

Plotting with ggplot()

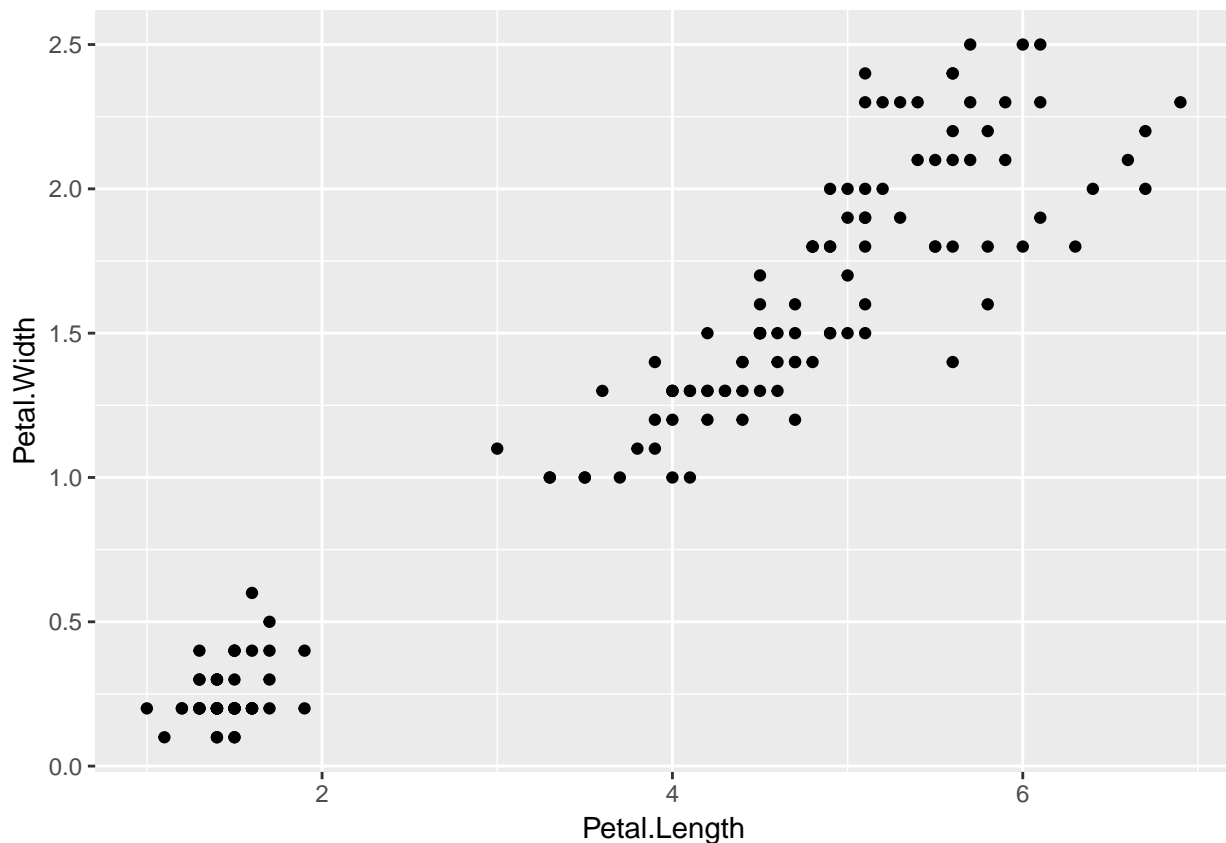
Plotting with ggplot()

The ultimate goal of any analysis is to *communicate* your results and understandings. Plotting your data is the most effect way to communicate a summary of it to an audience. Today will focus on the basics on how to do that in R, using the tidyverse package ggplot, with all examples with the built-in iris dataset.

Basic ggplot syntax

You have to specify a data table, at least one column from the table, and a geometry.

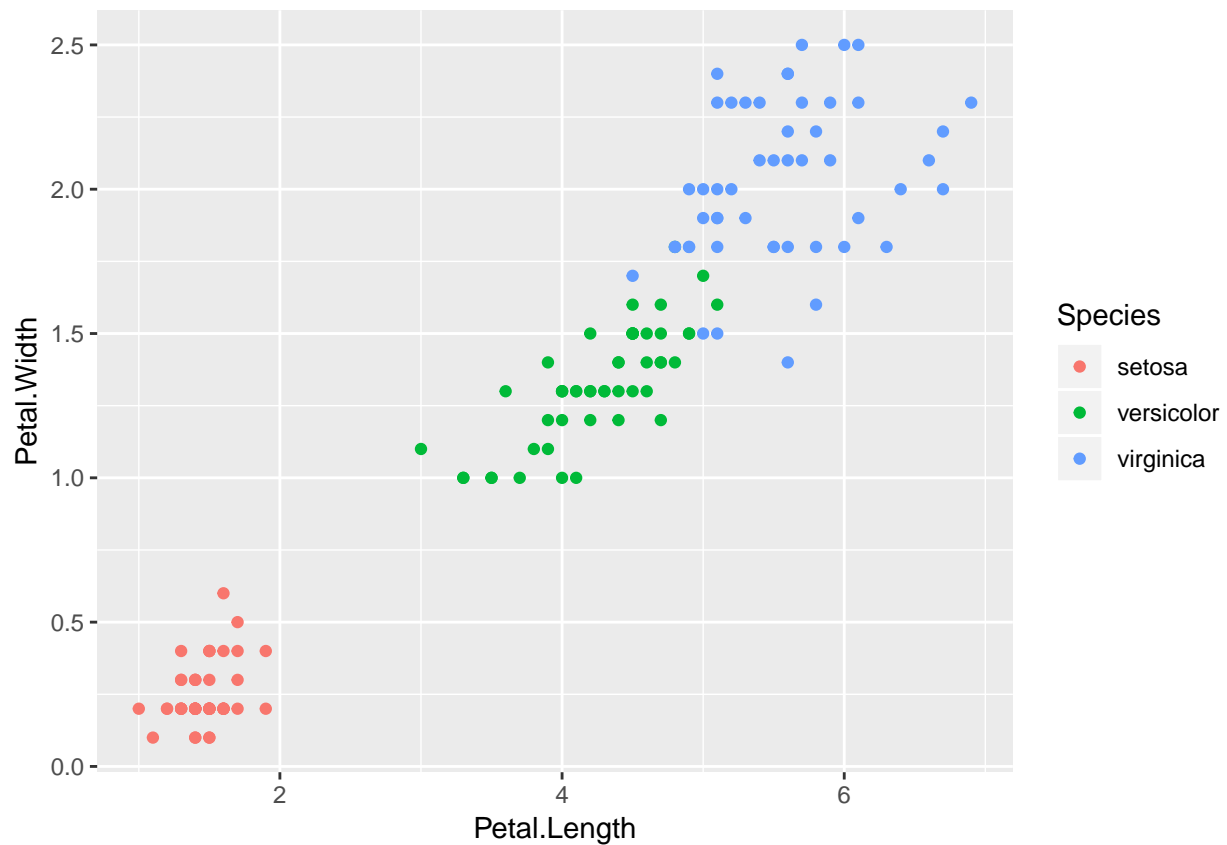
```
ggplot(iris, aes(x = Petal.Length, y = Petal.Width)) + geom_point()
```



The scale of the axes are automatically determined from the data and they're labelled witht the column name.

Adding an aesthetic like color modifies all the points and automatically adds a legend.

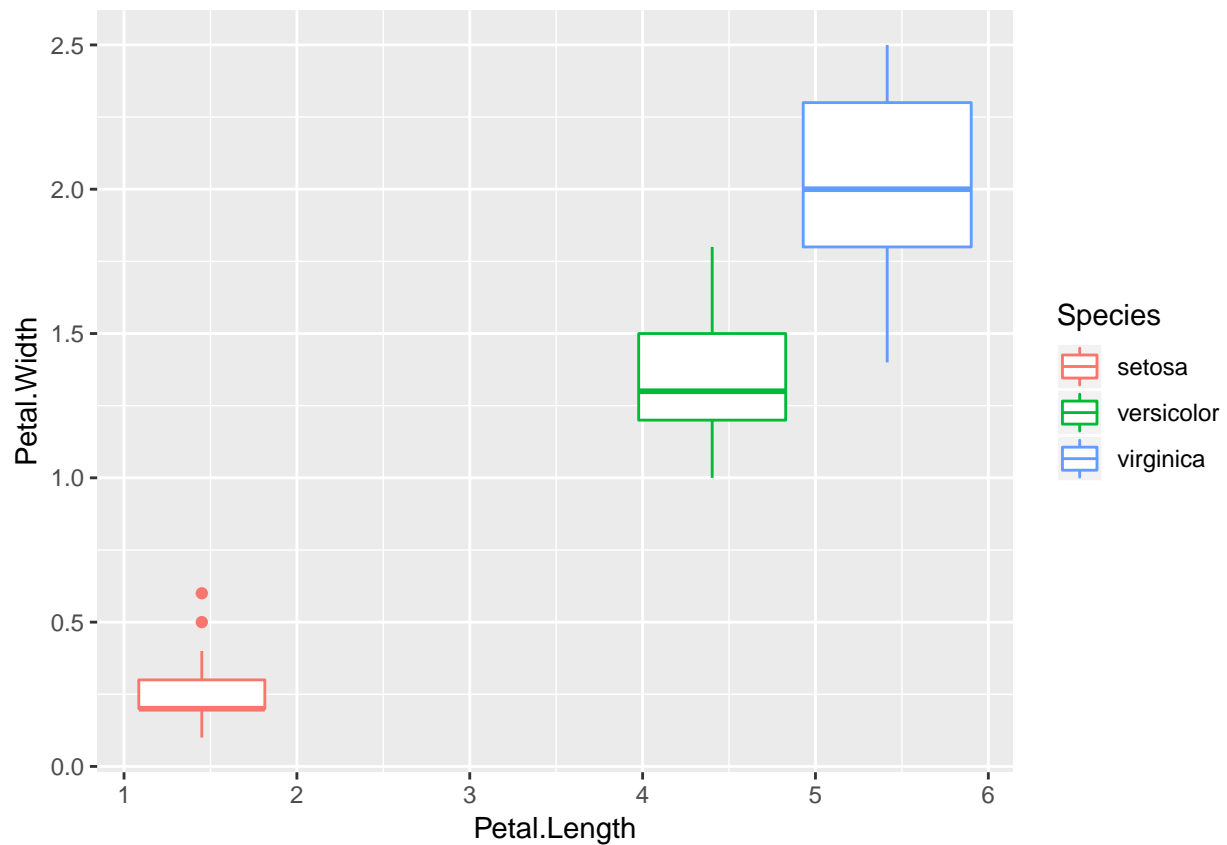
```
ggplot(iris, aes(x = Petal.Length, y = Petal.Width, color = Species)) + geom_point()
```



The legend title, like the axes labels, is the name of the column given to `color =` and the legend labels are whatever is in that column. For a continuous scale `ggplot()` will have a bar showing the range of colors and what values they correspond to.

You just have to change the geom to change the plot to another geometry/type

```
ggplot(iris, aes(x = Petal.Length, y = Petal.Width, color = Species)) + geom_boxplot()
```



Geometries (aka Types of Plots)

