

SOC-5811 Week 2: Introduction to R and ggplot2

Nick Graetz

University of Minnesota, Department of Sociology

9/10/2025

1/96

IN THIS LECTURE

1. What is R?
2. RStudio interface
3. Packages
4. R as calculator
5. Anatomy of a function
6. Help files
7. R scripts

WHAT IS R?

- ▶ R is a language for statistical computing and graphics
- ▶ Originally developed in 1992 by Robert Gentleman and Ross Ihaka based on the programming language S
- ▶ The core of the R language is maintained by the R Core Team
- ▶ A (very) large number of packages which add additional functionality are maintained by other contributors

WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities



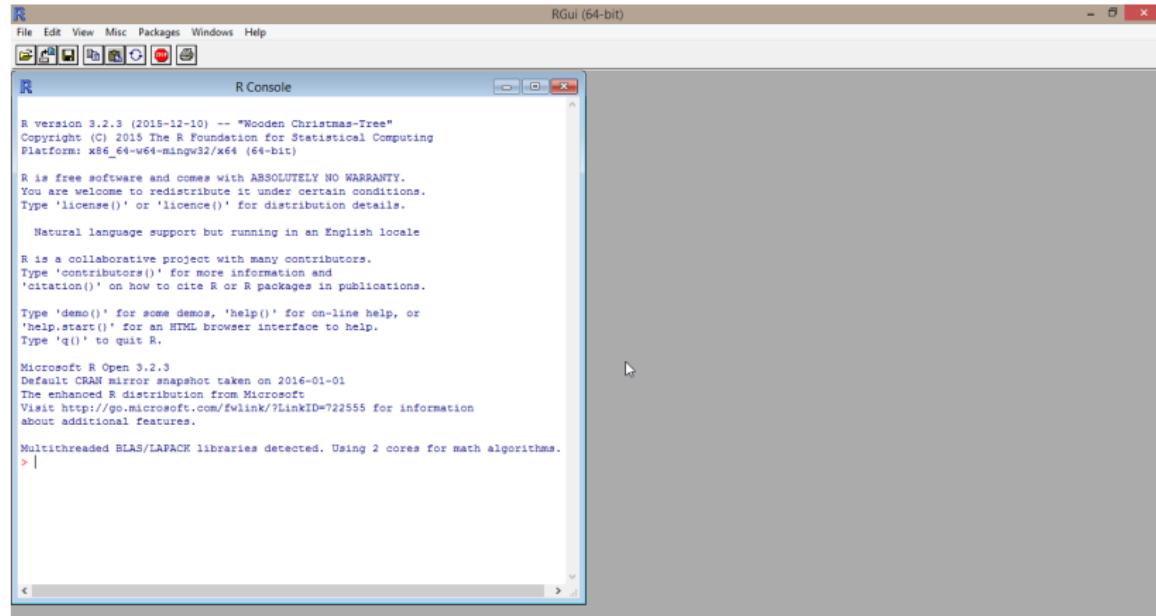
WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities
- ▶ R is free software
 - ▶ You don't have to pay for it (and you can share it with anyone)
 - ▶ You can use and modify it as you see fit

WHY USE R?

- ▶ R can do many useful things
 - ▶ Flexible data management
 - ▶ Powerful statistical capabilities, particularly for modeling
 - ▶ Extensive graphics capabilities
- ▶ R is free software
 - ▶ You don't have to pay for it (and you can share it with anyone)
 - ▶ You can use and modify it as you see fit
- ▶ R has a large (and enthusiastic) user base
 - ▶ This makes finding help relatively straightforward
 - ▶ New methods are often implemented in R very quickly

R (GUI) INTERFACE



WHAT IS RSTUDIO?

- ▶ “Integrated development environment” (IDE)

WHAT IS RSTUDIO?

- ▶ “Integrated development environment” (IDE)
- ▶ Convenient interface for R which incorporates a number of useful features for developing code
 - ▶ syntax highlighting
 - ▶ code completion
 - ▶ code navigation
 - ▶ debugging tools
 - ▶ etc.

WHAT IS RSTUDIO?

- ▶ “Integrated development environment” (IDE)
- ▶ Convenient interface for R which incorporates a number of useful features for developing code
 - ▶ syntax highlighting
 - ▶ code completion
 - ▶ code navigation
 - ▶ debugging tools
 - ▶ etc.
- ▶ Also provides integration with other useful tools
 - ▶ Shiny (for developing web apps)
 - ▶ R Markdown (for authoring documents and slides)
 - ▶ Git/Subversion (for version control)

RSTUDIO INTERFACE

The screenshot displays the RStudio interface with the following components:

- Script Editor (Left):** Shows the code for generating 1,000 draws from a normal distribution, summarizing the draws, and plotting a histogram.
- Console (Bottom Left):** Displays the R command history and the resulting output, including the generated data summary and histogram plot.
- Plots (Top Right):** A histogram titled "Histogram of x" showing the frequency distribution of the generated data.
- Environment (Top Right):** Shows the current environment variables.
- Help Viewer (Bottom Right):** Provides documentation for the `Normal` function, specifically the `stats` version.

```
example.r
1 # generate 1,000 draws from a normal distribution
2 x <- rnorm(1000)
3
4 # summarize these draws
5 summary(x)
6
7 # plot a histogram
8 hist(x)
9
```

```
> # generate 1,000 draws from a normal distribution
> x <- rnorm(1000)

> # summarize these draws
> summary(x)
   Min. 1st Qu. Median Mean 3rd Qu. Max.
-4.24700 -0.68000  0.03167 -0.01287  0.62250  3.08900

> # plot a histogram
> hist(x)
>
```

Histogram of x

Frequency

X

Normal (stats)

The Normal Distribution

Description

Density, distribution function, quantile function and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.

RSTUDIO INTERFACE

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help.
- Tools Menu:** Opened, showing options like Import Dataset, Install Packages..., Check for Package Updates..., Version Control, Shell..., Addins, Keyboard Shortcuts Help, Modify Keyboard Shortcuts..., Project Options..., and Global Options... (highlighted).
- Environment Tab:** Shows a histogram titled "Histogram of x". The x-axis is labeled "x" and ranges from -4 to 2. The y-axis is labeled "Frequency" and ranges from 0 to 200. The distribution is roughly symmetric and centered around 0.
- Console Tab:** Displays R code and its output. The code generates 1,000 draws from a normal distribution, summarizes the draws, and plots a histogram.
- R Script Tab:** Shows the same R script as the Console tab.
- R Documentation Tab:** Shows the "Normal (stats)" page under "R Documentation". It includes sections for "Description" and "Density, distribution function, quantile function and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd."

```
> # generate 1,000 draws from a normal distribution
> x <- rnorm(1000)

> # summarize these draws
> summary(x)
   Min. 1st Qu. Median  Mean 3rd Qu. Max.
-4.24700 -0.68000  0.03167 -0.01287  0.62250  3.08900

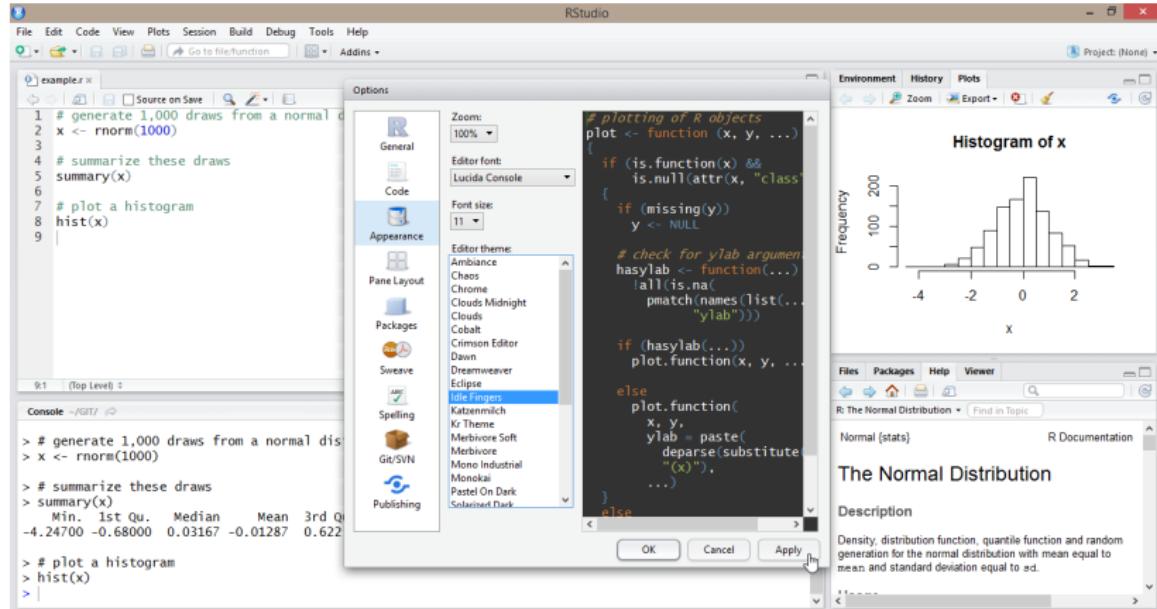
> # plot a histogram
> hist(x)
>
```

RSTUDIO INTERFACE

The screenshot displays the RStudio interface with the following components:

- Code Editor:** Shows a script named "example.r" containing R code to generate 1,000 draws from a normal distribution, summarize the draws, and plot a histogram.
- Console:** Shows the output of the R code run in the console.
- Options Dialog:** A modal dialog titled "Options" is open, showing settings for R version (64-bit), working directory (~/.R), and various startup behaviors. It includes sections for General, Code, Appearance, Pane Layout, Packages, Sweave, Spelling, Git/SVN, and Publishing.
- Plots:** A histogram titled "Histogram of x" is displayed, showing the frequency distribution of the generated draws.
- Help Viewer:** The "The Normal Distribution" page is open in the help viewer, providing documentation for the R function "Normal".

