

DATA VISUALIZATION

Gerko Vink 

g.vink@uu.nl

Methodology & Statistics @ Utrecht University

12 Jun 2025

DISCLAIMER

I owe a debt of gratitude to many people as the thoughts and code in these slides are the process of years-long development cycles and discussions with my team, friends, colleagues and peers. When someone has contributed to the content of the slides, I have credited their authorship.

Images are either directly linked, or generated with StableDiffusion or DALL-E. That said, there is no information in this presentation that exceeds legal use of copyright materials in academic settings, or that should not be part of the public domain.

Warning

You **may use** any and all content in this presentation - including my name - and submit it as input to generative AI tools, with the following **exception**:

- You must ensure that the content is not used for further training of the model

SLIDE MATERIALS AND SOURCE CODE

Materials

- lecture slides on Moodle
- course page: www.gerkovink.com/sur
- source: github.com/gerkovink/sur

RECAP

Gisteren hebben we deze onderwerpen behandeld:

- Ontbrekende waarden identificeren
- Synthetische imputaties maken

SOME ADDITION

I've showed you how to draw a *bootstrap* sample from a dataset in the following way:

```
1 which_rows <- sample(1:nrow(mice::boys), size = 100, replace = TRUE)
2 mice::boys[which_rows, ]
```

But, there is a much easier way:

```
1 mice::boys %>%
2   slice_sample(n = 100, replace = TRUE) %>%
3   glimpse()
```

Rows: 100

Columns: 9

```
$ age <dbl> 1.790, 1.086, 13.656, 13.300, 0.188, 17.911, 12.501, 2.502, 0.101, ...
$ hgt <dbl> 90.2, 74.5, 146.9, 165.5, 58.5, 181.2, 170.5, 91.1, 55.8, 59.5, 17...
$ wgt <dbl> 14.50, 8.92, 34.70, 41.90, 6.03, 86.80, 53.40, 15.80, 5.06, 5.13, ...
$ bmi <dbl> 17.82, 16.07, 16.07, 15.29, 17.61, 26.43, 18.36, 19.03, 16.25, 14.0...
$ hc <dbl> 51.2, 45.1, 55.1, 55.6, 41.5, 58.3, 56.4, 49.7, 38.5, 38.0, 57.0, ...
$ gen <ord> NA, NA, G2, G2, NA, G5, G2, NA, NA, NA, G5, G2, NA, NA, NA, G4, NA...
$ phb <ord> NA, NA, P3, P2, NA, P5, P1, NA, NA, NA, P4, P1, NA, NA, NA, P5, NA...
$ tv <int> NA, NA, NA, 5, NA, 15, 6, NA, NA, NA, 20, 2, NA, NA, NA, 25, NA, N...
$ reg <fct> east, west, city, east, east, north, south, south, east, west, sou...
```

TODAY

Vandaag behandelen we de volgende onderwerpen:

- Basisplots: histogrammen, scatterplots en boxplots
- Geavanceerde plots met ggplot2
- Aanpassen van grafieken voor publicatie
- Exporteren van grafieken en resultaten

WE USE THE FOLLOWING PACKAGES

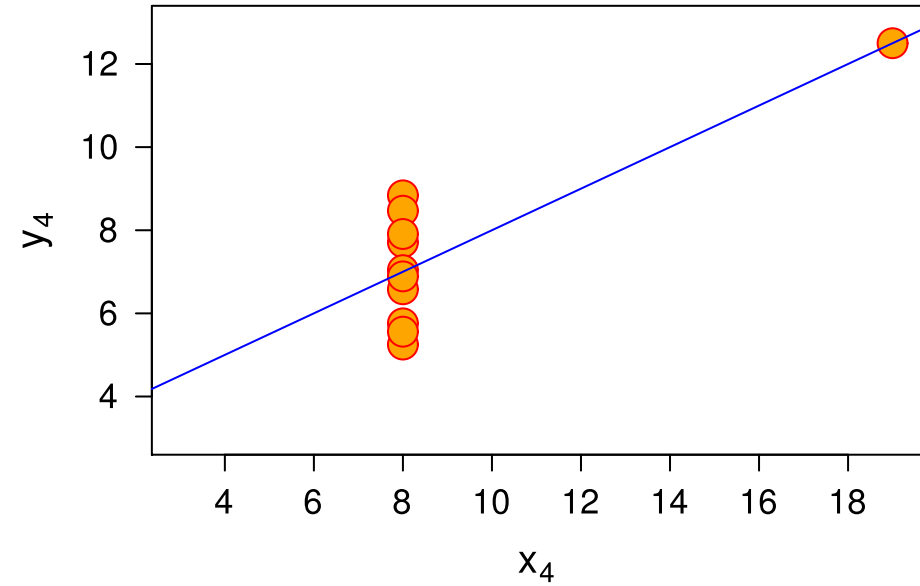
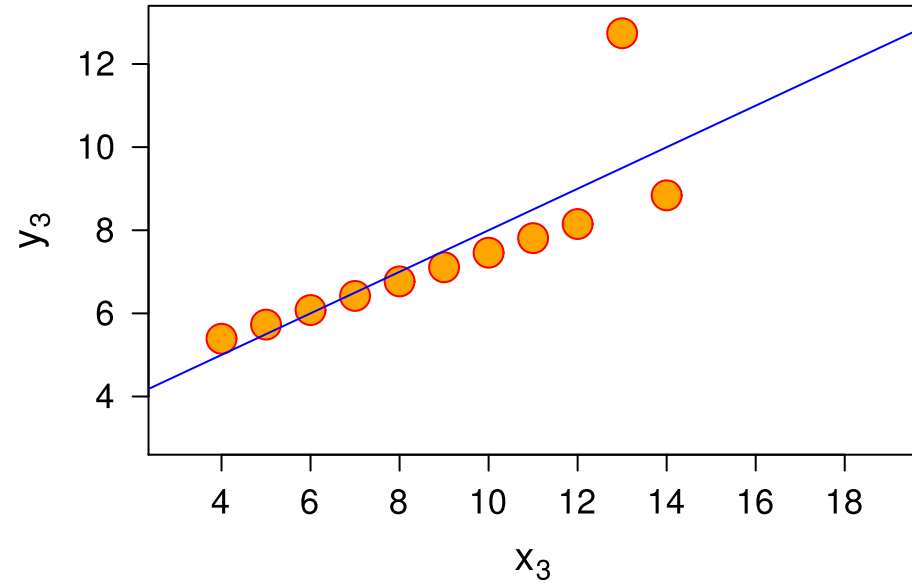
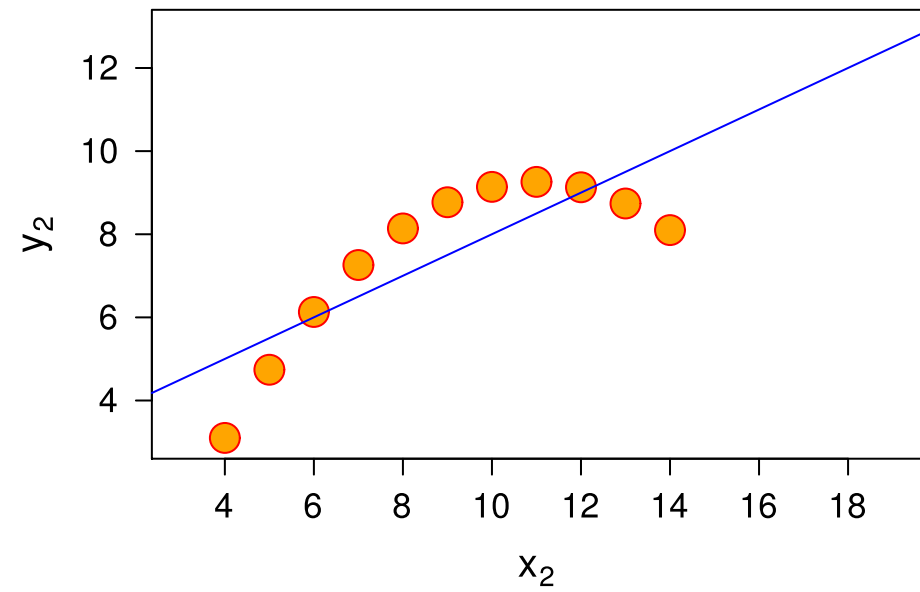
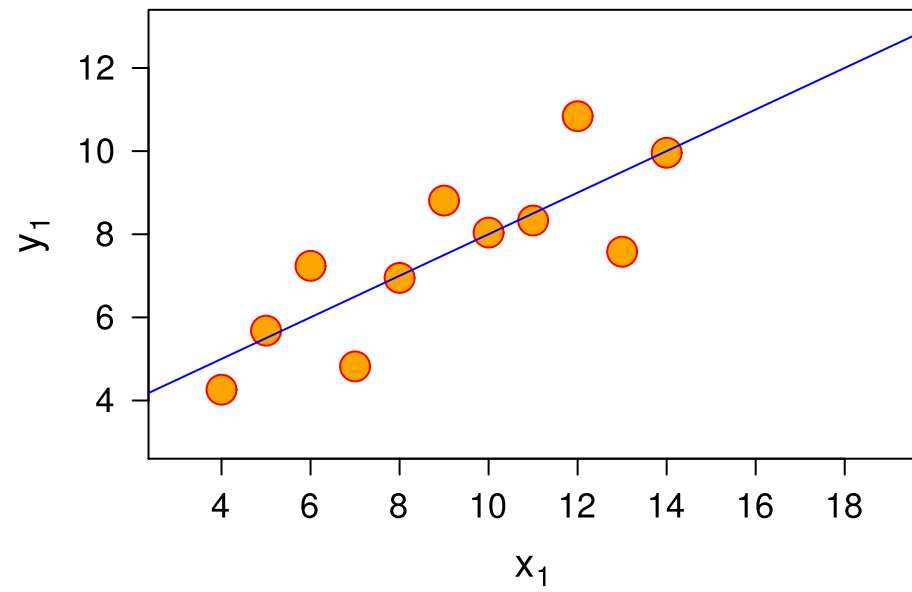
```
1 library(mice)      # Boys dataset
2 library(dplyr)     # Data manipulation
3 library(magrittr)  # Pipes
4 library(ggplot2)   # Plotting suite
```

WHY VISUALISE?

- We can process a lot of information quickly with our eyes
- Plots give us information about
 - Distribution / shape
 - Irregularities
 - Assumptions
 - Intuitions
- Summary statistics, correlations, parameters, model tests, p -values do not tell the whole story

ALWAYS PLOT YOUR DATA!

WHY VISUALISE?

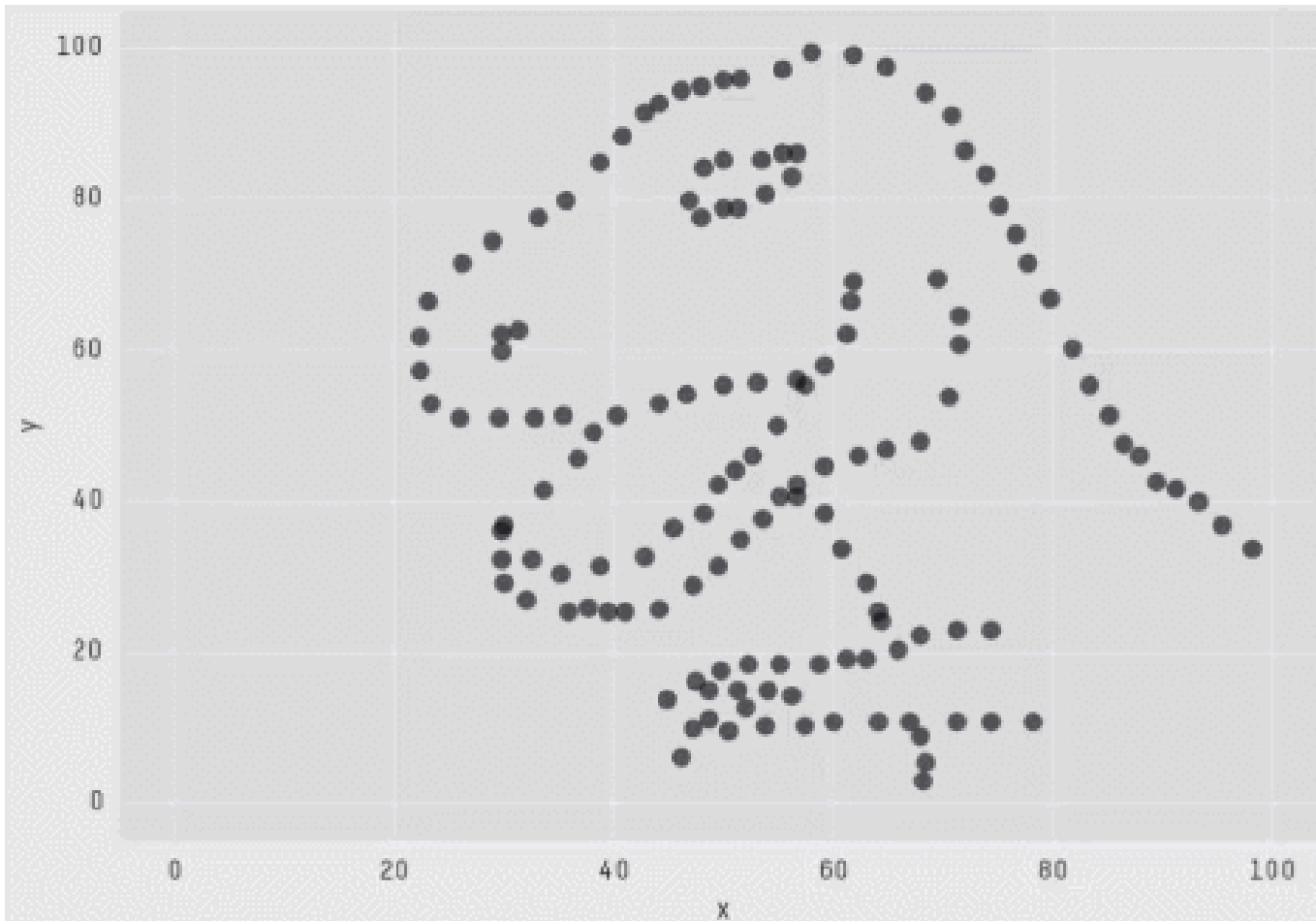


Source: Anscombe, F. J. (1973). "Graphs in Statistical Analysis". *American Statistician*. 27 (1): 17-21.