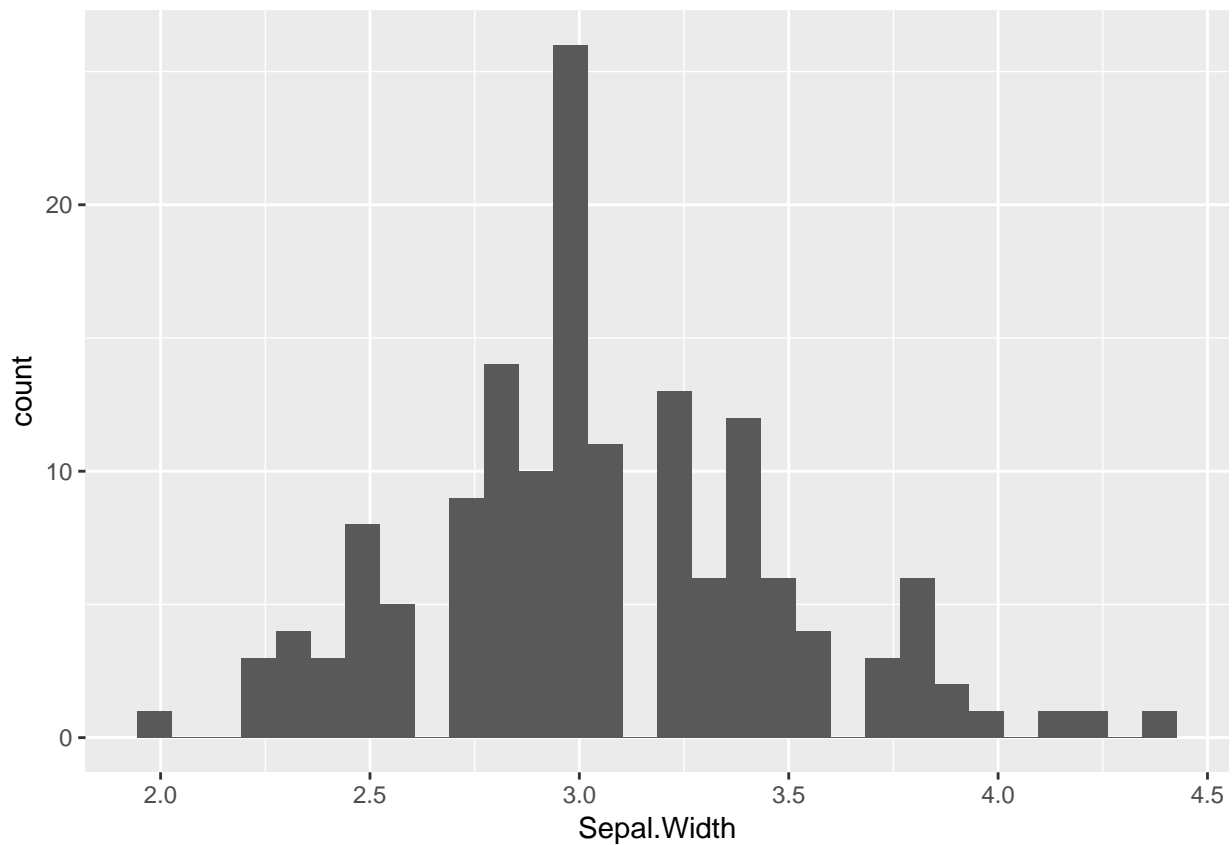
One Variable

Continuous

Histogram

```
ggplot(iris, aes(x = Sepal.Width)) + geom_histogram()
```

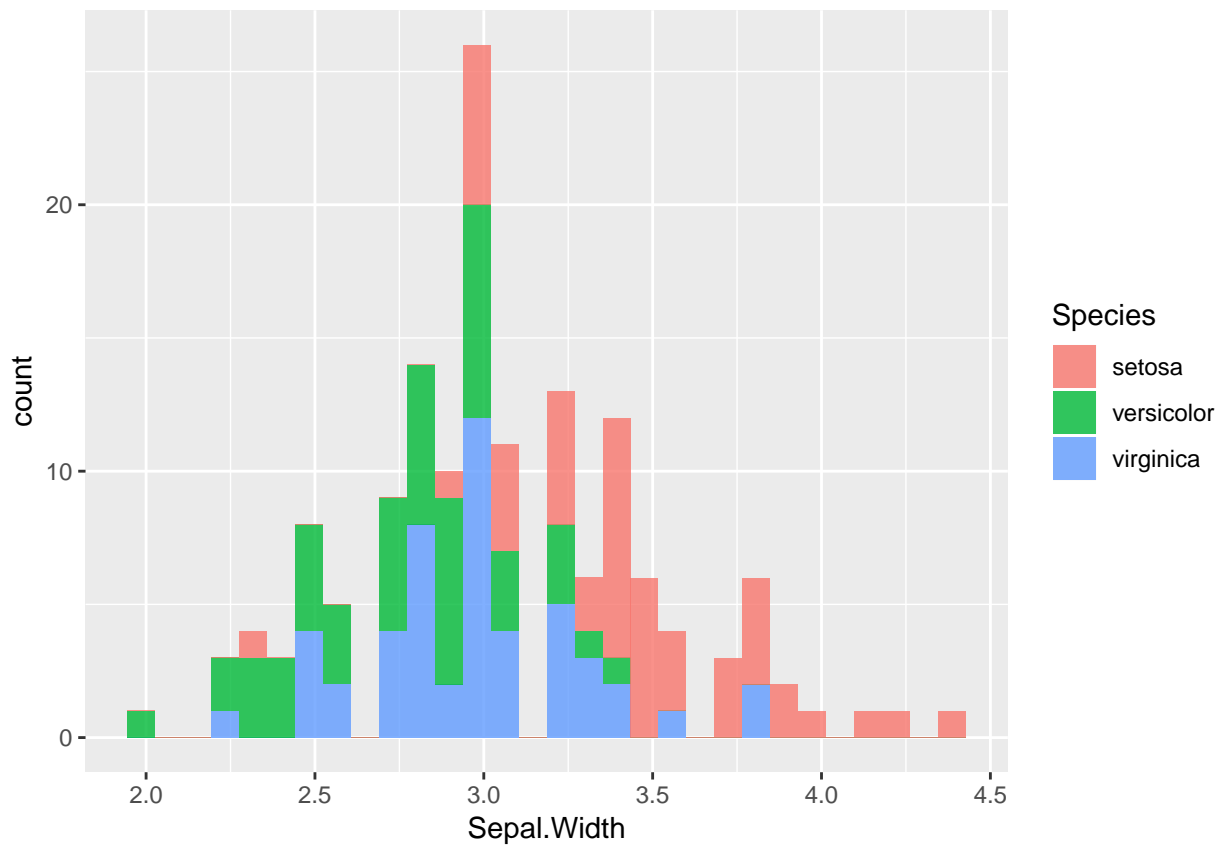
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Be careful with histograms! By default, when you add fill, color, etc, `geom_histogram()` stacks the different bars on top of each other. For *overlapping* histograms, which are more effective visually and what people are used to seeing, use `position = 'identity'` inside of `geom_histogram()`

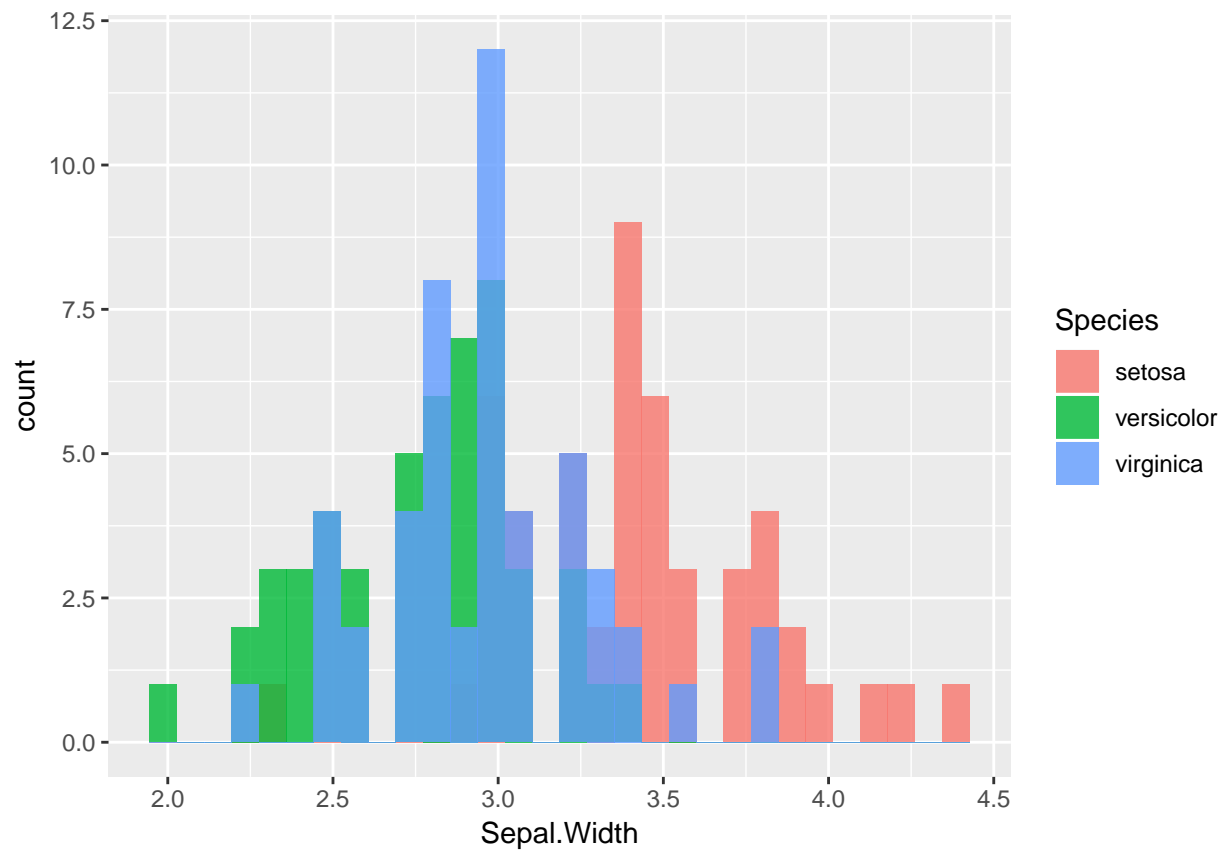
```
# default histogram
ggplot(iris, aes(x = Sepal.Width, fill = Species)) + geom_histogram(alpha = 0.8)
```

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



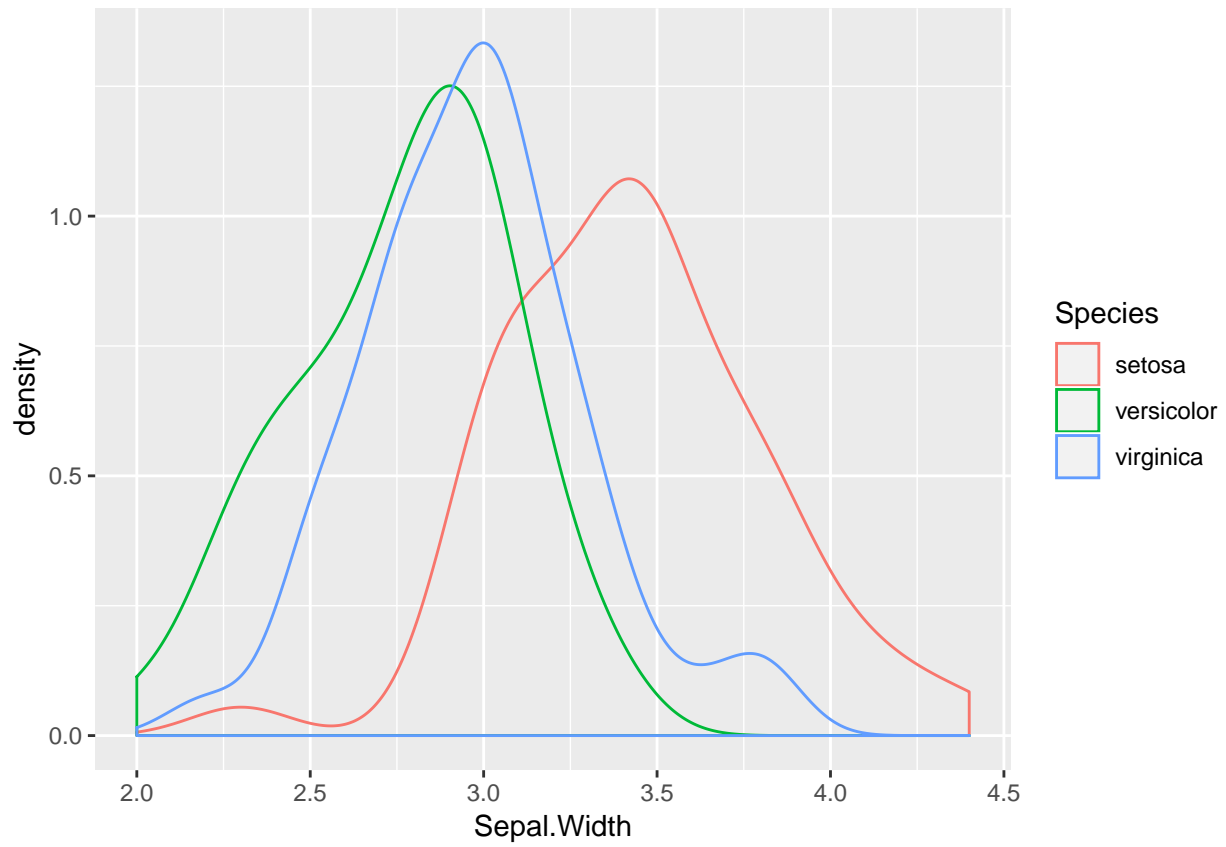
```
# use position = 'identity' for overlapping histograms
ggplot(iris, aes(x = Sepal.Width, fill = Species)) + geom_histogram(position = 'identity', alpha = 0.8)

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



Density plot

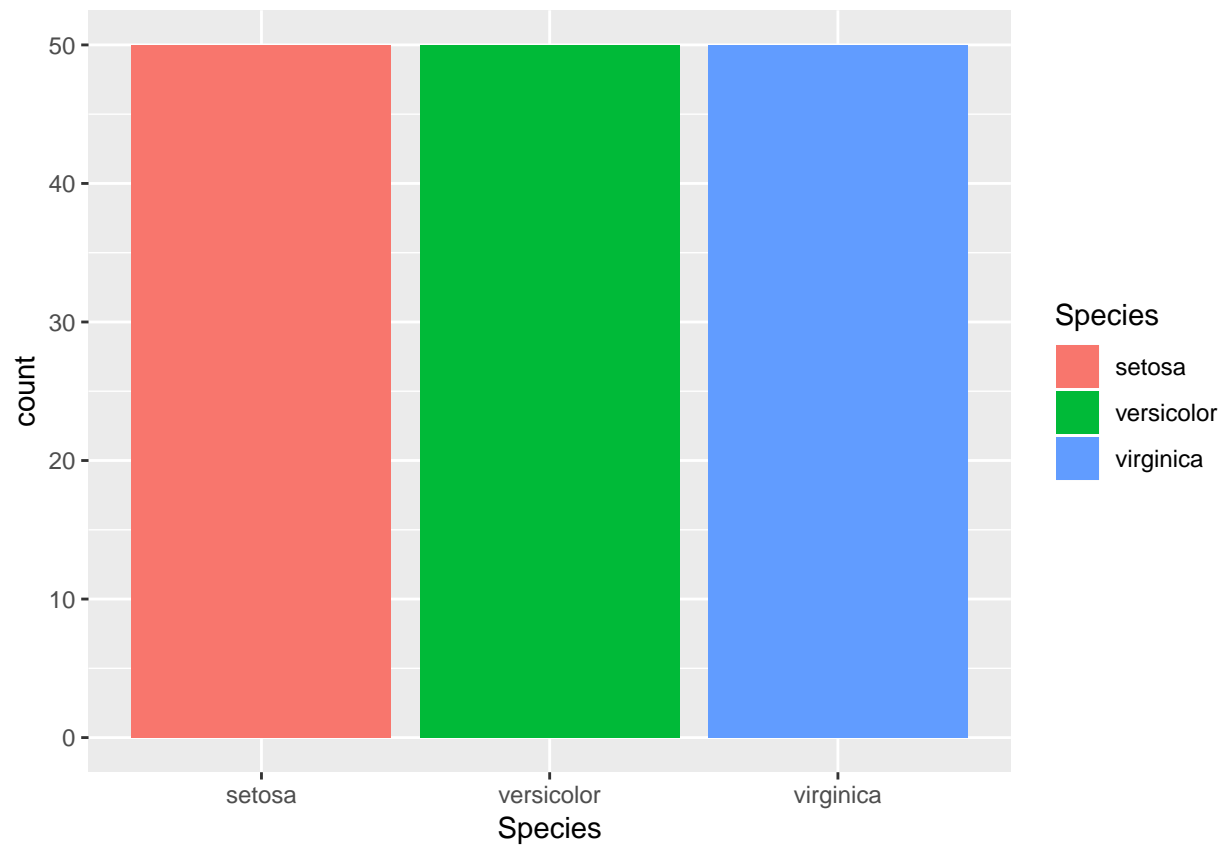
```
ggplot(iris, aes(x = Sepal.Width, color = Species)) + geom_density()
```



Discrete

Bar plot

