

R Data Visualization Primer

Load ggplot2

You will want to use the *ggplot2* package for all of your data visualization. Make sure to load it before you try to generate any graphs:

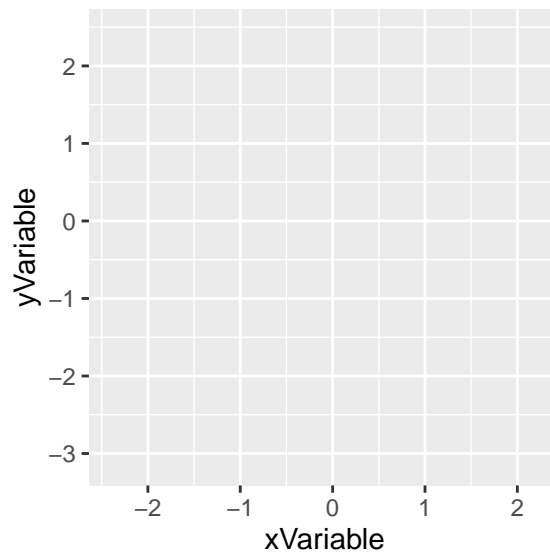
```
#load ggplot2  
#Make sure you have installed ggplot2 before you have run this  
#If not, first run: install.packages("ggplot2")  
library(ggplot2)
```

Generating a plot

ggplot2 works by generating a plot base and adding elements to that plot one option at a time.

To generate this plot base, use the *ggplot()* function:

```
plotName <- ggplot(data = dataframe,  
                  aes(x = xVariable,  
                     y = yVariable))  
  
#plotName = a data structure for storing your plot  
#dataFrame = the dataframe containing your data  
#xVariable = your x-axis variable  
#yVariable = your y-axis variable  
  
#To view your plot, execute:  
plotName
```



The *aes()* option inside the *ggplot2()* function sets the aesthetic options for your plot. You can use this to set, among other things:

- Which variables you plot and on which axes
- The colors and fills of your graph
- The size, transparency, and shape of graph elements

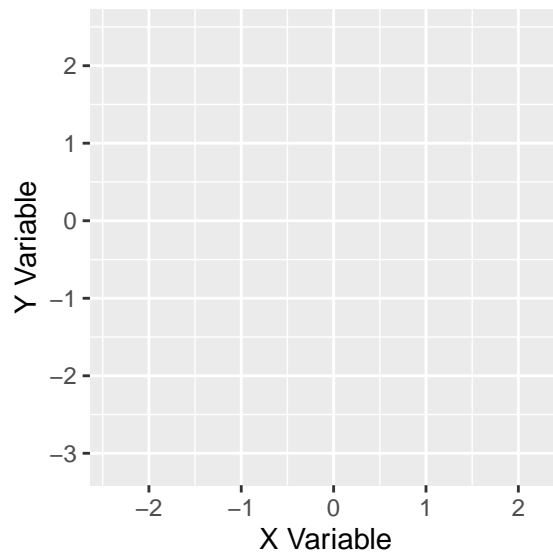
To add on to your base plot, you add additional functions using the `+`. For example, to label your axes, use the `labs()` function:

```
plotName <- plotName +
  labs(x = "X Variable", #<---
       y = "Y Variable") #<---
```

#labs() = labels elements of your graph
#You can use this to label the axes, legend, graph title, and more
#Run ?labs to see all of your options

#Notice that instead of regenerating the whole plot from scratch every time, you can add new elements o
#Both approaches work, but if you choose to modify an existing plot, you have to be mindful of what you