Gerko Vink @ Anton de Kom Universiteit, Paramaribo



WHY VISUALISE?

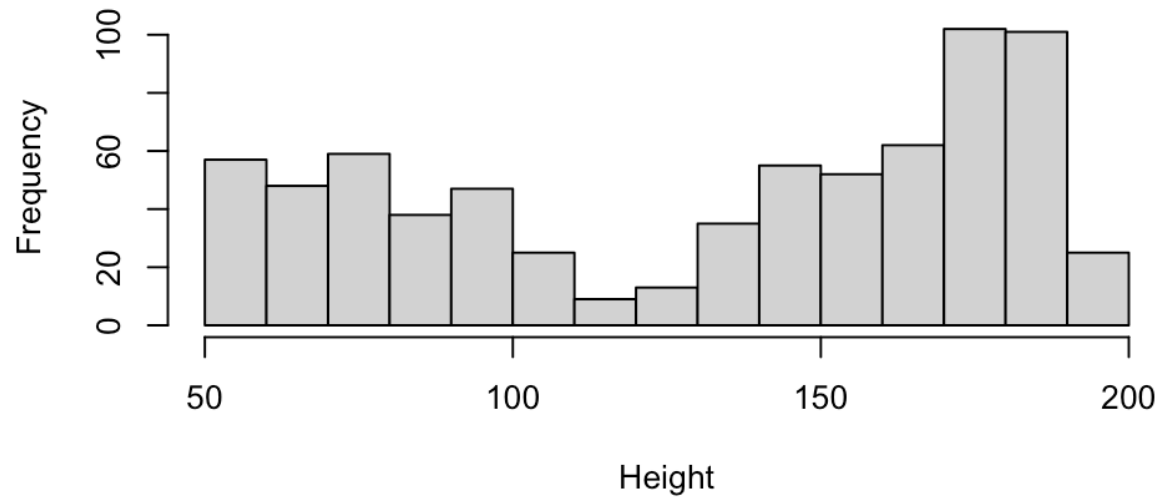


X Mean: 54.2659224
Y Mean: 47.8313999
X SD : 16.7649829
Y SD : 26.9342120
Corr. : -0.0642526

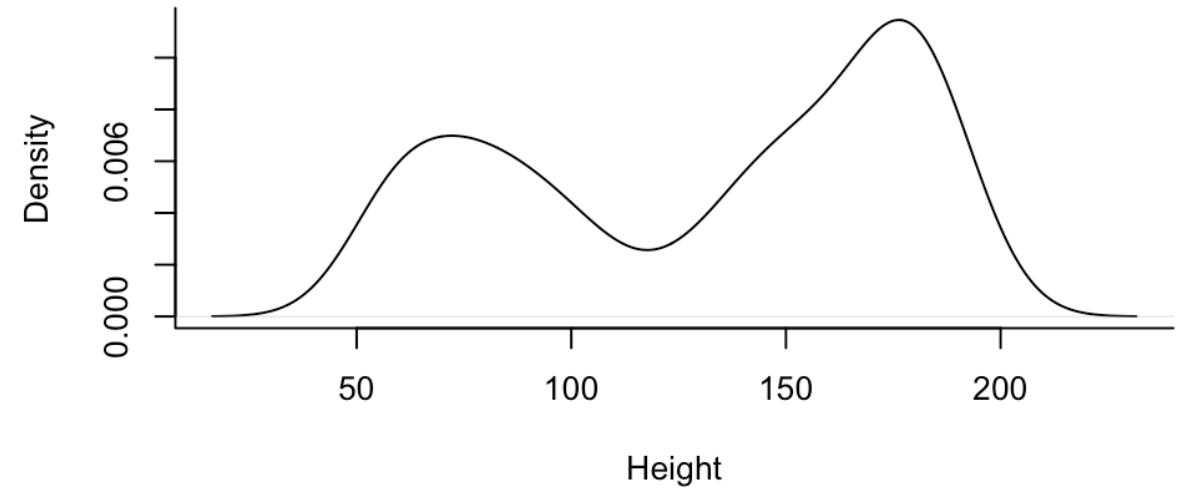
Source: <https://www.autodeskresearch.com/publications/samestats>