RSTUDIO INTERFACE



RSTUDIO INTERFACE

The screenshot shows the RStudio interface with the following components:

- File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help** menu bar.
- Addins** dropdown menu.
- example.r** file open in the Source pane.
- Console** pane showing R code execution and output.
- Options** dialog box open, titled "Choose the layout of the panes in RStudio by selecting from the controls in each quadrant." It has three tabs: **Source**, **Environment, History, Plots, Pre**, and **Console**.
 - Source** tab: Environment, History, Plots, Pre checkboxes are checked.
 - Environment** tab: Environment, History, Files, Plots, Packages, Help, Build, VCS, Viewer checkboxes are checked.
 - Console** tab: Environment, History, Files, Plots, Packages, Help, Build, VCS, Viewer checkboxes are checked.
- Environment**, **History**, **Plots** tabs are selected in the top right.
- Histogram of x** plot in the Plots pane.
- Environment**, **Packages**, **Help**, **Build**, **VCS**, **Viewer** tabs are selected in the bottom right.
- R: The Normal Distribution** documentation page is open in the bottom right.

PACKAGES

Most basic R functionality is part of base and is loaded automatically when you start R. Additional functionality can be added through packages.

PACKAGES

Most basic R functionality is part of base and is loaded automatically when you start R. Additional functionality can be added through packages.

The first time you use a package, it needs to be installed:

```
install.packages ("ggplot2")
```

After that, you just need to load the package using the library () command whenever you start a new instance of R:

```
library (ggplot2)
```

R AS CALCULATOR

R can be used as a calculator by just typing in the console.

All of the basic arithmetic operators (+, -, *, /, ^) do what you would expect them to do, following normal order of operations conventions:

```
230 + 97
```

```
## [1] 327
```

```
500/20
```

```
## [1] 25
```

R AS CALCULATOR

Parentheses can be used to alter the order of operations:

```
300/20^1/2
```

```
## [1] 7.5
```

```
(300/20)^(1/2)
```

```
## [1] 3.872983
```

R AS CALCULATOR: QUICK EXERCISE

1. How many seconds are in September?
2. What is 80 degrees Fahrenheit in degrees Celsius?
3. How much longer is 1 mile than 1600 meters (in feet)?



R AS CALCULATOR: QUICK EXERCISE

1. How many seconds are there in September?

```
30 * 24 * 60 * 60
```

```
## [1] 2592000
```

2. What is 80 degrees Fahrenheit in degrees Celsius?

```
(80 - 32) * (5/9)
```

```
## [1] 26.66667
```

3. How much longer is 1 mile than 1600 meters (in feet)?

```
5280 - 1600 * 3.28084
```

```
## [1] 30.656
```

16/96

FUNCTIONS

R functions are used to transform input into output in some way.

For example...

```
log(10)
```

```
## [1] 2.302585
```

```
exp(3)
```

```
## [1] 20.08554
```

```
sqrt(80)
```

```
## [1] 8.944272
```

FUNCTIONS: ANATOMY

```
log(x = 300, base = 10)
```

```
## [1] 2.477121
```

1. Function name: **log()**
2. Argument name(s): **x, base**
3. Argument value(s): **300, 10**
4. Output: **2.4771213**

FUNCTIONS: ARGUMENT ORDER

Arguments can be specified in any order *if they are named*:

```
log(x = 300, base = 10)
```

```
## [1] 2.477121
```

```
log(base = 10, x = 300)
```

```
## [1] 2.477121
```

FUNCTIONS: ARGUMENT NAMES

Arguments don't need to be named, but then *there is only one correct order*:

```
log(x = 300, base = 10)
```

```
## [1] 2.477121
```

```
log(base = 10, x = 300)
```

```
## [1] 2.477121
```

```
log(300, 10)
```

```
## [1] 2.477121
```

```
log(10, 300)
```

```
## [1] 0.4036944
```

20/96



UNIVERSITY OF MINNESOTA

Driven to DiscoverSM



MINNESOTA POPULATION CENTER

FUNCTIONS: DEFAULTS

Some (but not all) arguments have defaults and don't need to be specified, assuming you are happy with the default:

```
log(x = 300)
```

```
## [1] 5.703782
```

```
log(base = 10)
```

```
## Error: argument "x" is missing, with no default
```

FUNCTIONS: COMBINING

Functions can be combined or nested with other functions and operators:

```
exp(log(10) + log(10))
```

```
## [1] 100
```

```
log(x = (4 * 10) / 7, base = 10)
```

```
## [1] 0.756962
```

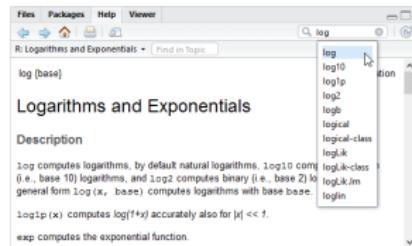
HELP FILES

Every function has a help file.

You can access a help file from the console:

```
help(log)
```

or from the help tab in RStudio:



HELP FILES

`log {base}`

R Documentation

Logarithms and Exponentials

Description

`log` computes logarithms, by default natural logarithms, `log10` computes common (i.e., base 10) logarithms, and `log2` computes binary (i.e., base 2) logarithms. The general form `log(x, base)` computes logarithms with base `base`.

`log1p(x)` computes $\log(1+x)$ accurately also for $|x| << 1$.

`exp` computes the exponential function.

`expm1(x)` computes $\exp(x) - 1$ accurately also for $|x| << 1$.

HELP FILES

Usage

```
log(x, base = exp(1))
logb(x, base = exp(1))
log10(x)
log2(x)

log1p(x)

exp(x)
expml(x)
```

Arguments

`x`
a numeric or complex vector.

`base`
a positive or complex number: the base with respect to which logarithms are computed. Defaults to `e=exp(1)`.

HELP FILES

Details

All except `logb` are generic functions: methods can be defined for them individually or via the [Math](#) group generic.

`log10` and `log2` are only convenience wrappers, but logs to bases 10 and 2 (whether computed *via* `log` or the wrappers) will be computed more efficiently and accurately where supported by the OS. Methods can be set for them individually (and otherwise methods for `log` will be used).

`logb` is a wrapper for `log` for compatibility with S. If (S3 or S4) methods are set for `log` they will be dispatched. Do not set S4 methods on `logb` itself.

All except `log` are [primitive](#) functions.

Value

A vector of the same length as `x` containing the transformed values. `log(0)` gives `-Inf`, and `log(x)` for negative values of `x` is `NaN`. `exp(-Inf)` is 0.

For complex inputs to the `log` functions, the value is a complex number with imaginary part in the range $[-\pi, \pi]$: which end of the range is used might be platform-specific.

HELP FILES

References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole. (for `log`, `log10` and `exp`)

Chambers, J. M. (1998) *Programming with Data. A Guide to the S Language*. Springer. (for `logb`)

See Also

[Trig](#), [sqrt](#), [Arithmetic](#).

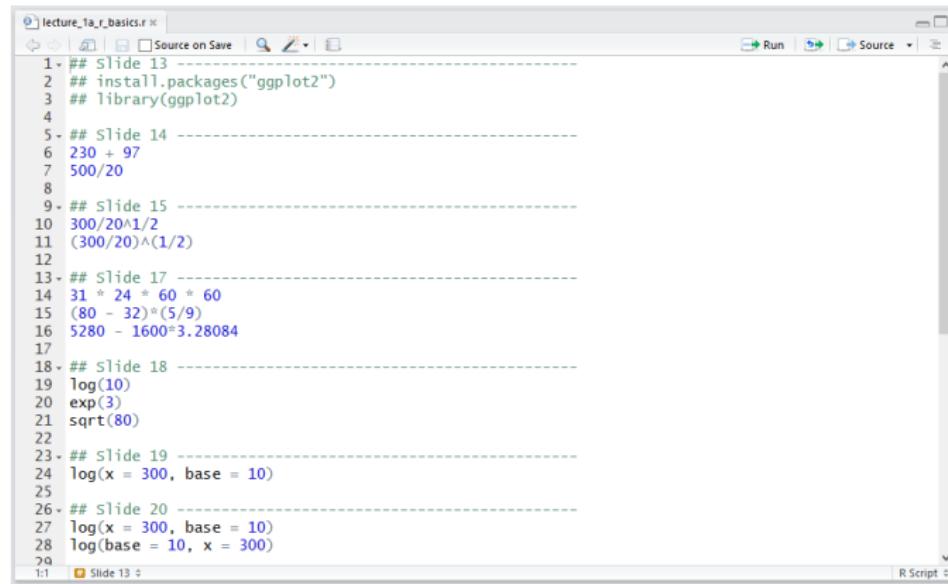
Examples

```
log(exp(3))
log10(1e7) # = 7

x <- 10^{-(1+2*1:9)}
cbind(x, log(1+x), logp(x), exp(x)-1, expml(x))
```

R SCRIPTS

An R script is a text file (.r extension) with a series of R commands and (ideally) some useful commentary.



```
lecture_1a_r_basics.r
1 ## Slide 13 -----
2 ## install.packages("ggplot2")
3 ## library(ggplot2)
4
5 ## Slide 14 -----
6 230 + 97
7 500/20
8
9 ## Slide 15 -----
10 300/20^(1/2)
11 (300/20)^(1/2)
12
13 ## Slide 17 -----
14 31 * 24 * 60 * 60
15 (80 - 32)^(5/9)
16 5280 - 1600^3.28084
17
18 ## Slide 18 -----
19 log(10)
20 exp(3)
21 sqrt(80)
22
23 ## Slide 19 -----
24 log(x = 300, base = 10)
25
26 ## Slide 20 -----
27 log(x = 300, base = 10)
28 log(base = 10, x = 300)
29
```



WHY USE A SCRIPT?

