# Introduction To R

Chris Chia

November 24, 2019

# What is R? Why R?

R is a programming language that allows you to perform statistics and visualisations. R is excellent because:

- Used widely among statisticians and mathematicians and gaining traction amongst economists, political scientists, social scientists
- Large community means - can easily troubleshoot any problems online
- Large collection of in-built libraries/tools to support your data analysis needs - beautiful data visualisation tools.
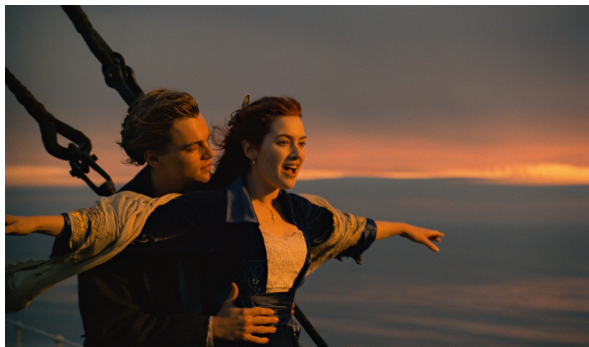- More flexible, versatile than Stata

# Outline

We're going to try and achieve these two outcomes in this 20 minute tutorial:

1. Introduction: What is R is and why use R?

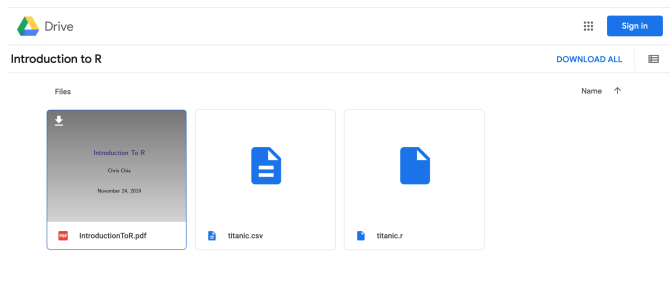2. Explore data storytelling by getting hands on with R and create data visualisations

# Let's get hands-on

- We want to ask questions, and use data to answer them - but also in a visual and interpretable way
- We'll use the "Titanic" dataset (from Stanford CS109), a table of passengers and their information (sex, age, family, passenger class) - and whether they survived.
- We'll explore whether **survival rate** is different for different **passenger classes**

# R and Data Setup

- Open RStudio. If you don't have it installed, download it at `http://rstudio.com`
- Download the titanic dataset (under the file name "titanic.csv"), and the file "titanic.R" from this url: `http://shorturl.at/fkSVW`

# Setup

Open the "titanic.r" file in RStudio, which will contain the commands we want to use to create data visualisations.

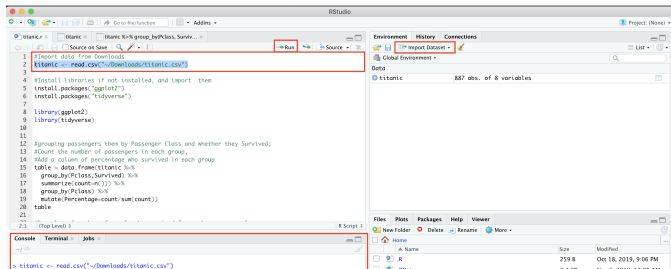from the menu File → Open File. (or Cmd+O on Mac, Ctrl + O on Windows)
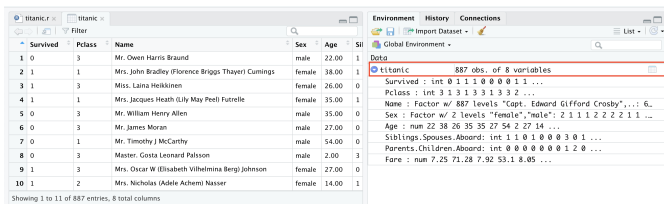
You should see a screen like this:

# Data Import

- In the code editor, select the line "**titanic** ← **read.csv("**∼**/Downloads/titanic.csv")**" . Make sure to change ∼**/Downloads/** to the folder where you saved "titanic.csv", and then press the "Run" button. This assigns the dataset to a variable (name) "titanic" by executing the command on the console.
- If that does not work, try importing it directly from the environment pane menu Import Dataset → From Text(base)
- This is why we use a ".r" file - so that others can reuse and replicate the steps we took. Also, all the commands need to be run again when we restart

# Exploring the Data

In the environment pane, press the right arrow icon, and click on the titanic rectangle. This is the dataset in raw table form. You should see something like this:



We can see that it has information about passengers, their Passenger Class (**PClass**) which is ordered from 1-3 (highest to lowest), **Sex** (male, female 1 or 2), **Age**, Family information, and whether they survived.
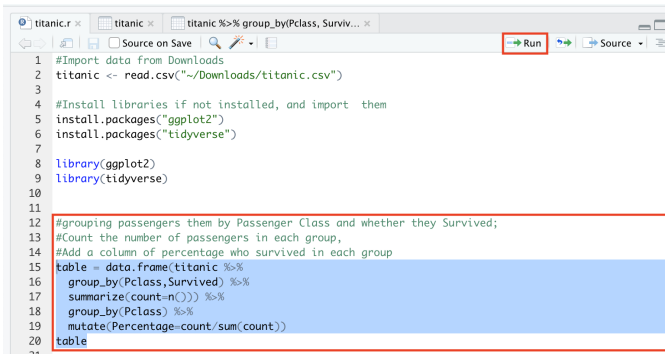
# Import libraries

We need to install and import libraries for the data visualisation. Select these lines in the code editor, then press 'Run'.

# Exploring Data

- Select the following lines, which groups the data by passenger class and survival, and counts how many are in each category, in the code editor, then press 'Run'. You should see this output in your console:

# First Data Visualisation - Grouped Bar Plot

Now select the following lines, which creates a bar plot of the number of
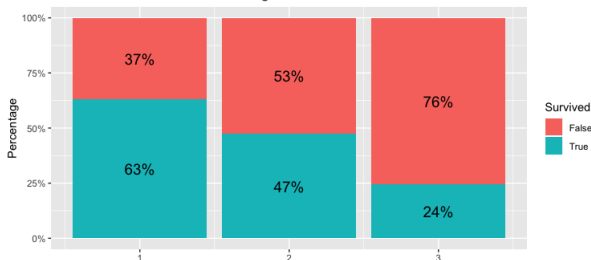each people grouped by each passenger class and whether they survived

# Second Data Visualisation - Stacked Percentage Bar Plot

We want to be able to make a comparison among ticket classes. Select the following lines, which will create a bar plot of survival in each passenger class as a percentage

# Extension/Challenge

See if survival rates are different for passengers of different Sex

# Further Development

- Explore relationships between survival and the other factors, e.g. Age and Family Members
- Look into building a logistic regression model with R to examine the relationship of each factor with survival
- Explore the Titanic dataset on Kaggle, a data science competition and learning platform