

- As we start with ggplot2 it is important to understand the structure of this.
 - The base graphics built into R require the use of many different functions and each of them seem to have their own method for how to use them.
 - ggplot2 will be more fluid and the more you learn about it the more amazing of graphics you can create.
- ggplot2 components**

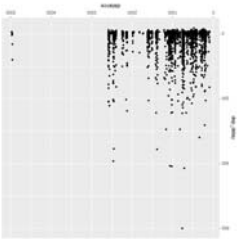
Intro to R Programming for Biostatistics
Day 3 - Graphics in R
Adam J Sullivan

<https://adamsullivan.com/teaching/biostats/day3/>

- What can't ggplot2 do?**
- A good place to start might be with what ggplot2 cannot do. From here we will introduce what it can do.
 - 3d graphs.
 - Interactive graphs, use ggvis
 - Diags, see lgraph

- ggplot2**
- We will begin our journey into statistical graphics with the package ggplot2.
 - This is another package by Hadley Wickham and is part of the tidyverse.
 - This means we can use piping or chaining to build our graphics.





ggplot2 Basics

9/12

12/12

Aesthetic Attributes

- We have focused on the data up until this point and now we will look at the aesthetic attributes.
- These consist of:
 - Color
 - size
 - shape

Aesthetic Attributes

8/12

11/12

ggplot2 Basics

- What the code first does is takes a random 0.5% sample of all of the flights data.
- Given that the original data has 336,776 flights, it can be hard to visualize this much data with any clarity so we will observe a sample for this.
- We then see that the aesthetic mapping is distance by departure delay.
- Finally we have a layer of points. This then leads to the following graph:

7/12

10/12

```
library(dplyr)
library(ggplot2)
library(tidyflights)
data = flights %>% sample_frac(.005)
ggplot(data, aes(distance, y= dep_delay)) +
  geom_point()
```

ggplot2 Basics

For example, we will create a simple scatter plot of distance by departure delay:

- We will get started with the components of every ggplot2 object:
 1. data
 2. aesthetic mappings between variables in the data and visual properties.
 3. At least one layer which describes how to render the data.
- Many of these are with the **geom()** function.

ggplot2 Basics

As we proceed through this section we will begin the graph things in the following pattern:

1. data
2. aesthetic mappings
3. geometric objects
4. statistical transformations
5. scales
6. coordinate systems
7. position adjustments
8. faceting

Color

- To start we will focus on the color. We will look at 1. Color by groups.
- 2. Color of points.

Color by Groups

- An important way to distinguish data can be to change the color of the groups.
- `ggplot2` has many default scales that convert your groups to color levels.
- These can be overridden but we will stick to the basics for now.

---- class #id

Color of points

Not only can we add color to the aesthetic portion but we can add it into the particular layers.

```
library(dplyr)
library(ggplot2)
library(tyfrflight13)

data = flights %>% sample_frac(.05)

ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point(color=blue)
```

Color by Groups

- We consider the same plot that we used before and now we will add color by the carrier brand:

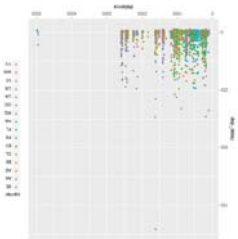
```
library(dplyr)
library(ggplot2)
library(tyfrflight13)

data = flights %>% sample_frac(.05)

ggplot(data, aes(x=distance, y= dep_delay, color=carrier)) +
  geom_point()
```

- The code shows that in the aesthetic portion (`aes()`) and we have added that carrier is associated with color.