Typing in the console is fine for quick calculations or experimentation with a command, but a script provides...

- ▶ a full record of all commands required to carry out an analysis
- ▶ a convenient mechanism for repeating an analysis without needing to retype everything (no need to reinvent the wheel)
- ▶ a starting point for writing new code
- ▶ a vehicle for providing context and commentary for your code

WHY USE A SCRIPT?

Any analysis you do should be saved as a script!

Without a script...

- ▶ you will forget what you've done
- ▶ you will forget why you did it
- ▶ no one else will ever know what you did or why you did it
- ▶ you will have do things over again for no reason

RUNNING A SCRIPT

If your script is open in RStudio, you can run the whole thing using `ctrl + shift + enter` or just a single line (or highlighted block) using `ctrl + enter`.

Or you can run a script from the command line using the `source()` function:

```
source ("C:/Users/ngraetz/Dropbox/my_script.R")
```

COMMENTING A SCRIPT

R will ignore any line in a script that starts with #, so you can use this to add comments to your code:

```
# add 1-5  
1 + 2 + 3 + 4 + 5
```

```
## [1] 15
```

```
# find the natural log of 10  
log(10)
```

```
## [1] 2.302585
```

COMMENTING A SCRIPT

Use comments to:

- ▶ Label blocks of code. This will help you navigate your code later
- ▶ Explain why you're doing something (if it's not self-evident)
- ▶ Write yourself (and other users) notes about particularly tricky lines of code

COMMENTING A SCRIPT

Use comments to:

- ▶ Label blocks of code. This will help you navigate your code later
- ▶ Explain why you're doing something (if it's not self-evident)
- ▶ Write yourself (and other users) notes about particularly tricky lines of code

You want to provide enough information so that your future self, or someone else, can quickly understand the structure and purpose of your code at a later date.

However, it is possible to provide too much information, making your code more cumbersome (e.g., writing out what each line of code does).



HEADERS

It's also good practice to use '#' to provide some sort of header at the top of your code:

```
#####
## Author:      John Doe
##
## Description: A short description of what this code does
##               and any important context for why.
##
## Output:       A list of files that are output by this
##               code.
##
## Notes:        Anything someone should know when running
##               this code.
#####
```

DATA VISUALIZATION

1. Understanding the ggplot approach
2. Aesthetics
3. Geoms
4. Facets
5. Options and customization
6. Reshaping
7. Saving plots
8. Additional packages



WHAT IS GGPLOT2?

- ▶ `ggplot2` contains functions that allow you to build complex graphics using a relatively small set of building blocks
- ▶ It's based on the book "Grammar of Graphics", which defined a fundamental theory of data visualization
- ▶ NOTE: the online documentation for `ggplot2` is fantastic, and lays all the functions out in terms of these building blocks:

<http://ggplot2.tidyverse.org/reference/>

<https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf>

LOAD LIBRARIES & DATA

```
library(ggplot2)
library(data.table)
mmr_data <- fread("data/mmr_data.csv")
head(mmr_data)
```

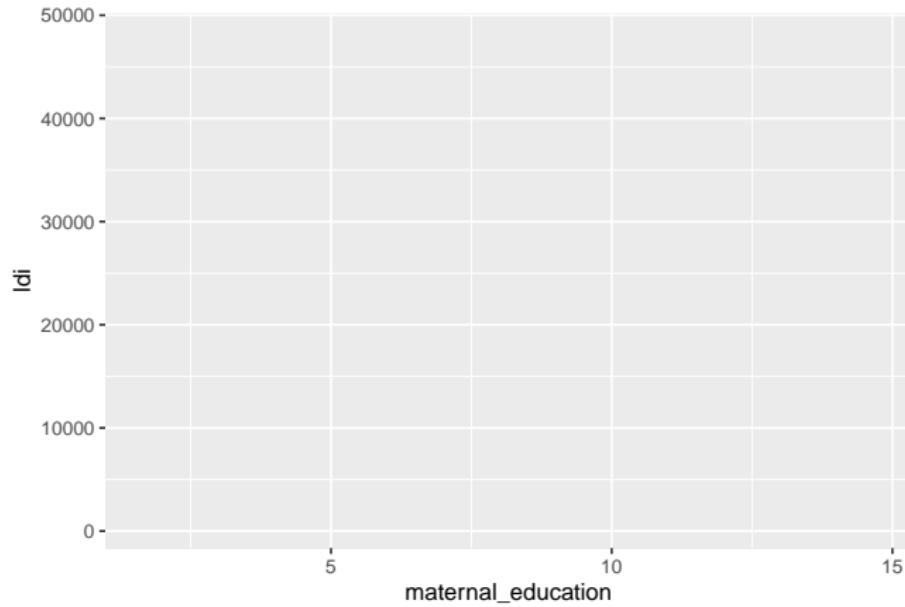
```
##   year_id      mmr maternal_education      ldi location_name
##   <int>    <num>          <num>    <num>    <char>
## 1: 2015 52.57428      9.900764 10593.983     China
## 2: 2015 50.41785      5.943825 2773.896   Cambodia
## 3: 2015 25.58855     11.535423 20782.643  Malaysia
## 4: 2015 61.25871     14.697507 33327.094     Japan
## 5: 2015 25.15193     14.635294 40454.078  Australia
## 6: 2015 33.02467     14.244808 35569.391 United Kingdom
##                     super_region_name      region_name
##                     <char>            <char>
## 1: Southeast Asia, East Asia, and Oceania      East Asia
## 2: Southeast Asia, East Asia, and Oceania Southeast Asia
## 3: Southeast Asia, East Asia, and Oceania Southeast Asia
## 4: High-income High-income Asia Pacific
## 5: High-income           Australasia
## 6: High-income        Western Europe
```

HOW DOES GGPLOT2 WORK?

First, you set up the graph:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi))
```

HOW DOES GGPLOT2 WORK?

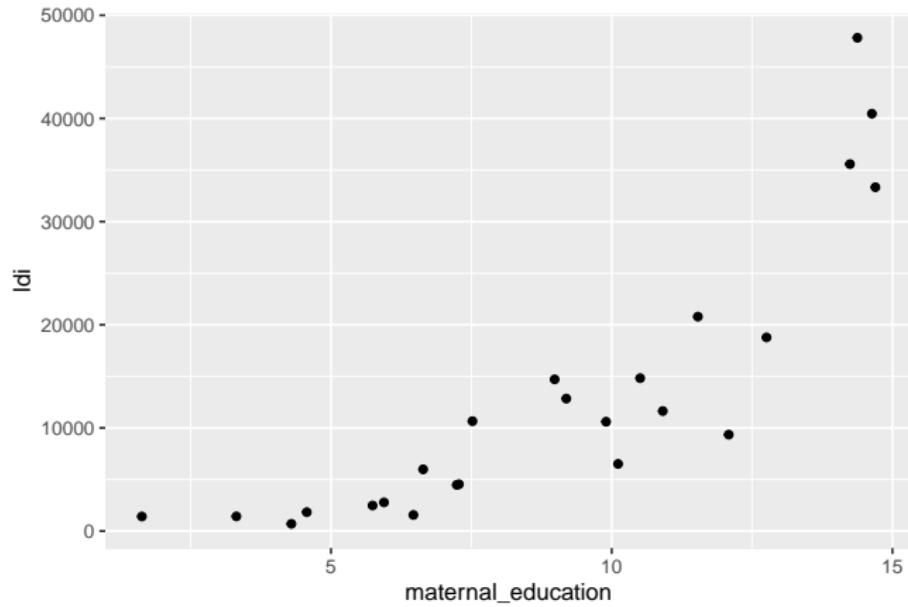


HOW DOES GGPLOT2 WORK?

Then, you add to it. Basically telling ggplot what type of graph to make:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi)) +  
  geom_point()
```

HOW DOES GGPLOT2 WORK?



WHAT ARE THE BUILDING BLOCKS OF A GGPLOT?