#To view your plot, execute:
 plotName



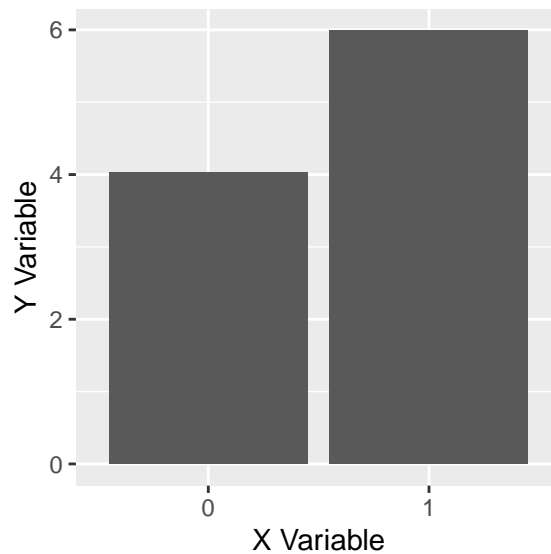
Bar Plots

To make a simple bar plot of group means across a categorical predictor, add the `stat_summary()` function:

```
plotName <- ggplot(data = dataFrame,
  aes(x = xVariable,
      y = yVariable)) +
  labs(x = "X Variable",
      y = "Y Variable") +
  stat_summary(fun = "mean", #<---
    geom = "bar", #<---
    position = "dodge") #<---

#stat_summary = computes and plots a summary statistic based on your data
#fun = the statistic you want to compute
#geom = the shape you want to plot
#position = "dodge" = place the bars side by size, rather than stack

#To view your plot, execute:
plotName
```



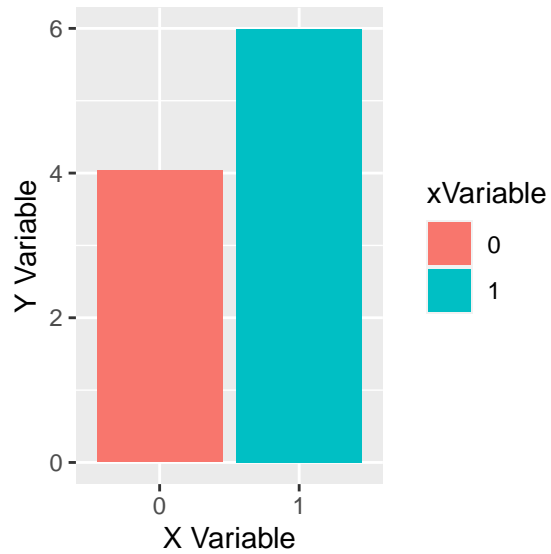
You can add color using the `fill =` option inside `aes()`:

```
plotName <- ggplot(data = dataFrame,
  aes(x = xVariable,
      y = yVariable,
      fill = xVariable)) + #<---
  labs(x = "X Variable",
      y = "Y Variable") +
  stat_summary(fun = "mean",
    geom = "bar",
    position = "dodge")
```

#fill = fills in geoms according to the variable specified

#To view your plot, execute:

plotName



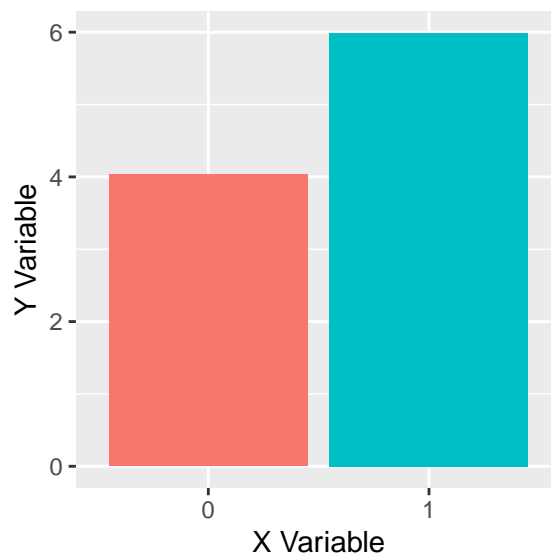
You can remove the legend using the *guides()* function:

```
plotName <- plotName +  
  guides(fill="none") #<---
```