BASE R PLOTS

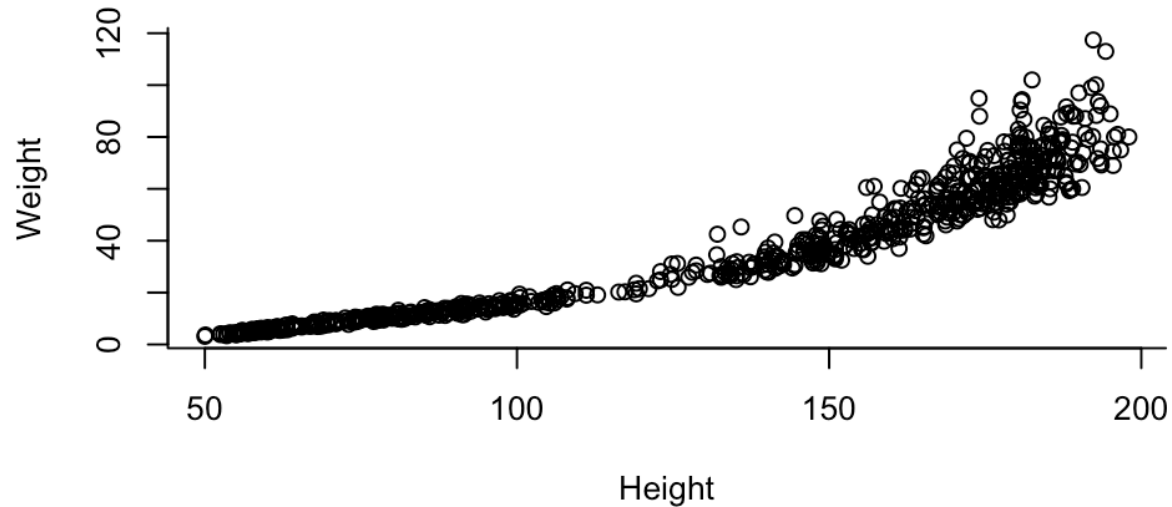
Histogram



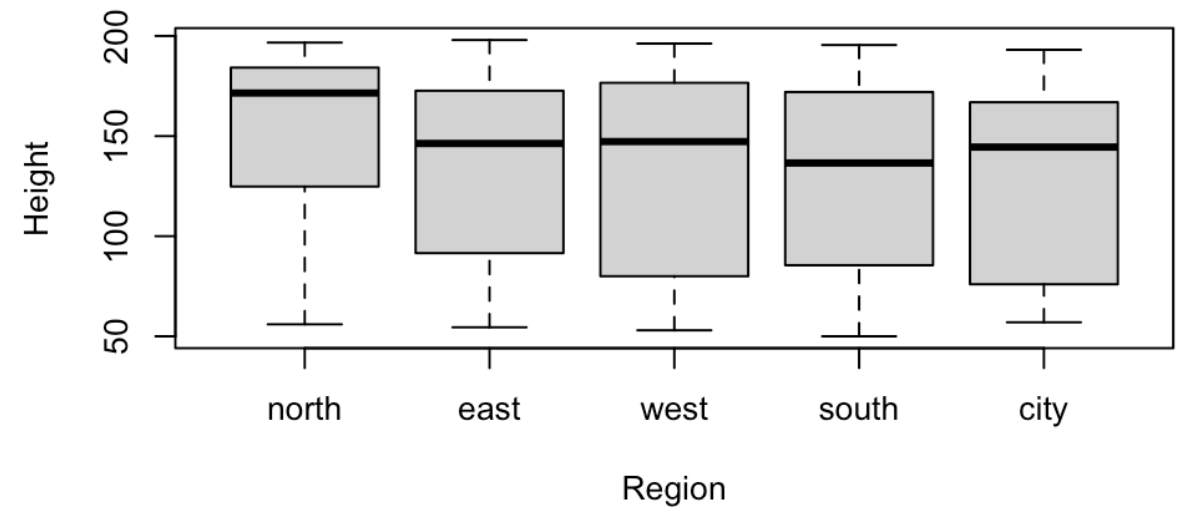
Density plot



Scatter plot

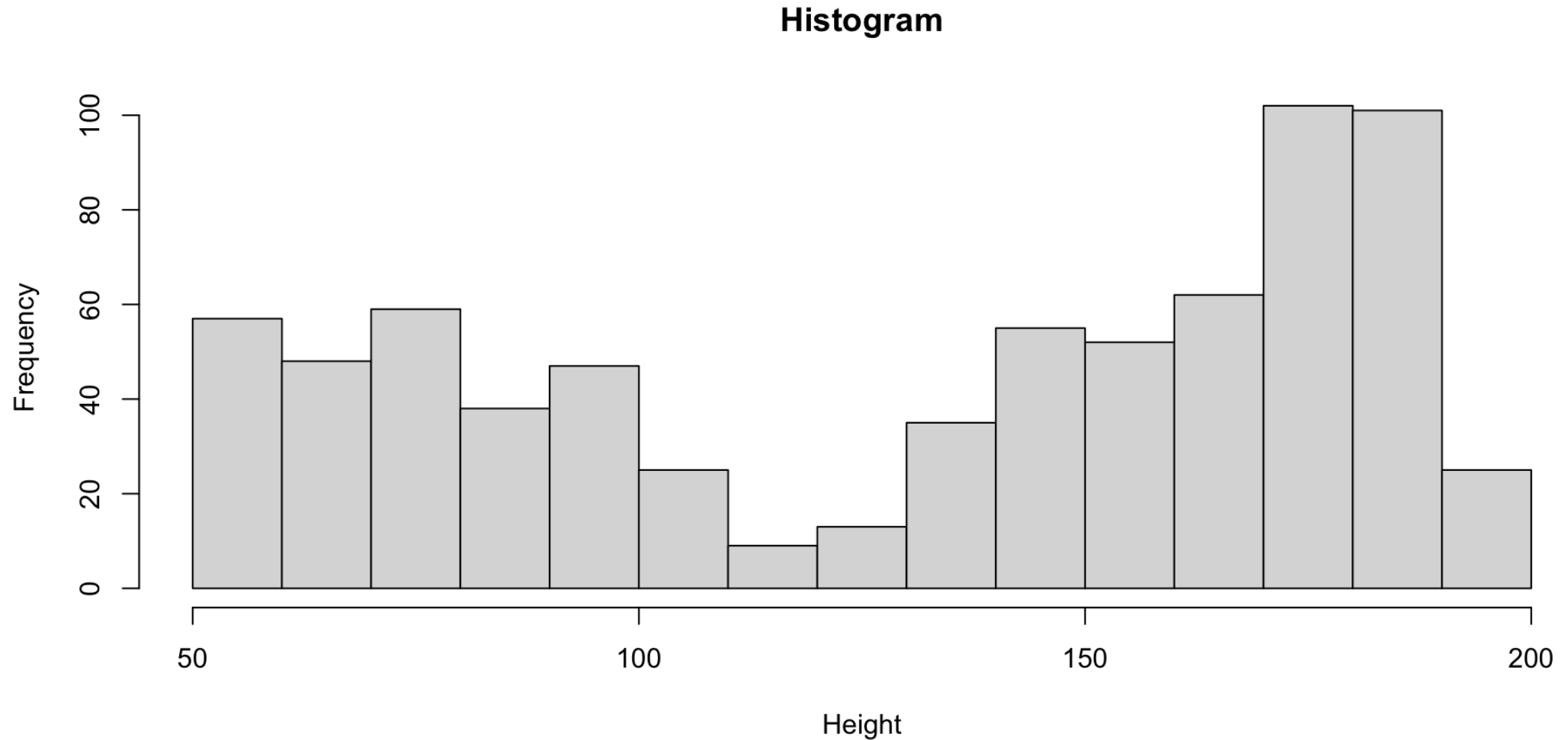


Boxplot



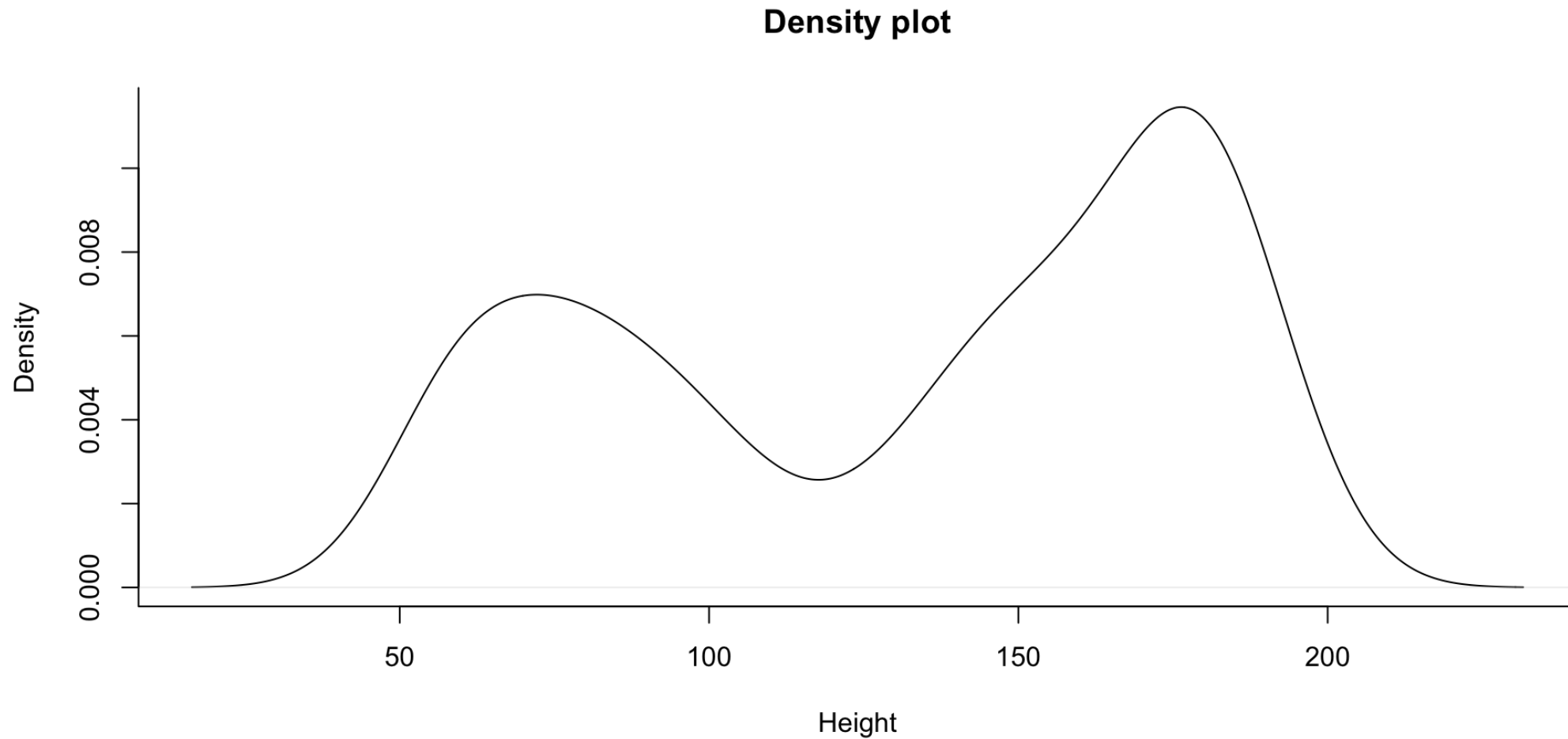
HISTOGRAM

```
1 hist(boys$hgt, main = "Histogram", xlab = "Height")
```



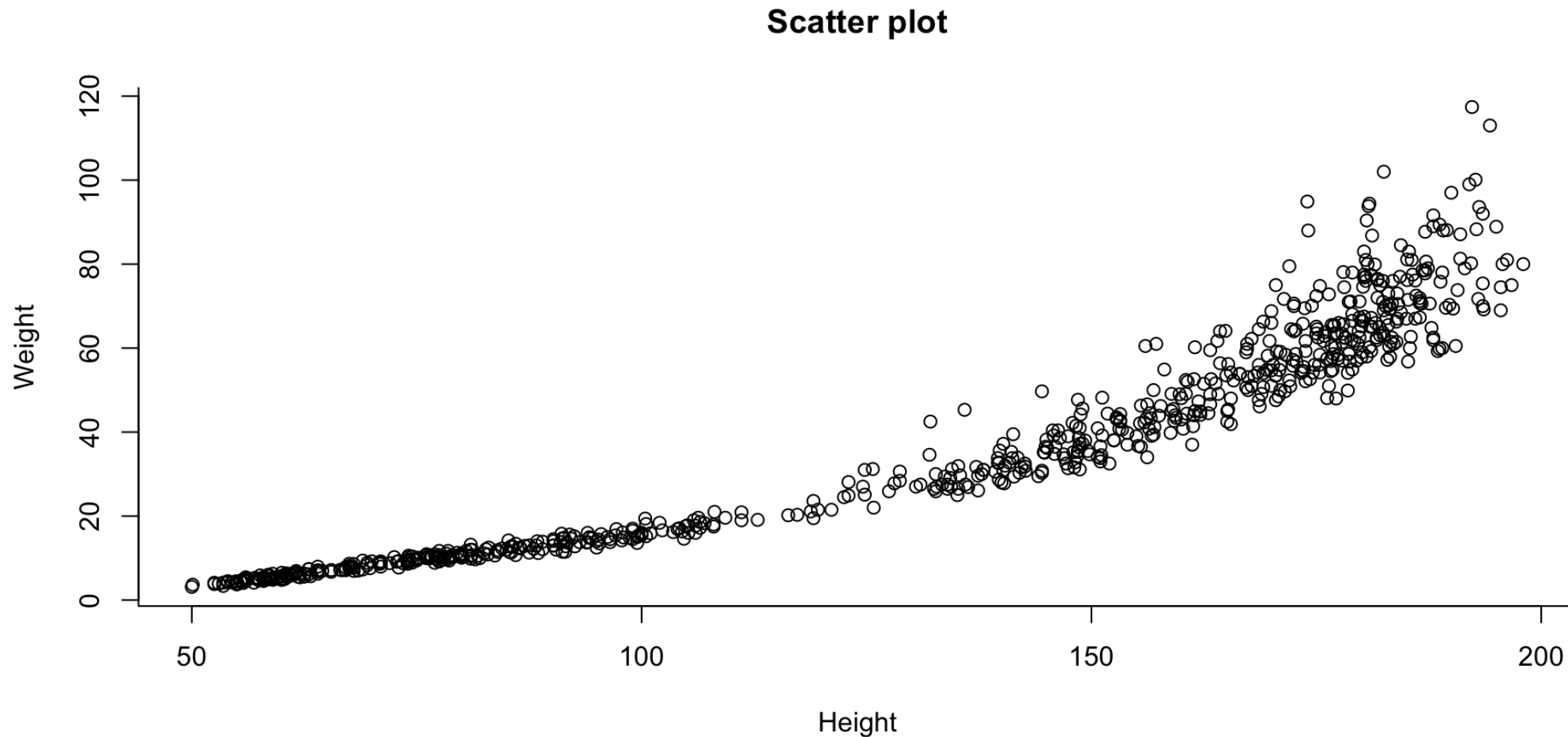
DENSITY

```
1 dens <- density(boys$hgt, na.rm = TRUE)  
2 plot(dens, main = "Density plot", xlab = "Height", bty = "L")
```



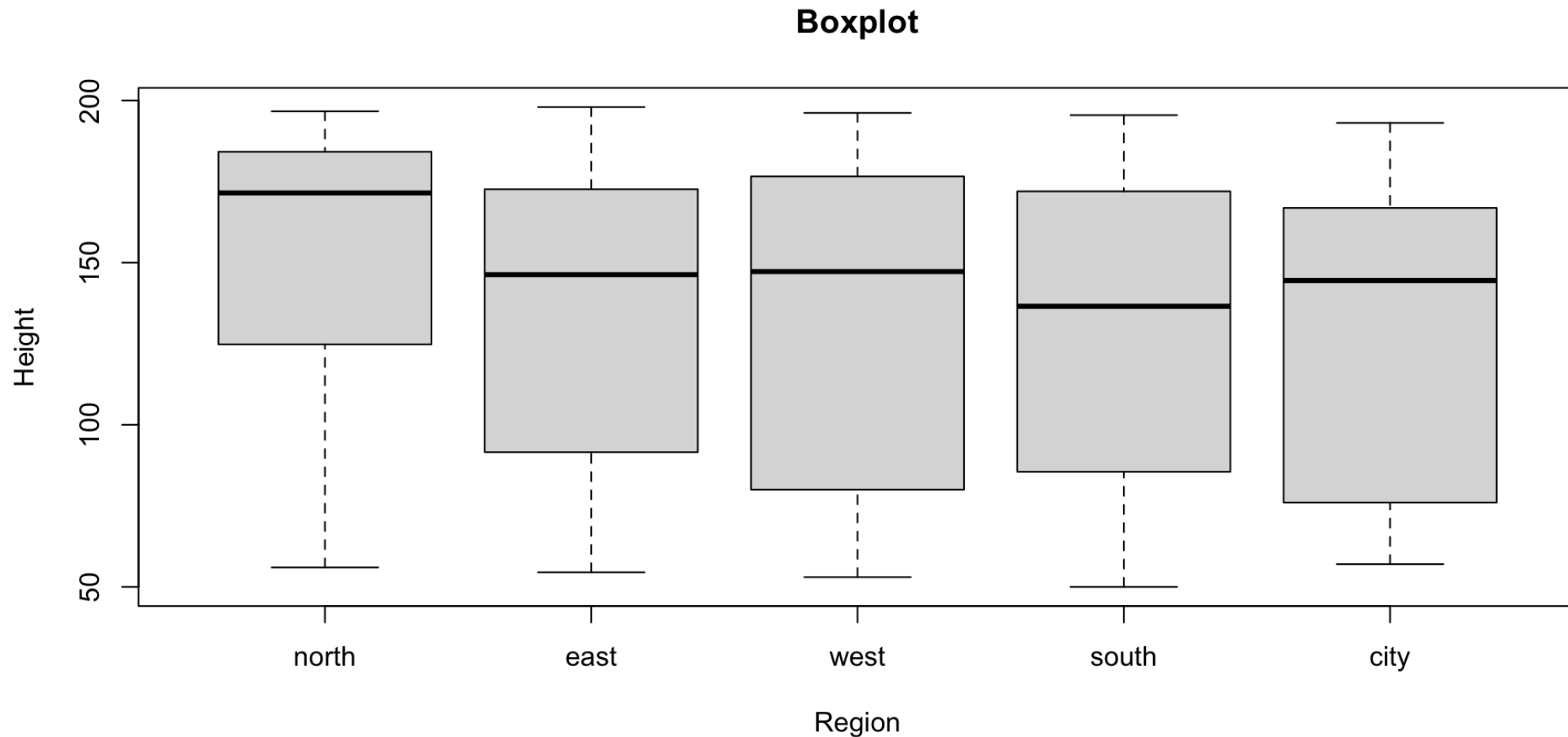
SCATTER PLOT

```
1 plot(x = boys$hgt, y = boys$wgt, main = "Scatter plot",  
2      xlab = "Height", ylab = "Weight", bty = "L")
```



