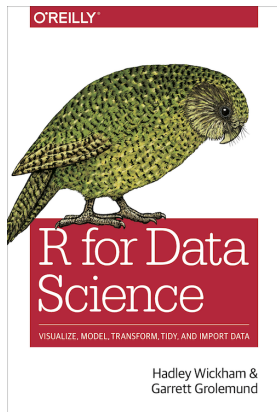# Socio-Informatics 348
## Data Visualisation with the Tidyverse
## Part 2

Dr Lisa Martin

Department of Information Science
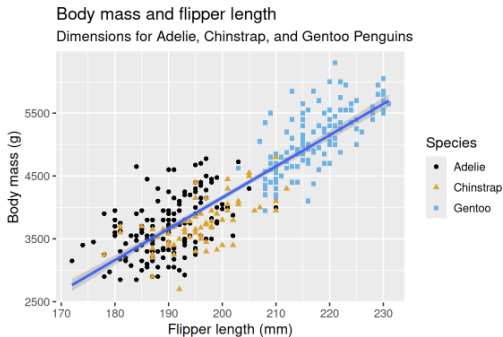Stellenbosch University

# Today's Reading



*R for Data Science, Wholegame, Visualisation*

# Side Notes

- **Dark Mode:** Tools ▷ Global Options ▷ Appearance ▷ Editor theme
- **Shortcuts:** Tools ▷ Keyboard Shortcuts Help (Alt+Shift+K)

# Where we left off...

- Scatterplot with palmerpenguins dataset
- ggplot layers
- Data, Aesthetics (aes) and Geometry (geom)

# Visualising Distributions

- Understanding the distribution of a variable is important
- Helps to identify patterns, outliers, and the overall shape of the data
- Geom used depends on the type of variable
- Continuous variables: geom_histogram(), geom_density()
- Categorical variables: geom_bar()

# Visualising Distributions

# Visualising Distributions

- glimpse()

# Visualising Distributions

- view()

# Visualising Distributions

**Categorical Variables: geom_bar()**

- Can only take one of a small set of values.
- The height of the bars displays how many observations occurred with each x value.
- geom_bar() uses the count of observations by default

# Visualising Distributions

## Categorical Variables: geom_bar()

# Visualising Distributions

**Categorical Variables: geom_bar() - ordered**

# Visualising Distributions

**Numerical Variables: geom_histogram()**



*Note: Warning about NAs*

# Visualising Distributions

## Numerical Variables: geom_histogram() - bin size

# Visualising Distributions

**Numerical Variables: geom_density()**

- Smoothed version of a histogram
- Good for continuous data

# Visualising Relationships

**Numerical and Categorical: geom_boxplot()**

- Visualise the distribution of a continuous variable across categories
- Displays the median, quartiles, and potential outliers
- Spread of the distribution - symmetric or skewed to one side