#guides(fill="none") = blanks out the legend for color fill

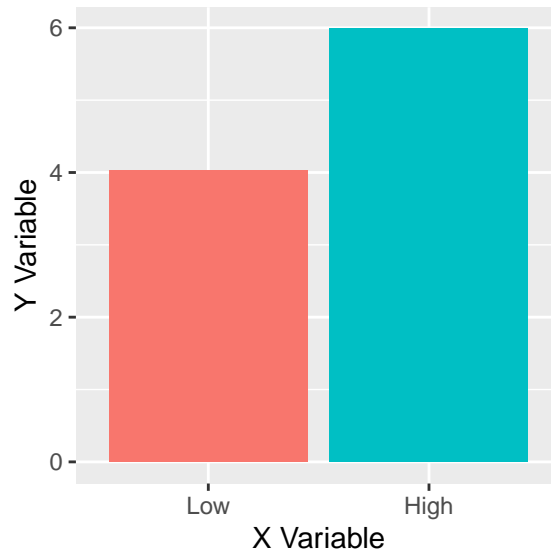
#To view your plot, execute:

plotName



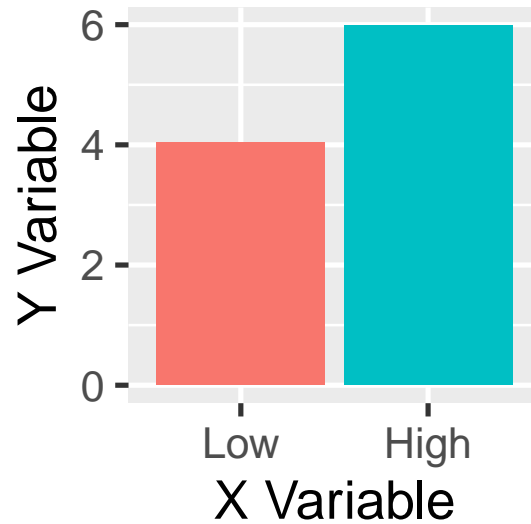
You can also relabel the x-axis ticks using `scale_x_discrete()`:

```
plotName <- plotName +  
  scale_x_discrete(labels = c("Low", "High")) #<---  
  
#scale_x_discrete = lets you set options for the x-axis, including the axis ticks  
#do ?scale_x_discrete to see all options  
#Use scale_y_discrete for a categorical y-axis  
#Use scale_x_continuous or scale_y_continuous for continuous variables  
  
#To view your plot, execute:  
plotName
```



You can set the text size for the whole graph using `base_size =` inside `theme_gray()`:

```
plotName <- plotName +  
  theme_gray(base_size = 20) #<---  
  
#theme_gray() = determines the overall look of the graph and allows you to set global parameters, like  
#Other themes (such as classic() and bw()) are available.  
  
#To view your plot, execute:  
plotName
```



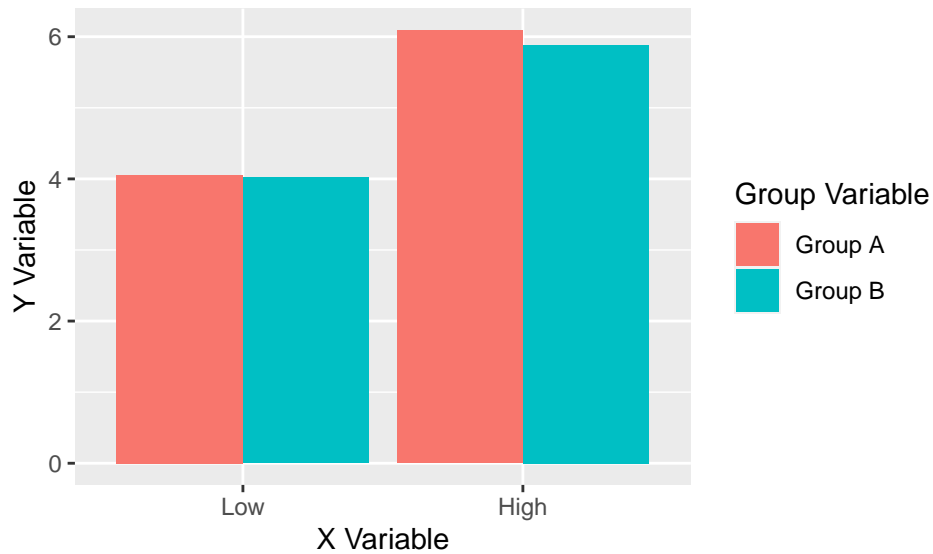
To add a grouping variable, use it as the *fill()* variable:

```
plotName <- ggplot(data = dataframe,
  aes(x = xVariable,
      y = yVariable,
      fill = groupVariable)) + #<---
  labs(x = "X Variable",
      y = "Y Variable",
      fill = "Group Variable") + #<---
  stat_summary(fun = "mean",
      geom = "bar",
      position = "dodge") +
  scale_x_discrete(labels = c("Low", "High")) +
  scale_fill_discrete(labels = c("Group A", "Group B")) #<---
```