- ▶ **Aesthetics**
- ▶ **Geoms**
- ▶ **Facets**
- ▶ Positions
- ▶ Scales
- ▶ Labels
- ▶ Themes

AESTHETICS

We map aesthetics with `aes` in the initial `ggplot()` call:

```
ggplot(data = mmr_data,  
       aes(x = maternal_education,  
            y = ldi)) +  
  geom_point()
```

AESTHETICS

The `aes` in the initial `ggplot()` call

“Aesthetic mapping” is how you tell `ggplot` which variable is `x`, which is `y`

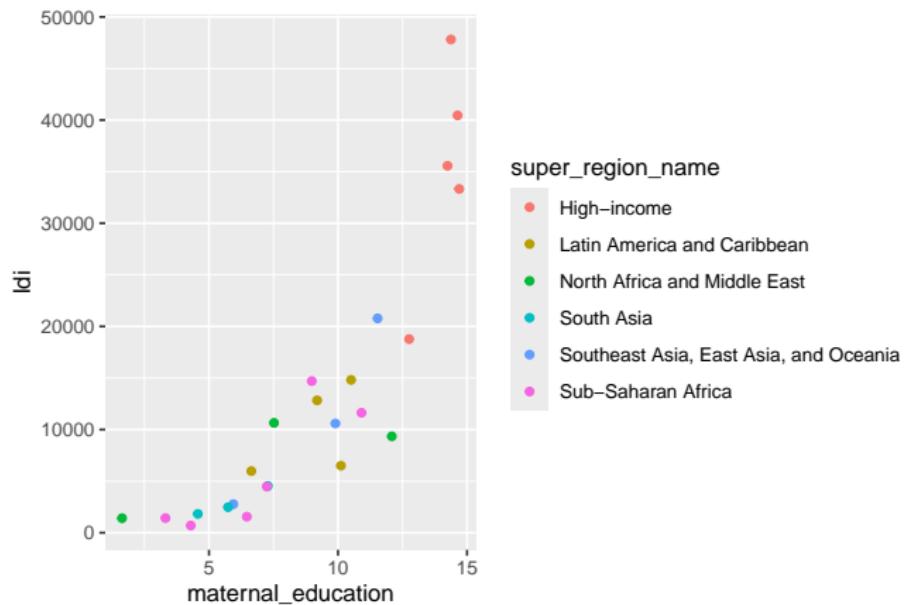
But, you can use them for more than just the axes:

- ▶ color (border color)
- ▶ fill (fill color)
- ▶ shape
- ▶ linetype (solid, dashed, dotted etc.)
- ▶ size
- ▶ alpha (transparency)
- ▶ labels

EXAMPLE OF AESTHETIC MAPPING

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi,  
                      color = super_region_name)) +  
  geom_point()
```

EXAMPLE OF AESTHETIC MAPPING



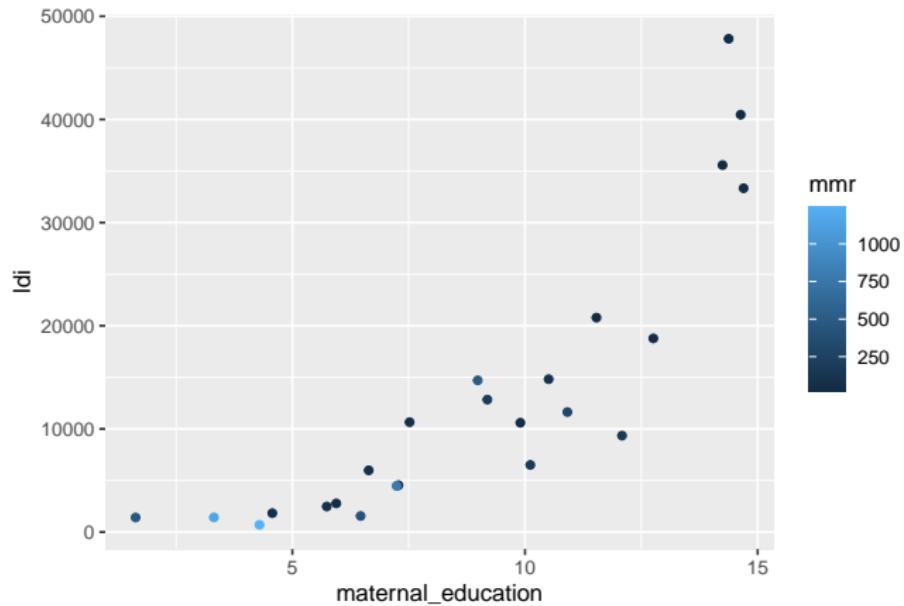
Note that ggplot conveniently makes a legend for you! In ggplot lingo, legends are called “scales”

EXAMPLE OF AESTHETIC MAPPING

In many cases, aesthetic mapping works for both continuous and categorical data:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi,  
                      color = mmr)) +  
  geom_point()
```

EXAMPLE OF AESTHETIC MAPPING

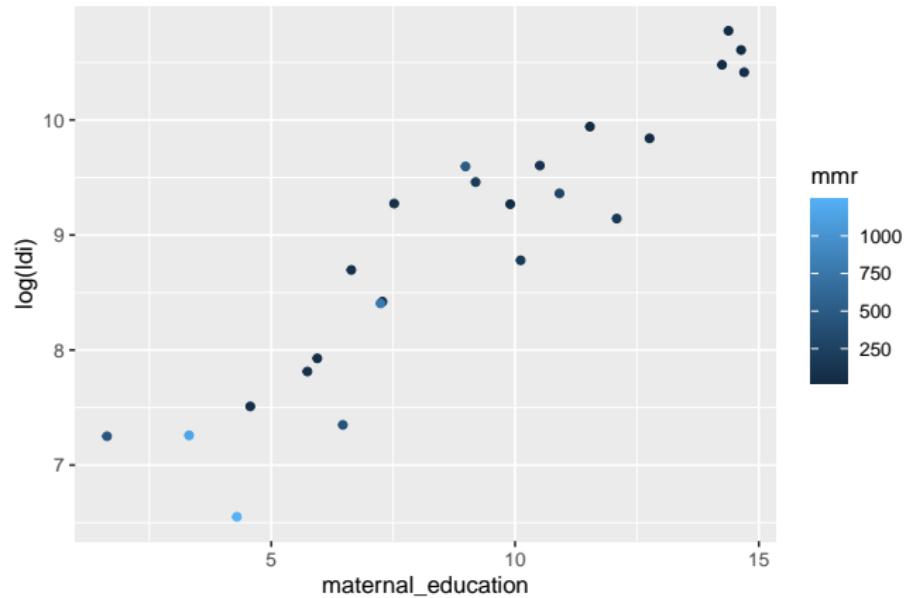


EXAMPLE OF AESTHETIC MAPPING

ggplot allows you to manipulate variables “on the fly”:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = log(ldi),  
                      color = mmr)) +  
  geom_point()
```

EXAMPLE OF AESTHETIC MAPPING



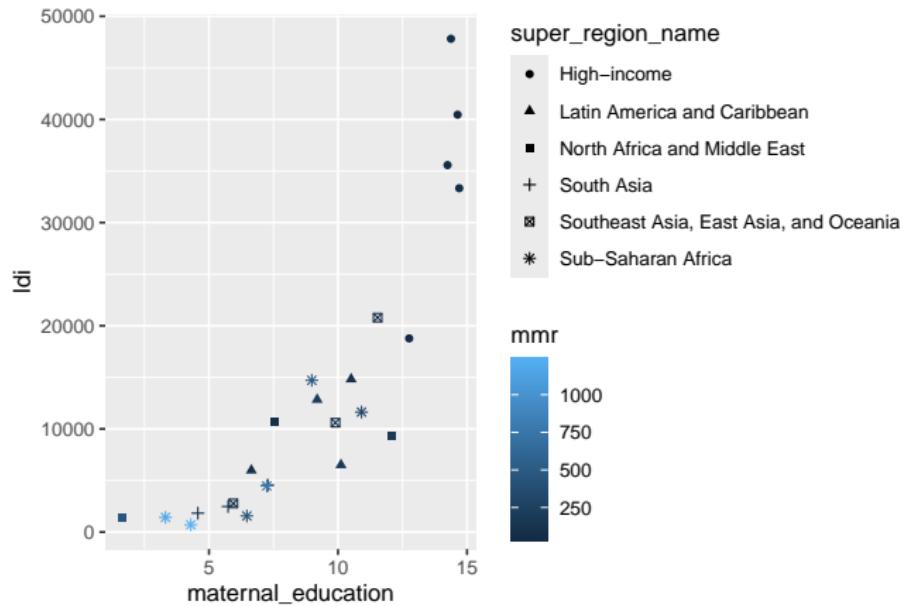
EXAMPLE OF AESTHETIC MAPPING

You can keep adding more aesthetics to add more information to your graph:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi,  
                      color = mmr,  
                      shape = super_region_name)) +  
  geom_point()
```

Note that not all aesthetics are meaningful for all geoms (e.g., linetype doesn't make sense if there are no lines in your graph)

EXAMPLE OF AESTHETIC MAPPING



WHAT ARE THE BUILDING BLOCKS OF A GGPLOT?