*R4DS, Figure 1.1*

# Visualising Relationships

**Numerical and Categorical: geom_boxplot()**

# Visualising Relationships

**Numerical and Categorical: geom_density()**



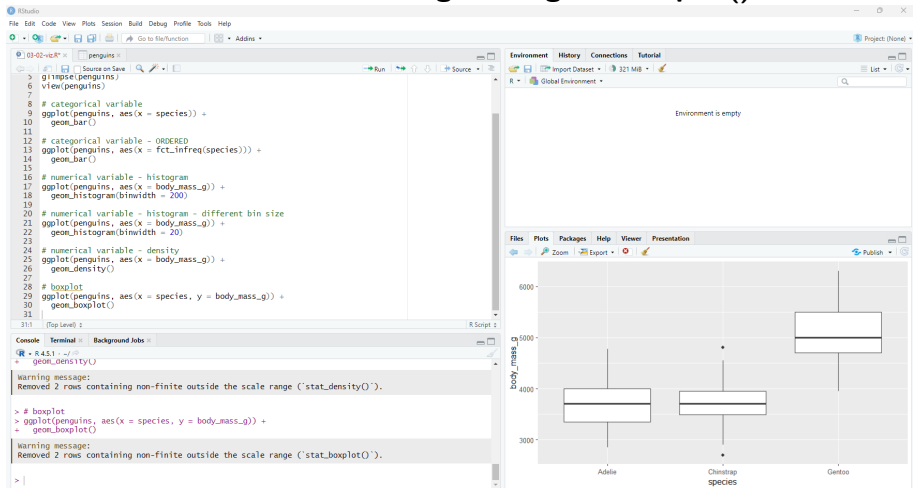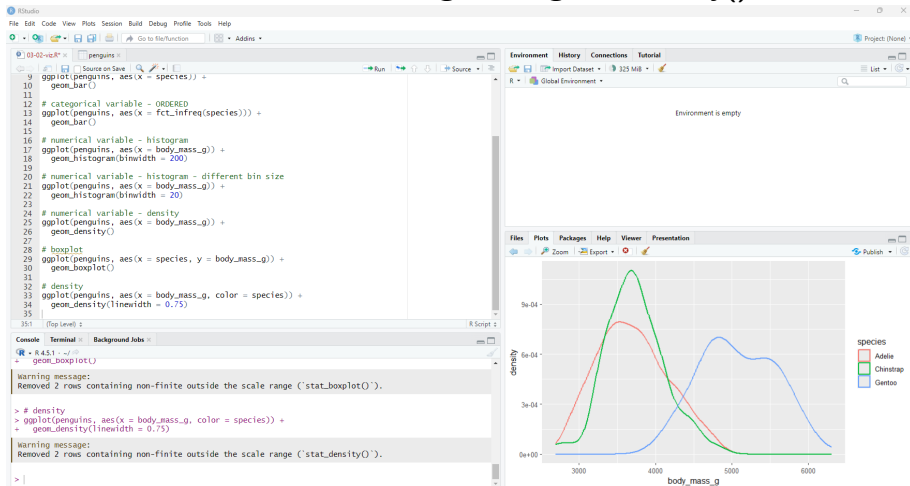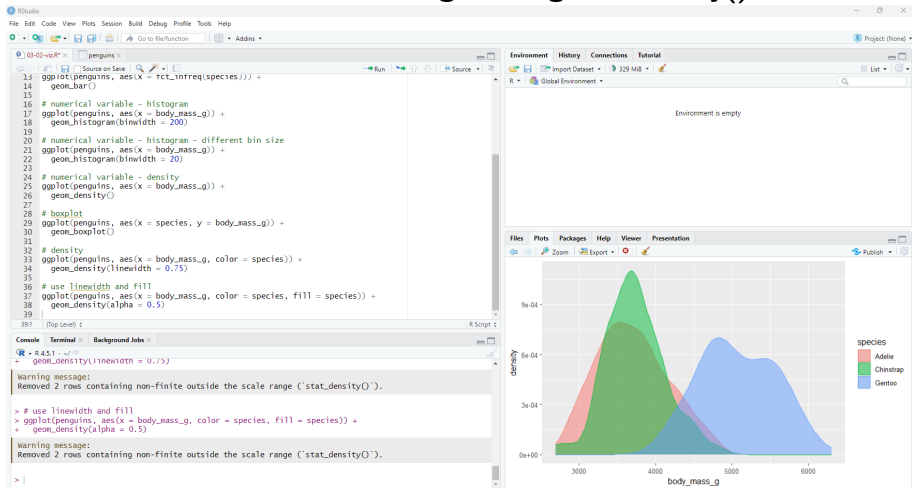*Note: linewidth to set the thickness of the line*

