessor",
                                "Associate\nProfessor",
                                "Full\nProfessor")),
          y = salary,
          fill=rank)) +
  geom_boxjitter(color="black",
                jitter.color = "darkgrey",
                errorbar.draw = TRUE) +
  scale_y_continuous(label = dollar) +
  labs(title = "Academic Salary by Rank",
       subtitle = "9-month salary for 2008-2009",
       x = "",
       y = "") +
  theme_minimal() +
  theme(legend.position = "none")

```

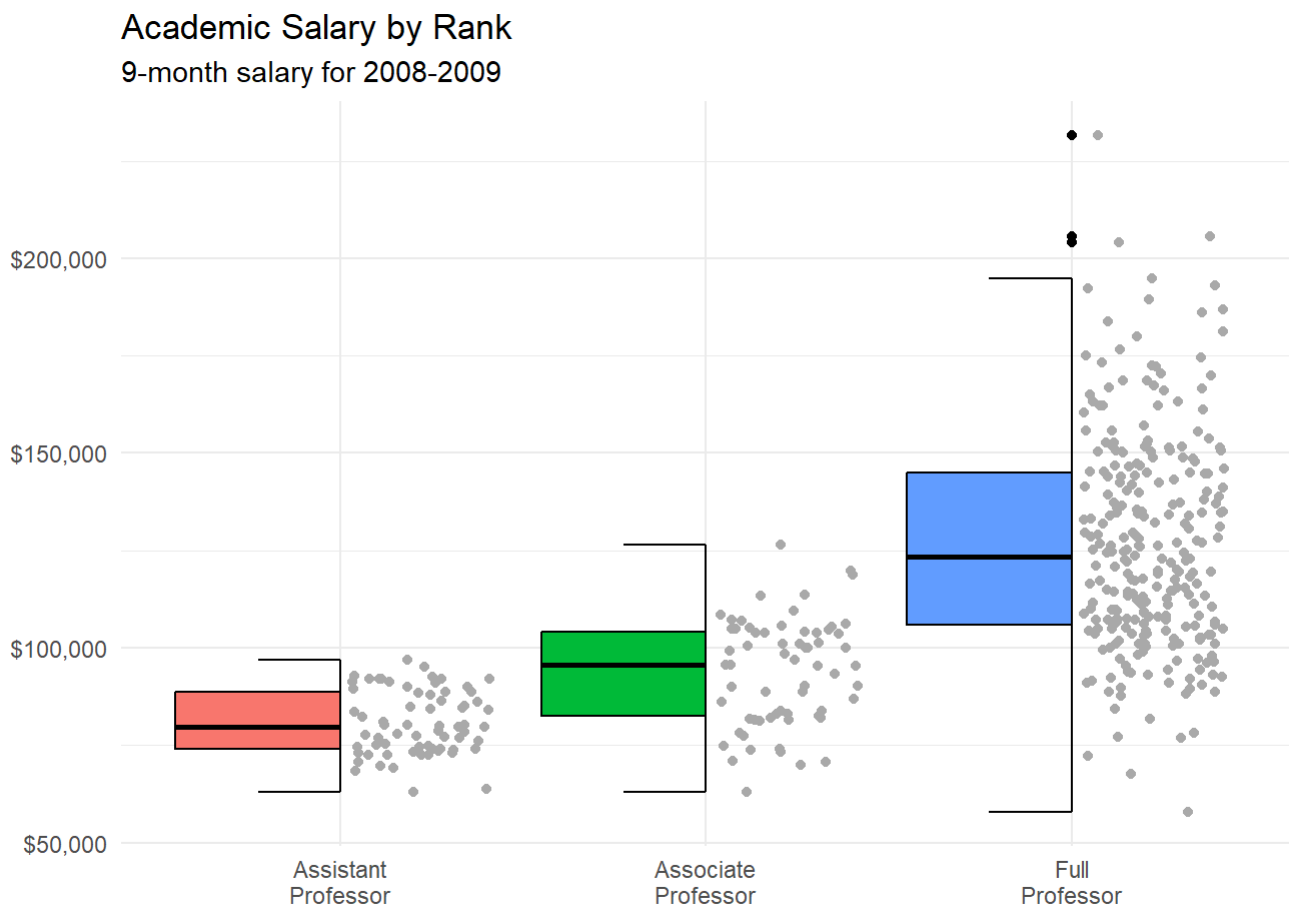


Figure 4.31: Using `geom_boxjitter`

4.3.8 Beeswarm Plots

Beeswarm plots (also called violin scatter plots) are similar to jittered scatterplots, in that they display the distribution of a quantitative variable by plotting points in way that reduces overlap. In addition, they also help display the density of the data at each point (in a manner that is similar to a [violin plot](#)). Continuing the previous example

```

# plot the distribution of salaries
# by rank using beewarm-style plots
library(ggbeeswarm)
library(scales)
ggplot(Salaries,
       aes(x = factor(rank,
                      labels = c("Assistant\nProfessor",
                                "Associate\nProfessor",
                                "Full\nProfessor")),
          y = salary,
          color = rank)) +
  geom_quasirandom(alpha = 0.7,
                  size = 1.5) +
  scale_y_continuous(label = dollar) +
  labs(title = "Academic Salary by Rank",
       subtitle = "9-month salary for 2008-2009",
       x = "",
       y = "") +
  theme_minimal() +
  theme(legend.position = "none")

```

Academic Salary by Rank
9-month salary for 2008-2009

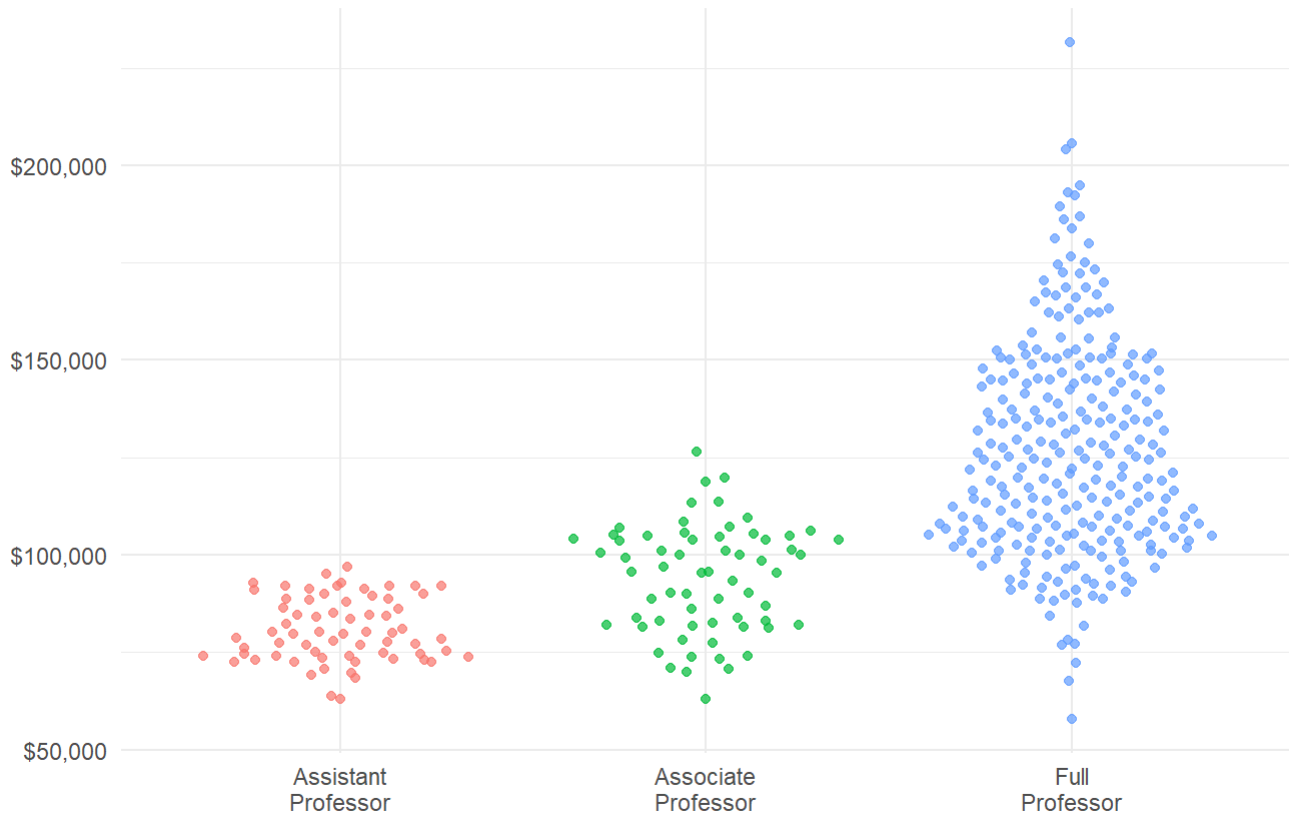


Figure 4.32: Beeswarm plot

The plots are create using the `geom_quasirandom` function. These plots can be easier to read than simple jittered strip plots. To learn more about these plots, see [Beeswarm-style plots with ggplot2](#).

4.3.9 Cleveland Dot Charts

Cleveland plots are useful when you want to compare a numeric statistic for a large number of groups. For example, say that you want to compare the 2007 life expectancy for Asian country using the [gapminder](#) dataset.