Color by Groups



Size

Now that we can color the points and color by groups we can also add the attributes of different sizes.

```
library(dplyr)
library(ggplot2)
library(tyfrflight13)

data = flights %>% sample_frac(.005)

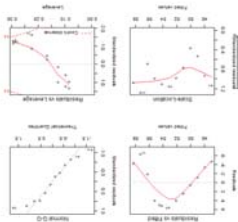
ggplot(data, aes(x=distance, y= dep_delay, size=air_time)) +
  geom_point()
```

- Note that whenever you see an attribute inside the `aes()` function it applies that attribute to a particular variable.
- In this case, the size of the points will be increased depending on the air time of the flight.


```
data = flights %>% sample_frac(.01)
ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point() +
  facet_grid(carrier)
```

We then will note a similar effect when we use `facet_grid()`:

facet_grid

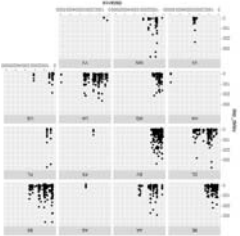


facet_wrap

- Notice that we are still working with the distance versus delay.
- We start out with our original scatter plot.
- Then we add yet another layer.
- This layer is the `facet_wrap()` where we wrap it based on carrier.

facet_grid

- This is where the language of graphs really helps.
- We first take the data and group it based on distance and departure delay.
- We state to place these as points on a graph.
- Finally we use the `facet_grid()` to take that plot and split it by the carrier.
- Each time you add a layer you can accomplish a little more towards your goal.



facet_wrap

- `geom_smooth` fits a smoothing line in data
- `geom_boxplot` box and whisker plot of data
- `geom_histogram` and `geom_freqpoly` distribution graphs
- `geom_bar` distribution of categorical data
- `geom_path` and `geom_line` lines between data points

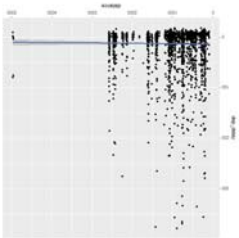
So far we have been focusing on scatter plots. As we continue to move through this section we will note that there are many other `geom` functions that can be used:

What about Other plots?

Smoothing

Smoothing

· Again, the smoothing line comes after our points which means it's another layer added onto our graph:



31/12

Smoothing

- It can be hard to view trends with just points alone.
- Many times we wish to add a smoothing line in order to see what the trends look like.
- This can be especially helpful when trying to understand regressions.

```
ggplot(data, aes(x=distance, y= dep_delay)) +  
  geom_point() +  
  geom_smooth(se=FALSE)
```

- Note that the `geom_smooth()` function adds confidence bands on the smooth as well.
- We can remove these by adding `se=FALSE` inside the `geom_smooth()` function:

Smoothing

30/12

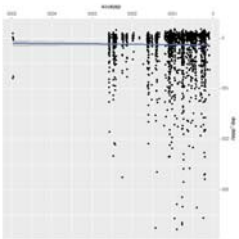
Smoothing

- We will take out scatter plot and apply a smoothing line to this:

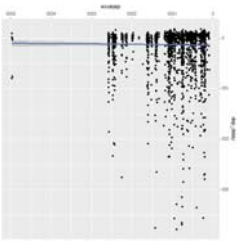
```
ggplot(data, aes(x=distance, y= dep_delay)) +  
  geom_point() +  
  geom_smooth()
```

30/12

Smoothing



30/12



Varying the Smooth

39/112

42/112

```
ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point() +
  geom_smooth(method="loess")
```

- Loess smoothing is a process by which many statistical software do smoothing.
- In ggplot2 this should be done when you have less than 1000 points, otherwise it can be time consuming.

Loess Smooths

39/112

41/112

```
ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point() +
  geom_smooth(span = 0.1)
ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point() +
  geom_smooth(span = 1)
```

- We have so far just seen how to add the smooth without being able to do anything but add or subtract the confidence bands.
- We now will change the smoothness of our smooth that we added.
- To do so we add span=___ inside the geom_smooth() layer:

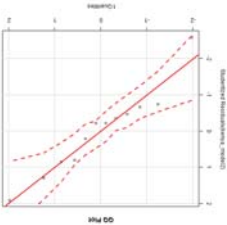
Varying the Smooth

Different Types of Smooths

- There are different types of smooths that we can do. We will consider:
- loess
- gam

37/112

40/112

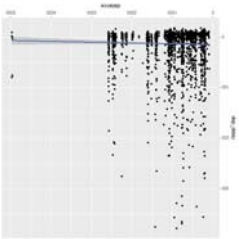


Varying the Smooth

- Note that with span = 0.1 we have a more rough smooth than we had previously.