```
ggplot(iris, aes(x = Species, fill = Species)) + geom_bar()
```

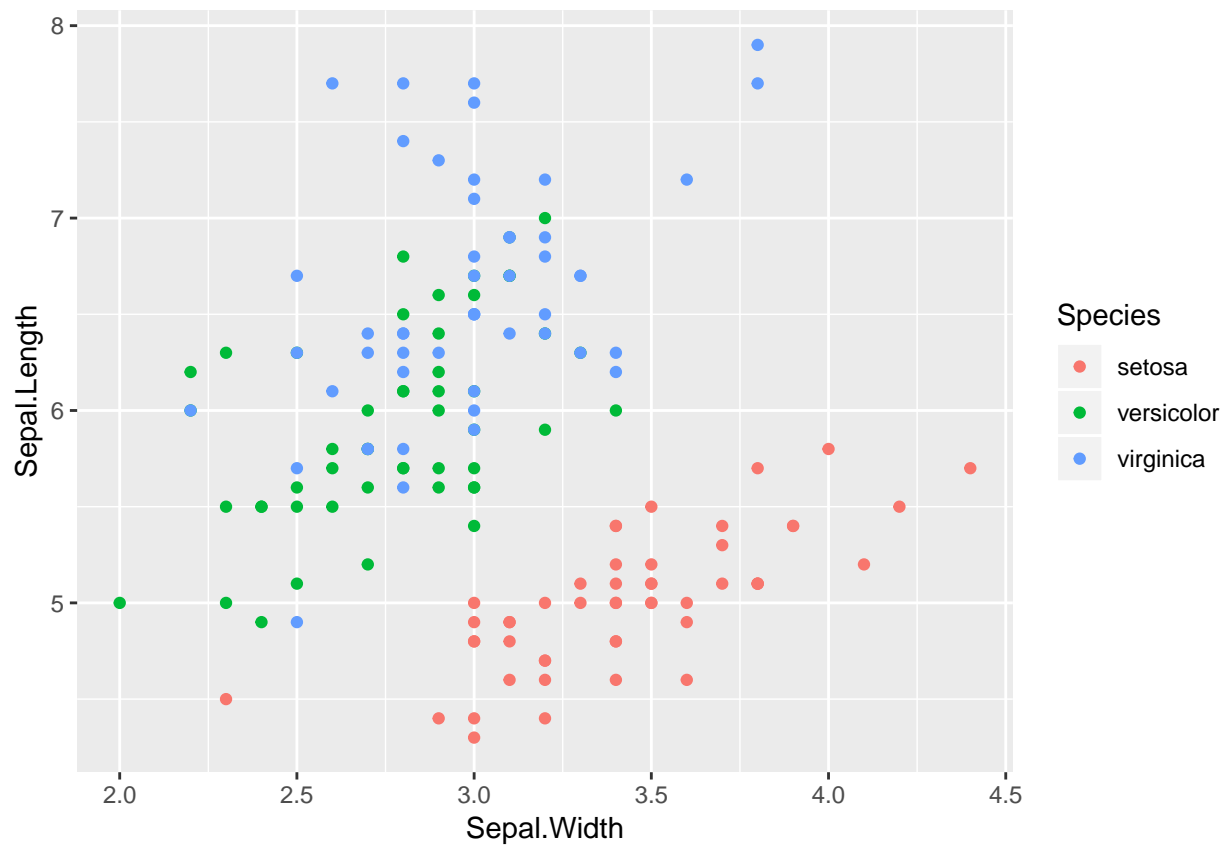


Two Variables

Continuous X, Continuous Y

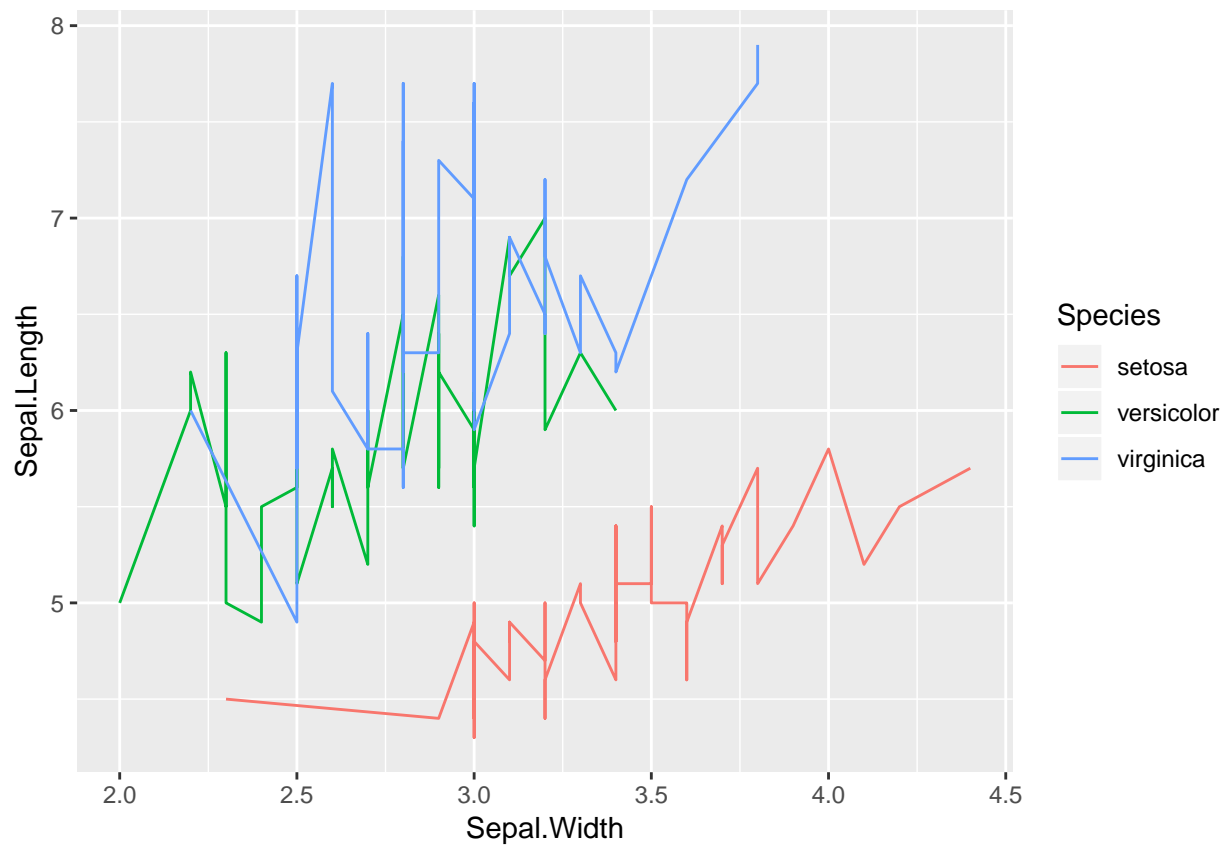
Scatter plot

```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point()
```

Line plot

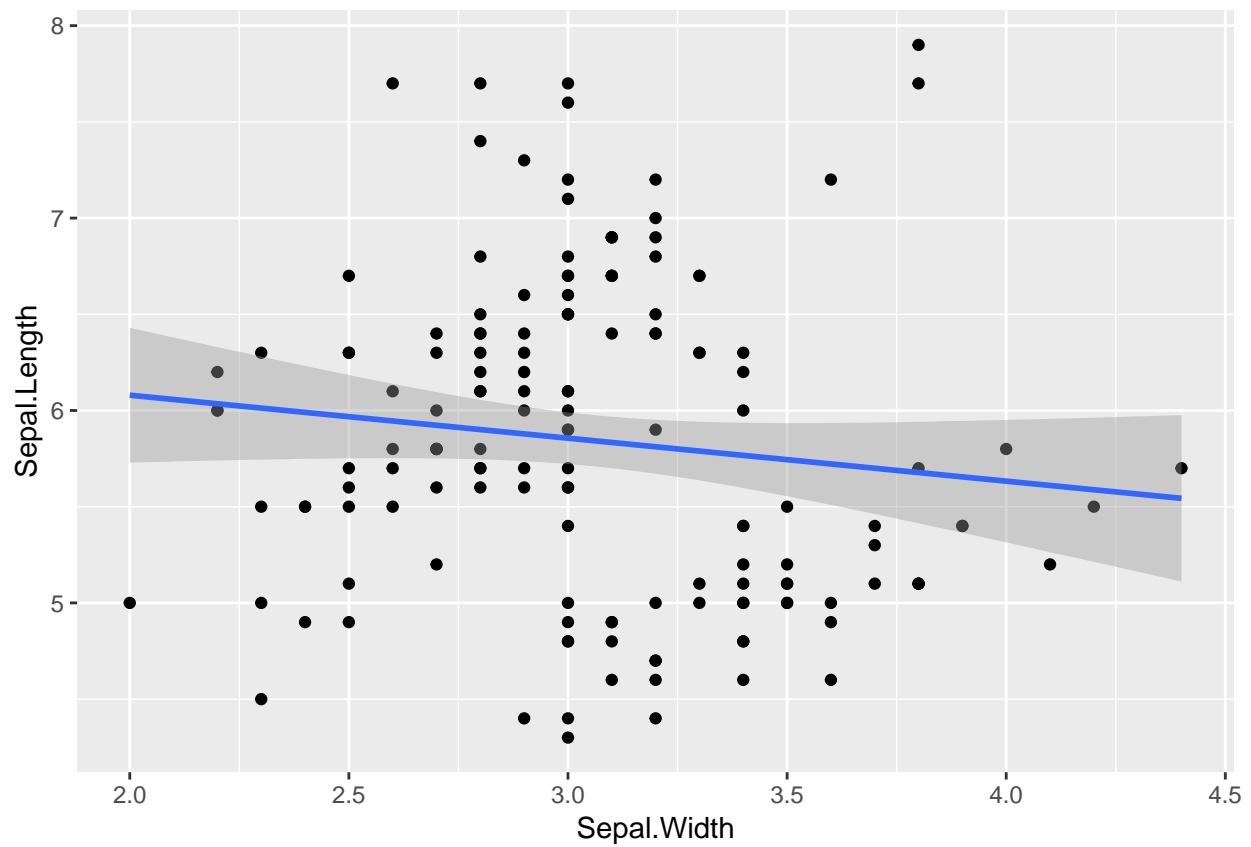
```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_line()
```



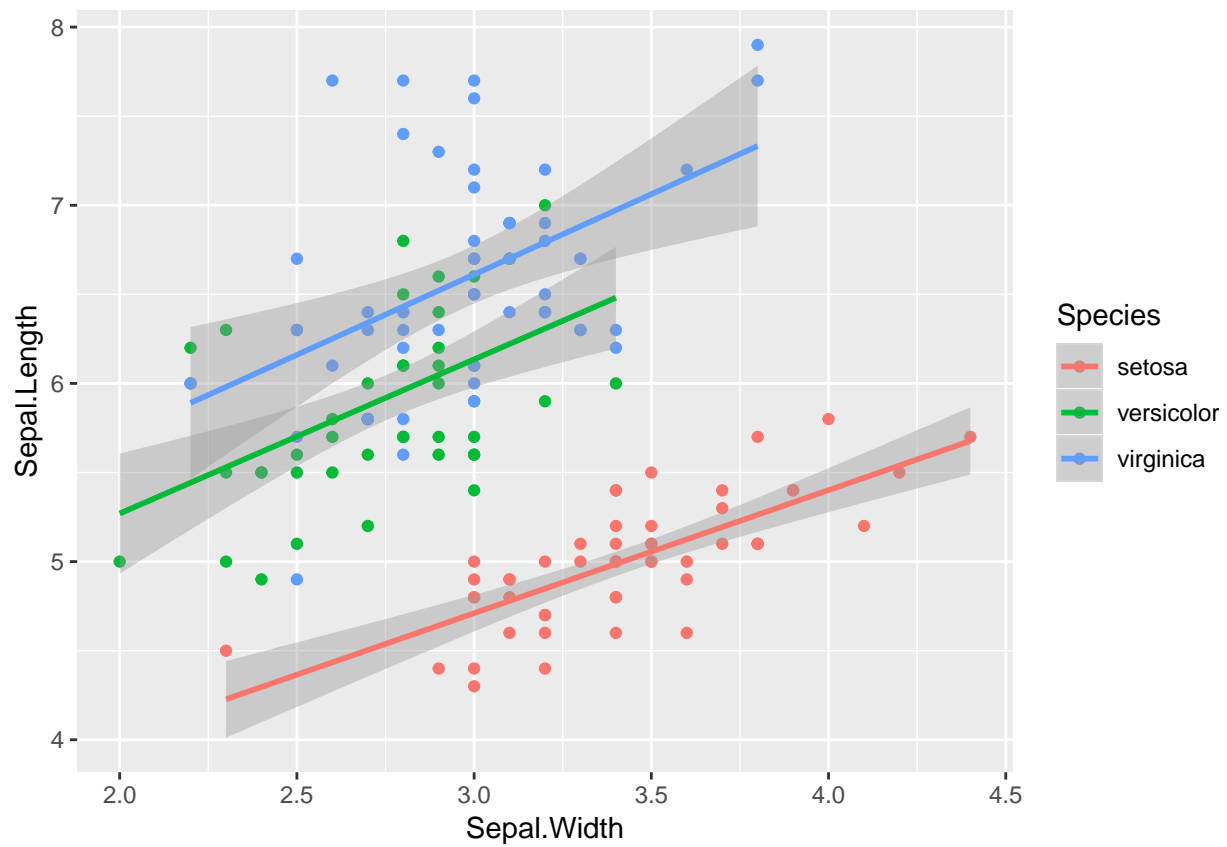
Adding a line of best fit

If you want to fit a straight line (or other type of fit) to your scatter plot, add on `geom_smooth()`. Use `method = lm` inside it to add a line (see the `geom_smooth()` documentation for other fit methods). Also, it automatically supplies confidence intervals with whatever is fit.