```
data(gapminder, package="gapminder")

# subset Asian countries in 2007
library(dplyr)
plotdata <- gapminder %>%
  filter(continent == "Asia" &
         year == 2007)

# basic Cleveland plot of life expectancy by country
ggplot(plotdata,
       aes(x= lifeExp, y = country)) +
  geom_point()
```

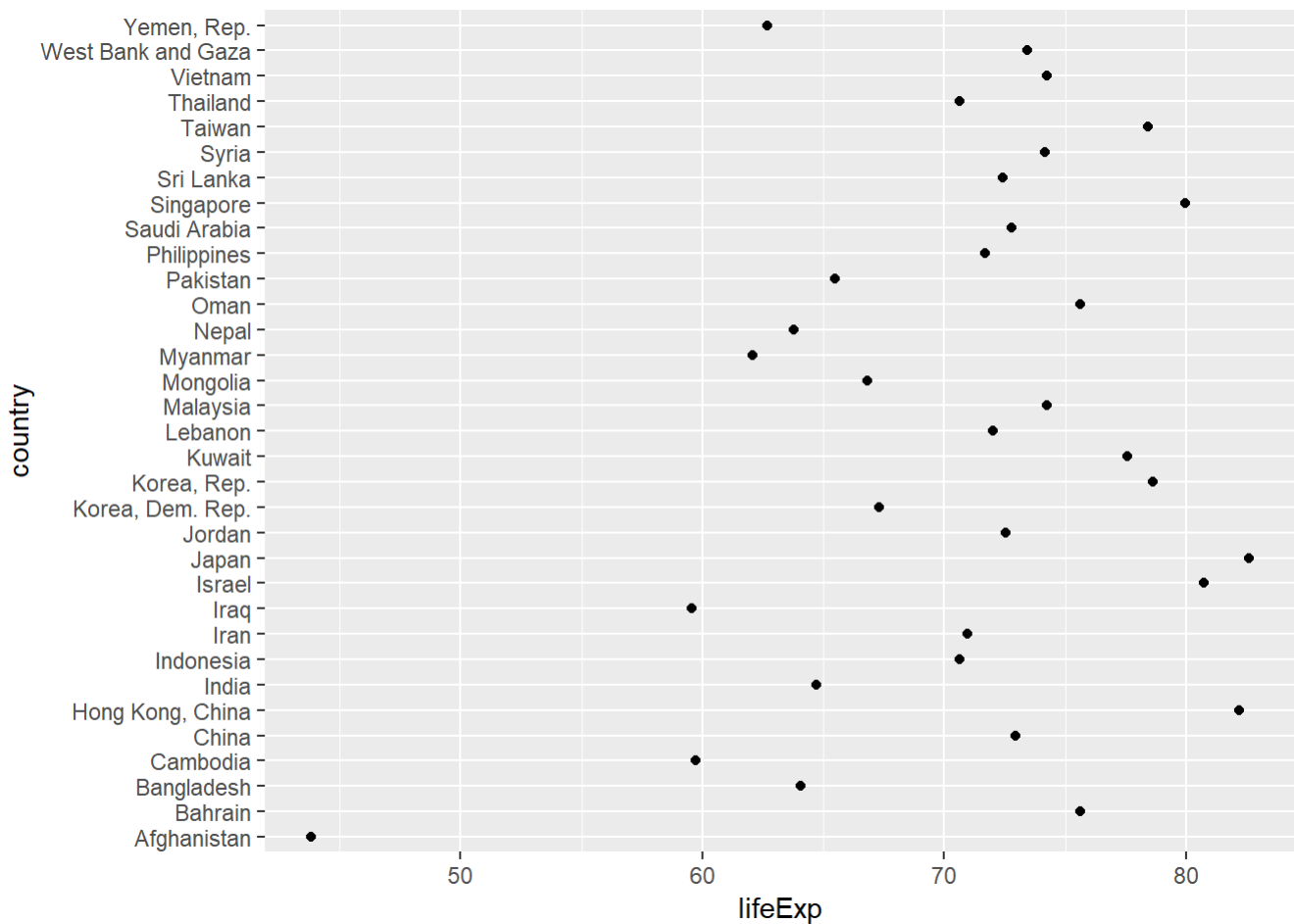


Figure 4.33: Basic Cleveland dot plot

Comparisons are usually easier if the y-axis is sorted.

```
# Sorted Cleveland plot
ggplot(plotdata,
  aes(x=lifeExp,
    y=reorder(country, lifeExp))) +
  geom_point()
```

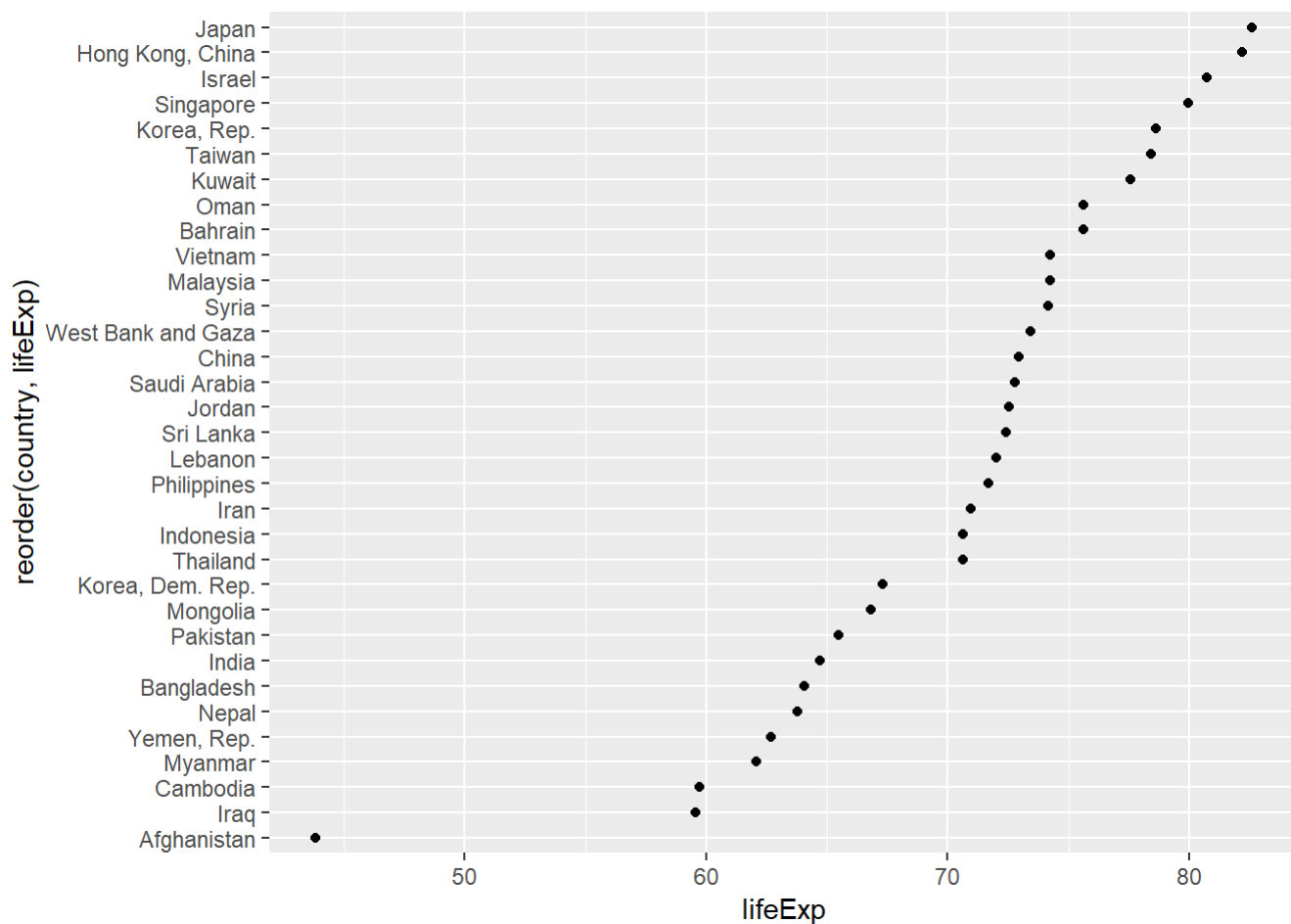


Figure 4.34: Sorted Cleveland dot plot

Finally, we can use options to make the graph more attractive.