#scale_fill_discrete() = sets options for the color fill and the legend, such as the legend labels

#To view your plot, execute:

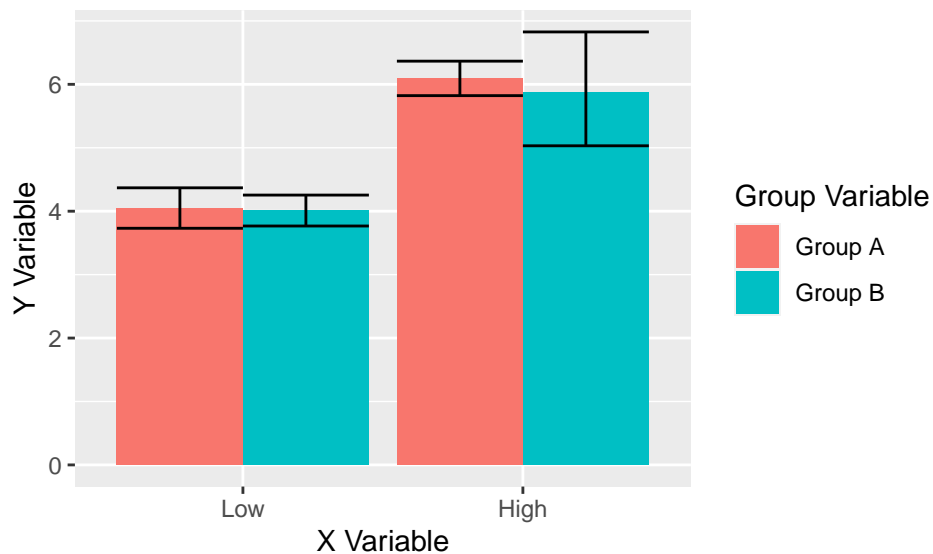
```
plotName
```



Finally, to add error bars representing 95% confidence intervals, you can use the *“mean_cl_boot”* option in *stat_summary*:

```
plotName <- plotName +
  stat_summary(fun.data = "mean_cl_boot", #<---
    geom = "errorbar", #<---
    position = "dodge") #<---
```