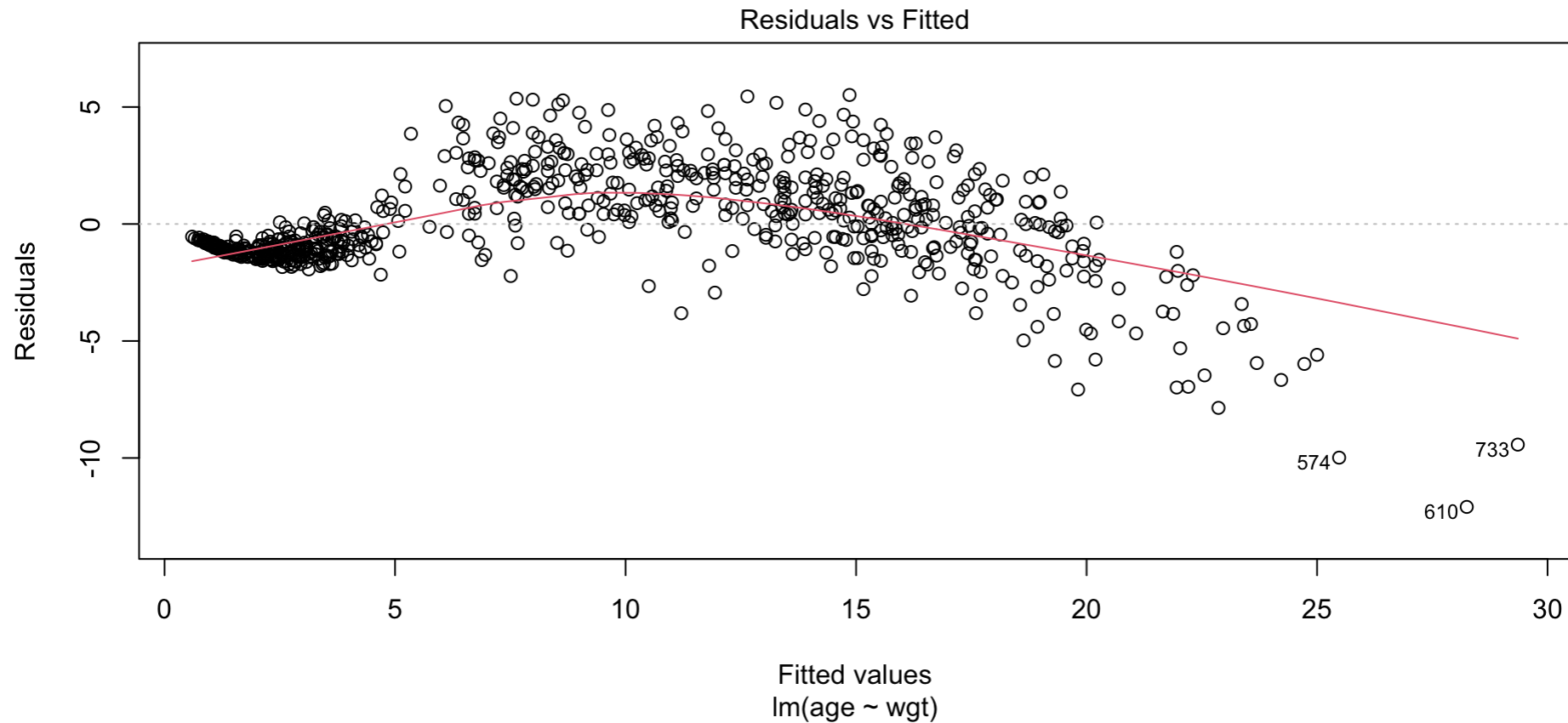
BOX PLOT

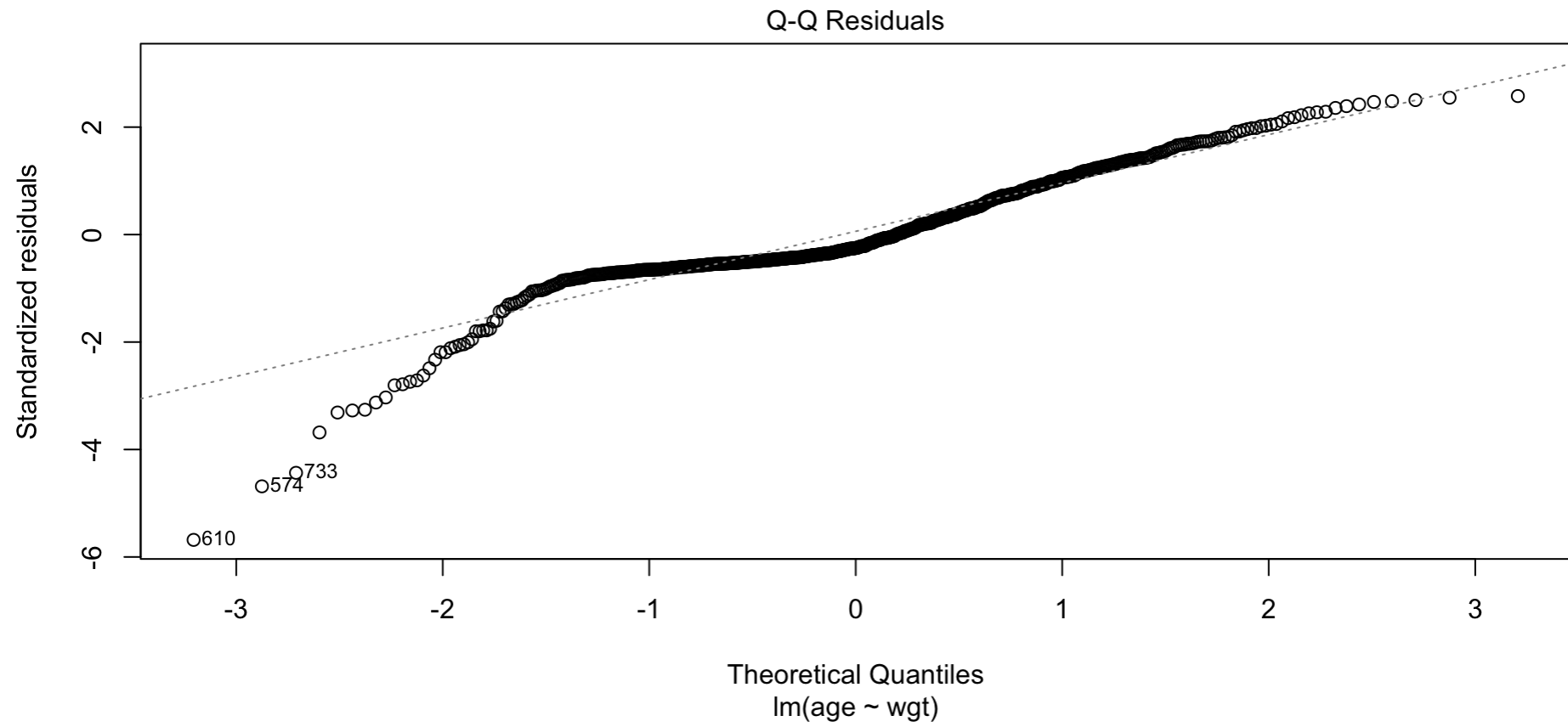
```
1 boxplot(boys$hgt ~ boys$reg, main = "Boxplot",  
2         xlab = "Region", ylab = "Height")
```

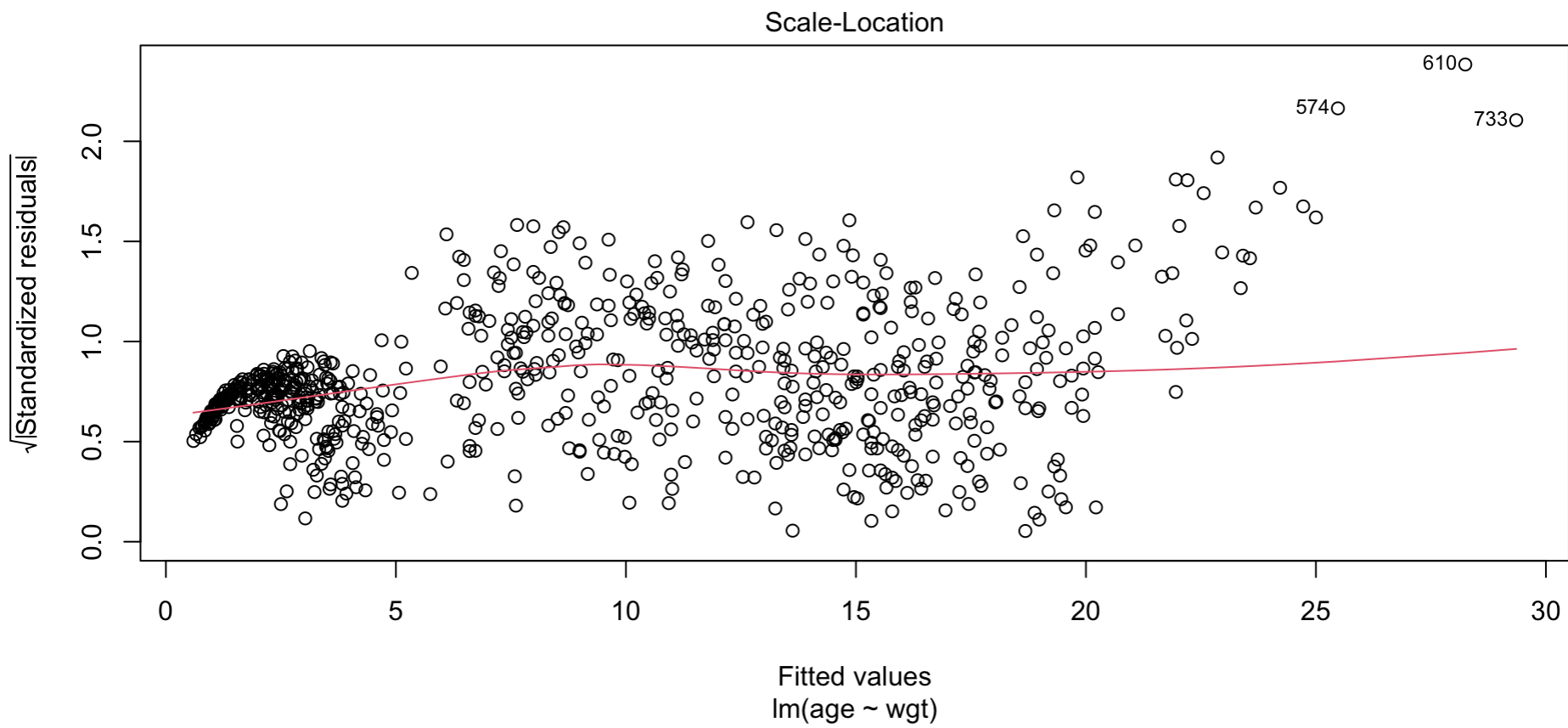


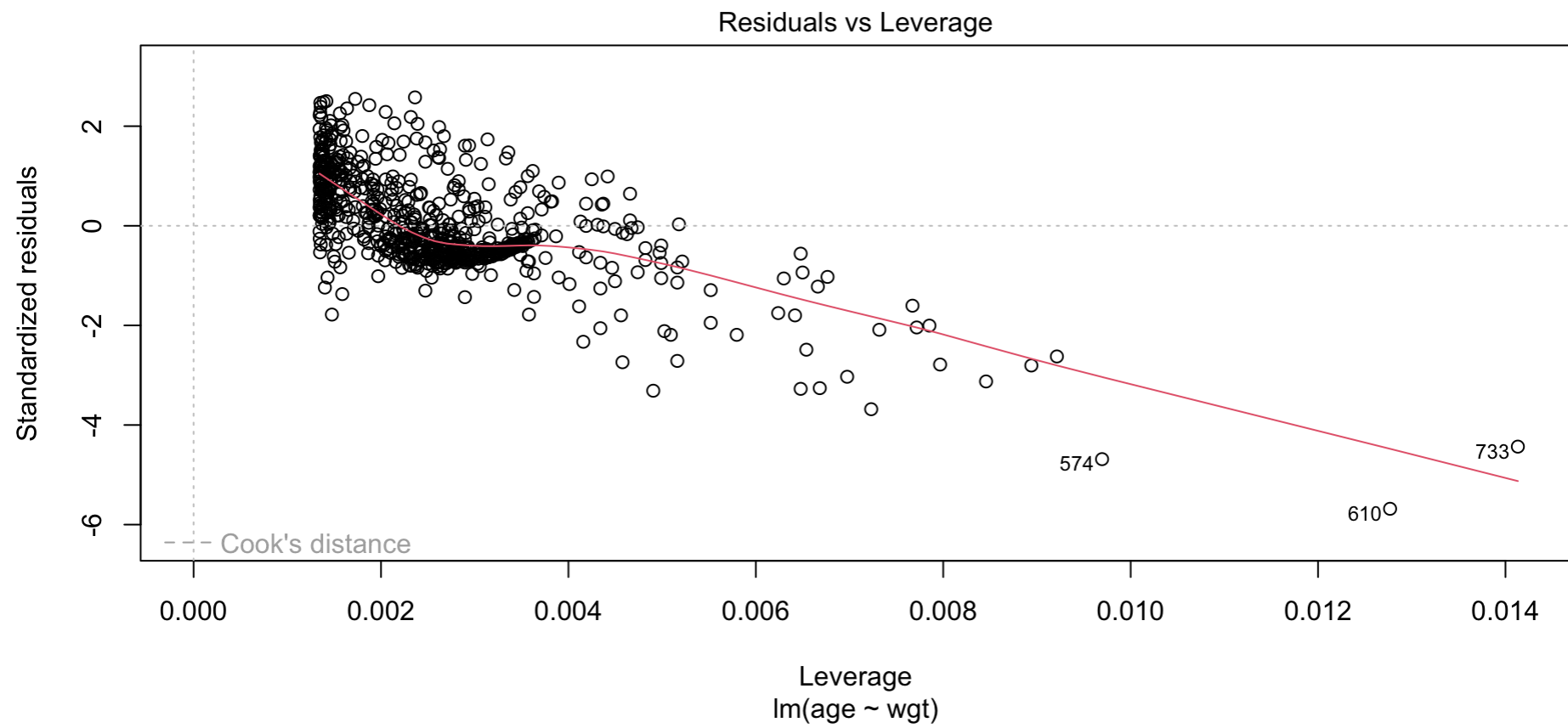
MANY R OBJECTS ALSO HAVE A `plot()` METHOD

```
1 boys %$% lm(age~wgt) %>% plot()
```









NEAT! BUT WHAT IF WE WANT MORE CONTROL?

GGPLOT2

WHAT IS `ggplot2`?

Layered plotting based on the book **The Grammar of Graphics** by Leland Wilkinson.