1. Graph just the data step: `ggplot(data, aes(x=distance, y= dep_delay))`
2. Then add just the smooth `ggplot(data, aes(x=distance, y= dep_delay)) + geom_smooth()`
3. Finally add the points in: `ggplot(data, aes(x=distance, y= dep_delay)) + geom_smooth() + geom_point()`

On Your Own: Rstudio Practice



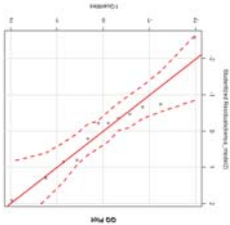
- The code for this is very similar and we can see how it looks below:

Gam Smoothing

```
library(ggvis)
ggplot(data, aes(x=distance, y= dep_delay)) +
  geom_point() +
  geom_smooth(method="gam", formula = y ~ s(x))
```

- We specify this by adding method="gam", formula = y~s(x) into the geom_smooth() layer.
- It works with a large number of points.
- gam smoothing is called generalized additive mode smoothing.

Gam Smoothing



- As you can see with the code we just add method="loess" into the geom_smooth() layer.

Loess Smooths

48/112

```
ggplot(data, aes(x=carrier, y= dep_delay)) +
  geom_jitter()
```

- We can add this as another layer just like we did with geom_point() .
- The jitter plot will add a small amount of random noise to the data and allow it to spread out and be more visible.
- In when you group continuous data into different categories, it can be hard to see where all of the data lies since many points can lie right on top of each other.

Jitter Plot

47/112

- geom_jitter adds random noise
- geom_boxplot boxplots
- geom_violin compact version of density

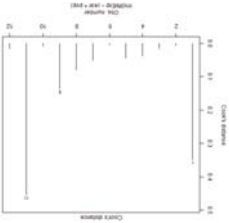
Graphs to Compare Categorical and Continuous Data

- Many times we need to compare categorical and continuous data.
- We will consider the following geom_ functions to do this:

46/112

Graphs to Compare Categorical and Continuous Data

Jitter Plot



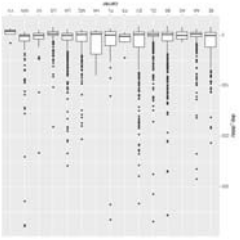
49/12

Boxplot

```
ggplot(data, aes(x=carrier, y=dep_delay)) +  
  geom_boxplot()
```

- Boxplots are one of the most commonly used statistics plots to display continuous data.
- It is extremely useful to evaluate the distribution of a continuous random variable across multiple groups.
- We can easily make this by adding a `geom_boxplot()` layer:

50/12



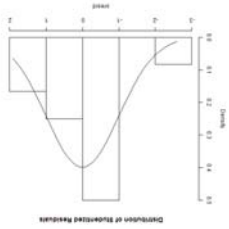
Boxplot

Violin

```
ggplot(data, aes(x=carrier, y=dep_delay)) +  
  geom_violin()
```

- Another plot to help display continuous data among different categories.
- In order to deal with multiple data points lying in a close area, the violin plot is wider at points where the data is bulted.
- We can simply code this with a `geom_violin()` layer:

52/12

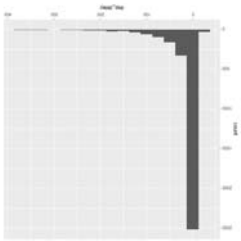


Violin

Continuous Distributions

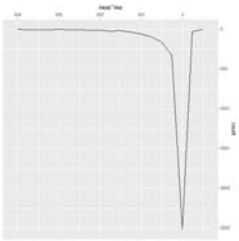
54/12

Histograms



53/112

Frequency Plots



60/112

Histograms

```
ggplot(data, aes(dep_delay)) +  
  geom_histogram(binwidth=25)  
ggplot(data, aes(dep_delay)) +  
  geom_histogram(binwidth=1)
```

blmwidth=...
geom_histogram(), if we would like to specify the width of bins we can do that by using
There are no gaps in a histogram. We can add a histogram layer simply by using the
They are not to be confused with bar charts though!
frequency of those bins is displayed.
This is common among continuous data where the data is split up into bins and the
Another very common graphic that most people have seen and used is the histogram.

56/112

Frequency Plots

```
ggplot(data, aes(dep_delay)) +  
  geom_freqpoly(binwidth=25)  
ggplot(data, aes(dep_delay)) +  
  geom_freqpoly(binwidth=1)
```

We again can use the binwidth=... command:
the geom_freqpoly() layer.
Instead of just having bars to display the frequency in a bin, the frequency plot would place a
Frequency plots are very similar to histograms.

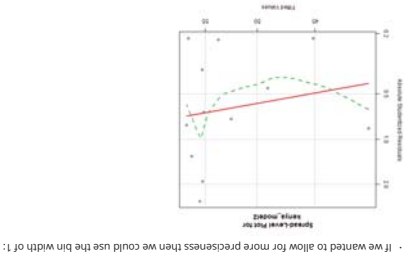
59/112

Continuous Distributions

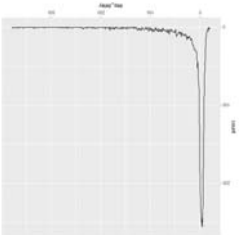
- If we consider just looking at continuous variables we become interested in understanding the distribution that this data takes on.
- We will explore continuous data using:
 - geom_histogram() shows us the distribution of one variable.
 - geom_freqplot uses lines rather than boxes to show the distribution.

55/112

Histograms



58/112



Frequency Plots

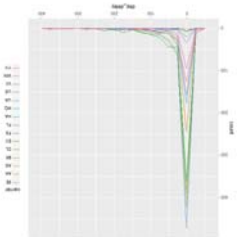
- We can also use a more precise bin width.

6/1/12

Adding Aesthetics

- Just like in the earlier part of this unit we saw that it was possible to add a great deal of aesthetics to plots.
- We will now view how these changes work on these geom_ functions:

```
ggplot(data, aes(dep_delay, color=carrier)) +
  geom_freqpoly(binwidth=25)
```



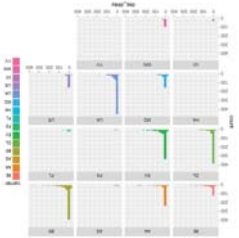
Adding Aesthetics

6/1/12

Adding Aesthetics

- Instead of just coloring the lines, we can use the fill=___ function in order to fill a color by carrier in this case.
- Then we create histograms and finally in order to separate the plots out so we can see things better we use faceting:

```
ggplot(data, aes(dep_delay, fill = carrier)) +
  geom_histogram(binwidth=20) +
  facet_wrap(~carrier)
```



Adding Aesthetics

6/1/12

Adding Aesthetics

- If we add grouping color by carrier we can see the plot below.
- Notice that we now have multiple frequency plots without having to use faceting.

6/1/12

6/1/12

Extra Graphs

6/7/12

Line and Path plots

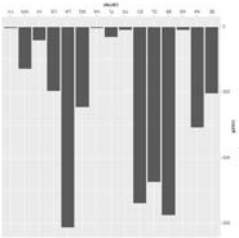
```
ggplot(data, aes(dep_delay, distance/air_time)) +  
  geom_line()
```

7/6/12

Bar Charts

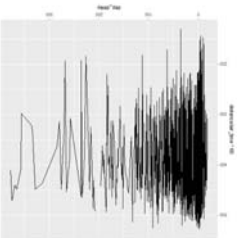
```
ggplot(data, aes(carrier)) +  
  geom_bar()
```

6/6/12



6/6/12

#Labeling



7/1/12

7/2/12

Labelling

- It is very important when making graphs to be able to label our features.
- We will look at various ways in which we can label our graphics now.