#To view your plot, execute:
plotName



Box and Violin Plots

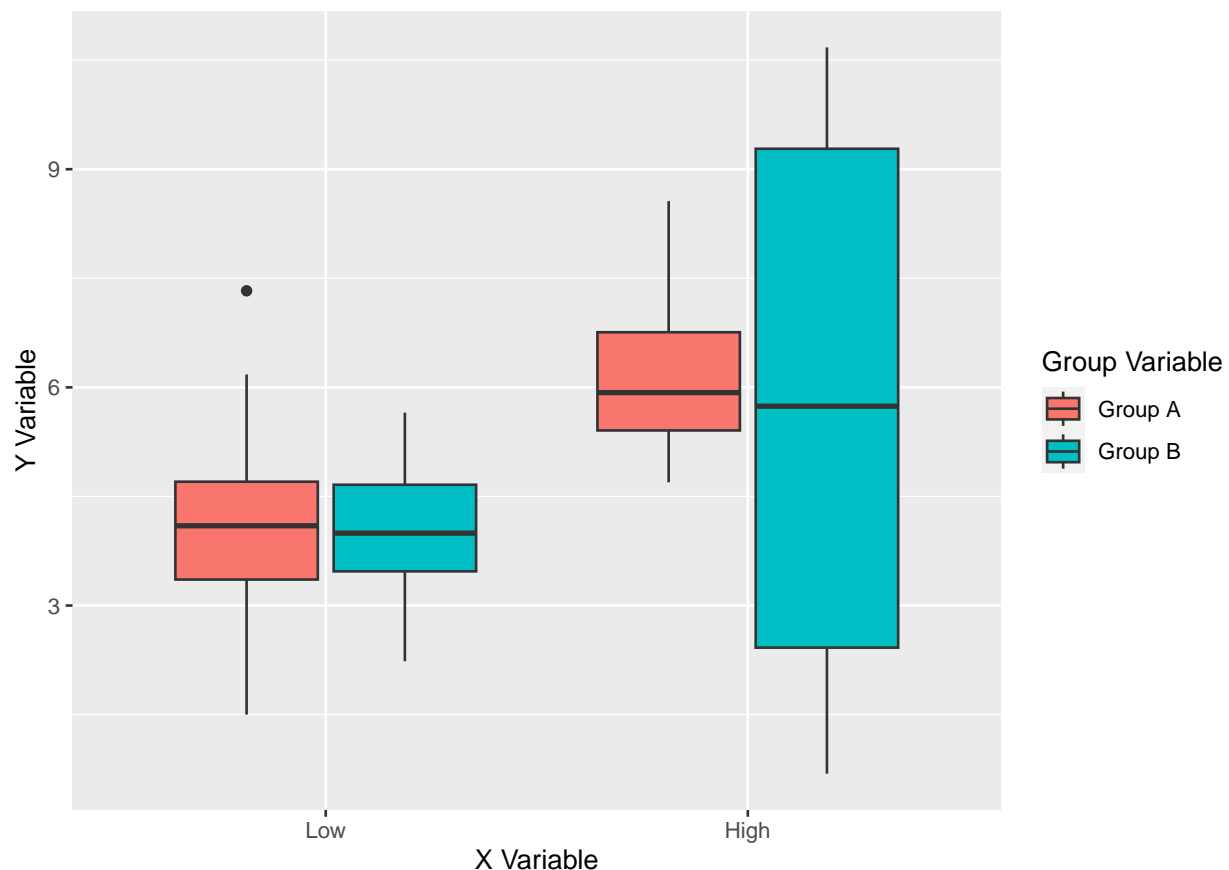
Bar plots are popular. But unless you are working with count data, bar plots are often not very good. Plots that actually show the distribution of the data are often more appropriate.

One example is the boxplot, which you can add with the `geom_boxplot()` function. In a box plot, the thick line represents the median, the edges of the box represent the IQR (bottom 25th percentile to top 75th percentile), and the “whiskers” extend out from the box up to $1.5 * \text{IQR}$. Any dots are potential outliers that sit very far from the rest of the data.

```
plotName <- ggplot(data = dataFrame,
  aes(x = xVariable,
      y = yVariable,
      fill = groupVariable)) +
  labs(x = "X Variable",
      y = "Y Variable",
      fill = "Group Variable") +
  scale_x_discrete(labels = c("Low", "High")) +
  scale_fill_discrete(labels = c("Group A", "Group B")) +
  geom_boxplot() #<---

#geom_boxplot() = plots data as a boxplot

#To view your plot, execute:
plotName
```