```
# with line of best fit
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length)) + geom_point() + geom_smooth(method = lm)
```



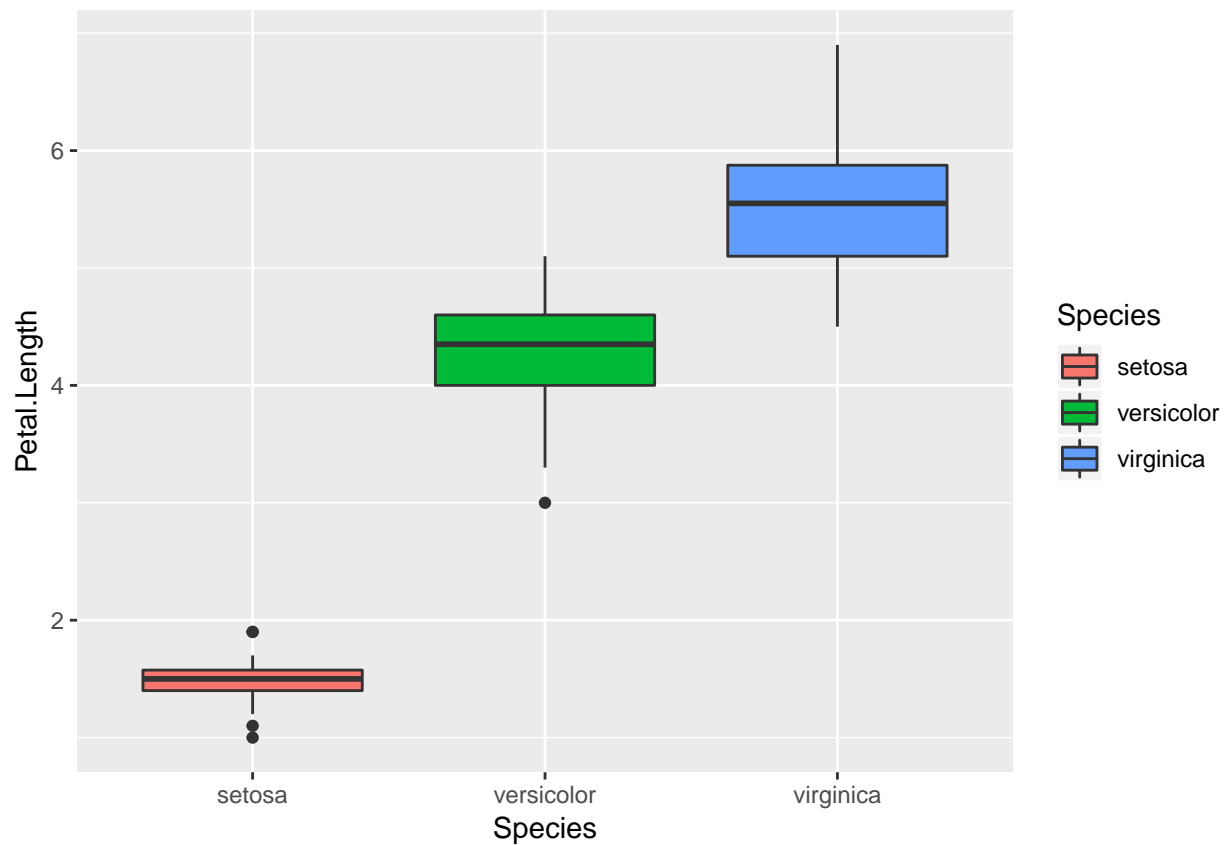
```
# if you color by something, each subgroup will automatically get it's own line  
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point() + geom_smooth(method = "lm")
```



Discrete X, Continuous Y

Boxplot

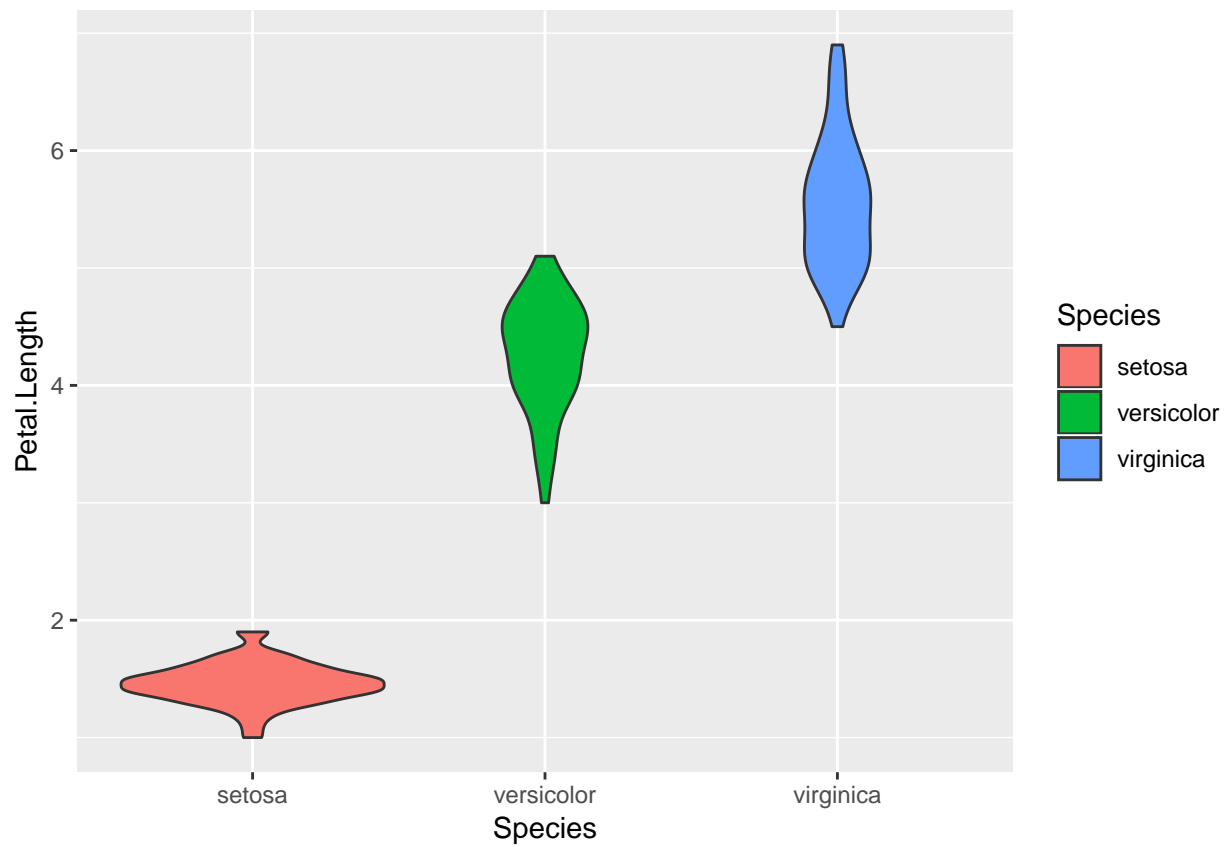
```
ggplot(iris, aes(x = Species, y = Petal.Length, fill = Species)) + geom_boxplot()
```



Violin plot

A violin plot is a mirrored density plot displayed like a boxplot. It gives you a better sense of the distribution of the underlying data than the five number summary of the boxplot.

```
ggplot(iris, aes(x = Species, y = Petal.Length, fill = Species)) + geom_violin()
```

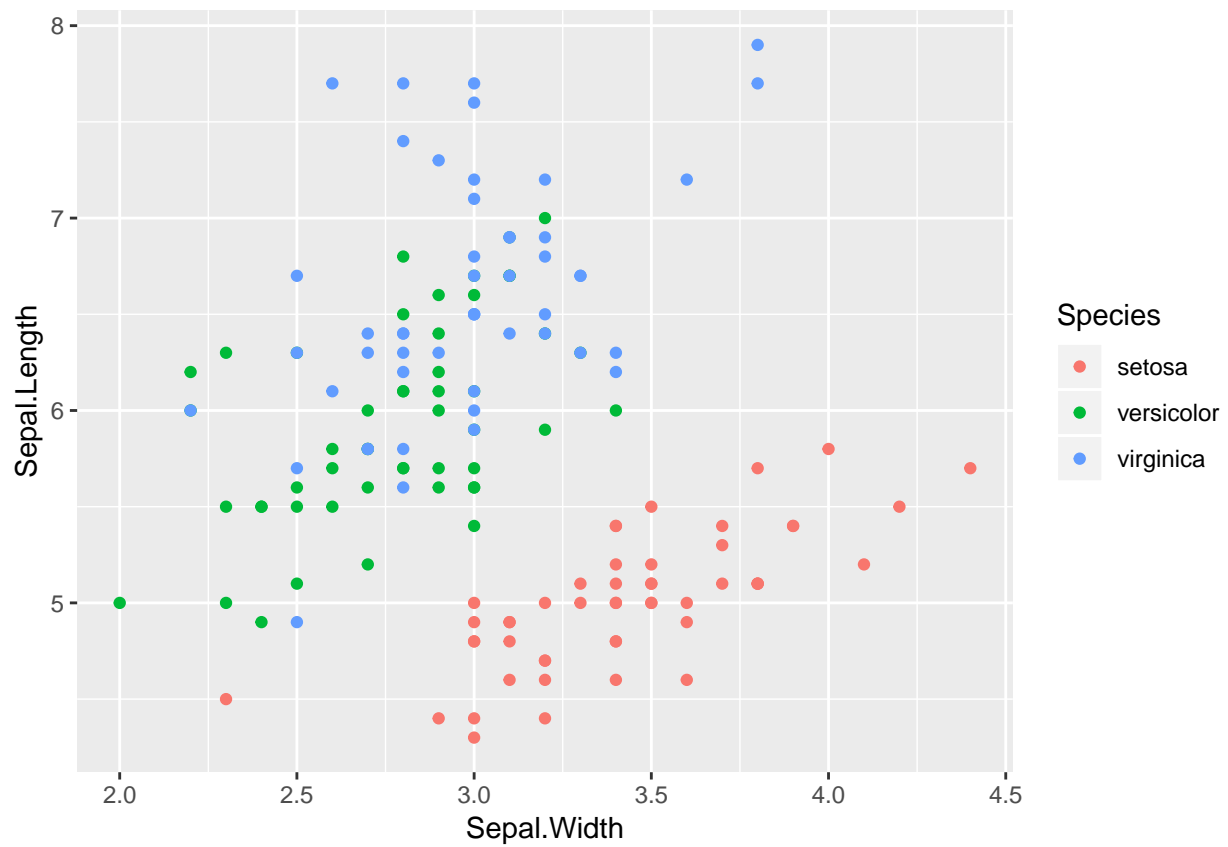


Aesthetics

Basic three

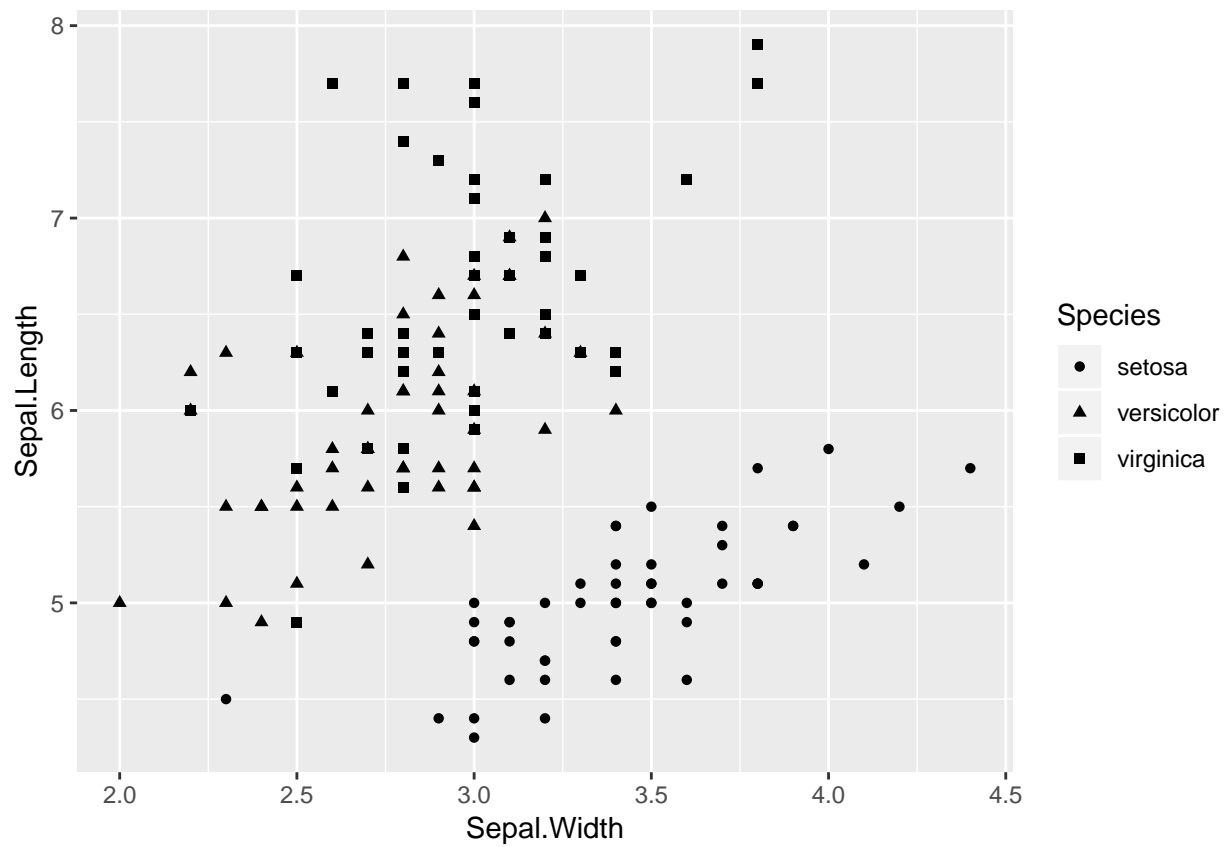
Color