Other Text Labels

- Aside from labelling the axes, many times we want to add other text into our graphics.
 - `geom_text` will allow a user to add text to a *graph*.
 - We simply add `geom_text()` as a layer and this layer has the following options:
 - the option `family` allows a user to specify font.
 - the option `fontface` allows a user to specify: plain, bold or italic.
 - `hjust`, `vjust` allows a user to specify location of the text.
 - `size` allows the user to adjust the size of a graph.

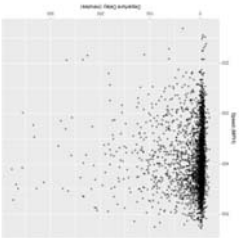
Labelling the Axes

- We can add a lot of features to the axes but for now we will just change labels.
- We use `xlab` and `ylab` for this, if we set them to `NULL` we have blank axes labels.
- For example we can make a graph based on departure delay and speed:

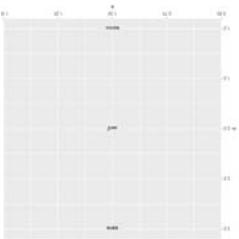
```
ggplot(data, aes(dep_delay, distance/air_time)) +  
  geom_point(alpha = 1 / 3) +  
  xlab("Departure delay (minutes)") +  
  ylab("Speed (mph)")
```

Labelling the Axes

- We can see that the `xlab()` and the `ylab()` functions are just added in as layers.
- They produce the graph below.



Font Families



7/6/12

7/6/12

7/6/12

7/6/12

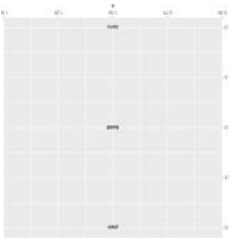
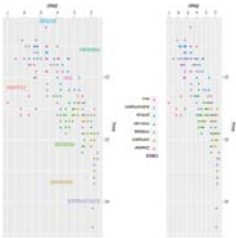
7/6/12

7/6/12

```
df <- data.frame(virt = c("a", "b", "c"), resp = c(1.2, 3.4, 2.5))
ggplot(df, aes(resp, virt)) +
  geom_text(aes(label = paste(" ", resp, "-")), nudge_y = -0.25) +
  xlim(1, 3.6)
```

- We use `nudge_x` and `nudge_y` in order to move the text in the x and y direction respectively.
- If we did not do this, our text would lie directly on top of the point.
- The nudge allows us to move the text horizontally or vertically to label points.

Nudge to label existing points



Font Face Styles

```
library(directlabels)
p1 = ggplot(mpg, aes(displ, hwy, colour = class)) +
  geom_point()
p2 = ggplot(mpg, aes(displ, hwy, colour = class)) +
  geom_point(aes(label = class), method = "smart", gp1d =
  "grid", arrange(p1, p2, ncol=2))
```

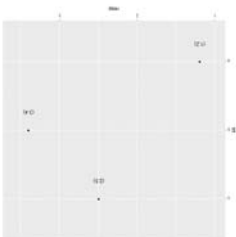
- 1. We created 2 plots `p1` and `p2`.
- 2. `p2` uses `geom_d1()` to add labels as a legend.
- 3. `grid.arrange()` allows us to place the graphs side by side.
- The code is a bit more advanced for this phase, below you can see that we have done a few different things.
- One other common feature of a graph can be to use labeling rather than a legend.