With `ggplot2` you

1. provide the *data*
2. define how to map variables to *aesthetics*
3. state which *geometric object* to display
4. (optional) edit the overall *theme* of the plot

`ggplot2` then takes care of the details

AN EXAMPLE: SCATTERPLOT

1: Provide the data

```
1 boys %>%  
2   ggplot()
```

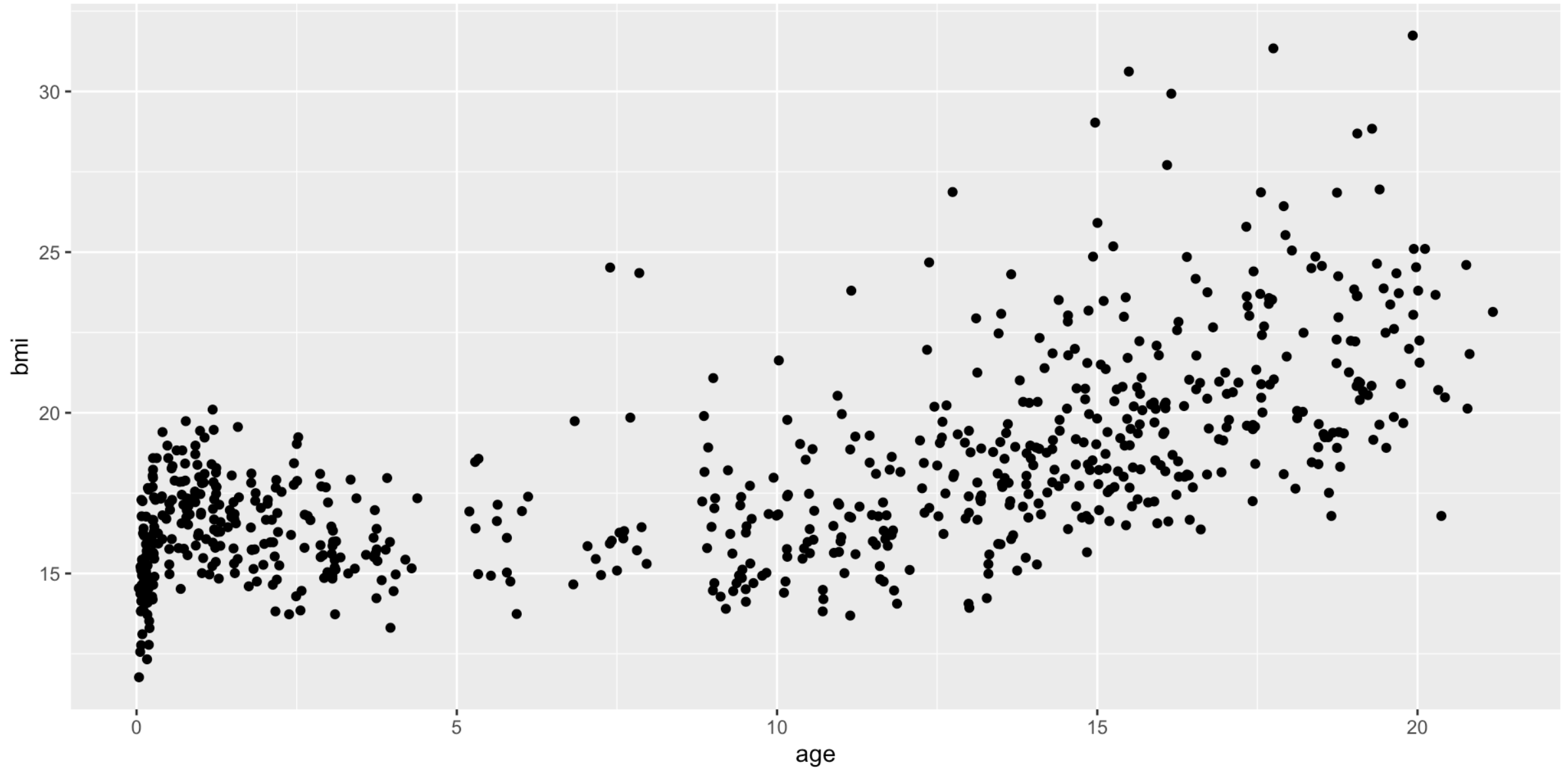
2: map variable to aesthetics

```
1 boys %>%  
2   ggplot(aes(x = age, y = bmi))
```

3: state which geometric object to display

```
1 boys %>%  
2   ggplot(aes(x = age, y = bmi)) +  
3   geom_point()
```


AN EXAMPLE: SCATTERPLOT



WHY THIS SYNTAX?

Create the plot

```
1 gg <-  
2   boys %>%  
3   ggplot(aes(x = age, y = bmi)) +  
4   geom_point(col = "dark green")
```

Add another layer (smooth fit line)

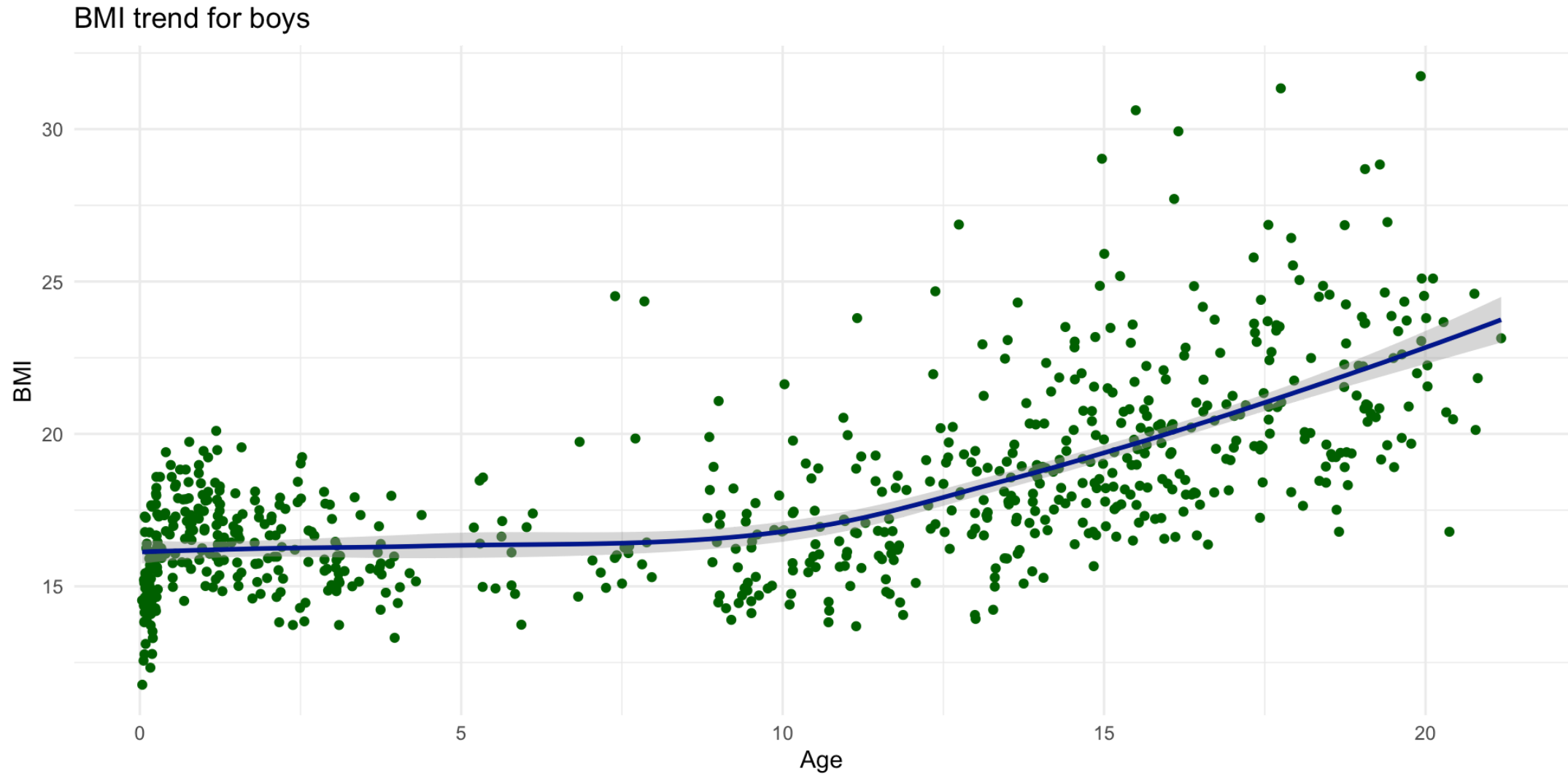
```
1 gg <- gg +  
2   geom_smooth(col = "dark blue")
```

Give it some labels and a nice look

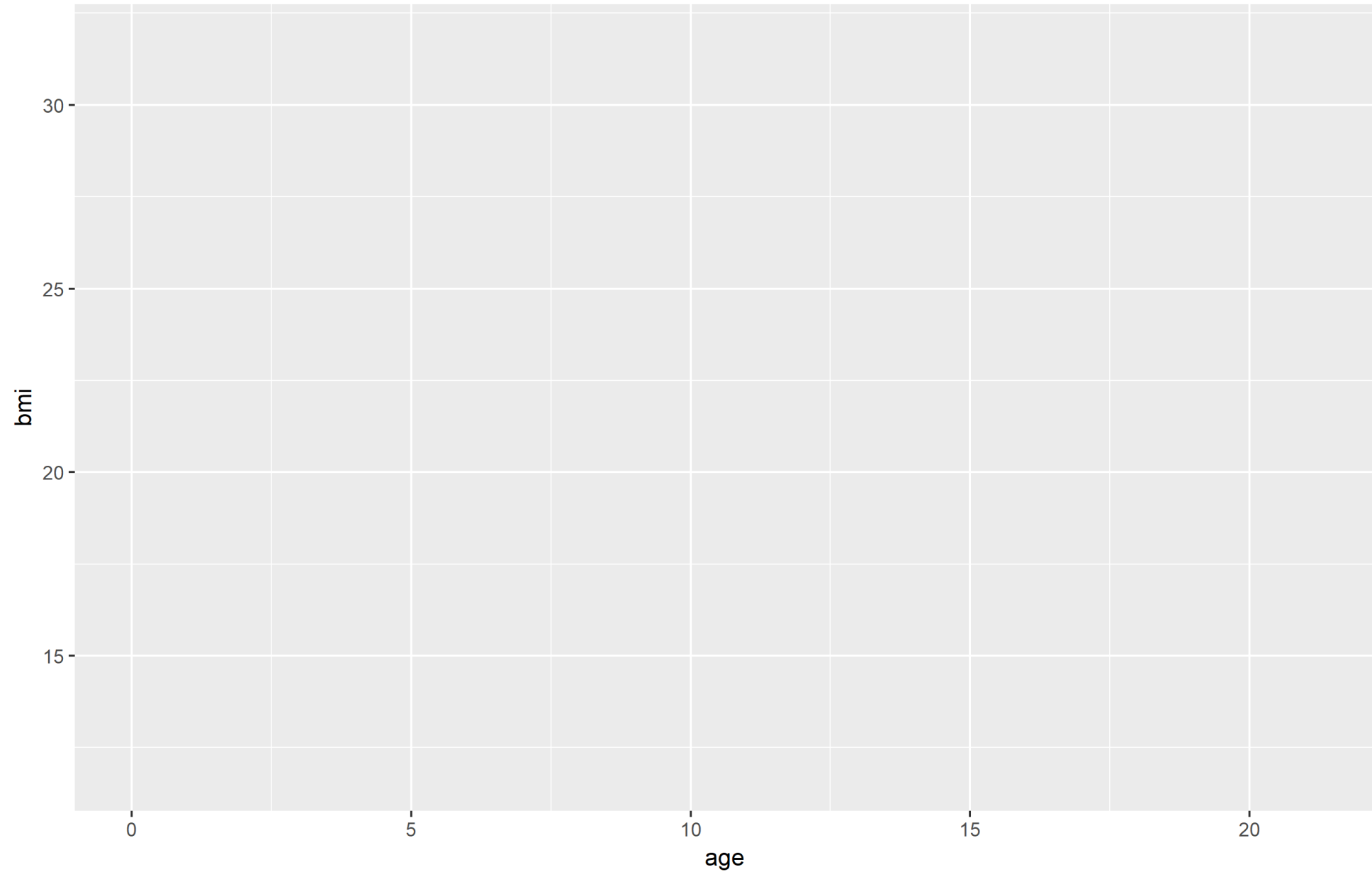
```
1 gg <- gg +  
2   labs(x = "Age", y = "BMI", title = "BMI trend for boys") +  
3   theme_minimal()
```

WHY THIS SYNTAX?

```
1 plot(gg)
```



WHY THIS SYNTAX?



AESTHETICS

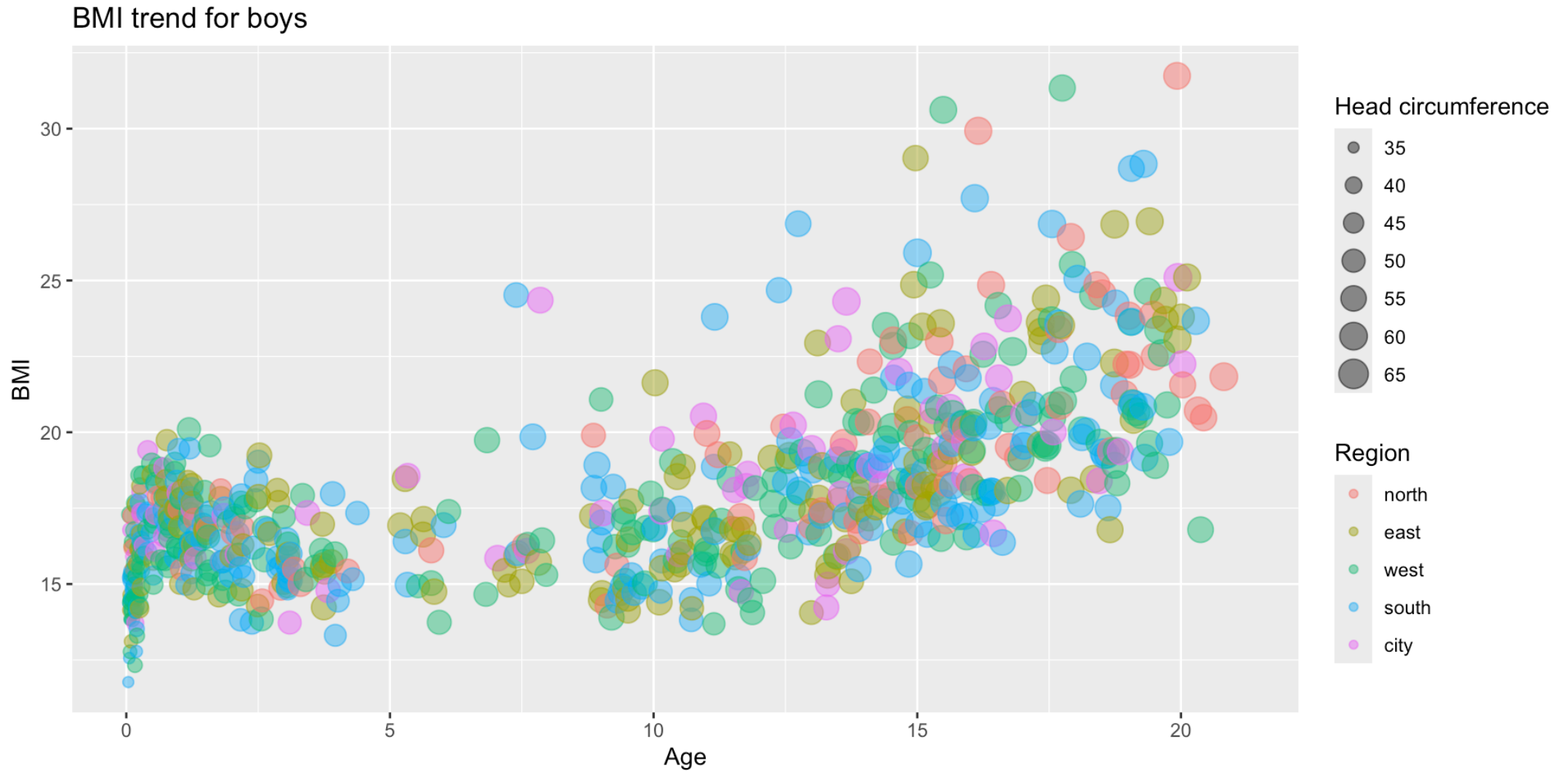
- x
- y
- size
- colour
- fill
- opacity (alpha)
- linetype
- ...