```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point()
```



Shape

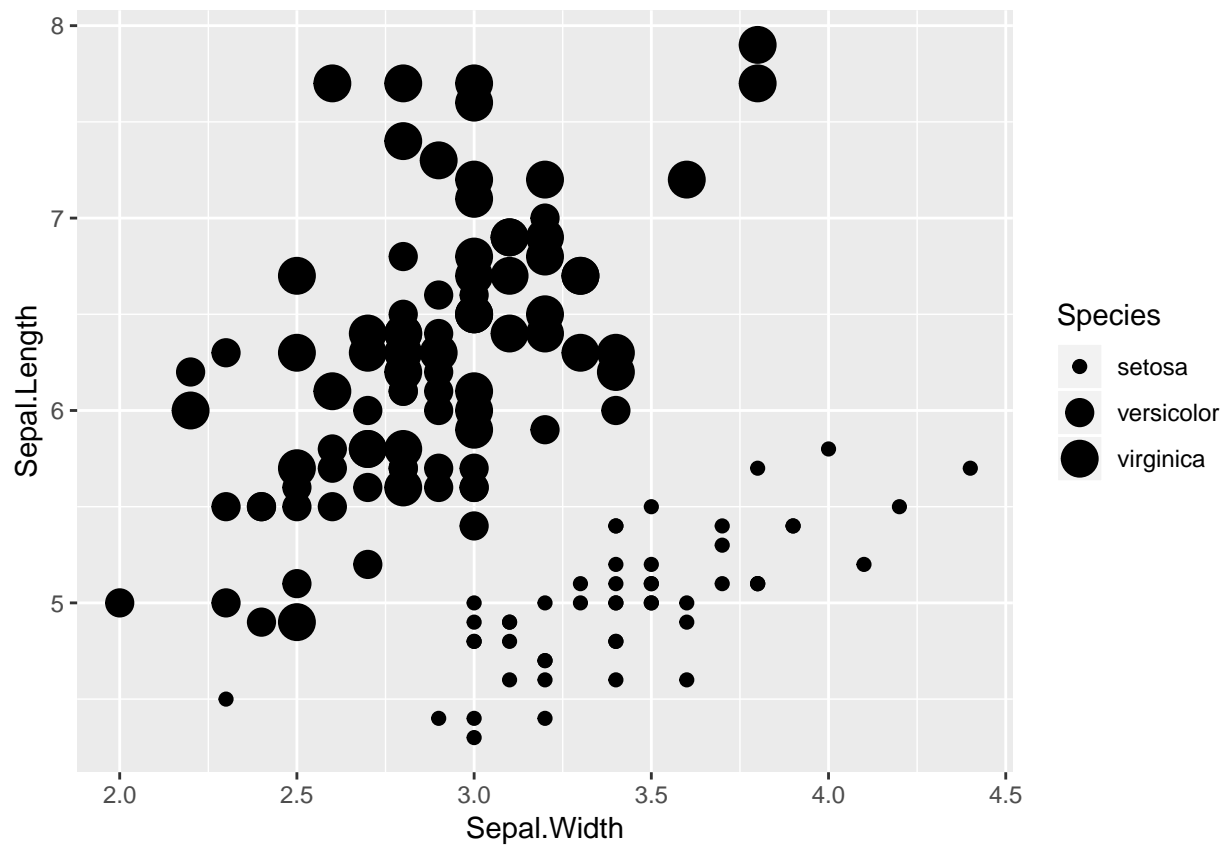
```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, shape = Species)) + geom_point()
```



Size

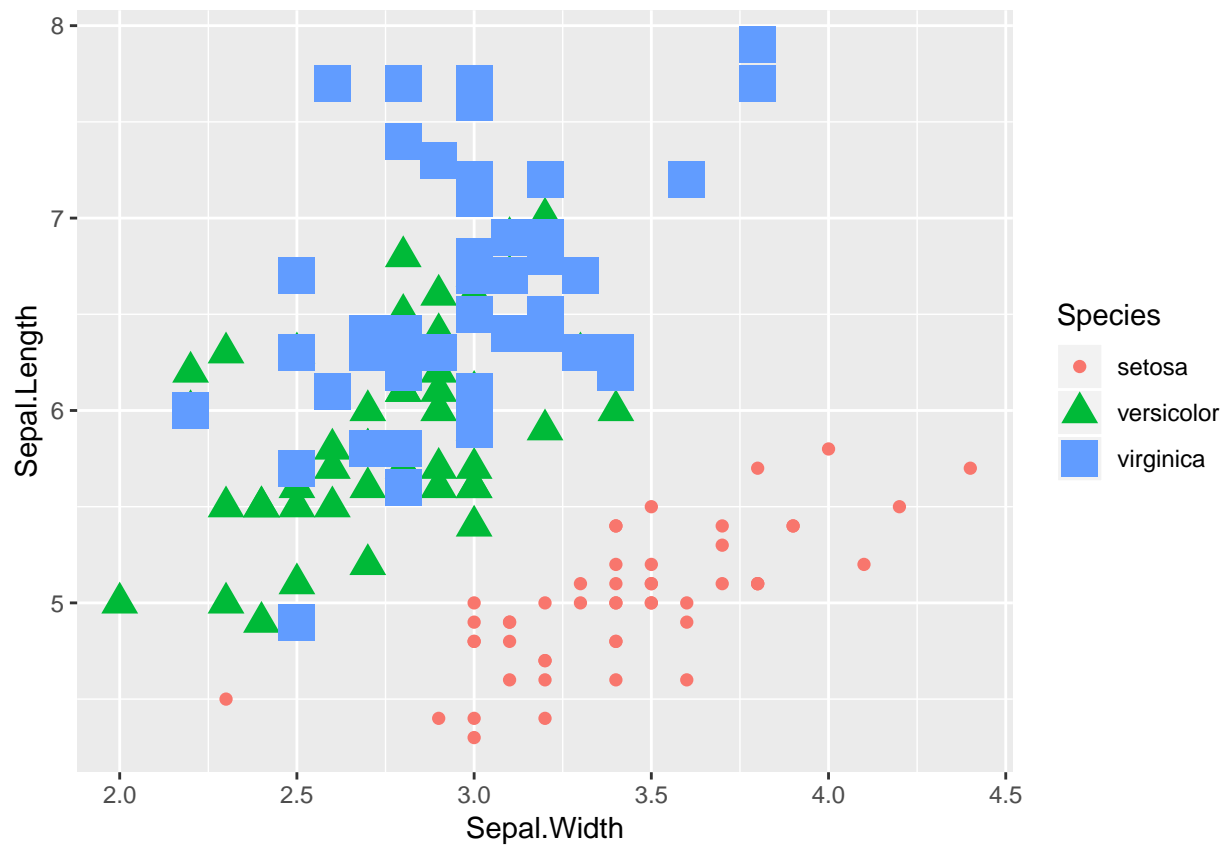
```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, size = Species)) + geom_point()
```

```
## Warning: Using size for a discrete variable is not advised.
```

Combine as many as you want

```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species, shape = Species, size = Species))  
## Warning: Using size for a discrete variable is not advised.
```

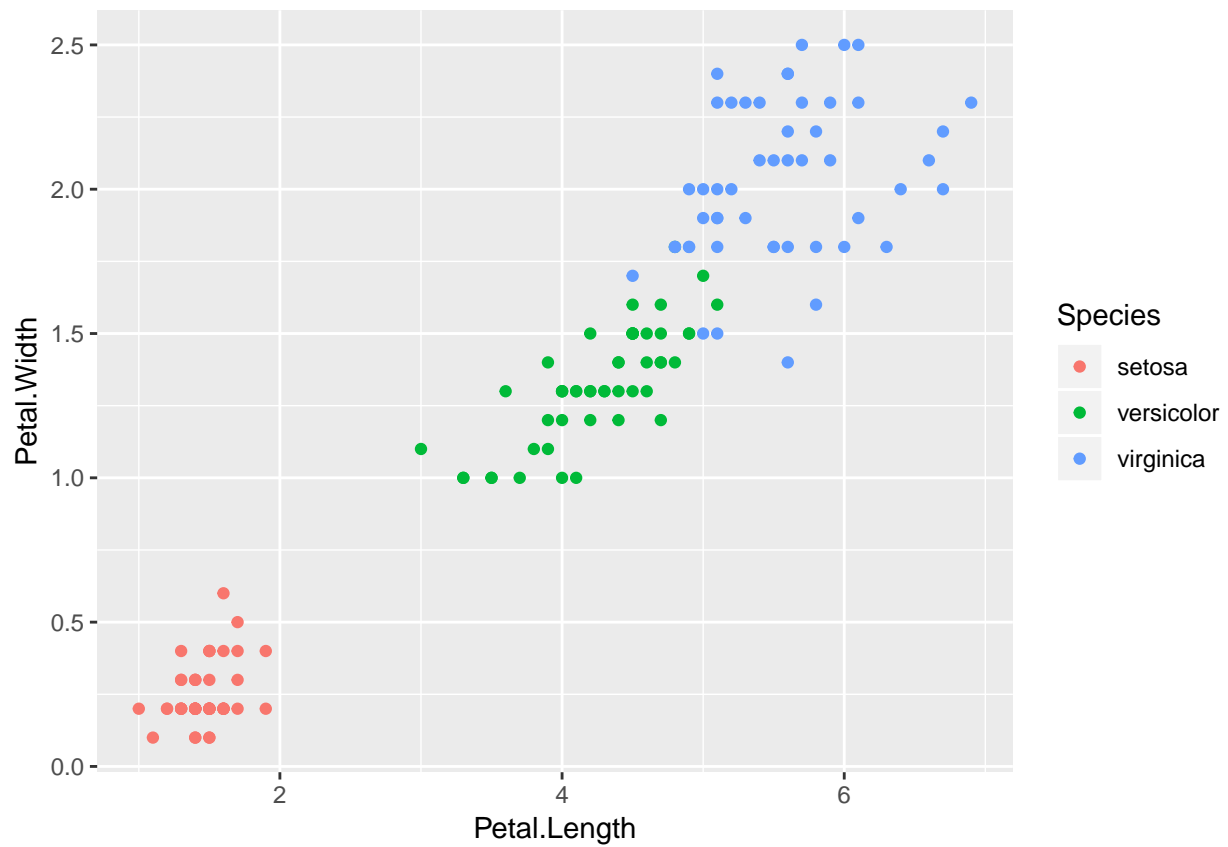


Miscellaneous Useful Features

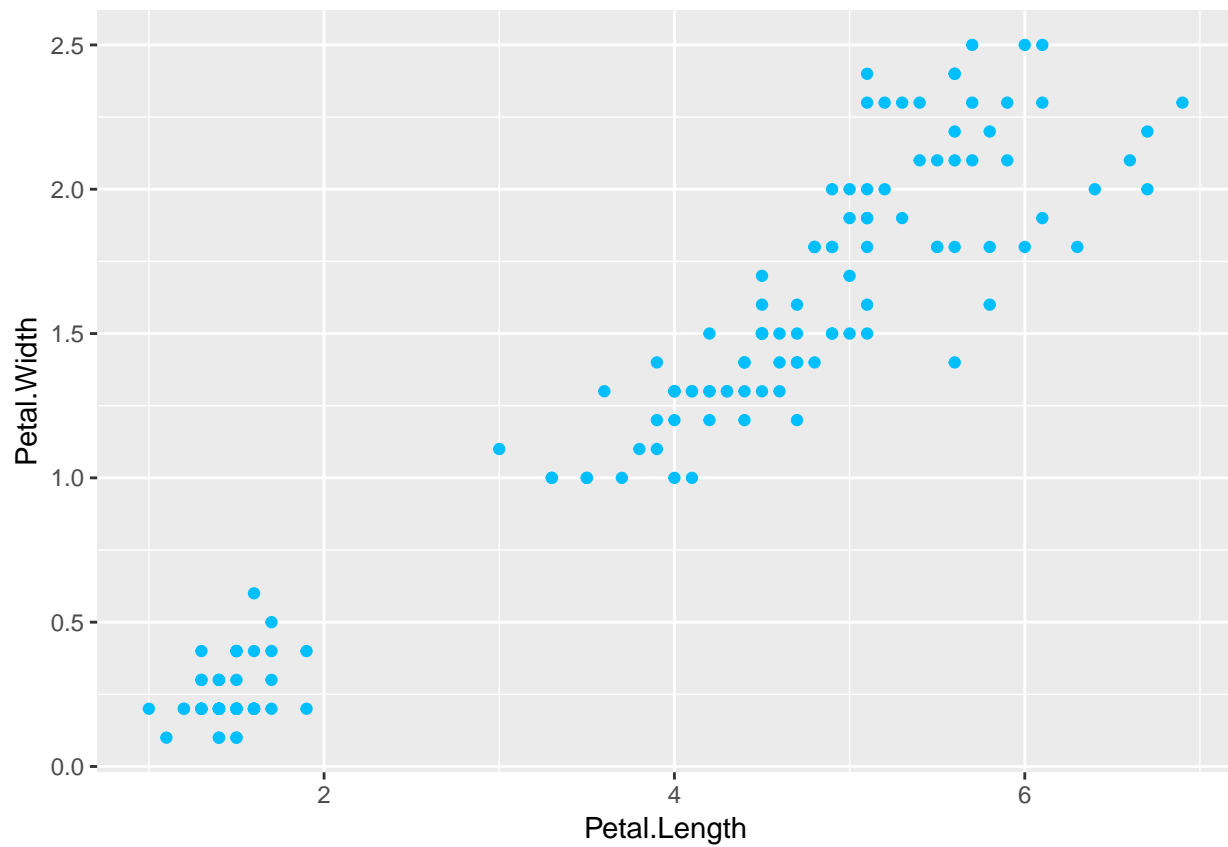
More about `aes()`

Whether you add color (or shape, size, etc.) inside or outside of `'aes()'` has a different outcome. Inside `aes()` ggplot modifies points differently, but outside `aes()` it applies the same thing to all points.

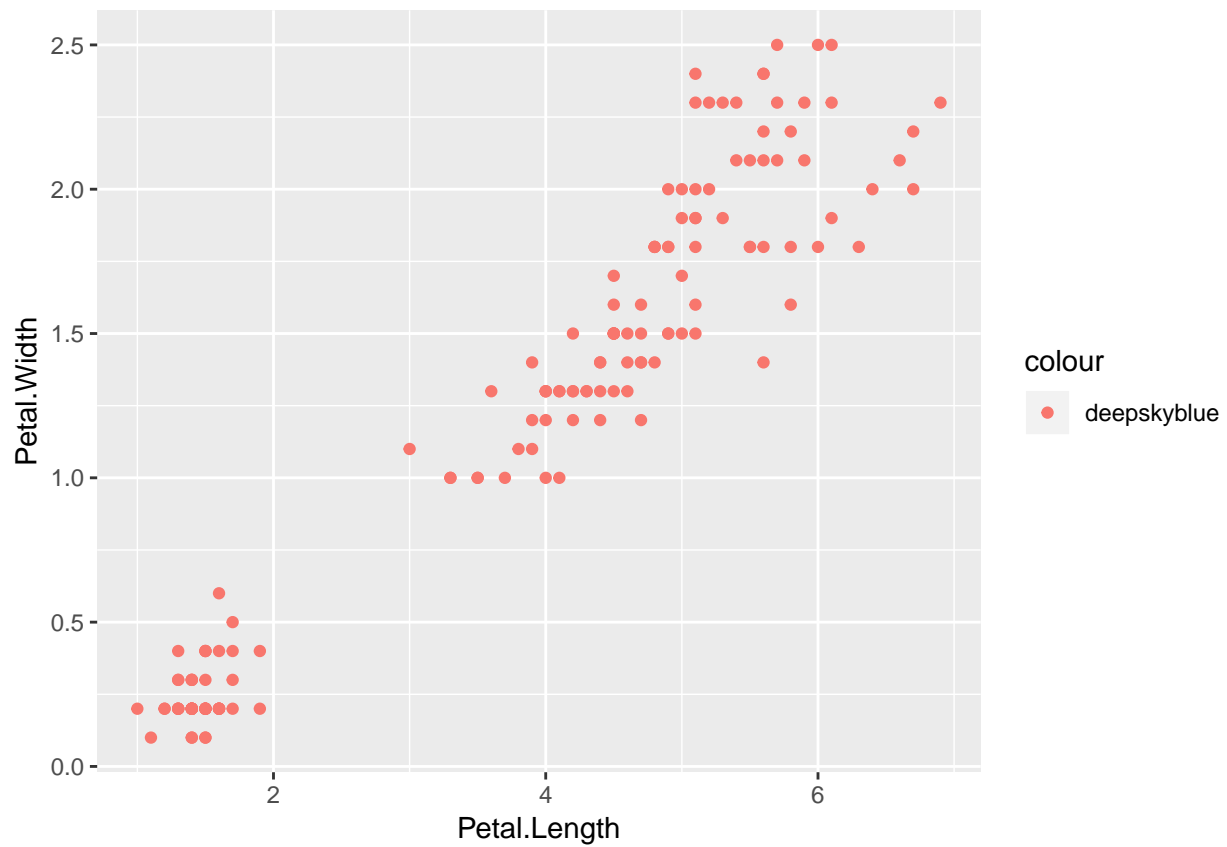
```
# where we've had it
ggplot(iris, aes(x = Petal.Length, y = Petal.Width, color = Species)) + geom_point()
```