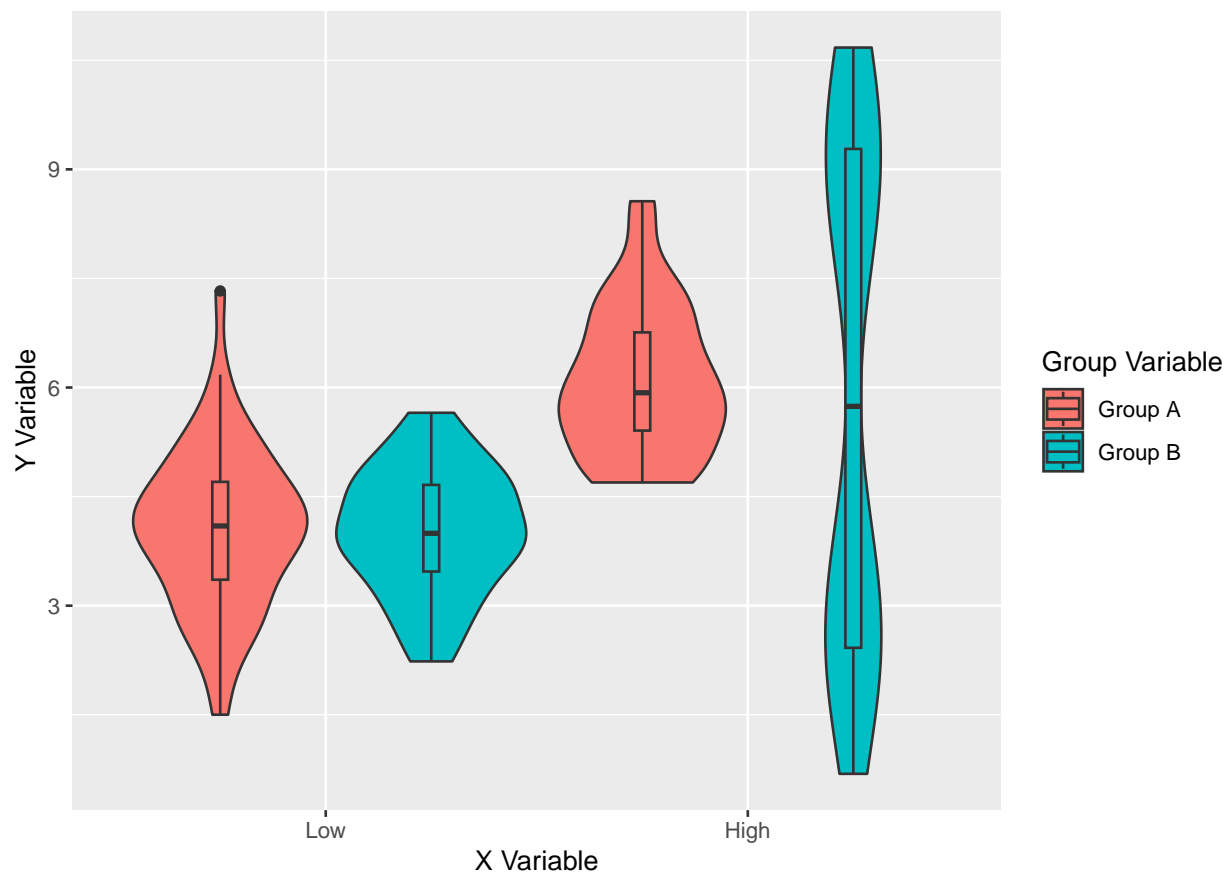


Another is the violin plot, which you can add with `geom_violin()`. A violin plot is like a sideways density distribution (think a smoothed histogram turned on its side). It is thicker where you have more data and thinner where you have less data. This can combine well with a boxplot:

```
plotName <- ggplot(data = dataFrame,
  aes(x = xVariable,
      y = yVariable,
      fill = groupVariable)) +
  labs(x = "X Variable",
      y = "Y Variable",
      fill = "Group Variable") +
  scale_x_discrete(labels = c("Low", "High")) +
  scale_fill_discrete(labels = c("Group A", "Group B")) +
  geom_violin(position = position_dodge(1)) + #<---
  geom_boxplot(position = position_dodge(1), #<---
    width = .075) #<---
```

#geom_violin() = plots data as a violin plot
#position_dodge() = sets the horizontal positions of the plot elements (useful if you want them to overlap)
#width = sets the width of the plot element

#To view your plot, execute:
 plotName



Note: both of these graphs make clear in this case that the distribution of the 4th group is different from the others. This is obscured in the bar plot. In fact, it's particularly clear from the violin plot that the 4th group is bimodal. This is a key advantage of a violin plot over a bar plot.