AESTHETICS

```
1 gg <-  
2   boys %>%  
3   filter(!is.na(reg)) %>%  
4  
5   ggplot(aes(x      = age,  
6             y      = bmi,  
7             size    = hc,  
8             colour  = reg)) +  
9  
10  geom_point(alpha = 0.5) +  
11  
12  labs(title  = "BMI trend for boys",  
13       x      = "Age",  
14       y      = "BMI",  
15       size    = "Head circumference",  
16       colour  = "Region")
```

AESTHETICS

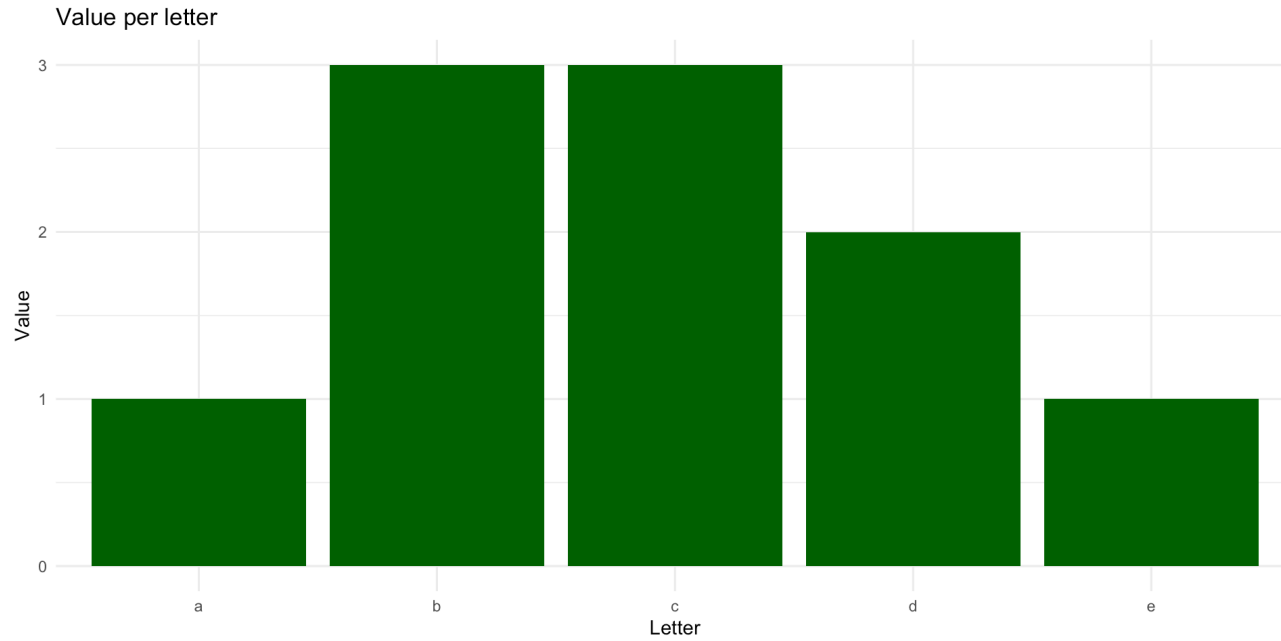
```
1 plot(gg)
```



GEOMS

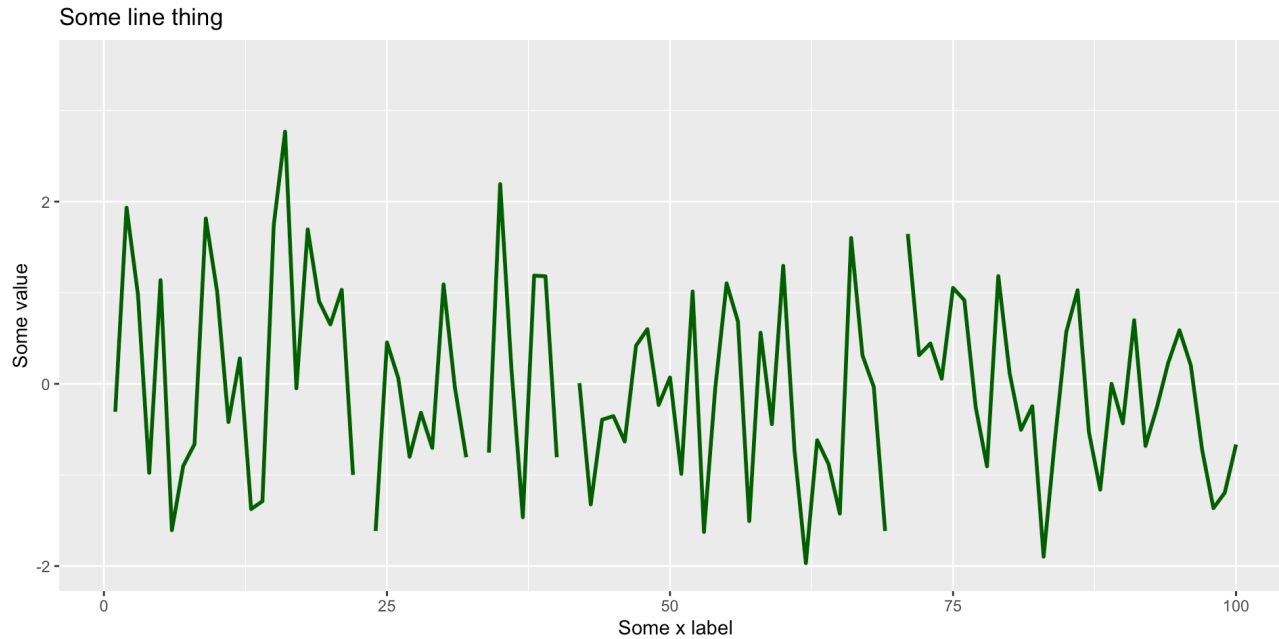
- `geom_point`
- `geom_bar`
- `geom_line`
- `geom_smooth`
- `geom_histogram`
- `geom_boxplot`
- `geom_density`

GEOMS: BAR



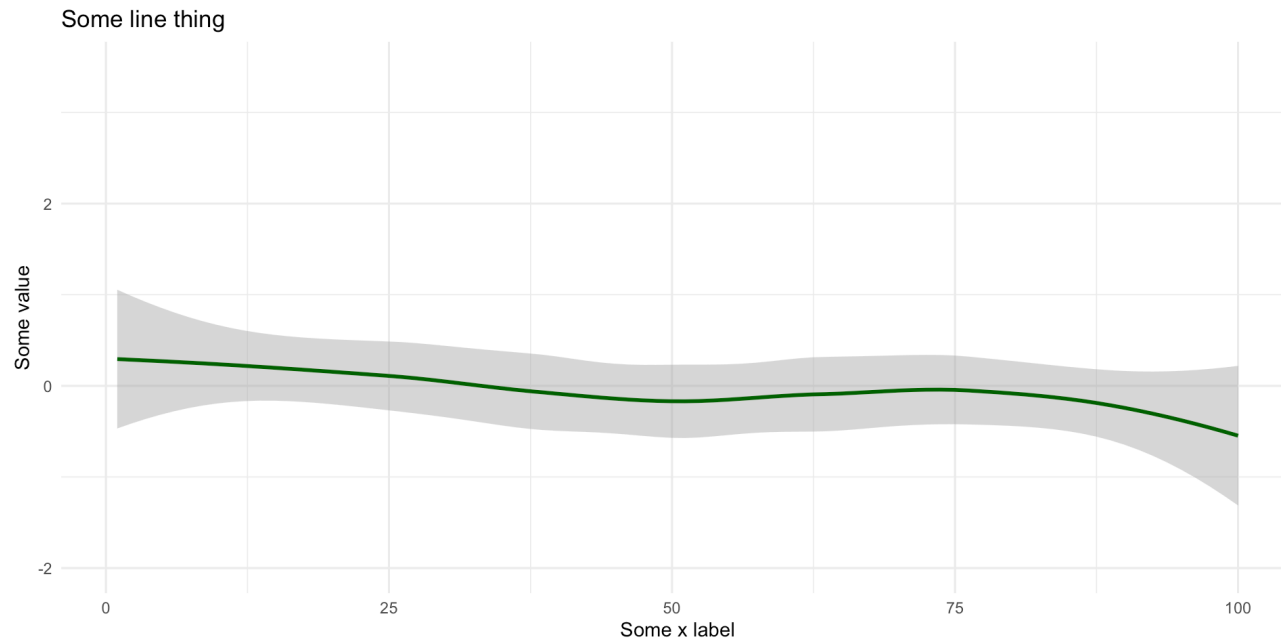
```
1 data.frame(x = letters[1:5],  
2           y = c(1, 3, 3, 2, 1)) %>%  
3   ggplot(aes(x = x, y = y)) +  
4   geom_bar(fill = "dark green",  
5           stat = "identity") +  
6   labs(title = "Value per letter",  
7        x     = "Letter",  
8        y     = "Value")
```

GEOMS: LINE



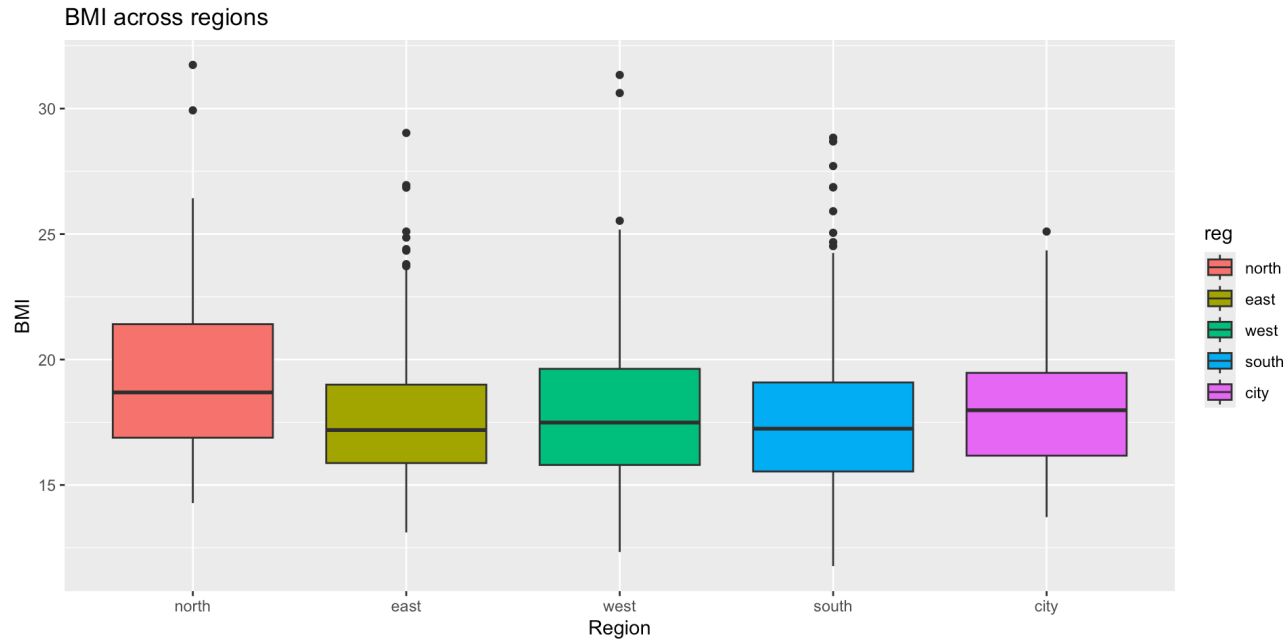
```
1 ggdat <- data.frame(x = 1:100,  
2                       y = rnorm(100))  
3 ggdat %>%  
4   ggplot(aes(x = x, y = y)) +  
5   geom_line(colour = "dark green",  
6             lwd = 1) +  
7   ylim(-2, 3.5) +  
8   labs(title = "Some line thing",  
9         x     = "Some x label",  
10        y     = "Some value")
```