```

# Fancy Cleveland plot
ggplot(plotdata,
       aes(x=lifeExp,
           y=reorder(country, lifeExp))) +
geom_point(color="blue",
           size = 2) +
geom_segment(aes(x = 40,
                 xend = lifeExp,
                 y = reorder(country, lifeExp),
                 yend = reorder(country, lifeExp)),
            color = "lightgrey") +
labs (x = "Life Expectancy (years)",
      y = "",
      title = "Life Expectancy by Country",
      subtitle = "GapMinder data for Asia - 2007") +
theme_minimal() +
theme(panel.grid.major = element_blank(),
      panel.grid.minor = element_blank())

```


Life Expectancy by Country

GapMinder data for Asia - 2007

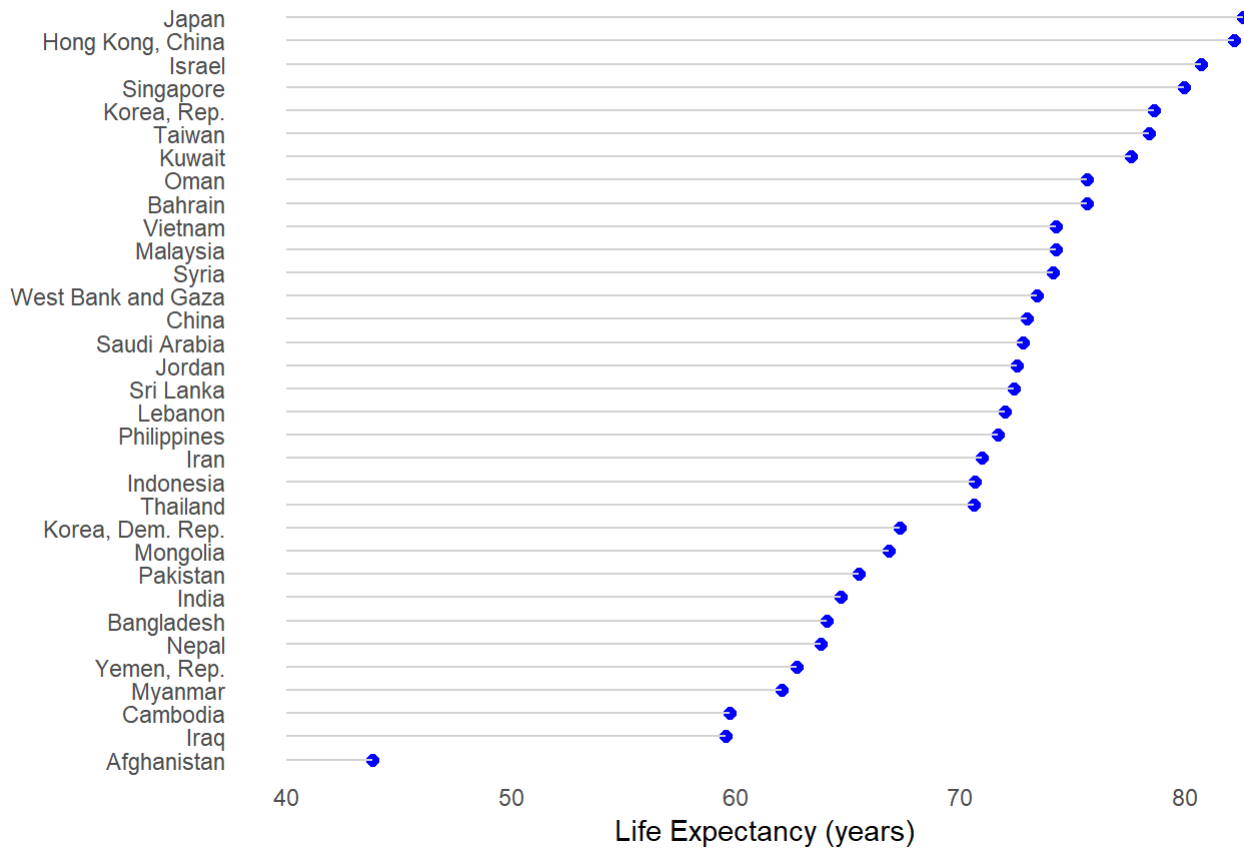


Figure 4.35: Fancy Cleveland plot

Japan clearly has the highest life expectancy, while Afghanistan has the lowest by far. This last plot is also called a lollipop graph (you can see why).

Chapter 5 Multivariate Graphs

Multivariate graphs display the relationships among three or more variables. There are two common methods for accommodating multiple variables: grouping and faceting.

5.1 Grouping

In grouping, the values of the first two variables are mapped to the x and y axes. Then additional variables are mapped to other visual characteristics such as color, shape, size, line type, and transparency. Grouping allows you to plot the data for multiple groups in a single graph.

Using the [Salaries](#) dataset, let's display the relationship between *yrs.since.phd* and *salary*.

```
library(ggplot2)
data(Salaries, package="carData")