```
# I just want my points to be another color  
ggplot(iris, aes(x = Petal.Length, y = Petal.Width)) + geom_point(color = 'deepskyblue')
```

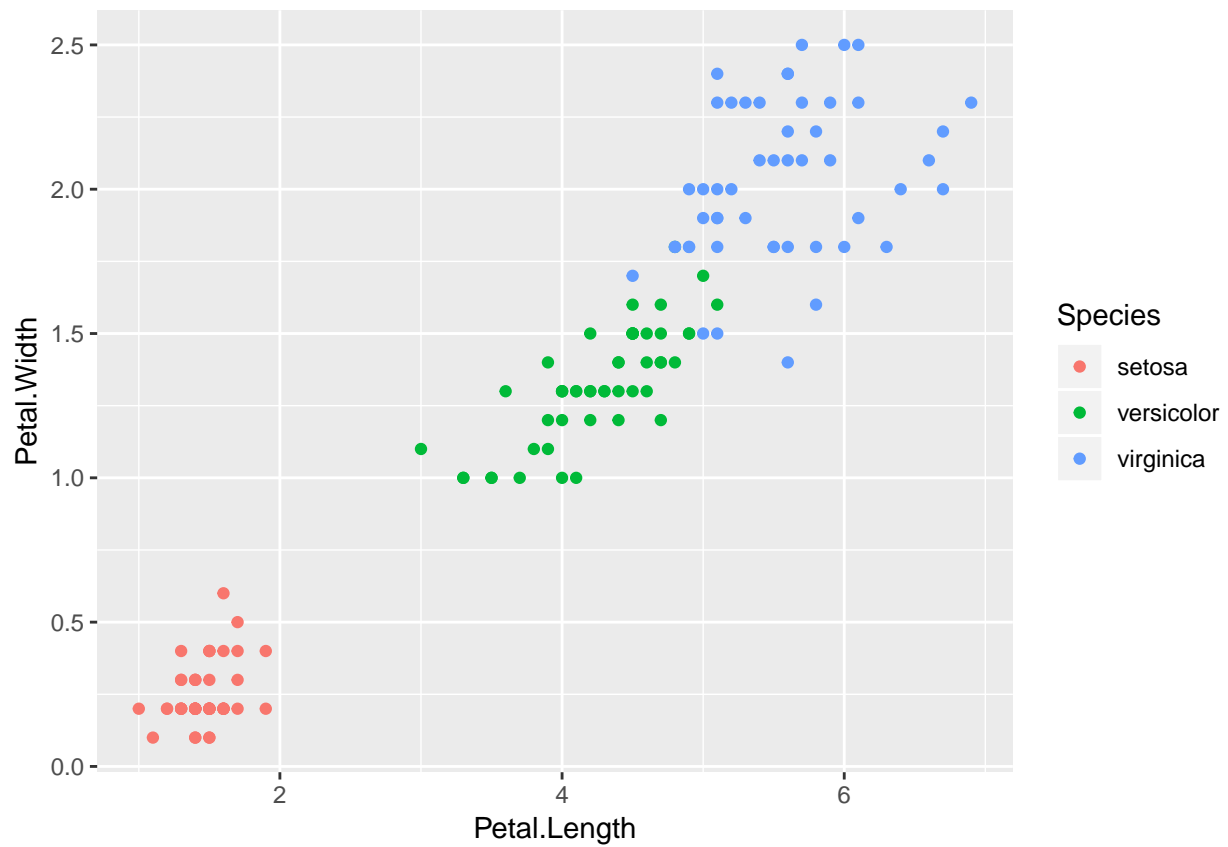


BUT if I put color = 'deepskyblue' instead an aes() it won't work; ggplot thinks it's a data feature
`ggplot(iris, aes(x = Petal.Length, y = Petal.Width, color = 'deepskyblue')) + geom_point()`

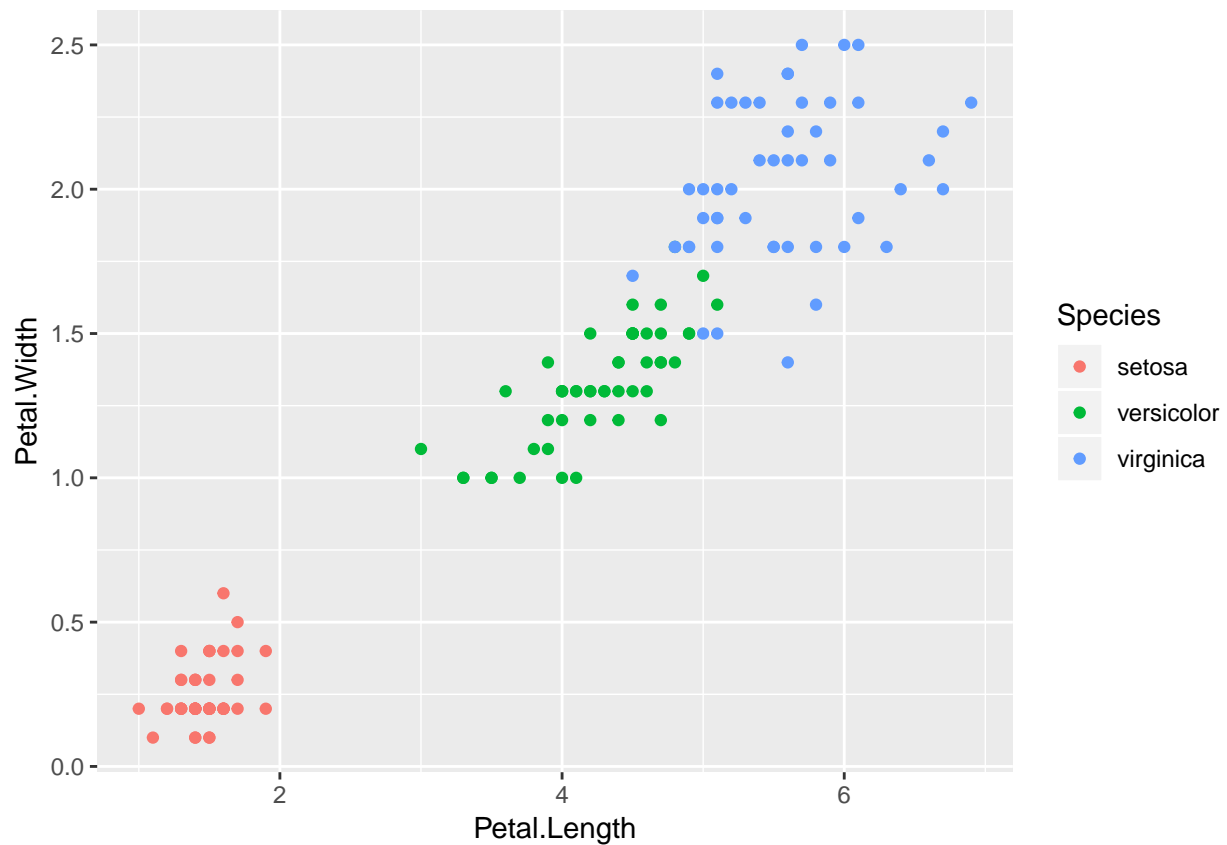


You can also put `aes()` inside either `ggplot()` or whatever geom you pick