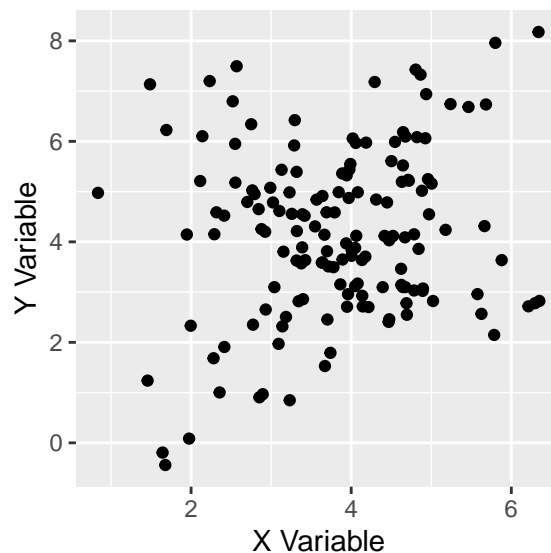
Scatter Plots

To make a scatter plot showing the relationship between two continuous variables, use `geom_point()`:

```
plotName <- ggplot(data = dataframe,
                   aes(x = xVariable,
                       y = yVariable)) +
  labs(x = "X Variable",
       y = "Y Variable") +
  geom_point() #<---

#geom_point() = plot data as points

#To view your plot, execute:
plotName
```



To add a trend line, use the `geom_smooth()` function:

```
plotName <- ggplot(data = dataframe,
                   aes(x = xVariable,
                       y = yVariable)) +
  labs(x = "X Variable",
       y = "Y Variable") +
  geom_point() +
  geom_smooth(method = "lm") #<---

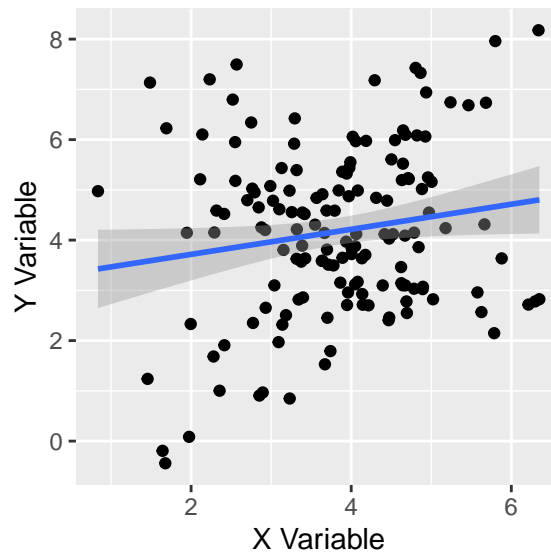
#geom_smooth() = add a "smoothed" trendline
```

```
#'method = "lm"' = add a regression line
```

```
#To view your plot, execute:
```

```
plotName
```

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



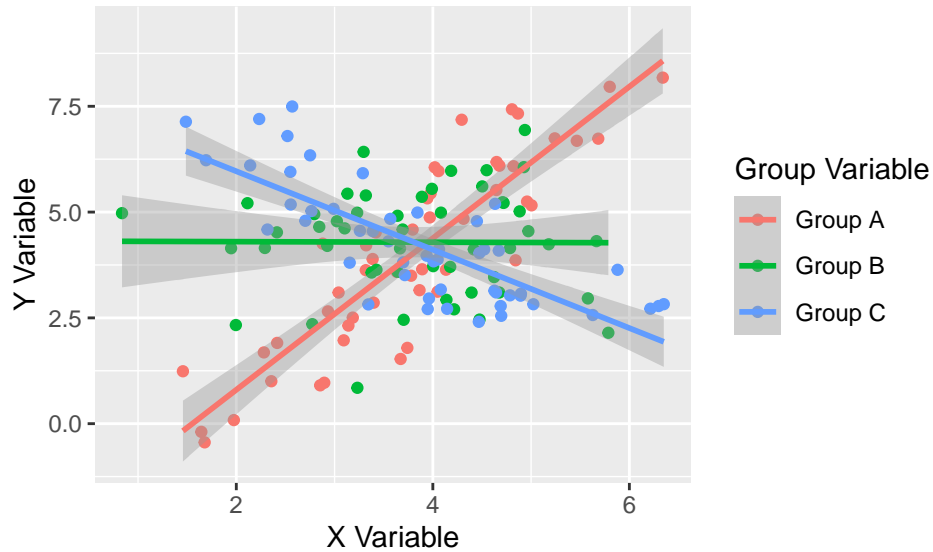
To add a grouping variable, use the `color =` option inside `aes()`:

```
plotName <- ggplot(data = dataframe,
  aes(x = xVariable,
      y = yVariable,
      color = groupVariable)) + #<---
  labs(x = "X Variable",
      y = "Y Variable",
      color = "Group Variable") + #<---
  geom_point() +
  geom_smooth(method = "lm") +
  scale_color_discrete(labels = c("Group A", "Group B", "Group C")) #<---

#geom_smooth() = add a "smoothed" trendline
#'method = "lm"' = add a regression line
#Notice that scale_color_discrete() is the color equivalent of scale_fill_discrete()

#To view your plot, execute:
plotName
```