# Visualising Relationships
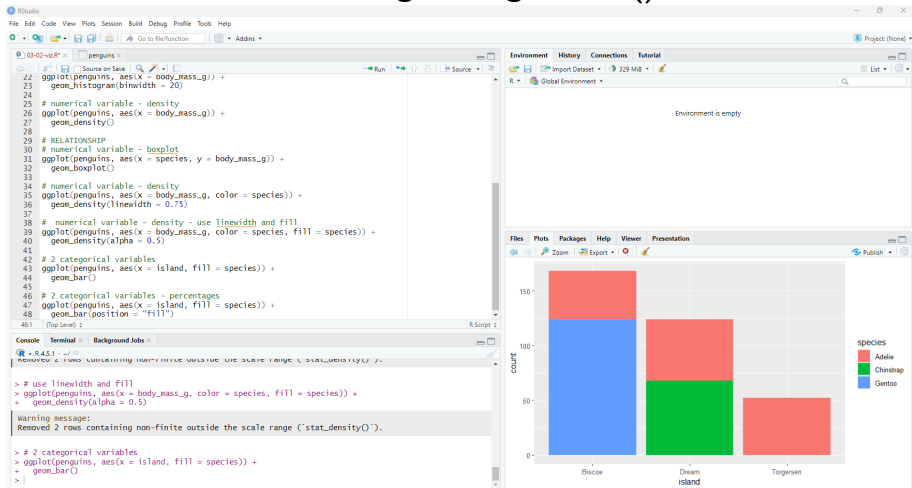
## Numerical and Categorical: geom_density()



*Note: linewidth and fill to set the thickness and colour of the area under the curve alpha to set the transparency of the fill*
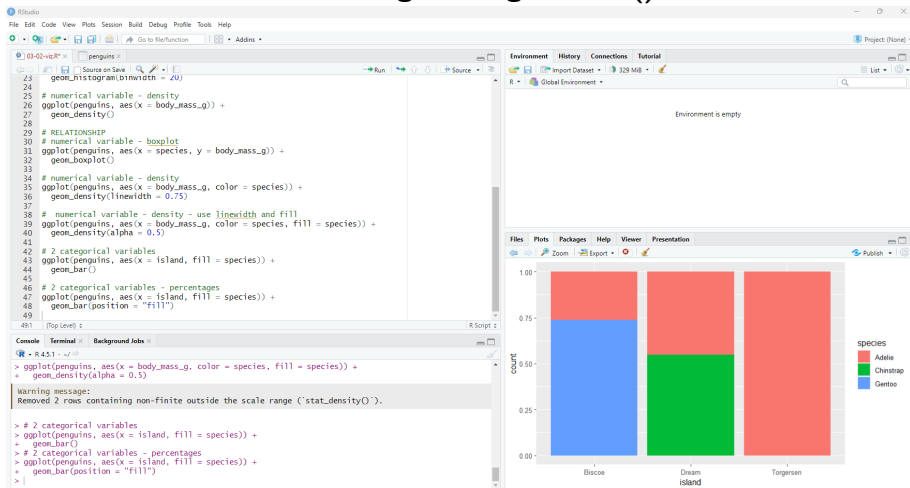
# Visualising Relationships

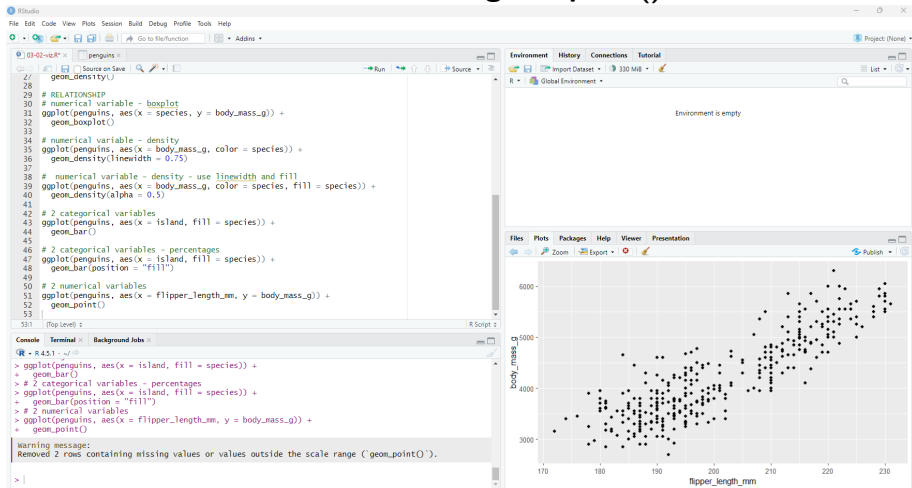## Two Categorical: geom_bar()