Labels Rather than a Legend

```
df <- data.frame(x = 1, y = 3.1, face = c("plain", "bold", "italic"))
ggplot(df, aes(x, y)) +
  geom_text(aes(label = face, forname = face))
```

- Many times we also wish to add other attributes to our text.
- The font face allows for plain text, **bold text**, and *italic text*.

Font Face Styles

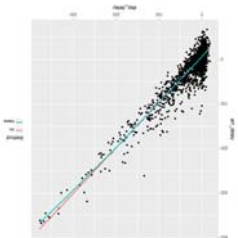


Nudge to label existing points

Other Aesthetic Mappings

Sometimes we map aesthetics to constant values. This allows us to distinguish between layers.

```
ggplot(data, aes(dep_delay, arr_delay)) +  
  geom_point() +  
  geom_smooth(aes(color="lan", method="lm", se=0) +  
    geom_smooth(aes(color="losses", method="loess", se=1) +  
    labs(color = "method")
```

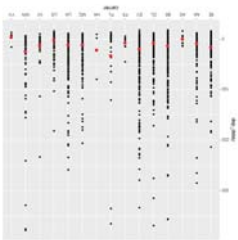


88/112

Statistical Transforms

- Many times we wish to do more than what we have seen at this point.
- We wish to add different statistical features to the graph, stat transforms the data.
- This is typically just a summary of some sort.
- Useful ones are smoothing or identity.
- You typically do not call them directly but the geom does.

```
ggplot(data, aes(carrier, dep_delay)) +  
  geom_point() +  
  stat_summary(geom = "point", fun.y = "mean", color = "red", size = 3)  
ggplot(data, aes(carrier, dep_delay)) +  
  geom_point() +  
  stat_summary(fun.y = "summary", fun.y = "mean", color = "red", size = 3)
```



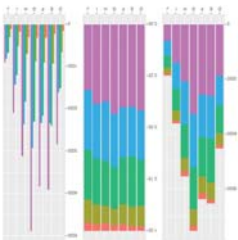
88/112

Position Adjustments

- We can use position adjustments to tweak the position of elements.
- For example with bars:
 - position_stack() stack overlapping bars
 - position_fill() stack overlapping bars and scale to 1
 - position_dodge() place overlapping bars next to each other.

```
dp1ot <- ggplot(diamonds, aes(color, fill = cut)) +  
  xlab(NULL) + ylab(NULL) + theme(legend.position = "none")  
# position stack is the default for bars, so geom_bar()  
# is equivalent to geom_bar(position = "stack").  
p1 = dp1ot + geom_bar()  
p2 = dp1ot + geom_bar(position = "fill")  
p3 = dp1ot + geom_bar(position = "dodge")  
grid.arrange(p1,p2,p3, ncol=3)
```

89/112



90/112

87/112

Modifying Axes and Scales

```
geom(data, aes(dep_delay, arr_delay)) +  
  geom_point(aes(color="carrier")) +  
  scale_y_continuous() +  
  scale_color_discrete()
```

- Scales are required and included in every plot.
 - If we do not specify them, `ggplot2` includes them in the background.
 - For example: `geom(data, aes(dep_delay, arr_delay)) + geom_point(aes(color="carrier"))`
- is read in by `ggplot2` as :

Scales

Axes and Legends

- We can modify both axes and legends. `ggplot2` actually considers these objects to be the same type of object.
- This means if we learn the tools to work with a legend then we can change the Axes in the same way and vice-versa.