- ▶ **Aesthetics**
- ▶ **Geoms**
- ▶ **Facets**
- ▶ Positions
- ▶ Scales
- ▶ Labels
- ▶ Themes

GEOMS

ggplot “geoms” (geometries) are the different types of graphs you can make:

- ▶ `geom_point()` for scatter plots
- ▶ `geom_line()` for line graphs
- ▶ `geom_bar()` for bar graphs
- ▶ And more: `geom_histogram()`, `geom_violin()`,
`geom_boxplot()`, `geom_errorbar()`, `geom_ribbon()`,
`geom_segment()`, `geom_path()`, `geom_tile()`,
`geom_polygon()`, etc.

There are dozens of different geometries you can use for ggplot.

See the ggplot cheat sheet for the whole list:

<https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf>

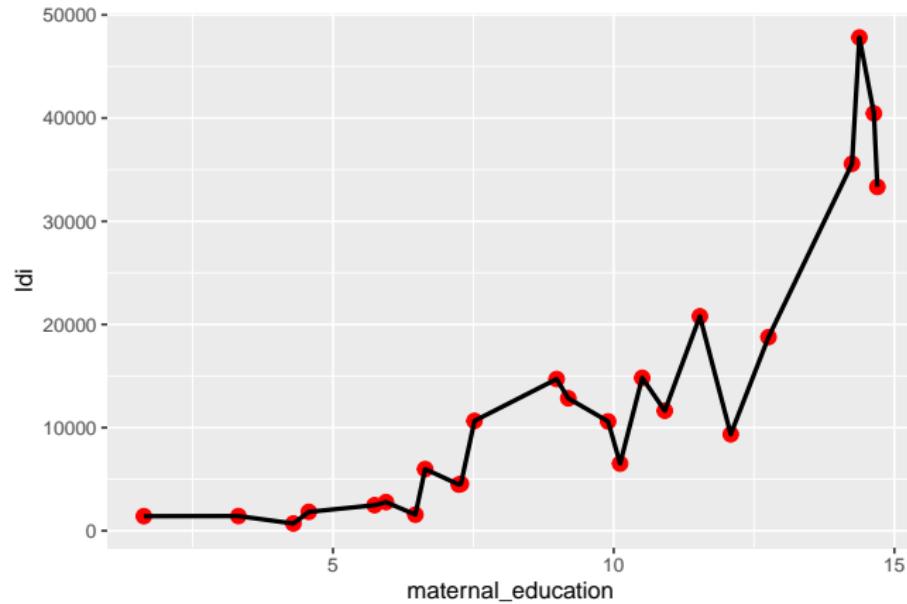


GEOMS

If you specify more than one geom, it “layers” them on top of each other

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi)) +  
  geom_point(size=3, color='red') +  
  geom_line(lwd=1)
```

GEOMS



Note: the order matters, it will layer geoms in order that they are written

GEOMS

Aesthetic arguments can also be provided directly to a geom in cases where you don't want them to map to some variable:

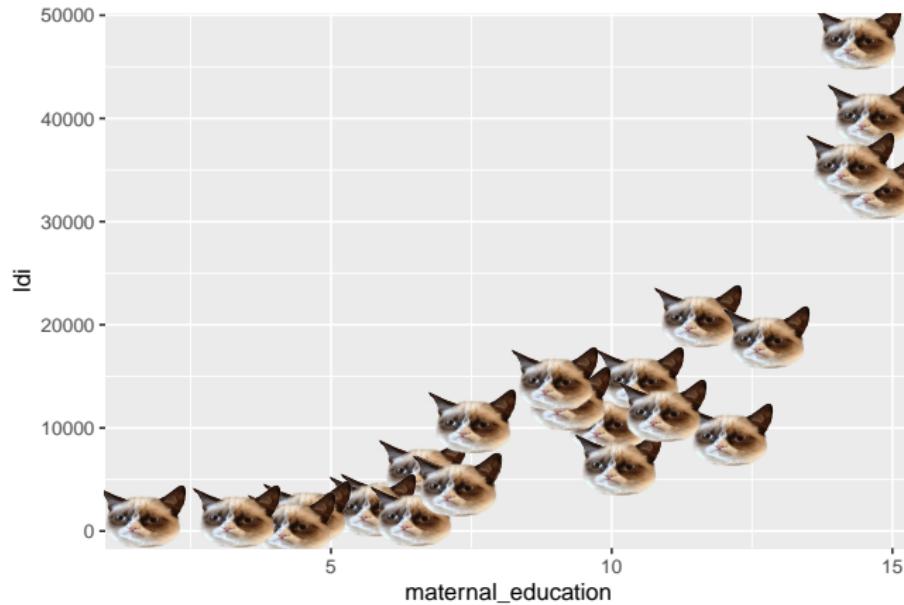
```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi)) +  
  geom_point(color = 'red',  
             size = 2,  
             alpha = .5)
```

GEOMS

Some libraries provide additional geoms:

```
library(ggcats)
plot <- ggplot(data = mmr_data,
                 aes(x = maternal_education,
                     y = ldi)) +
  geom_cat(cat='grumpy', size=3)
```

GEOMS



WHAT ARE THE BUILDING BLOCKS OF A GGPLOT?

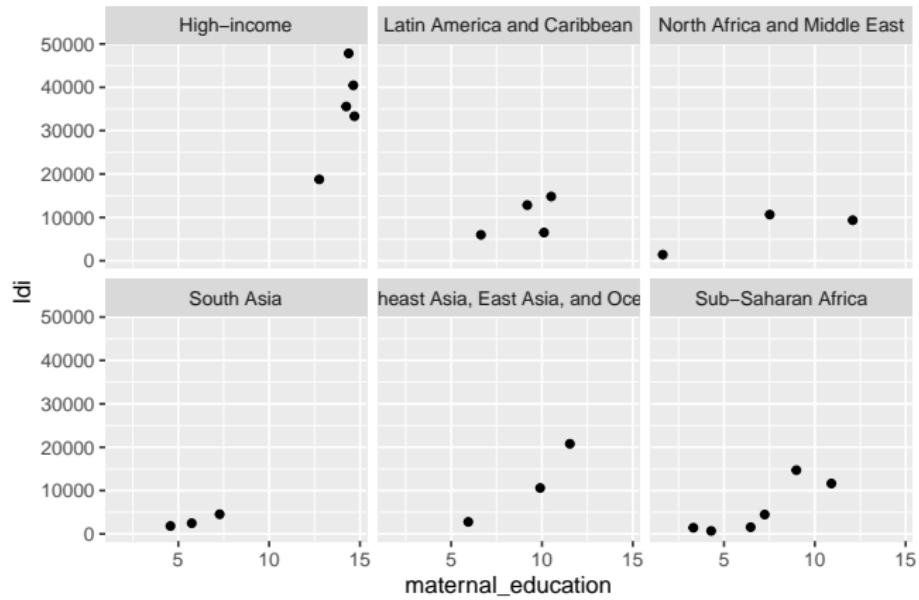
- ▶ **Aesthetics**
- ▶ **Geoms**
- ▶ **Facets**
- ▶ Positions
- ▶ Scales
- ▶ Labels
- ▶ Themes

FACETS

Facets allow you to incorporate more complexity into your graphs by adding multiple panels:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi)) +  
  geom_point() +  
  facet_wrap(~super_region_name)
```

FACETS



FACETS

`facet_wrap` stratifies and wraps. You can specify `nrow` and `ncol` to modify dimensions.

`facet_grid` forms a grid of panels based on row and column facetting variables. Example: `facet_grid(sex ~ age_group)` will create rows of panels based on sex and columns of panels based on age group.

WHAT ARE THE BUILDING BLOCKS OF A GGPLOT?

- ▶ **Aesthetics**
- ▶ **Geoms**
- ▶ **Facets**
- ▶ Positions
- ▶ Scales
- ▶ Labels
- ▶ Themes

POSITIONS

ggplot lets you modify where geoms appear relative to each other, using position functions:

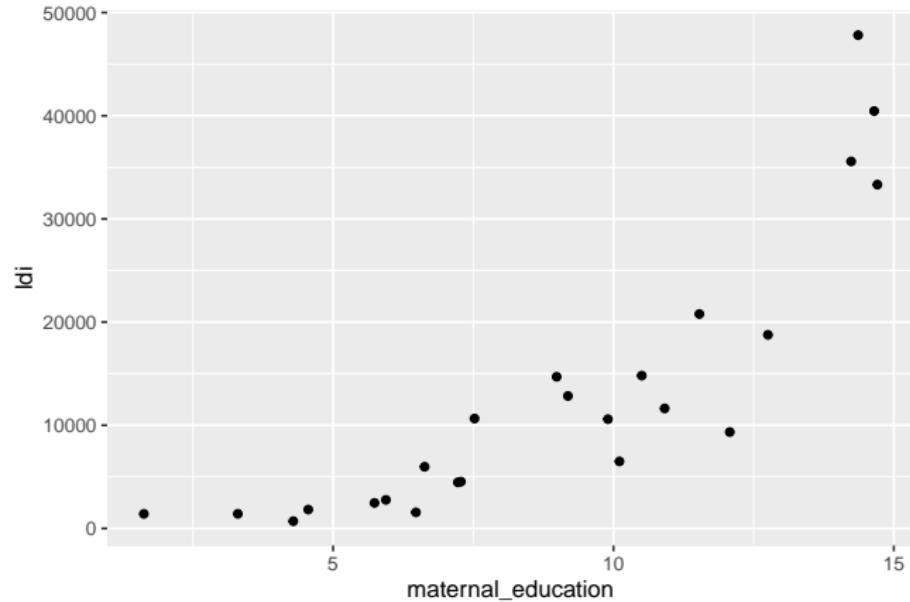
- ▶ `position_jitter()` randomly displaces points (usually just for `geom_point`)
- ▶ `position_dodge()` automatically (tries to) shift to avoid overlap
- ▶ `position_stack()` stack, or add together geoms (usually just for `geom_bar`)
- ▶ `position_fill()` rescale the y-axis so the geoms sum to 100% (usually just for `geom_bar`)