GEOMS: SMOOTH



```
1 ggdat %>%  
2   ggplot(aes(x = x, y = y)) +  
3   geom_smooth(colour = "dark green",  
4               lwd = 1,  
5               se = TRUE) +  
6   ylim(-2, 3.5) +  
7   labs(title = "Some line thing",  
8         x     = "Some x label",  
9         y     = "Some value")
```

GEOMS: BOXPLOT

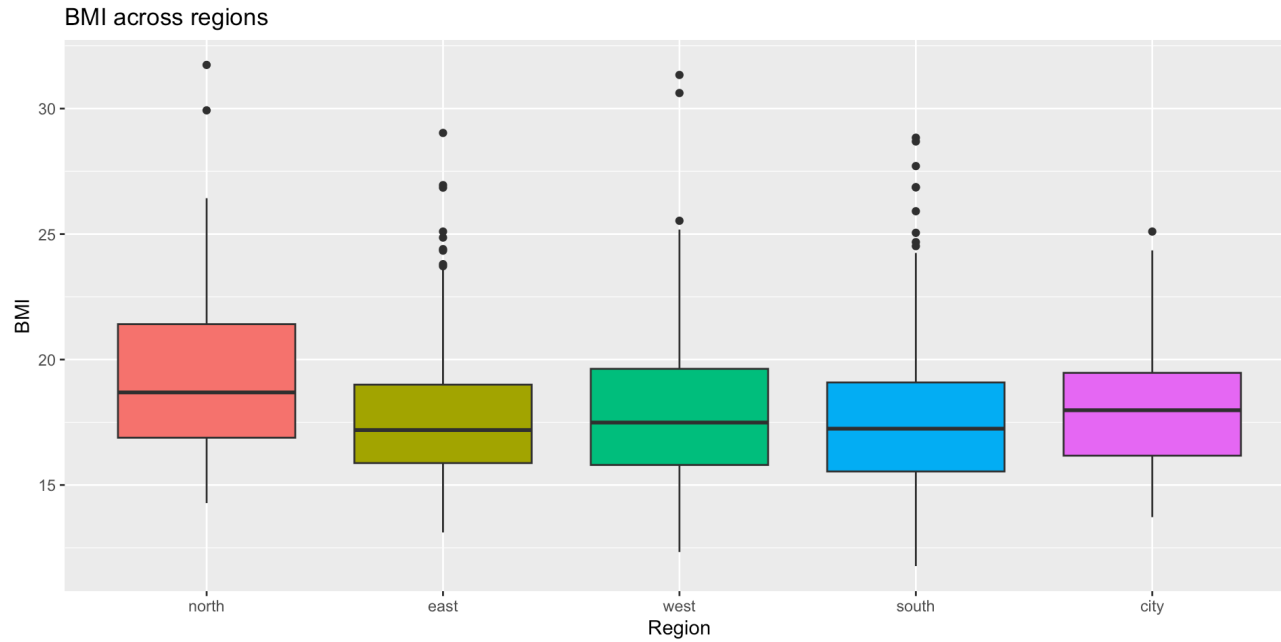


```

1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = reg,
4             y = bmi,
5             fill = reg)) +
6   geom_boxplot() +
7   labs(title = "BMI across regions",
8        x = "Region",
9        y = "BMI")

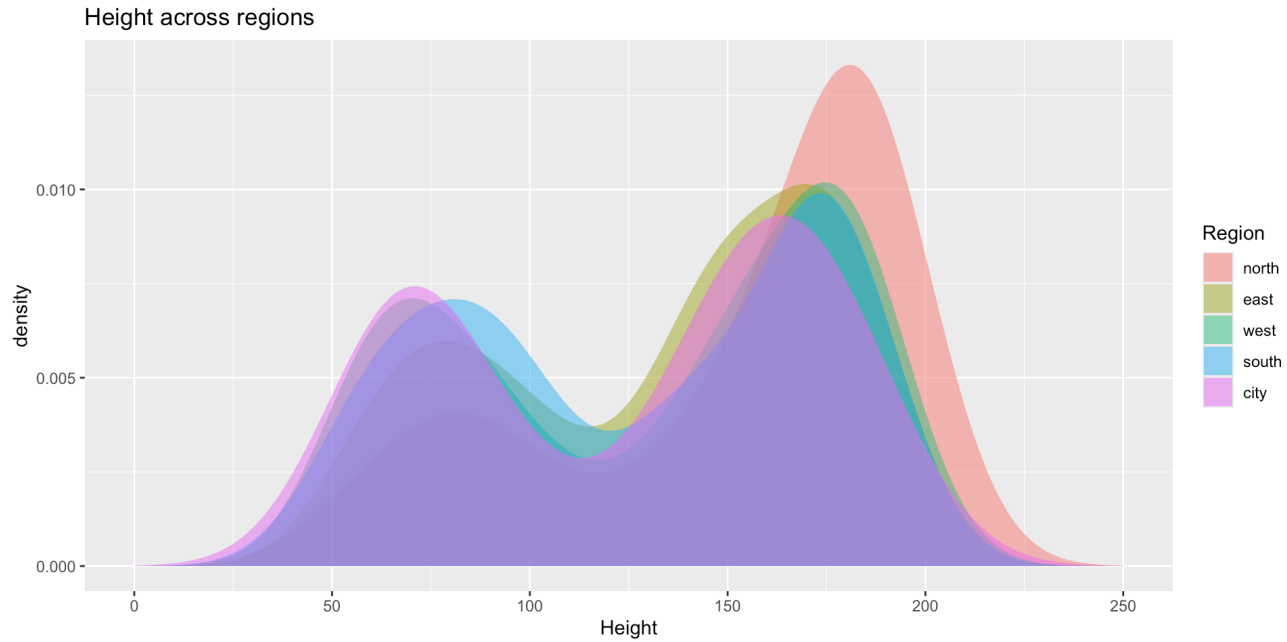
```

GEOMS: BOXPLOT WITHOUT LEGEND



```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = reg,
4             y = bmi,
5             fill = reg)) +
6   geom_boxplot() +
7   labs(title = "BMI across regions",
8        x = "Region",
9        y = "BMI") +
10  theme(legend.position = "none")
```

GEOMS: DENSITY

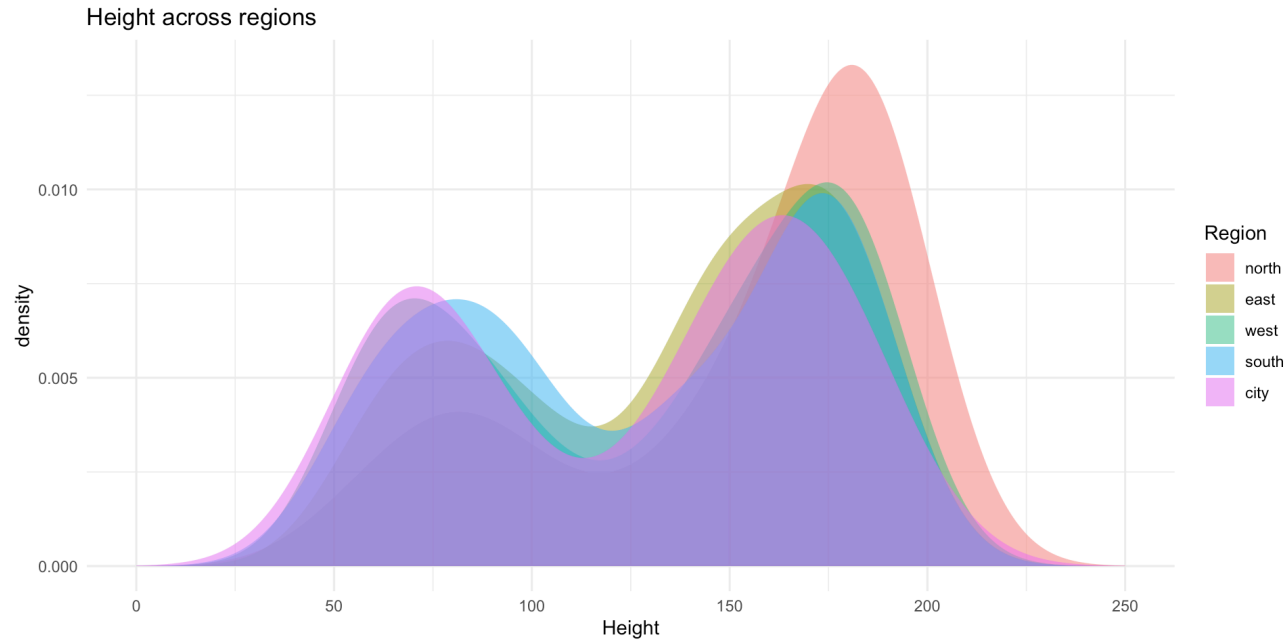


```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = hgt, fill = reg)) +
4   geom_density(alpha = 0.5,
5                 colour = "transparent") +
6   xlim(0, 250) +
7   labs(title = "Height across regions",
8         x      = "Height",
9         fill   = "Region")
```

CHANGING THE STYLE: THEMES

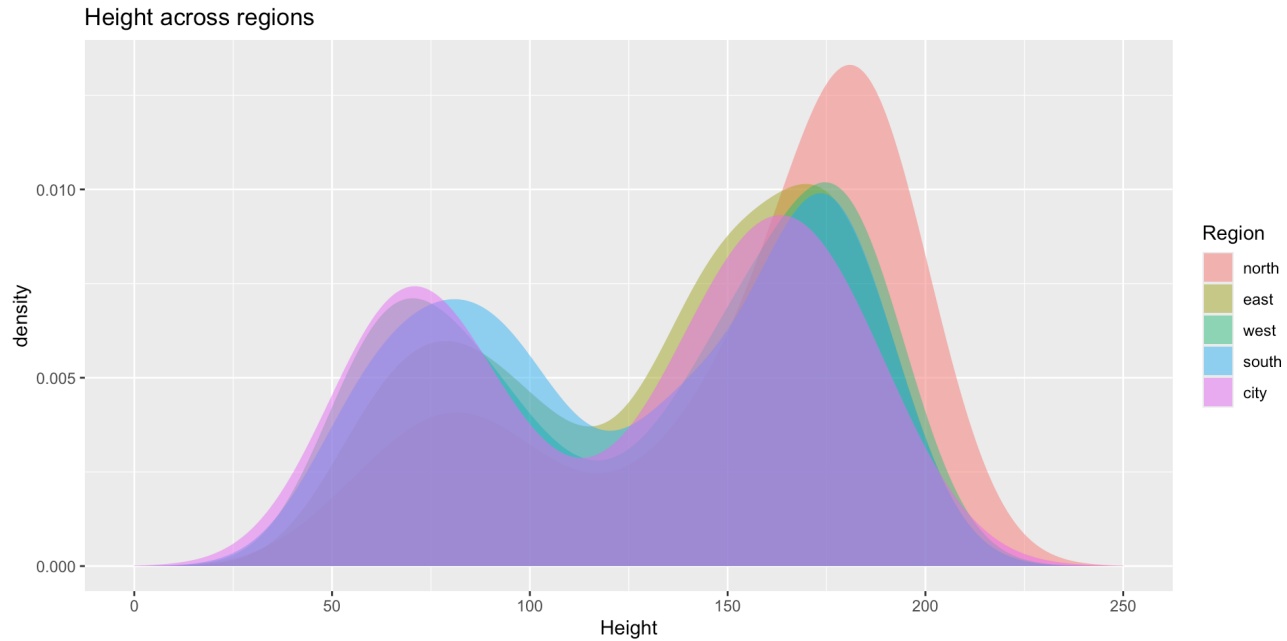
- Themes determine the overall appearance of your plot
- standard themes: e.g., `theme_minimal()`, `theme_classic()`, `theme_bw()`, ...
- extra libraries with additional themes: e.g., `ggthemes`
- customize own theme using options of `theme()`