# plot experience vs. salary
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary)) +
  geom_point() +
  labs(title = "Academic salary by years since degree")
```

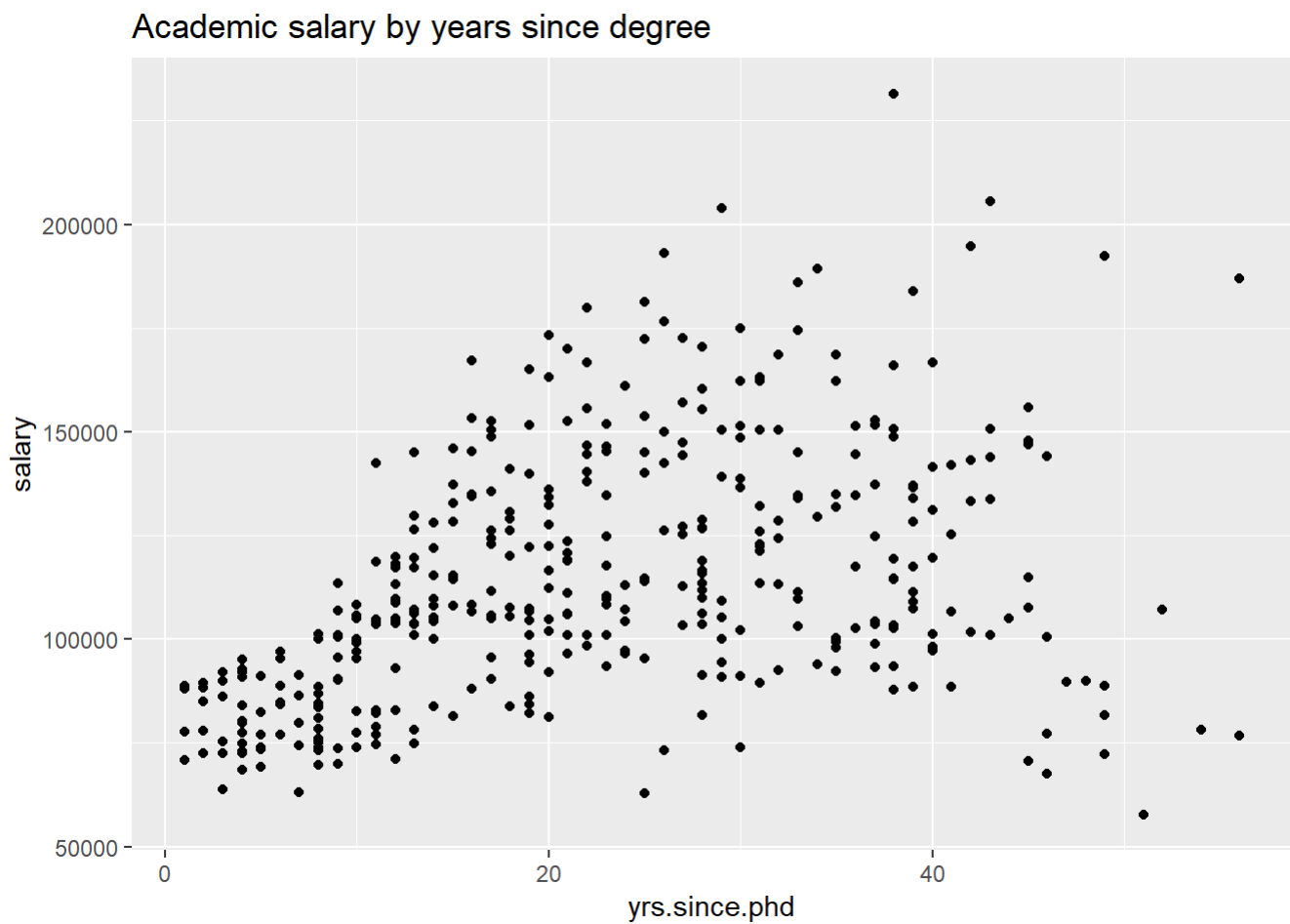


Figure 5.1: Simple scatterplot

Next, let's include the rank of the professor, using color.

```
# plot experience vs. salary (color represents rank)
ggplot(Salaries, aes(x = yrs.since.phd,
                    y = salary,
                    color=rank)) +
  geom_point() +
  labs(title = "Academic salary by rank and years since degree")
```

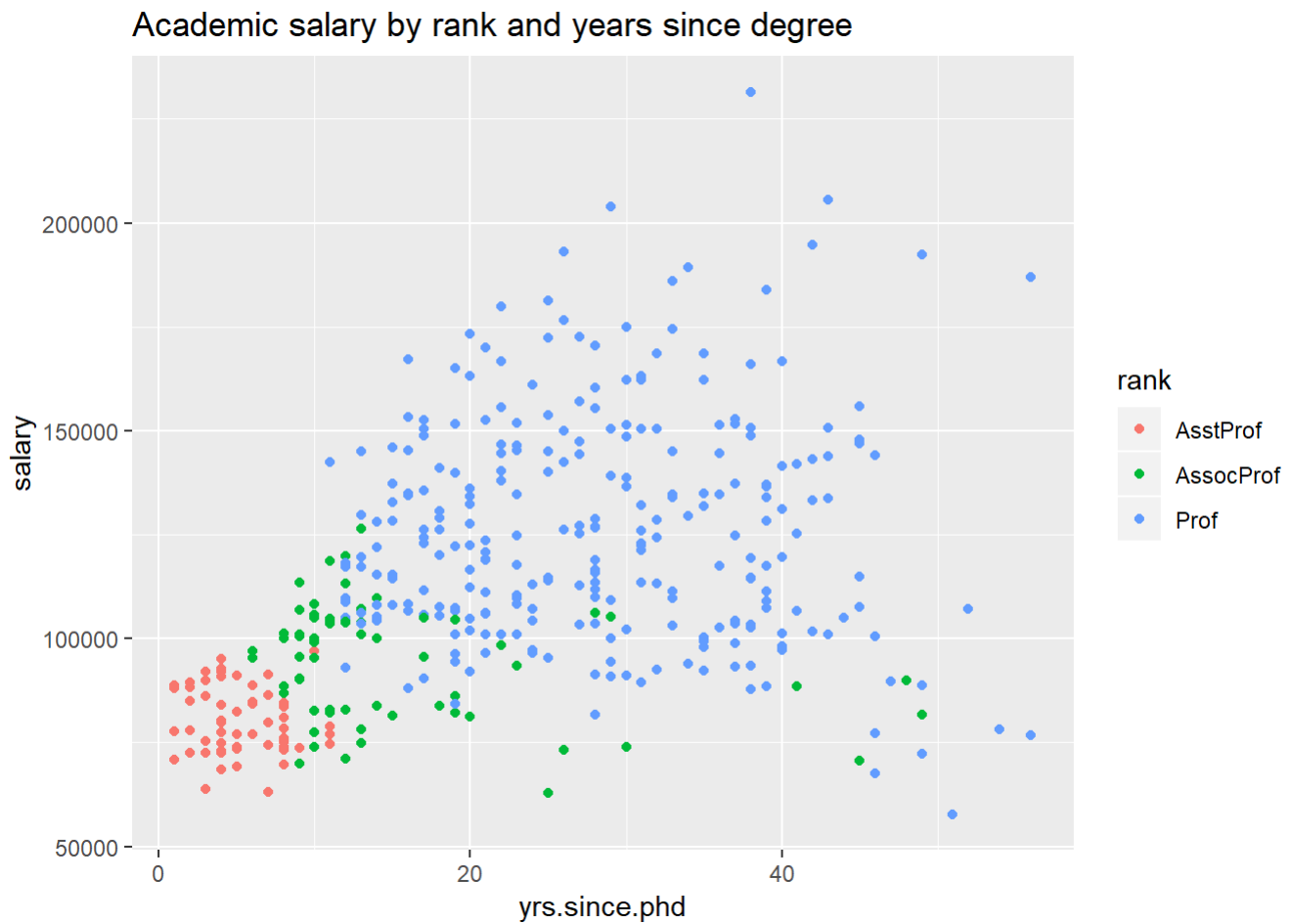


Figure 5.2: Scatterplot with color mapping

Finally, let's add the gender of professor, using the shape of the points to indicate sex. We'll increase the point size and add transparency to make the individual points clearer.