```
# inside ggplot()  
ggplot(iris, aes(x = Petal.Length, y = Petal.Width, color = Species)) + geom_point()
```



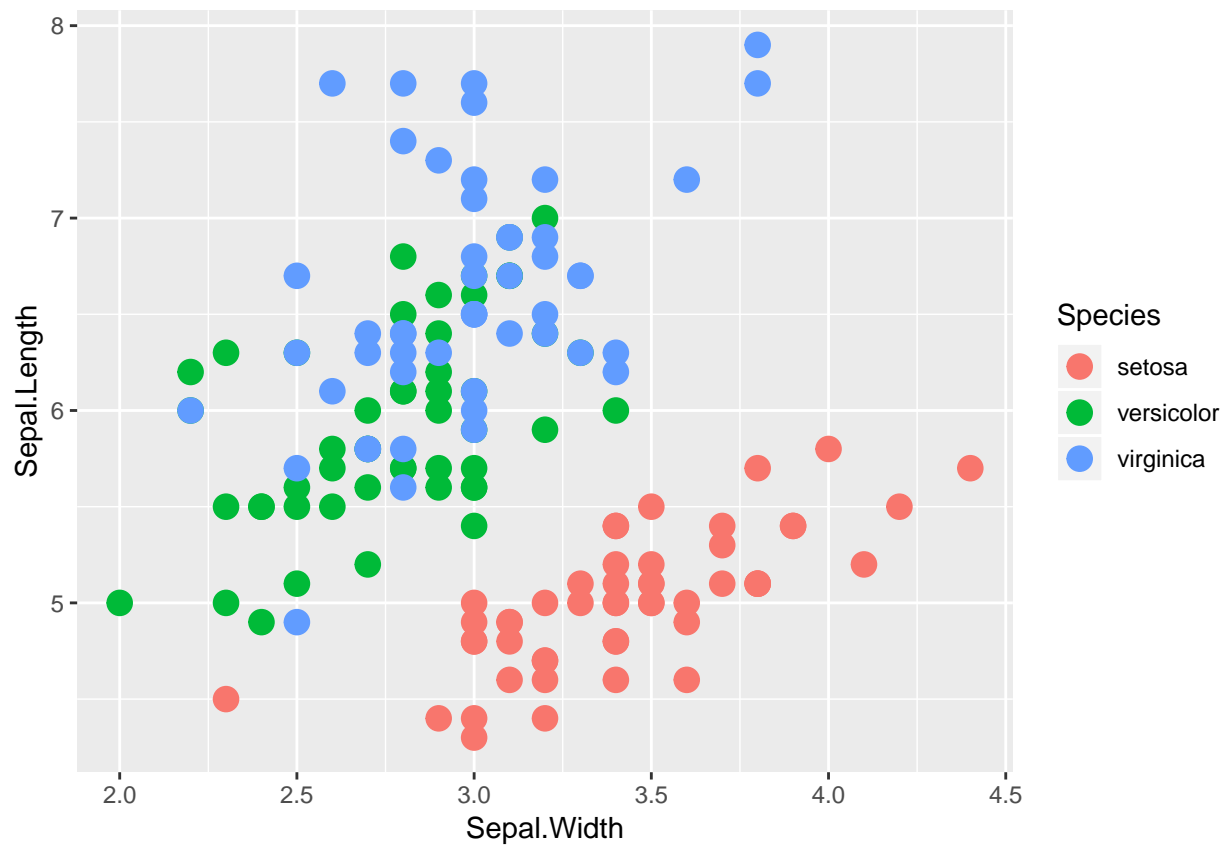
```
# inside geom_point()  
ggplot(iris) + geom_point(aes(x = Petal.Length, y = Petal.Width, color = Species))
```



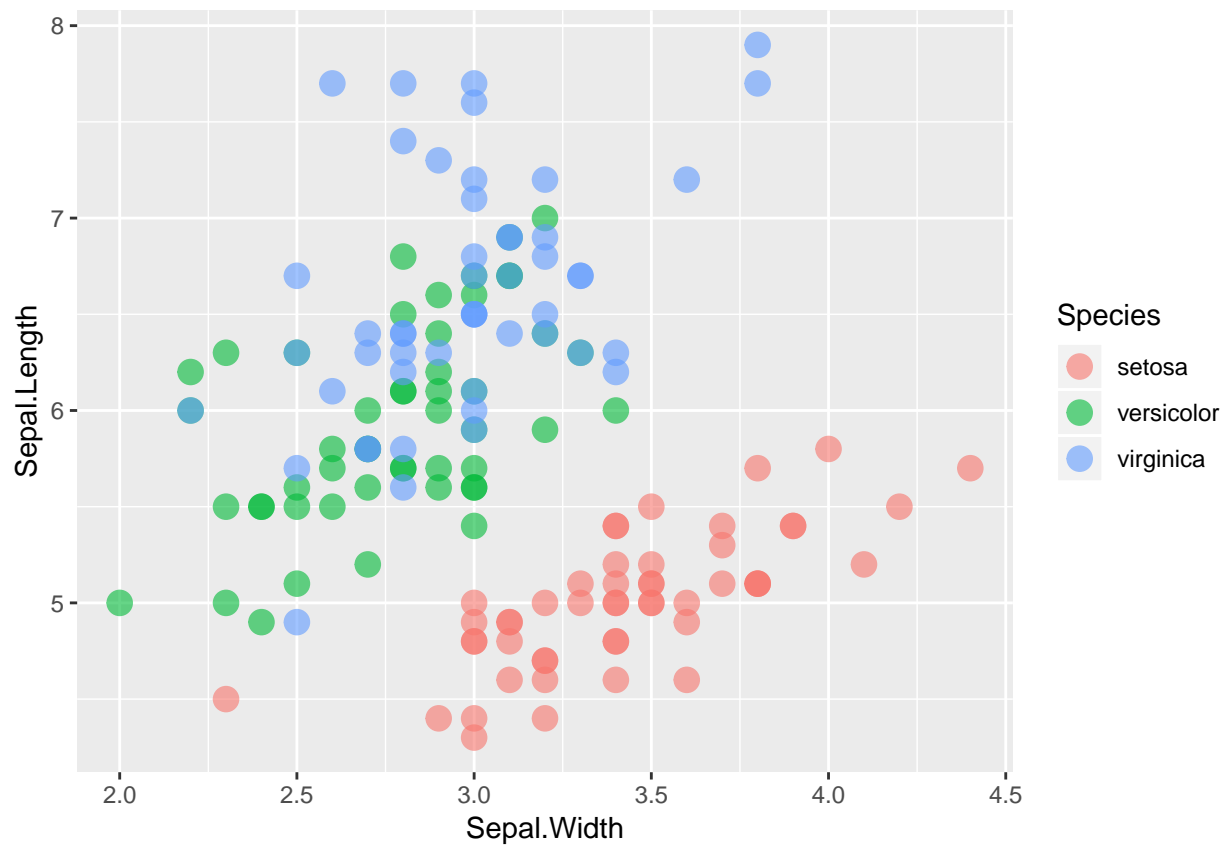
Alpha

The `alpha` = parameter sets the transparency of the plot. Alpha ranges from 0 to 1, where 0 is completely transparent and 1 is completely solid. Making your plots partially transparent is helpful when you have overplotting or any overlapping. (I made the points bigger to make the difference in alpha easier to see)

```
# without transparency
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point(size = 4)
```



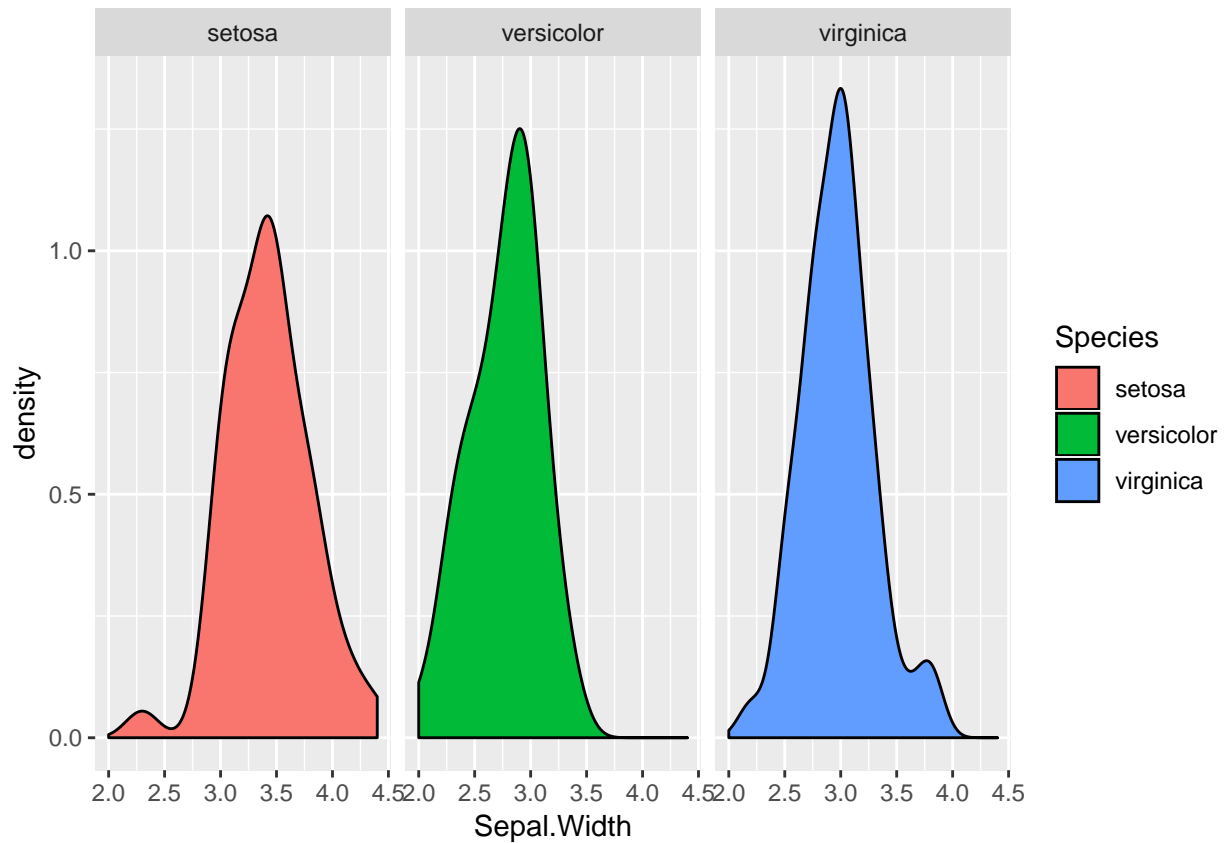
```
# with transparency  
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point(size = 4, alpha = 0.5)
```

facets

Sometimes you might want to make multiple plots based on an element in your data like: significant/not significant, sample, phenotype, etc. If it's a label in your table, you can add a facet to automatically split it.

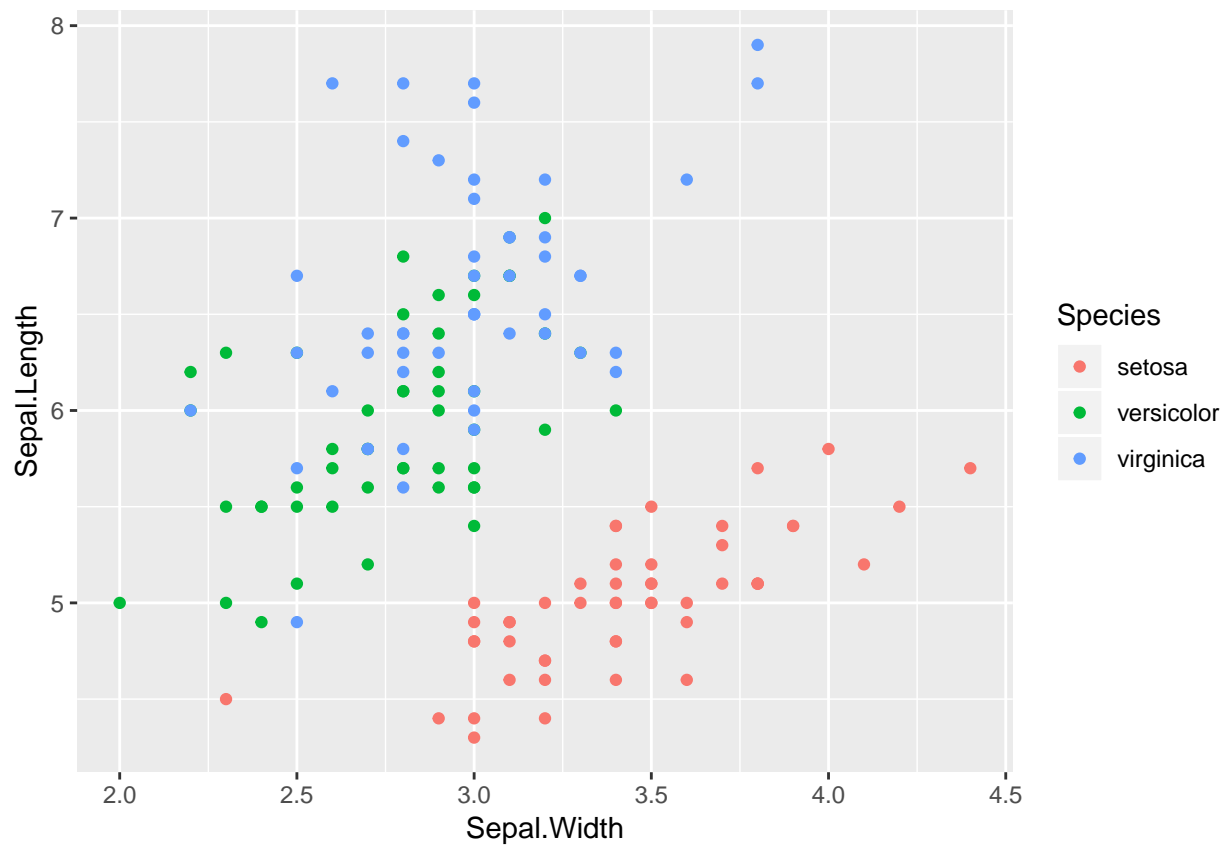
```
ggplot(iris, aes(x = Sepal.Width, fill = Species)) + geom_density() + facet_wrap(~ Species)
```



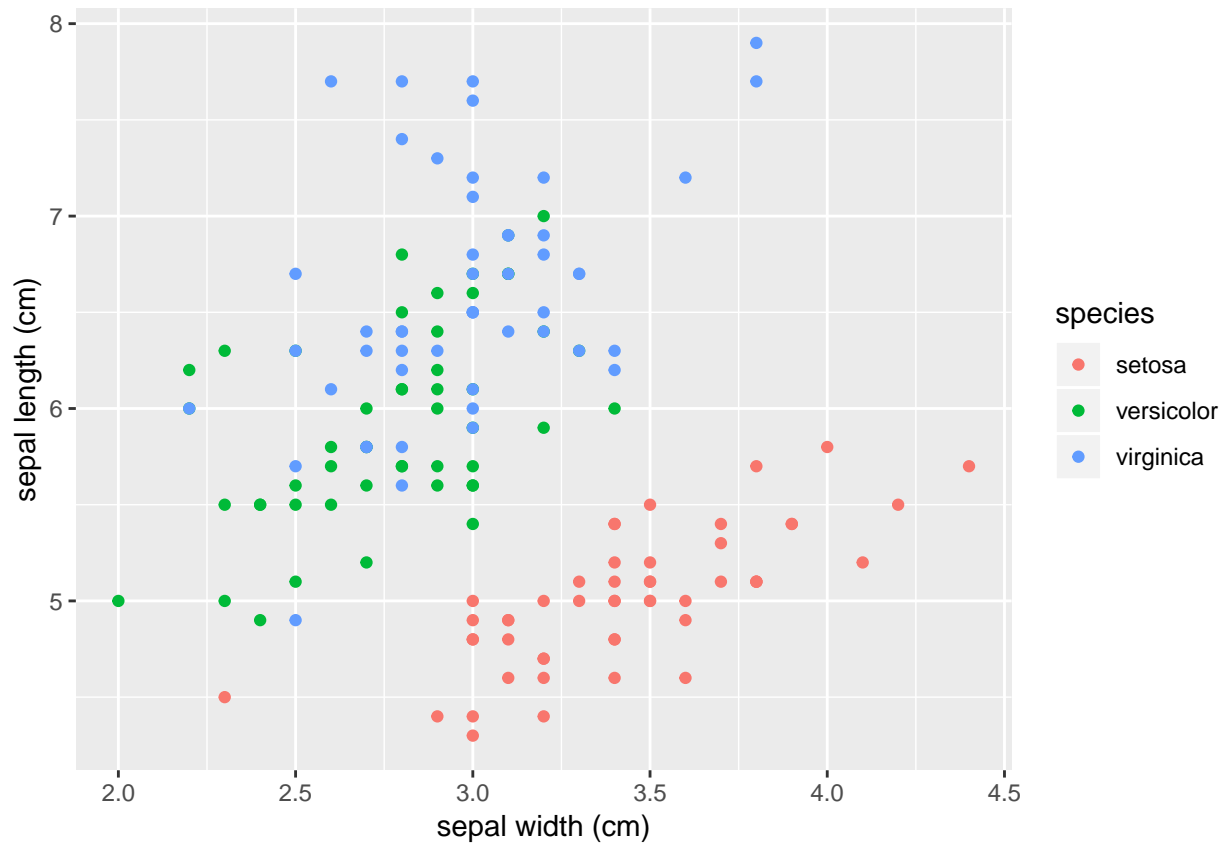
labels

The best practice is to have your column names and labels in your data table formatted nicely so you can plot and not think about it. But sometimes that isn't possible or you don't think about it and you need to rename axis, legend, etc. The easiest, with `labs()`