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



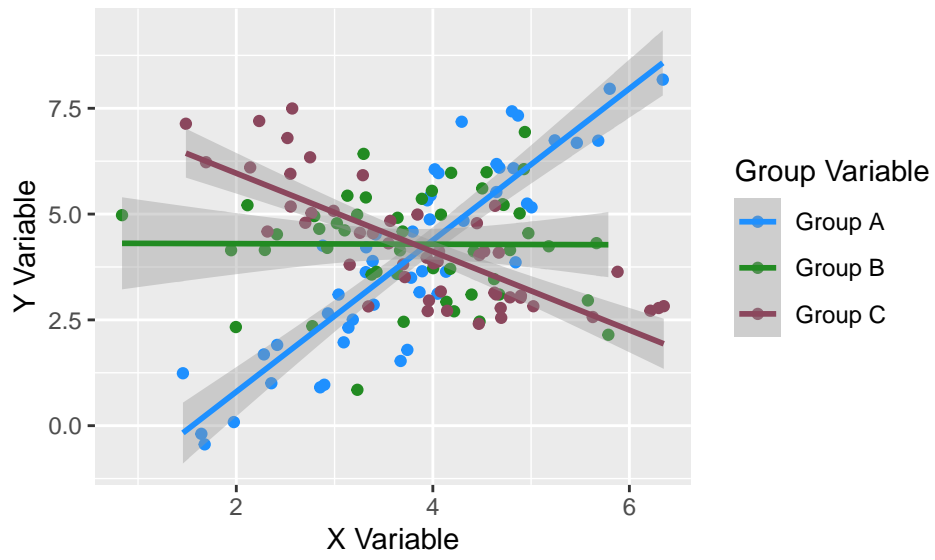
You can choose the color palette for the grouping using `scale_color_manual()`:

```
plotName <- ggplot(data = dataframe,
  aes(x = xVariable,
      y = yVariable,
      color = groupVariable)) +
  labs(x = "X Variable",
      y = "Y Variable",
      color = "Group Variable") +
  geom_point() +
  geom_smooth(method = "lm") +
  scale_color_manual( #<---
    values = c("dodgerblue", "forestgreen", "palevioletred4"), #<---
    labels = c("Group A", "Group B", "Group C") #<---
  )

#geom_smooth() = add a "smoothed" trendline
#'method = "lm"' = add a regression line

#To view your plot, execute:
plotName
```

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



Data Visualization Cheat Sheet

Plotting option	Code
Base plot	<code>plotName <- ggplot(data = dataFrame, aes(x = xVariable, y = yVariable, color = groupVariable))</code>
Bar plot	<code>+ stat_summary(fun="mean", geom="bar", position="dodge")</code>
95% CI Error Bars	<code>+ stat_summary(fun.data="mean_cl_boot", geom="errorbar", position="dodge")</code>
Box plot	<code>+ geom_boxplot(position = position_dodge(1), width=.5)</code>
Violin plot	<code>+ geom_violin(position = position_dodge(1))</code>
Scatter plot	<code>+ geom_point()</code>
Add trend line	<code>+ geom_smooth(method="lm")</code>
Add labels	<code>+ labs(x="X Label", y = "Y Label", color = "Color Label", fill = "Fill Label")</code>
Change axis labels	<code>+ scale_x_discrete(labels = c("Low", "High"))</code>
Change legend labels	<code>+ scale_color_discrete(labels = c("Low", "High"))</code>
Change legend colors	<code>+ scale_color_manual(values = c("red", "black"))</code>
Change font size	<code>+ theme_gray(base_size=20)</code>
Hide a legend	<code>+ guides(color="none",fill="none")</code>