```
# plot experience vs. salary
# (color represents rank, shape represents sex)
ggplot(Salaries,
  aes(x = yrs.since.phd,
    y = salary,
    color = rank,
    shape = sex)) +
  geom_point(size = 3,
    alpha = .6) +
  labs(title = "Academic salary by rank, sex, and years since degree")
```

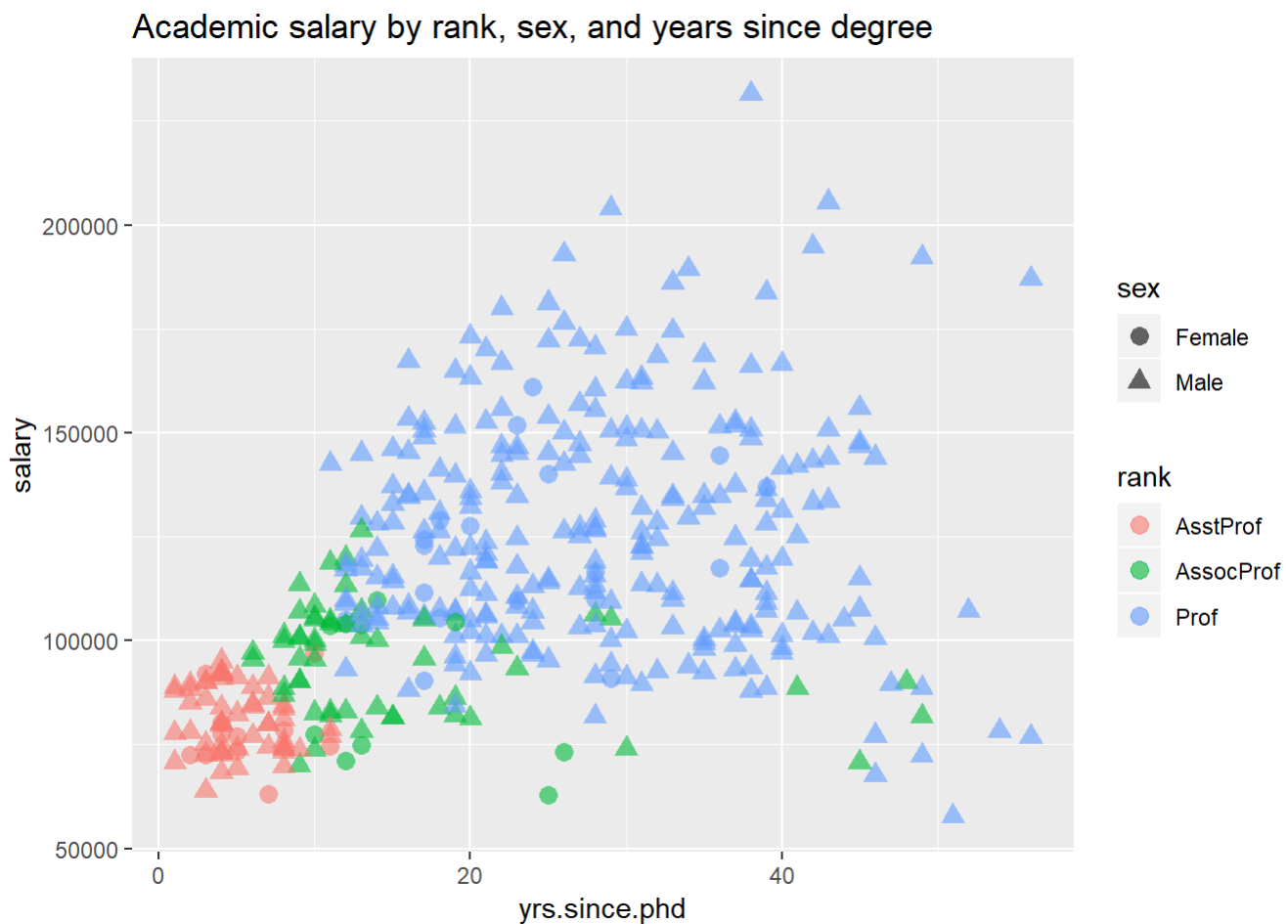


Figure 5.3: Scatterplot with color and shape mapping

I can't say that this is a great graphic. It is very busy, and it can be difficult to distinguish male from female professors. Faceting (described in the next section) would probably be a better approach.

Notice the difference between specifying a constant value (such as `size = 3`) and a mapping of a variable to a visual characteristic (e.g., `color = rank`). Mappings are always placed within the `aes` function, while the assignment of a constant value always appear outside of the `aes` function.

Here is a cleaner example. We'll graph the relationship between years since Ph.D. and salary using the size of the points to indicate years of service. This is called a bubble plot.

```
library(ggplot2)
data(Salaries, package="carData")

# plot experience vs. salary
# (color represents rank and size represents service)
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary,
           color = rank,
           size = yrs.service)) +
  geom_point(alpha = .6) +
  labs(title = "Academic salary by rank, years of service, and years since degree")
```

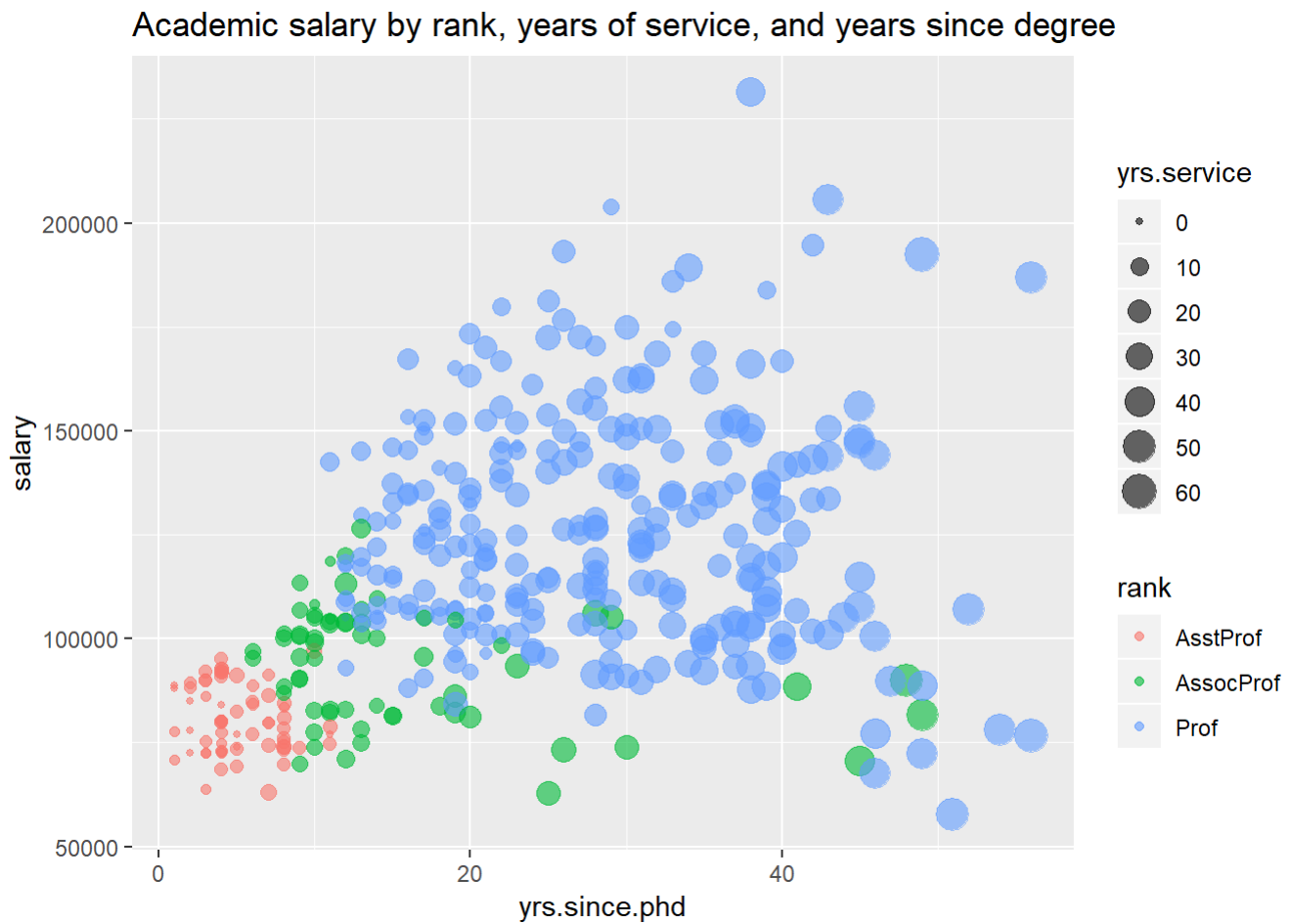


Figure 5.4: Scatterplot with size and color mapping

There is obviously a strong positive relationship between years since Ph.D. and year of service. Assistant Professors fall in the 0-11 years since Ph.D. and 0-10 years of service range. Clearly highly experienced professionals don't stay at the Assistant Professor level (they are probably promoted or leave the University). We don't find the same time demarcation between Associate and Full Professors.

Bubble plots are described in more detail in a later chapter.

As a final example, let's look at the *yrs.since.phd* vs *salary* and add *sex* using color and quadratic best fit lines.