```
# plot before labelling  
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point()
```



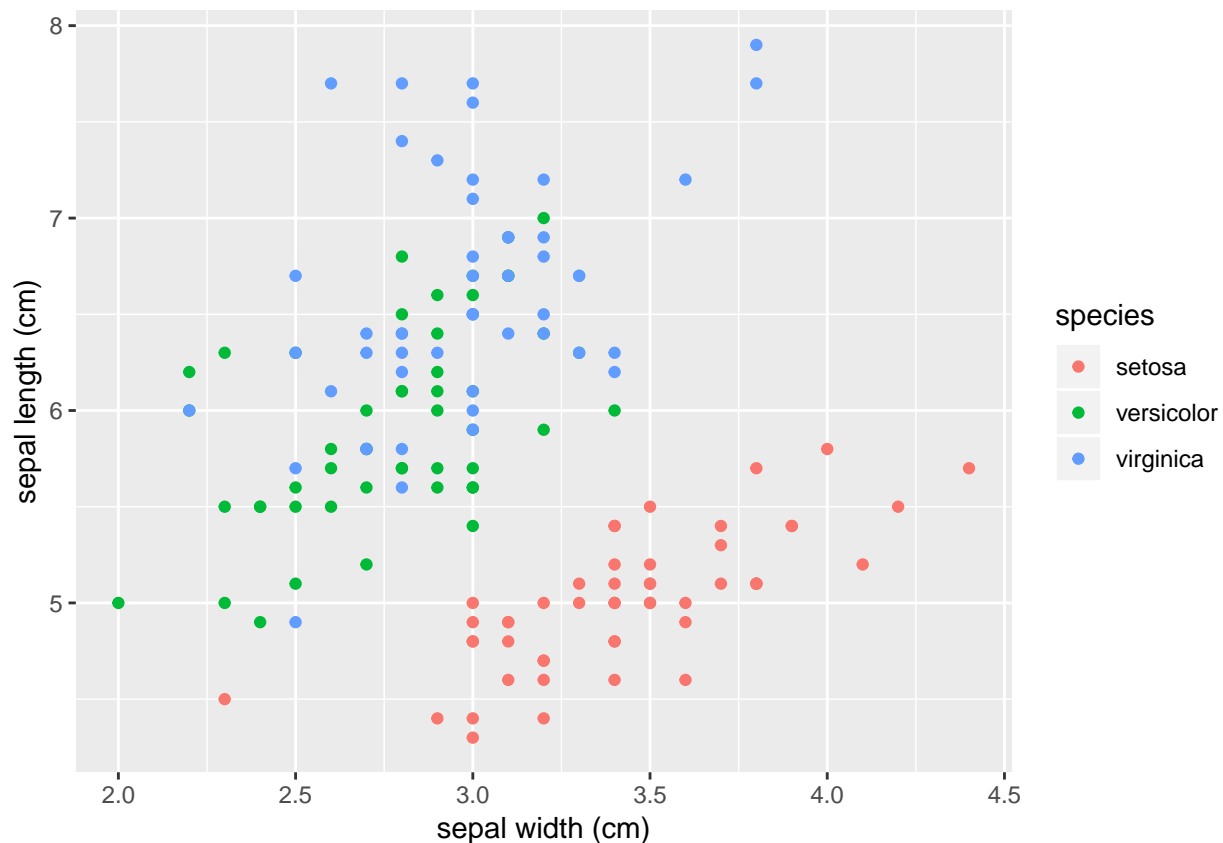
```
# with labels added  
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point() +  
  labs(x = 'sepal width (cm)', y = 'sepal length (cm)', color = 'species')
```



You can also modify the labels (and make more extensive modifications but only labels are shown here) with 'scale_?()?'. The syntax is `scale + + plot part to modify + type of data (discrete or continuous mainly)`. These are the scale modifiers you'll use most often:

- `scale_x_discrete()`
- `scale_x_continuous()`
- `scale_y_discrete()`
- `scale_y_continuous()`
- `scale_color_discrete()`
- `scale_color_continuous()`
- `scale_color_manual()`

```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point() +
  scale_x_continuous(name = 'sepal width (cm)') +
  scale_y_continuous('sepal length (cm)') +
  scale_color_discrete('species')
```



viridis

Viridis is a package of color scales <https://cran.r-project.org/web/packages/viridis/vignettes/intro-to-viridis.html> designed to address several problems. Viridis is

- **colorful:** It has the largest difference possible between the starting color value and the ending color value for to make differences easy to see
- **perceptually uniform:** The change from the starting color value to the ending color happens at the same rate so the differences in similar appearing colors are the same across the color scale
- **color blind friendly:** 8.5% of people (at least in the US) have some form of colorblindness. Chances are there's a colorblind person viewing your paper. Viridis is designed so that color blind people will still see contrasts. Also, viridis still has contrast in greyscale, so it works if a figure is printed in black and white or for someone completely colorblind.
- It looks **pretty**

Get started

If it's not already installed, install the package

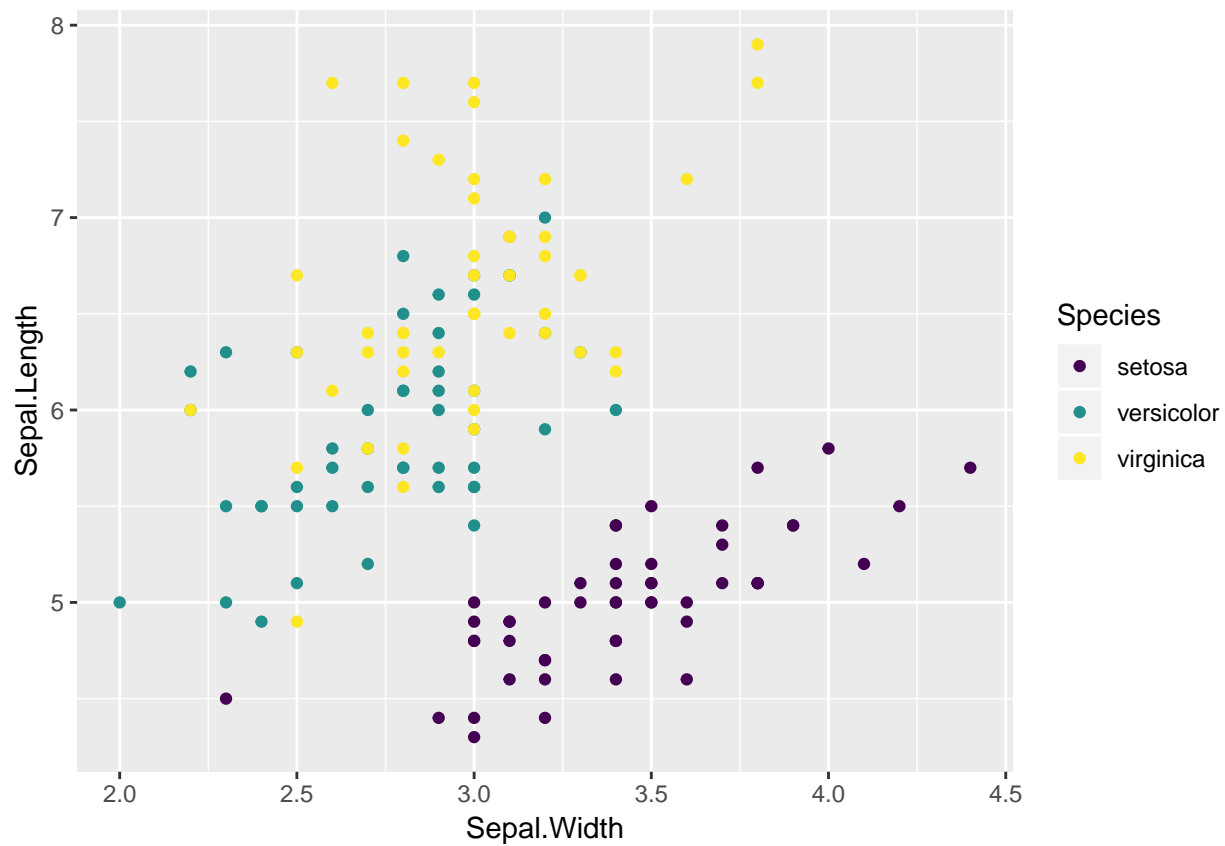
```
#install.packages('viridis')
```

You have to load the package before you can use the color scale

```
library(viridis)
```

Viridis coloring by discrete variable

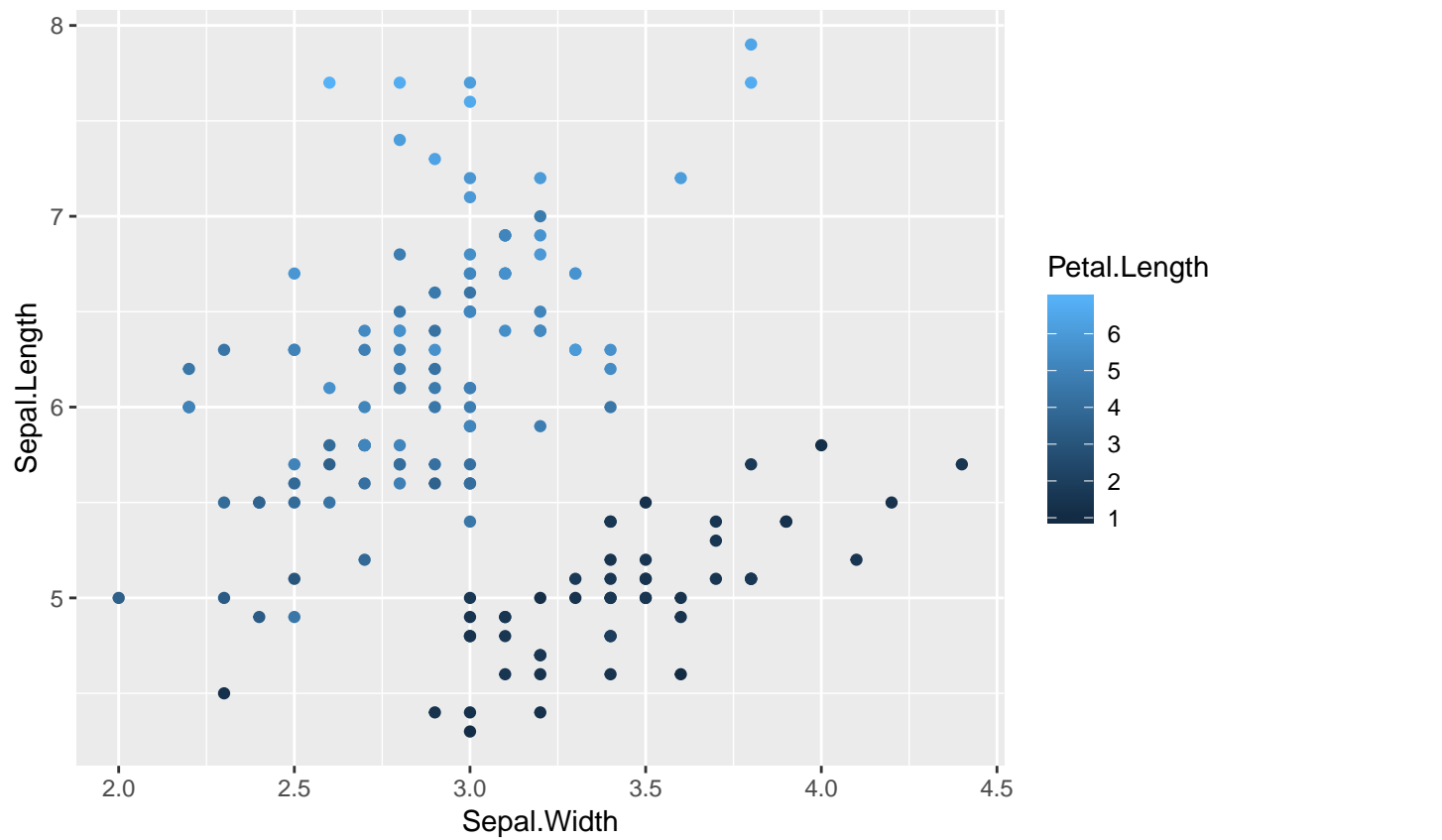
```
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Species)) + geom_point() + scale_color_viridis()
```



Viridis coloring by a continuous variable

The viridis scale has better contrast than the default ggplot color scale for continuous coloring

```
# default ggplot() continuous color scale
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Petal.Length)) + geom_point()
```



```
# viridis continuous color scale  
ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color = Petal.Length)) + geom_point() + scale_color_viridis()
```