Scale Title

- The first argument in a scale function is the axes/legend title.
- We can use 2 types of text:
 - Strings
 - Mathematical Expressions
- For example we will create 2 plots below.
- They will be the same plot but we will allow the first one to just be a string and the second to be a mathematical expression.

tick Label Key Label Labels

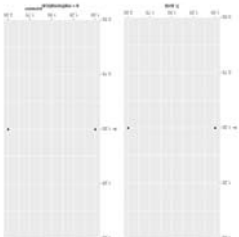
Ticks, grid line key breaks

Label Title name

Axis Legend Argument Name

```
df <- data.frame(x = 1:2, y = 1, z = "a")  
p <- ggplot(df, aes(x, y)) + geom_point()  
p1 = p + scale_x_continuous("X axis")  
p2 = p + scale_x_continuous(quote(a + mathematical ~ expression))  
grid.arrange(p1, p2, ncol=2)
```

Scale Title



Scale Title

9/7/12

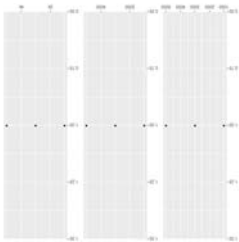
Earlier we learned about common labeling functions such as:

- xlab
- ylab
- labs

We can also use common text notations in order to add further details:

```
p <- ggplot(df, aes(x, y)) + geom_point(aes(colour = z))
p1 = p + xlab("x axis") + ylab("y axis")
p2 = p + labs(x = "x axis", y = "y axis",
  colour = "Colour in legend")
grid.arrange(p1, p2, ncol=2)
```

Labeling a Scale



10/7/12

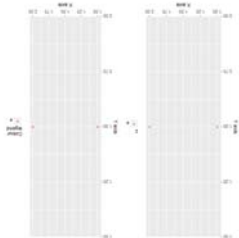
```
df <- data.frame(x = c(1, 3, 5) * 1000, y = 1)
aes <- ggplot(df, aes(x, y)) +
  geom_point() +
  labs(x = NULL, y = NULL)
aes + scale_x_continuous(breaks = c(2000, 4000))
aes + scale_y_continuous(breaks = c(2000, 4000))
labels = c("2k", "4k")
```

- The breaks argument controls what values appear as the tick marks on axes and keys.
- We not only like to be able to change the labels of scales but it can be helpful to choose the tick marks as well.

Breaks and Labels

```
leg <- ggplot(df, aes(y, x, fill = x)) +
  geom_tile() +
  labs(x = NULL, y = NULL)
leg + scale_fill_continuous(breaks = c(2000, 4000))
leg + scale_fill_continuous(breaks = c(2000, 4000), labels = c("2k", "4k"))
```

- As it was state before `ggplot2` considers axes and legends to be the same type.
- This means if we are creating a continuous scale with a bar graph coloring or even a heat map we can change the tick marks on the legend as well.
- We see that just like the axes above we now have three different legends with the tick marks and labels of them changed.



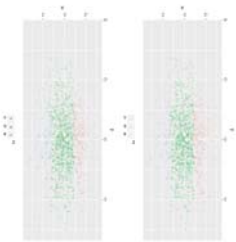
9/7/12

10/7/12


```
df <- data.frame(x = 1:3, y = 1:3, z = c("a", "b", "c"))
geom_point(aes(colour = z), size = 3) +
  xlab(NULL) +
  ylab(NULL)
base + theme(legend.position = "right") # the default
base + theme(legend.position = "none")
```

- We can choose the specific layout of the legend.
- For example we can place it at the "top", "bottom", "right" or not even have a legend at all.

Legend Layouts

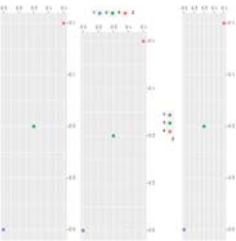


Legends

```
norm <- data.frame(x = norm(1000), y = norm(1000))
norm$z <- cut(norm$x, 3, labels = c("a", "b", "c"))
geom_point(aes(x), size = 0.1) +
  geom_point(aes(colour = z), alpha = 0.1)
geom_point(aes(colour = z), alpha = 0.1) +
  guides(colour = guide_legend(override.aes = list(alpha = 1)))
```

- Sometimes we wish to have the legends display different things.
- For example if we use transparent colors in the plot we may want solid colors in the legend.
- The alpha command will make the colors more transparent.
- However we can override the legend colors and set the alpha differently.

Legends



Legends