```
# plot experience vs. salary with
# fit lines (color represents sex)
ggplot(Salaries,
       aes(x = yrs.since.phd,
           y = salary,
           color = sex)) +
  geom_point(alpha = .4,
             size = 3) +
  geom_smooth(se=FALSE,
             method = "lm",
             formula = y~poly(x,2),
             size = 1.5) +
  labs(x = "Years Since Ph.D.",
       title = "Academic Salary by Sex and Years Experience",
       subtitle = "9-month salary for 2008-2009",
       y = "",
       color = "Sex") +
  scale_y_continuous(label = scales::dollar) +
  scale_color_brewer(palette = "Set1") +
  theme_minimal()
```

Academic Salary by Sex and Years Experience

9-month salary for 2008-2009

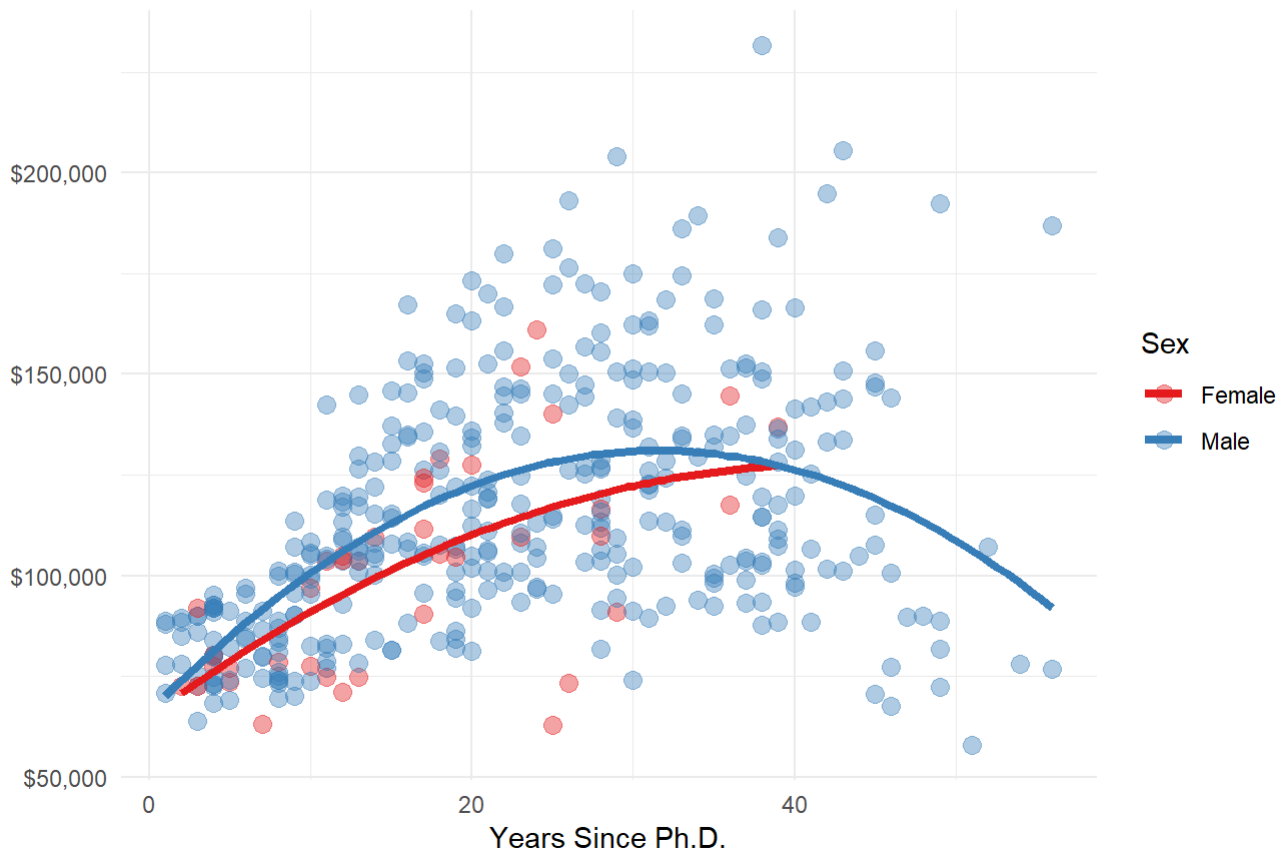


Figure 5.5: Scatterplot with color mapping and quadratic fit lines

5.2 Faceting

Grouping allows you to plot multiple variables in a single graph, using visual characteristics such as color, shape, and size.

In faceting, a graph consists of several separate plots or *small multiples*, one for each level of a third variable, or combination of variables. It is easiest to understand this with an example.

```
# plot salary histograms by rank
ggplot(Salaries, aes(x = salary)) +
  geom_histogram(fill = "cornflowerblue",
                 color = "white") +
  facet_wrap(~rank, ncol = 1) +
  labs(title = "Salary histograms by rank")
```


Salary histograms by rank

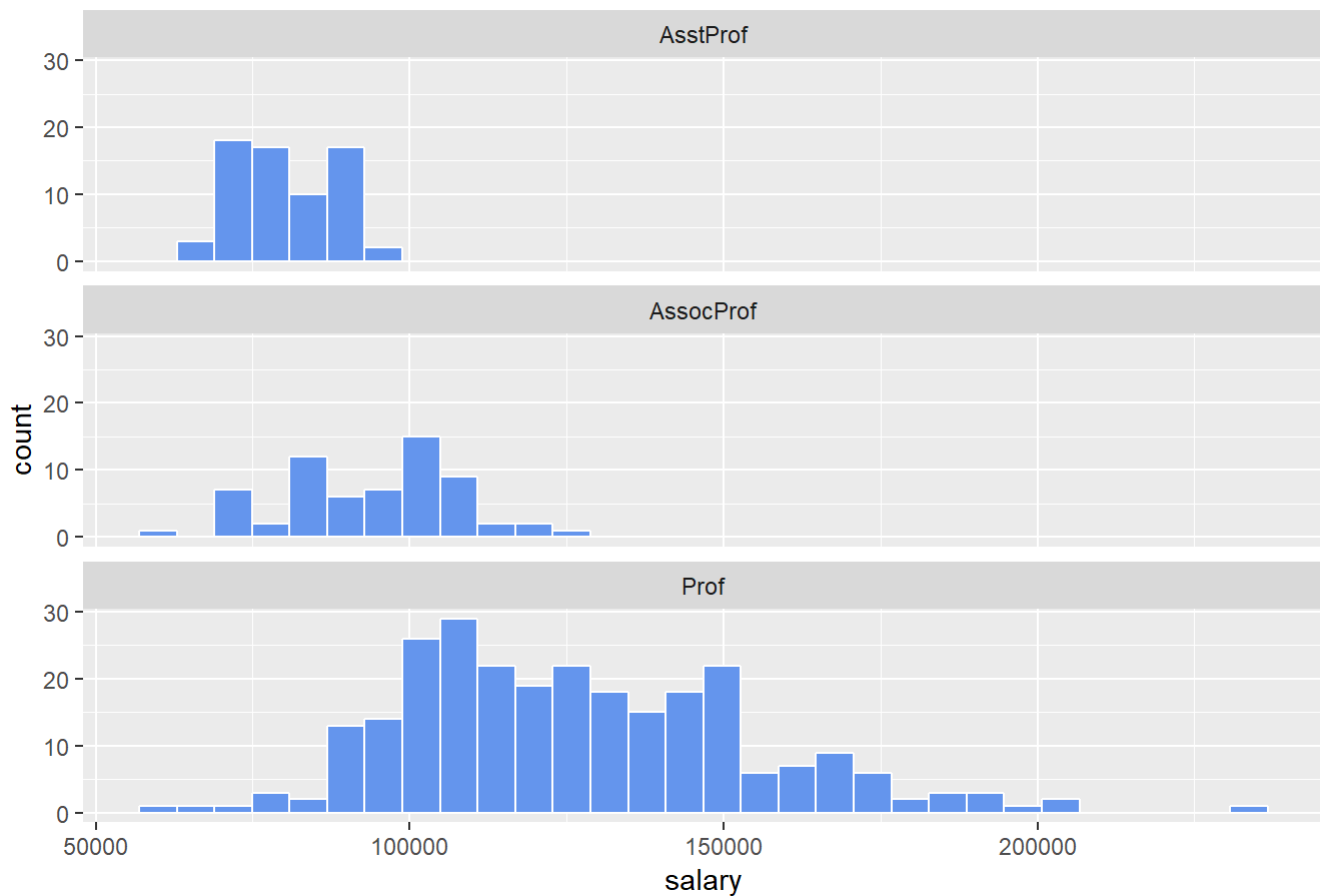


Figure 5.6: Salary distribution by rank

The `facet_wrap` function creates a separate graph for each level of `rank`. The `ncol` option controls the number of columns.

In the next example, two variables are used to define the facets.