CHANGING THE STYLE: THEMES



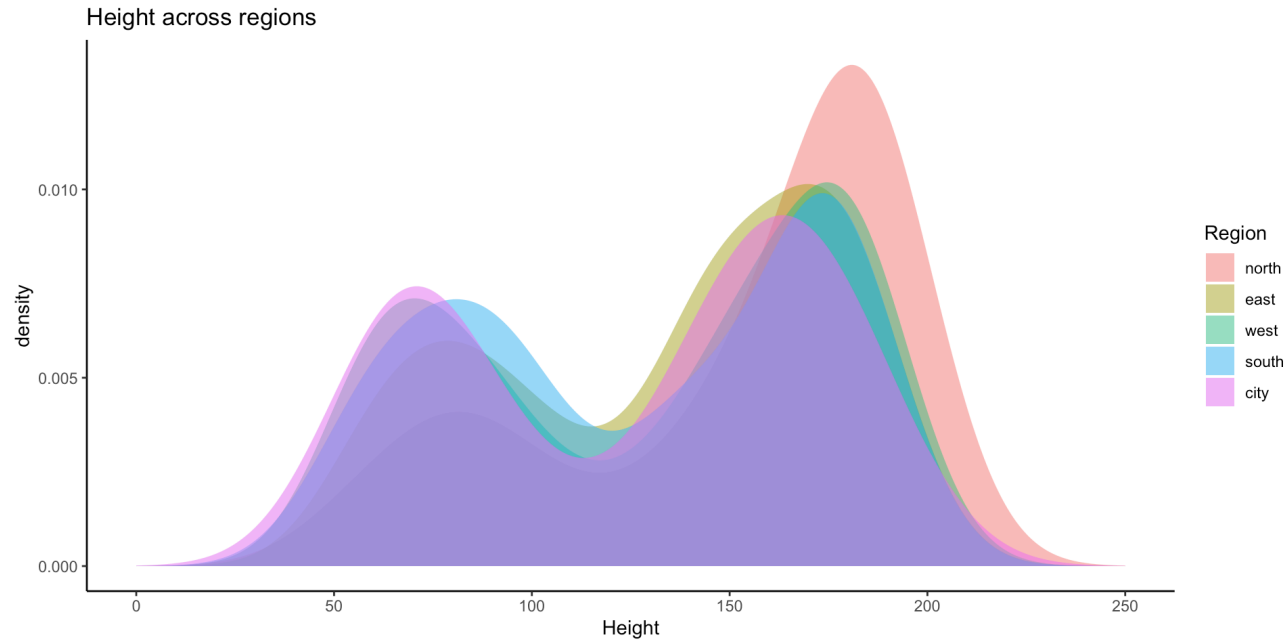
```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = hgt, fill = reg)) +
4   geom_density(alpha = 0.5,
5                 colour = "transparent") +
6   xlim(0, 250) +
7   labs(title = "Height across regions",
8         x      = "Height",
9         fill   = "Region") +
10  theme_minimal()
```


CHANGING THE STYLE: THEMES



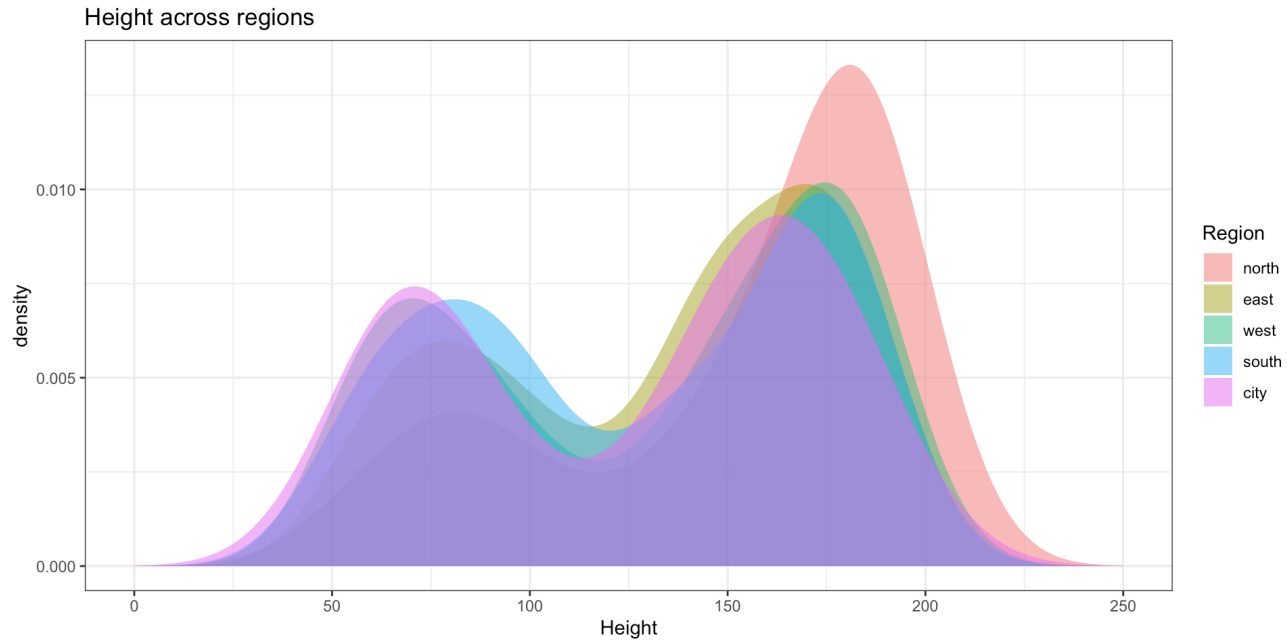
```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = hgt, fill = reg)) +
4   geom_density(alpha = 0.5,
5                 colour = "transparent") +
6   xlim(0, 250) +
7   labs(title = "Height across regions",
8         x      = "Height",
9         fill   = "Region") +
10  theme_gray()
```

CHANGING THE STYLE: THEMES



```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = hgt, fill = reg)) +
4   geom_density(alpha = 0.5,
5                 colour = "transparent") +
6   xlim(0, 250) +
7   labs(title = "Height across regions",
8         x      = "Height",
9         fill   = "Region") +
10  theme_classic()
```

CHANGING THE STYLE: THEMES



```
1 boys %>%
2   filter(!is.na(reg)) %>%
3   ggplot(aes(x = hgt, fill = reg)) +
4   geom_density(alpha = 0.5,
5                 colour = "transparent") +
6   xlim(0, 250) +
7   labs(title = "Height across regions",
8         x     = "Height",
9         fill  = "Region") +
10  theme_bw()
```

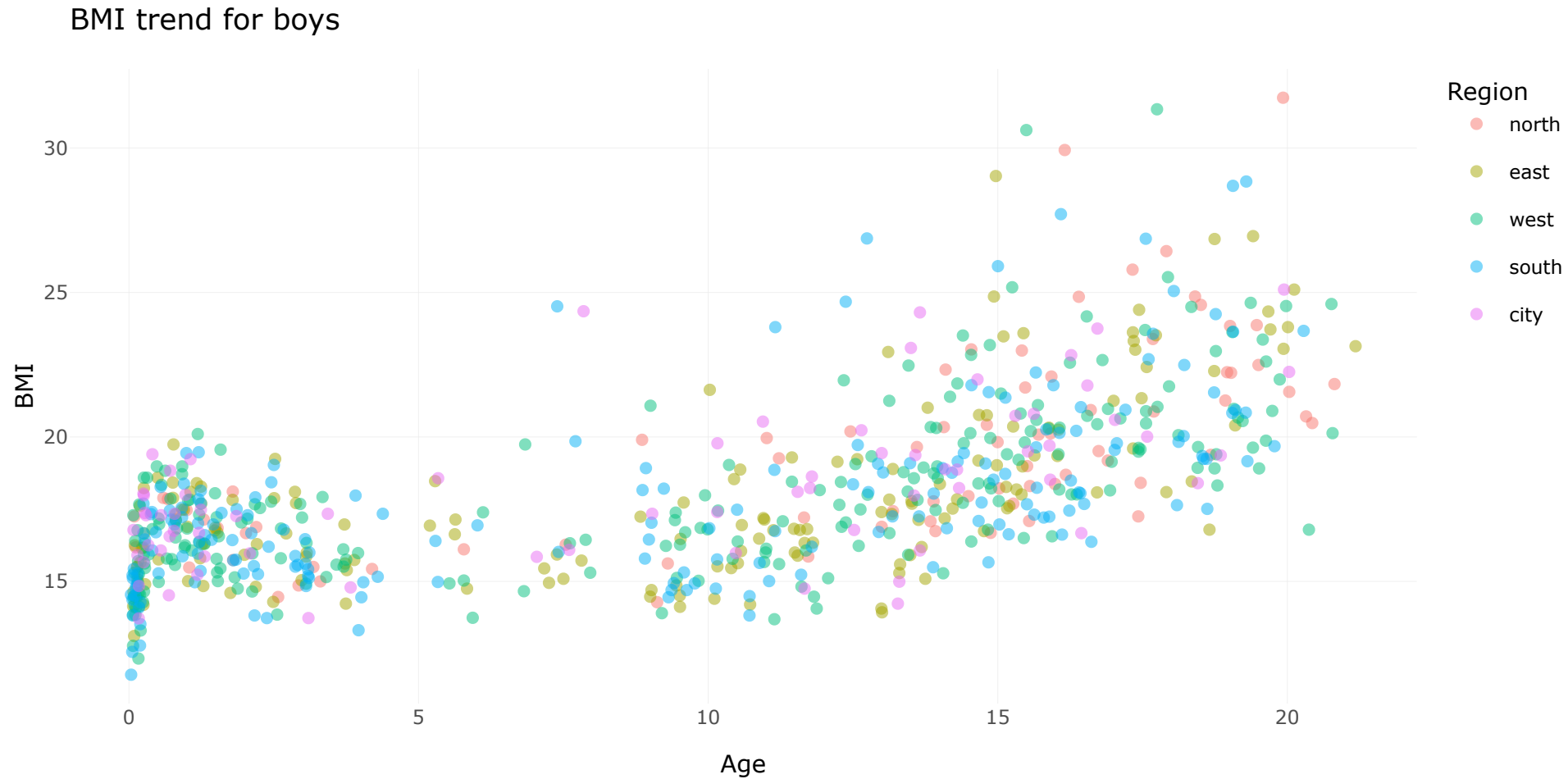
INTERACTIVE PLOTS

Use `plotly::ggplotly()` to make any ggplot interactive

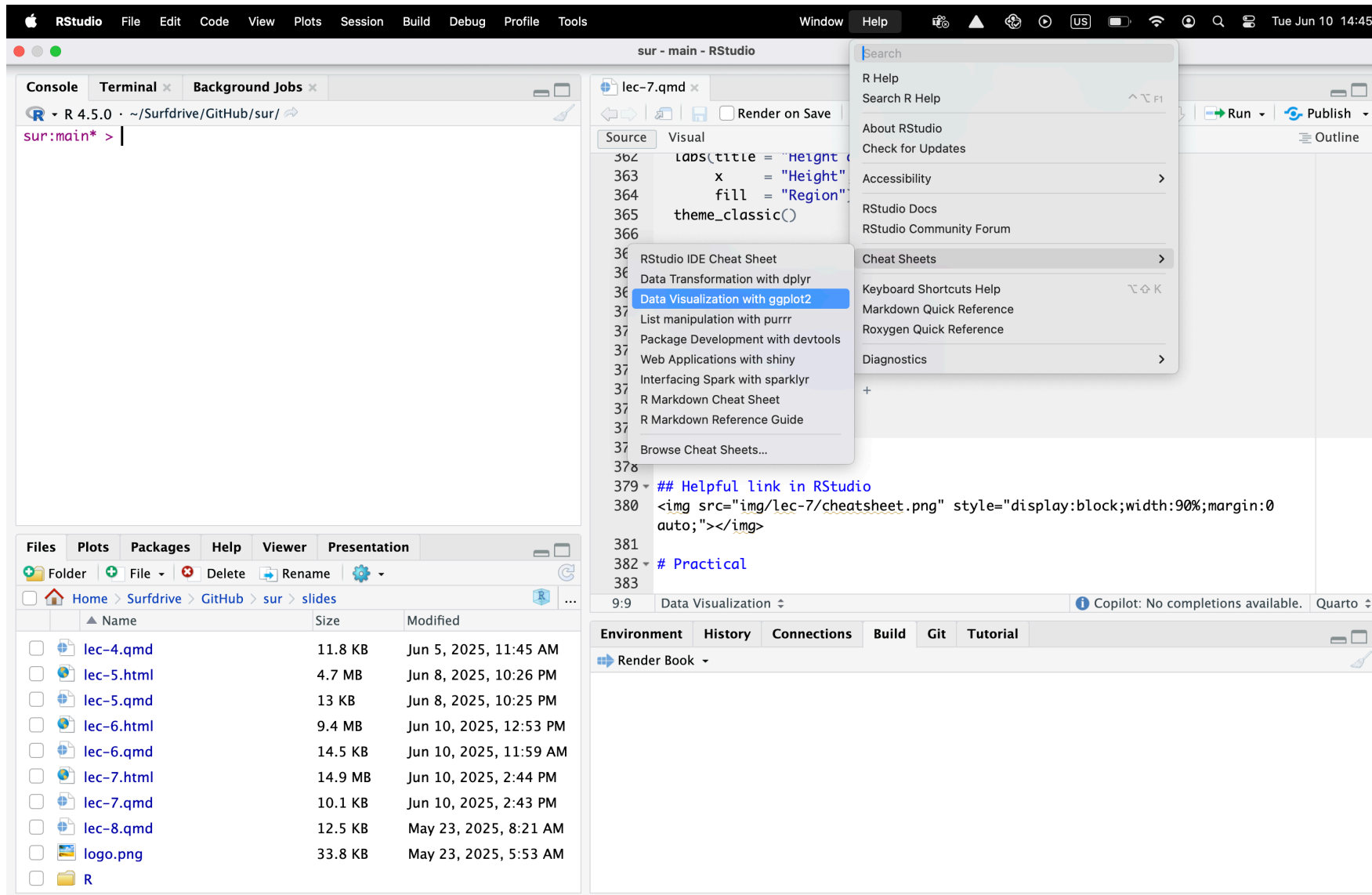
```
1 library(plotly)
2 gg <- boys %>%
3   filter(!is.na(reg)) %>%
4
5   ggplot(aes(x      = age,
6              y      = bmi,
7              colour = reg)) +
8
9   geom_point(alpha = 0.5) +
10
11   labs(title = "BMI trend for boys",
12        x     = "Age",
13        y     = "BMI",
14        colour = "Region") +
15   theme_minimal()
16
17 ggplotly(gg)
```

INTERACTIVE PLOTS

Use `plotly::ggplotly()` to make any ggplot interactive



HELPFUL LINK IN RSTUDIO



PRACTICAL