# Visualising Relationships

**Two Categorical: geom_bar()**



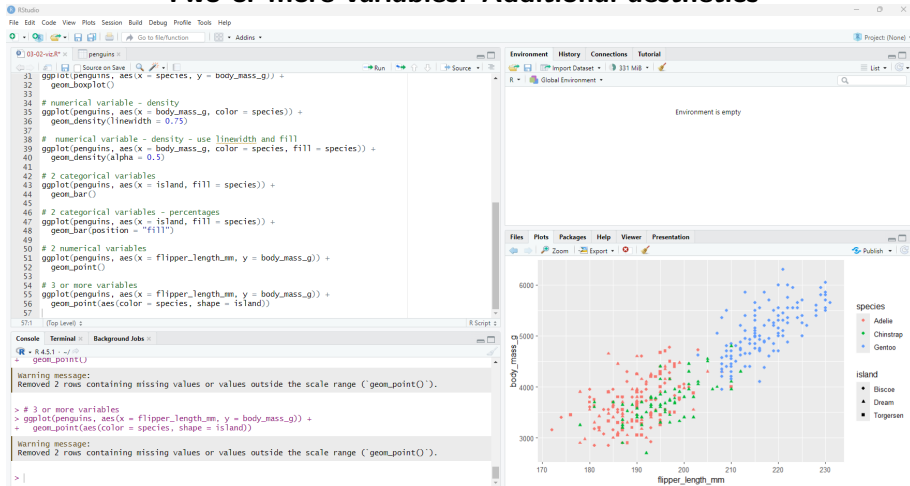*Note: Use of position = "fill" to show proportions instead of counts*

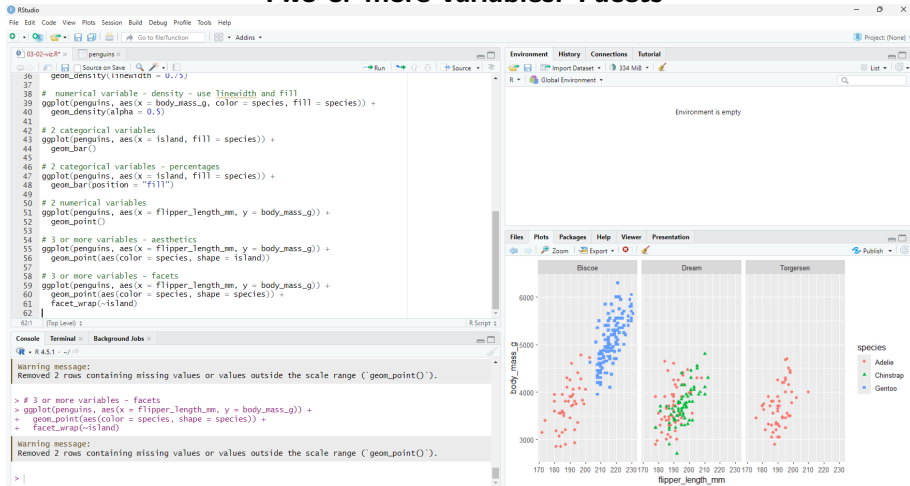# Visualising Relationships

## Two Numerical: geom_point()

# Visualising Relationships

## Two or more variables: Additional aesthetics

# Visualising Relationships

## Two or more variables: Facets

# Save your plots

**ggsave()**

```r
ggplot(penguins, aes(x = flipper_length_mm, y = body_mass_g)) +
  geom_point()
ggsave(filename = "penguin-plot.png")
```