```
# plot salary histograms by rank and sex
ggplot(Salaries, aes(x = salary / 1000)) +
  geom_histogram(color = "white",
                 fill = "cornflowerblue") +
  facet_grid(sex ~ rank) +
  labs(title = "Salary histograms by sex and rank",
       x = "Salary ($1000)")
```

Salary histograms by sex and rank

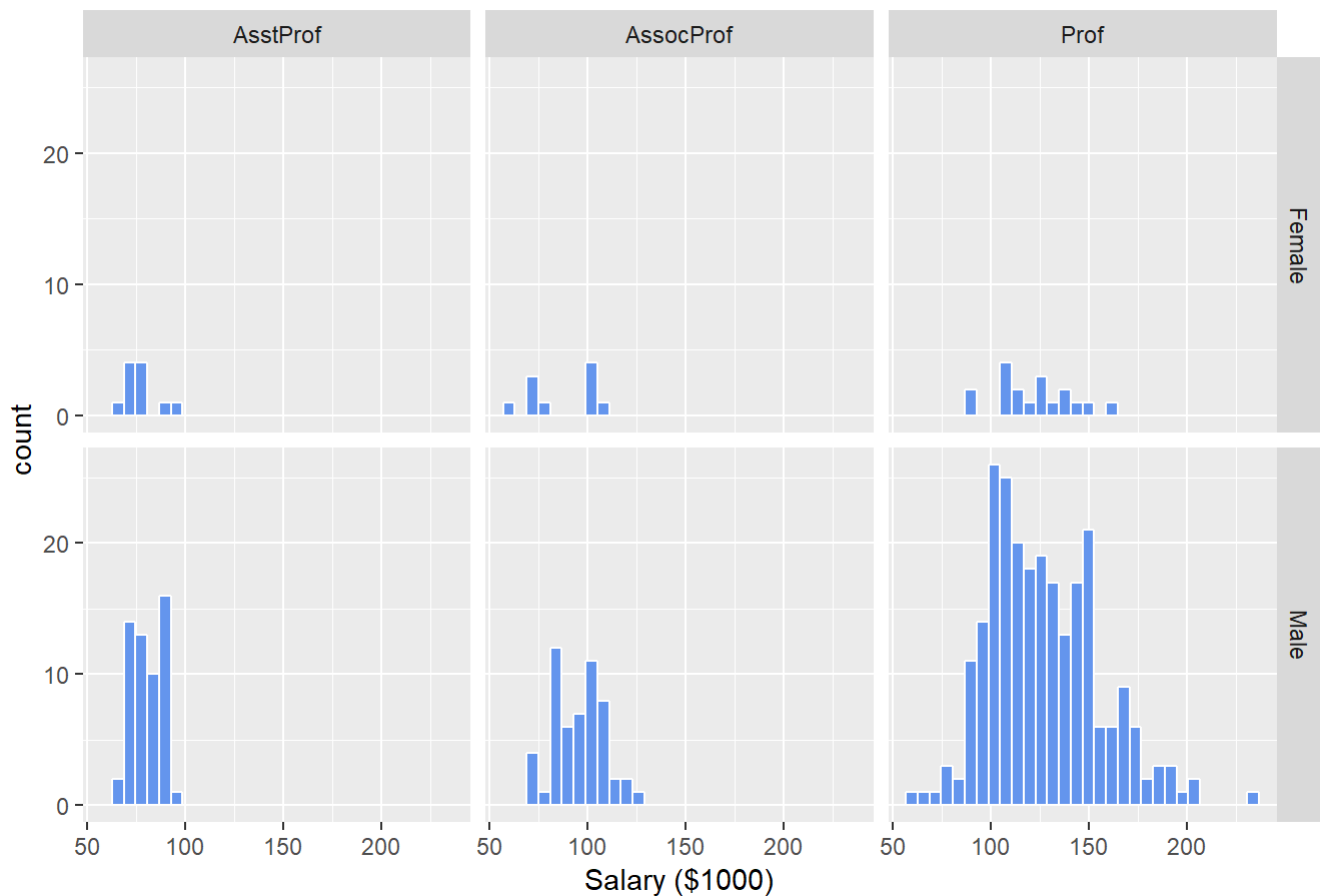


Figure 5.7: Salary distribution by rank and sex

The format of the `facet_grid` function is

```
facet_grid( row variable(s) ~ column variable(s))
```

Here, the function assigns `sex` to the rows and `rank` to the columns, creating a matrix of 6 plots in one graph.

We can also combine grouping and faceting. Let's use [Mean/SE plots](#) and faceting to compare the salaries of male and female professors, within rank and discipline. We'll use color to distinguish sex and faceting to create plots for rank by discipline combinations.

```

# calculate means and standard erroes by sex,
# rank and discipline

library(dplyr)

plotdata <- Salaries %>%
  group_by(sex, rank, discipline) %>%
  summarize(n = n(),
            mean = mean(salary),
            sd = sd(salary),
            se = sd / sqrt(n))

# create better labels for discipline
plotdata$discipline <- factor(plotdata$discipline,
                              labels = c("Theoretical",
                                          "Applied"))

# create plot
ggplot(plotdata,
       aes(x = sex,
           y = mean,
           color = sex)) +
  geom_point(size = 3) +
  geom_errorbar(aes(ymin = mean - se,
                   ymax = mean + se),
               width = .1) +
  scale_y_continuous(breaks = seq(70000, 140000, 10000),
                    label = scales::dollar) +
  facet_grid(. ~ rank + discipline) +
  theme_bw() +
  theme(legend.position = "none",
        panel.grid.major.x = element_blank(),
        panel.grid.minor.y = element_blank()) +
  labs(x="",
       y="",
       title="Nine month academic salaries by gender, discipline, and rank",

```