POSITIONS

`position_jitter` randomly displaces points (usually just for `geom_point`)

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi)) +  
  geom_point(position='jitter')
```

POSITIONS



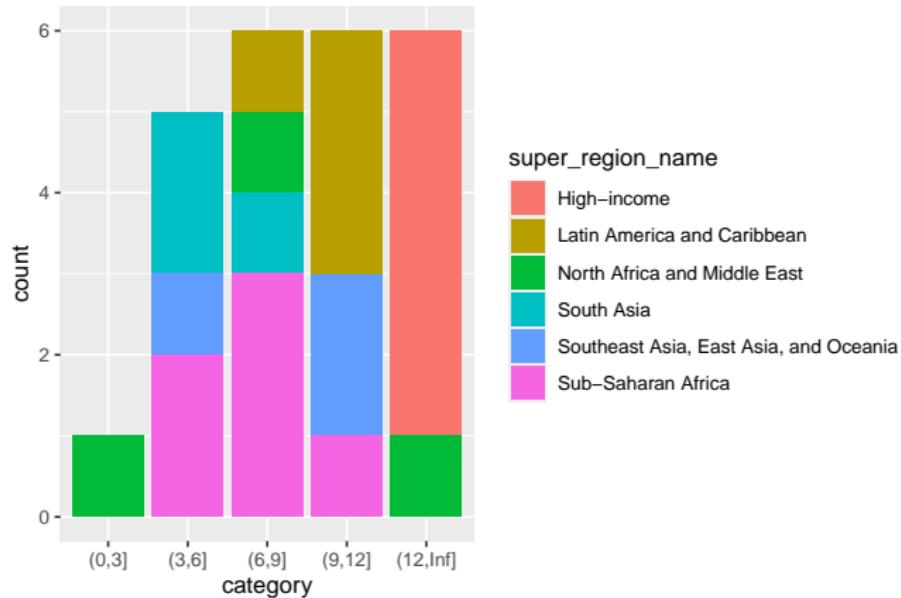
It's built right into the `geom_point()` function for convenience.

POSITIONS

position_stack is the default for geom_bar for factor variables:

```
mmr_data$category <- cut(mmr_data$maternal_education,  
                           breaks=c(0, 3, 6, 9, 12, Inf))  
plot <- ggplot(data = mmr_data,  
                 aes(x = category,  
                     fill = super_region_name)) +  
  geom_bar()
```

POSITIONS

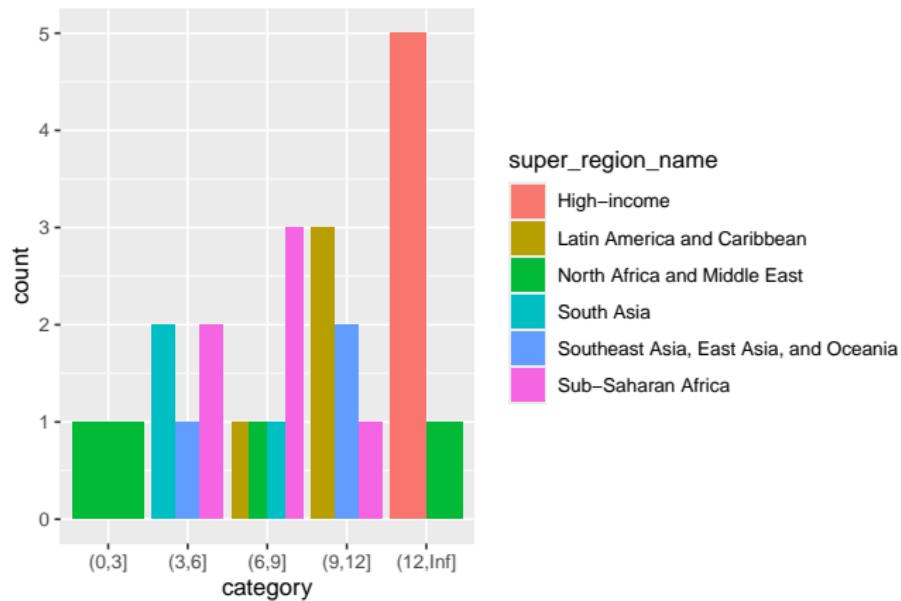


POSITIONS

position_dodge would put the bars side-by-side:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = category,  
                     fill = super_region_name)) +  
      geom_bar(position = 'dodge')
```

POSITIONS



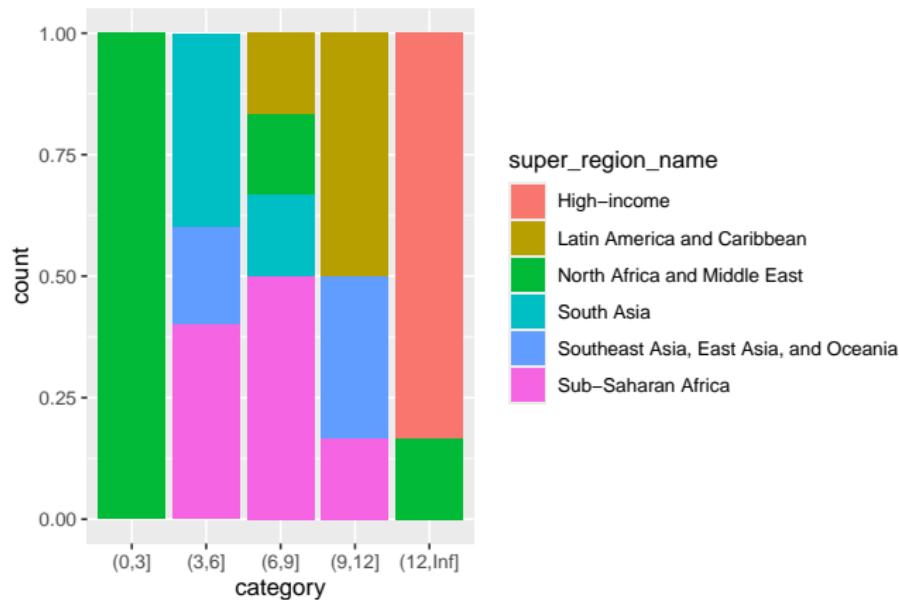
It's built right into the `geom_bar()` function for convenience.

POSITIONS

`position_fill` makes the bars sum to 100%:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = category,  
                      fill=super_region_name)) +  
  geom_bar(position='fill')
```

POSITIONS



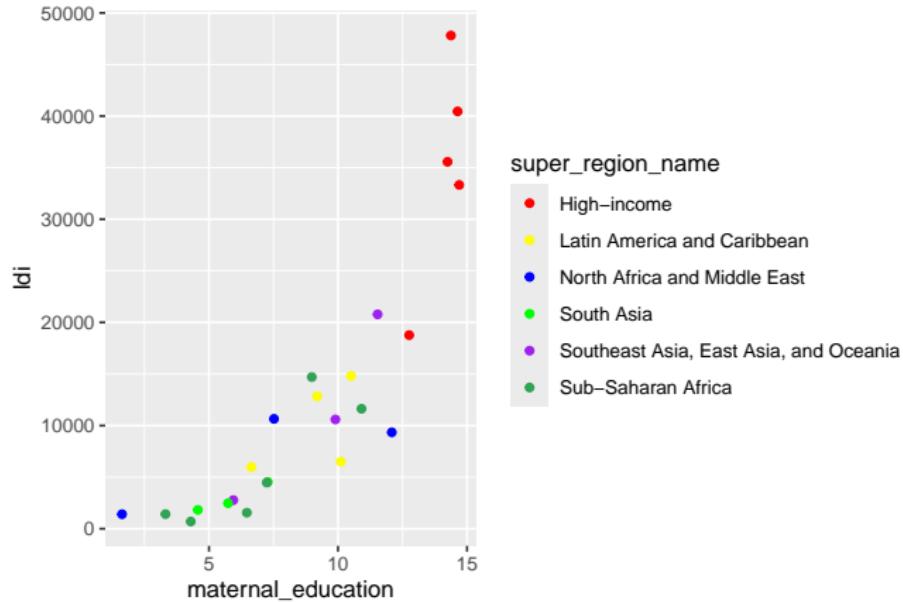
It's built right into the `geom_bar()` function for convenience

SCALES

You can also modify the “scales” (i.e., legends) to customize aesthetic mapping:

```
plot <- ggplot(data = mmr_data,
                 aes(x = maternal_education,
                     y = ldi,
                     color = super_region_name)) +
  geom_point() +
  scale_color_manual(values = c('red', 'yellow',
                                'blue', 'green',
                                'purple', '#31a354'))
```

SCALES



Every aesthetic (fill, color, shape, linetype) has corresponding
scale_* function (`scale_fill_manual`, `scale_color_manual`)
etc.)

SCALES

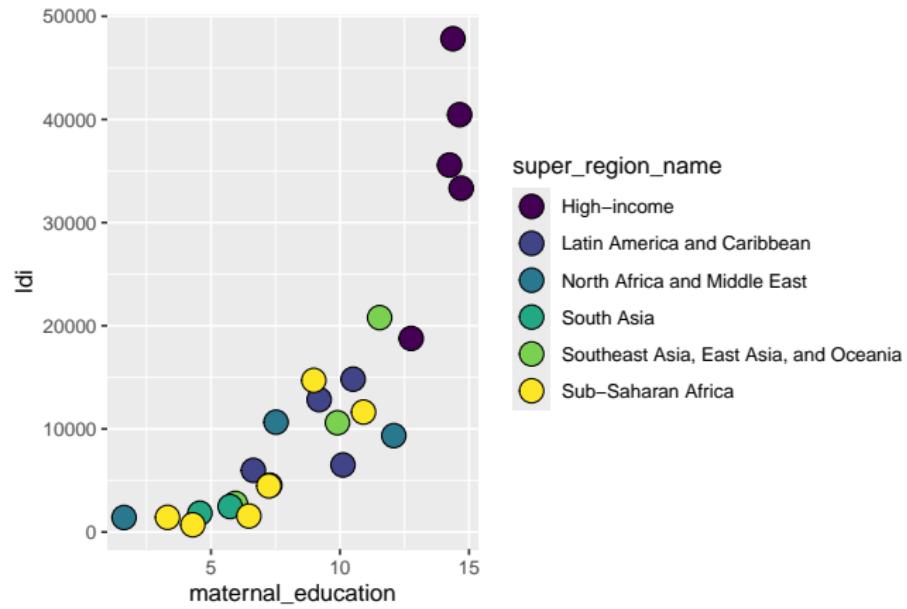
I recommend using the viridis library for colors:

```
library(viridis)

## Loading required package: viridisLite

plot <- ggplot(data = mmr_data,
                aes(x = maternal_education,
                    y = ldi,
                    fill = super_region_name)) +
  geom_point(shape=21, size=5) +
  scale_fill_viridis_d()
```

SCALES

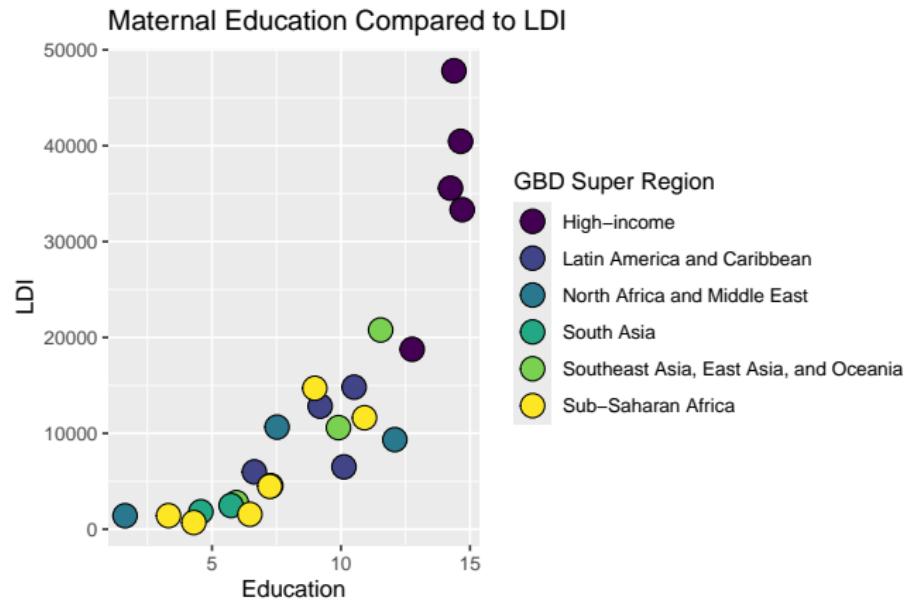


LABELS

Titles for everything can be added with the `labs()` function:

```
plot <- ggplot(data = mmr_data,
                 aes(x = maternal_education,
                     y = ldi,
                     fill = super_region_name)) +
  geom_point(shape=21, size=5) +
  scale_fill_viridis_d() +
  labs(title = 'Maternal Education Compared to LDI',
       y = 'LDI',
       x = 'Education',
       fill = 'GBD Super Region')
```

LABELS

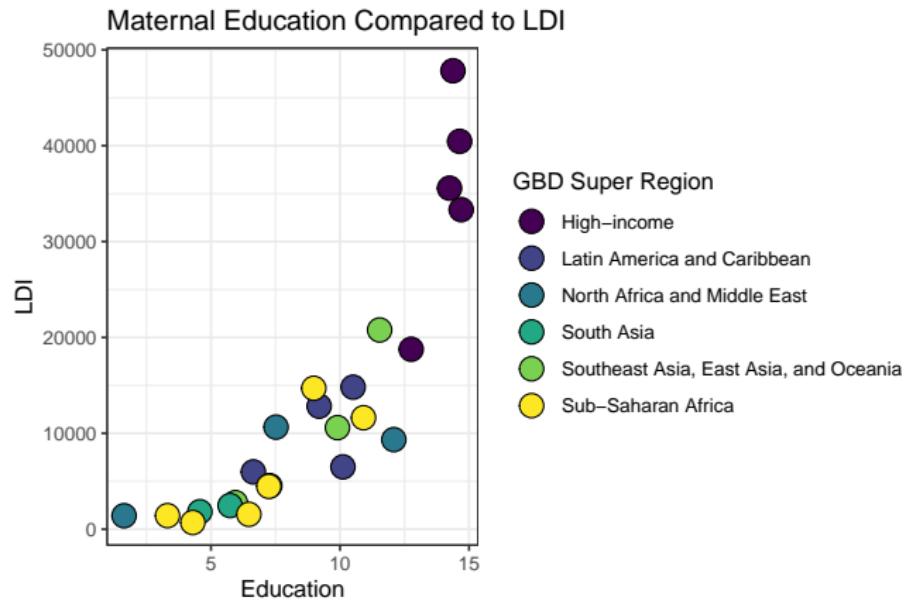


THEMES

ggplot also comes with handy “themes”, or preset options:

```
plot <- ggplot(data = mmr_data,  
                 aes(x = maternal_education,  
                      y = ldi,  
                      fill = super_region_name)) +  
  geom_point(shape=21, size=5) +  
  scale_fill_viridis_d() +  
  labs(title = 'Maternal Education Compared to LDI',  
       y = 'LDI',  
       x = 'Education',  
       fill = 'GBD Super Region') +  
  theme_bw()
```

THEMES

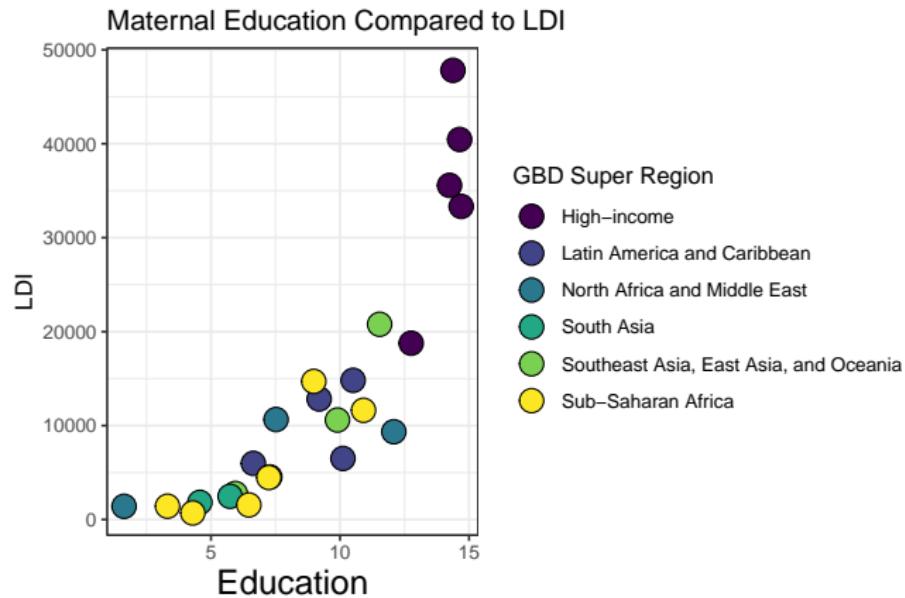


THEMES

Themes also allow you to adjust visuals that are not aesthetics:

```
plot <- ggplot(data = mmr_data,
                 aes(x = maternal_education,
                     y = ldi,
                     fill = super_region_name)) +
  geom_point(shape=21, size=5) +
  scale_fill_viridis_d() +
  labs(title = 'Maternal Education Compared to LDI',
       y = 'LDI',
       x = 'Education',
       fill = 'GBD Super Region') +
  theme_bw() +
  theme(axis.title.x = element_text(size=16))
```

THEMES



RESHAPING

`ggplot2` is designed to work with data shaped such that each desired aesthetic is mapped to **one** variable. If your data is not shaped this way, it's almost always easier to reshape the data than to try and make `ggplot2` work with original data structure.

RESHAPING

For example, if you want to plot the number of Ebola deaths by age group for both males and females, this is an inconvenient data structure since there are separate columns for deaths among males and females:

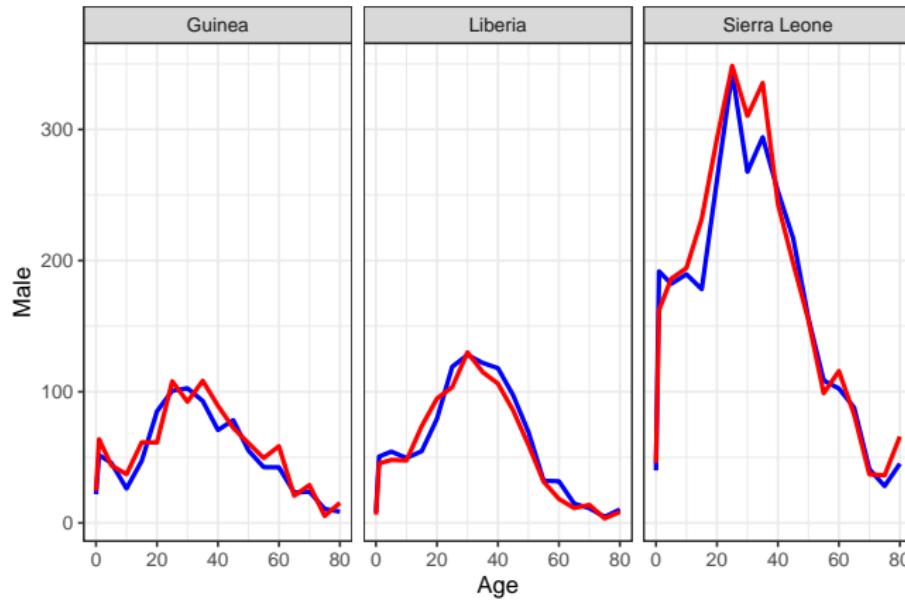
```
## Key: <Country, Age>
##      Country    Age Female   Male
##      <char> <int>  <num> <num>
## 1: Guinea     0    24.5  21.9
## 2: Guinea     1    63.8  51.7
## 3: Guinea     5    44.0  45.8
## 4: Guinea    10    37.1  26.2
## 5: Guinea    15    61.5  47.4
## 6: Guinea    20    61.3  85.1
```

RESHAPING

One option is to just add different geoms for each variable:

```
plot <- ggplot(data = wide_data,
                 aes(x = Age)) +
  geom_line(aes(y = Male), color = 'blue', lwd=1) +
  geom_line(aes(y = Female), color = 'red', lwd=1) +
  facet_wrap(~Country) +
  theme_bw()
```

RESHAPING



RESHAPING

A better option is to reshape **long** before attempting to plot these data:

```
long_data <- melt(wide_data,
                    id.vars = c("Country", "Age"),
                    value.name = "Deaths",
                    variable.name = "Sex")
head(long_data, 3)
```

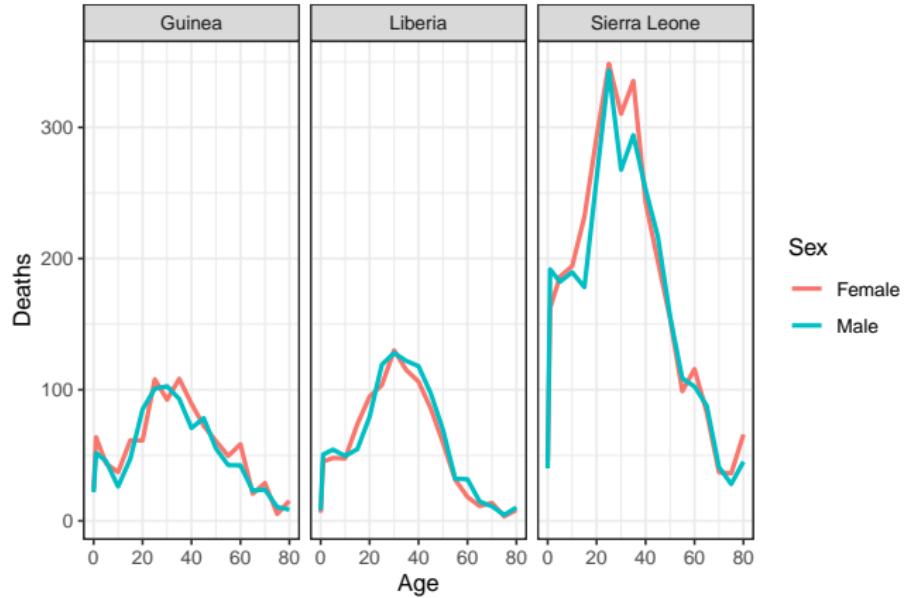
```
##      Country    Age     Sex Deaths
##      <char> <int> <fctr>  <num>
## 1: Guinea      0 Female   24.5
## 2: Guinea      1 Female   63.8
## 3: Guinea      5 Female   44.0
```

RESHAPING

```
plot <- ggplot(data = long_data,
                 aes(x = Age,
                     y = Deaths,
                     color = Sex)) +
  geom_line(lwd=1) +
  facet_wrap(~Country) +
  theme_bw()
```



RESHAPING



SAVING PLOTS

You can save your plot directly into a pdf or image file.

Use “`ggsave`” and specify an output file type:

```
ggsave('images/my_plot.pdf')
```

```
## Saving 6 x 4 in image
```

ADDITIONAL PACKAGES

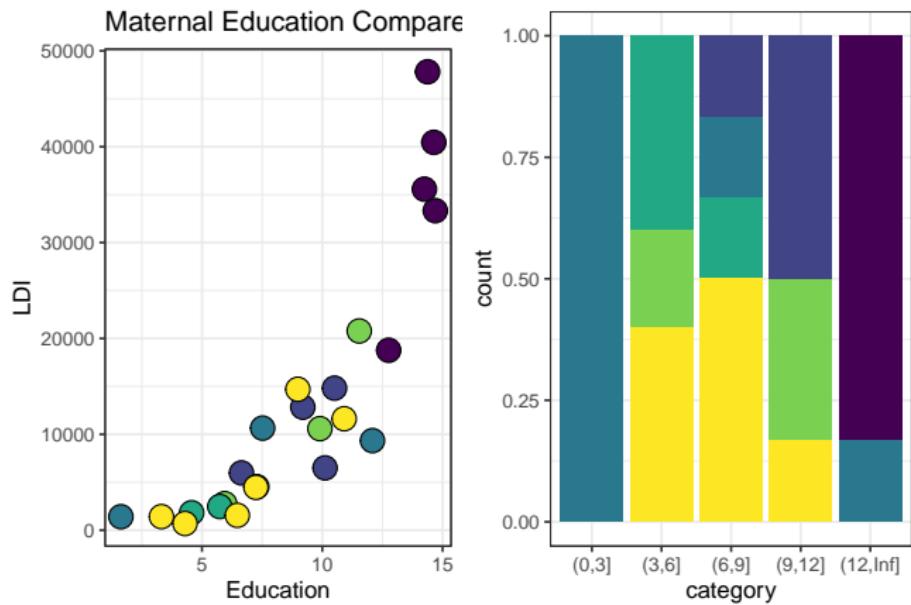
`ggplot2` has become so popular that other users have started writing add-ons to it:

- ▶ **cowplot** - plot tables and arrange multiple plots together
- ▶ **ggrepel** - label points nicely
- ▶ **viridis** - easy-to-use color schemes
- ▶ **ggthemes** - more themes, preset colors

ADDITIONAL PACKAGE: COWPLOT

```
library(cowplot)
plot <- ggplot(data = mmr_data,
    aes(x = maternal_education,
        y = ldi,
        fill = super_region_name)) +
  geom_point(shape=21, size=5) +
  scale_fill_viridis_d(guide='none') +
  labs(title = 'Maternal Education Compared to LDI',
      y = 'LDI',
      x = 'Education',
      fill = 'GBD Super Region') +
  theme_bw()
plot2 <- ggplot(data = mmr_data,
    aes(x = category,
        fill = super_region_name)) +
  geom_bar(position='fill') +
  scale_fill_viridis_d(guide='none') +
  theme_bw()
```

ADDITIONAL PACKAGE: COWPLOT



ADDITIONAL PACKAGE: GGREPEL

```
library(ggrepel)
plot <- ggplot(data = mmr_data,
                aes(x = maternal_education,
                    y = ldi,
                    fill = super_region_name,
                    label = location_name)) +
  geom_point(shape=21, size=5) +
  geom_text_repel() +
  scale_fill_viridis_d(guide='none') +
  labs(title = 'Maternal Education Compared to LDI',
       y = 'LDI',
       x = 'Education',
       fill = 'GBD Super Region') +
  theme_bw()
```

ADDITIONAL PACKAGE: GGREPEL