```

  subtitle = "(Means and standard errors)") +
  scale_color_brewer(palette="Set1")

```

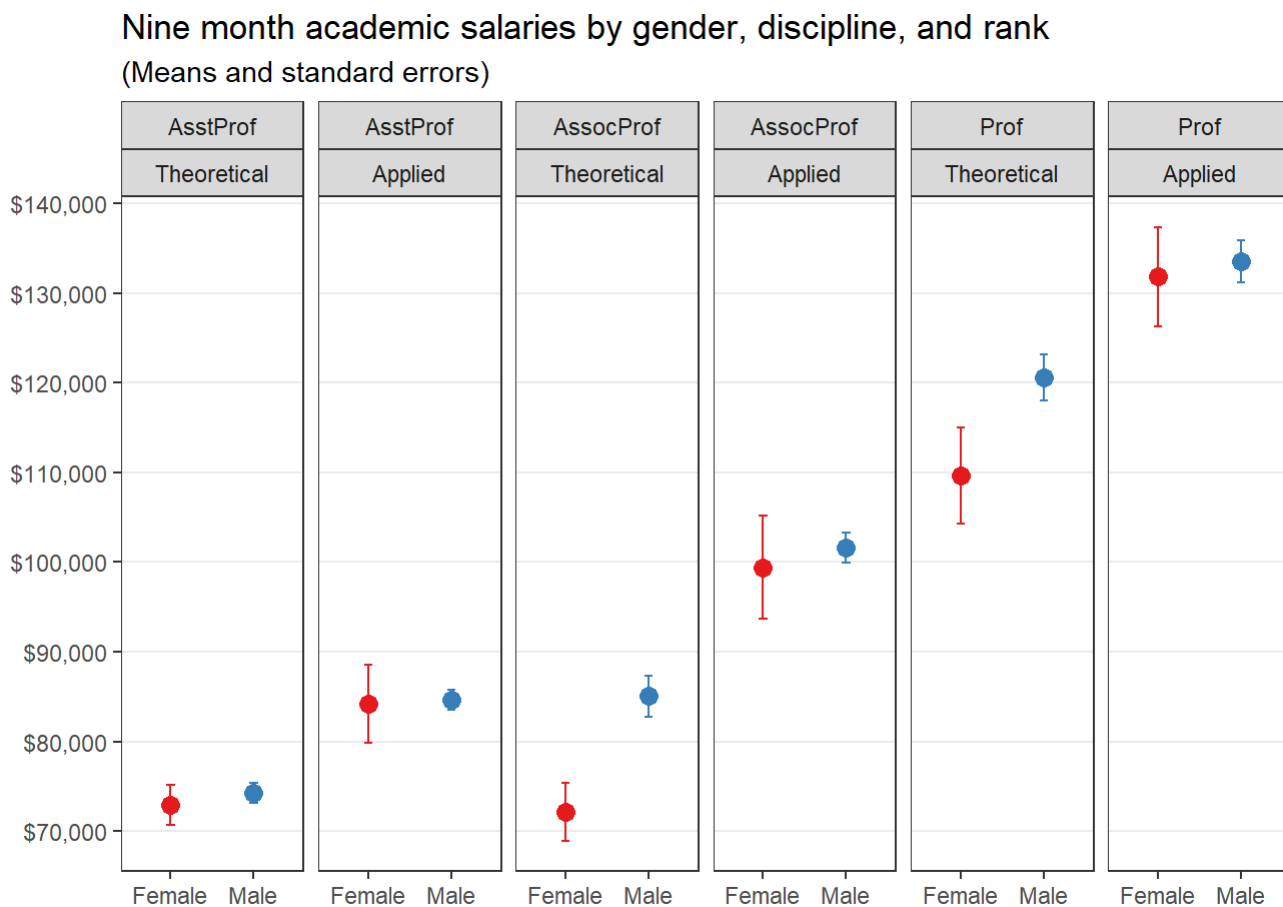


Figure 5.8: Salary by sex, rank, and discipline

The statement `facet_grid(. ~ rank + discipline)` specifies no row variable (.) and columns defined by the combination of *rank* and *discipline*.

The `theme_` functions create a black and white theme and eliminates vertical grid lines and minor horizontal grid lines. The `scale_color_brewer` function changes the color scheme for the points and error bars.

At first glance, it appears that there might be gender differences in salaries for associate and full professors in theoretical fields. I say “might” because we haven’t done any formal hypothesis testing yet (ANCOVA in this case).

See the [Customizing](#) section to learn more about customizing the appearance of a graph.

As a final example, we'll shift to a new dataset and plot the change in life expectancy over time for countries in the "Americas". The data comes from the [gapminder](#) dataset in the `gapminder` package. Each country appears in its own facet. The theme functions are used to simplify the background color, rotate the x-axis text, and make the font size smaller.

```
# plot life expectancy by year separately
# for each country in the Americas
data(gapminder, package = "gapminder")

# Select the Americas data
plotdata <- dplyr::filter(gapminder,
                          continent == "Americas")

# plot life expectancy by year, for each country
ggplot(plotdata, aes(x=year, y = lifeExp)) +
  geom_line(color="grey") +
  geom_point(color="blue") +
  facet_wrap(~country) +
  theme_minimal(base_size = 9) +
  theme(axis.text.x = element_text(angle = 45,
                                    hjust = 1)) +
  labs(title = "Changes in Life Expectancy",
        x = "Year",
        y = "Life Expectancy")
```

Changes in Life Expectancy

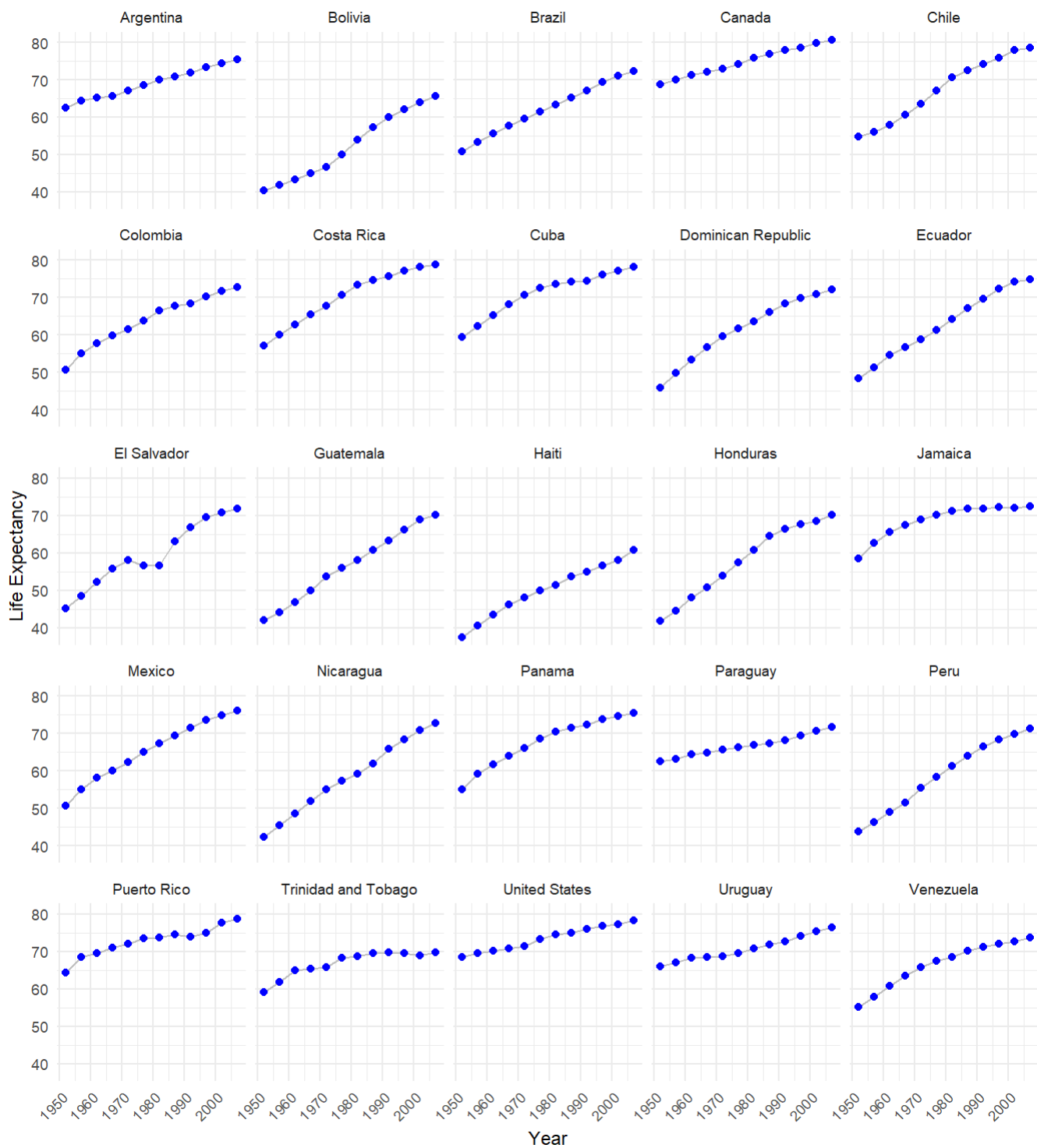


Figure 5.9: Changes in life expectancy by country

We can see that life expectancy is increasing in each country, but that Haiti is lagging behind.