# Data Science for Actuaries

Paul King, ATRC 28 June 2019

# Outline

- What do I mean by data science

- Data skills

- Case study

- Tools and infrastructure

# What do I mean by data science?

When you train as an actuary you'll learn how to analyse data, evaluate financial risks, and communicate this data to non-specialists.

Source: Institute and Faculty of Actuaries

UNIVERSITY OF LEICESTER

# Curriculum 2019, CM1

- Describe the possible aims of a data analysis (e.g. descriptive, inferential, and predictive).

- Describe the stages of conducting a data analysis to solve real-world problems in a scientific manner and describe tools suitable for each stage.

- Describe sources of data and explain the characteristics of different data sources, including extremely large data sets.

- **Explain the meaning and value of reproducible research and describe the elements required to ensure a data analysis is reproducible.**
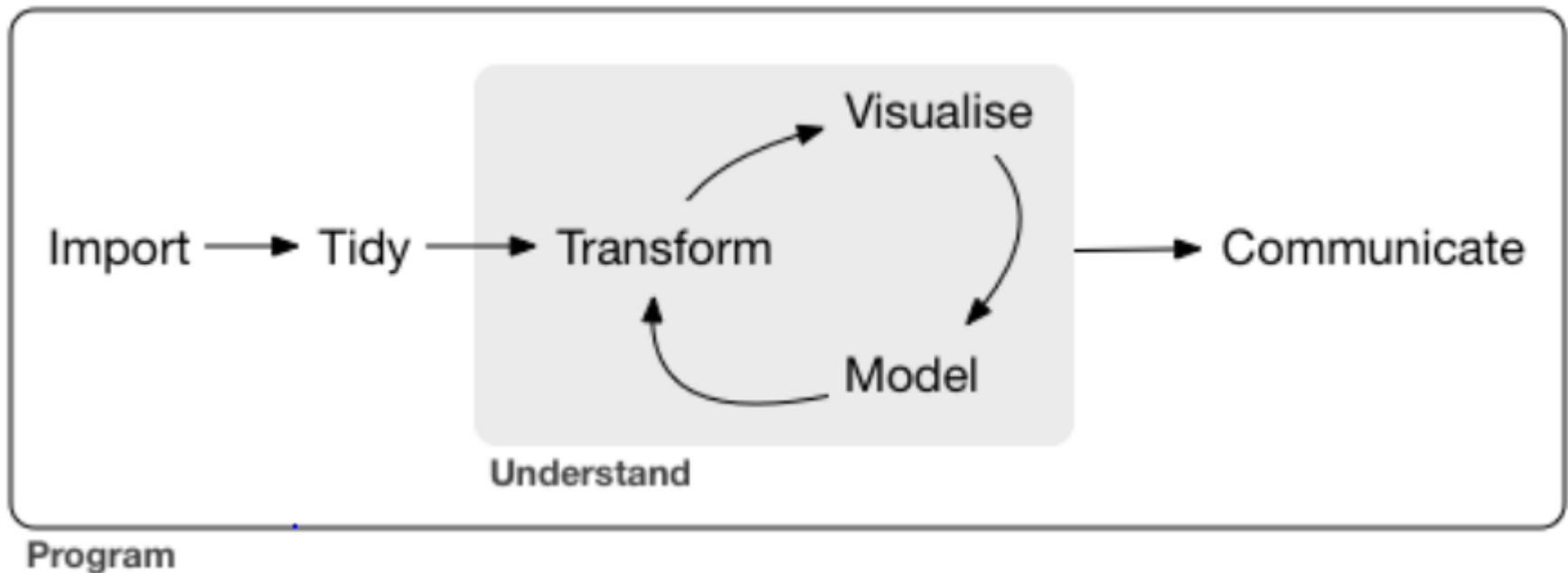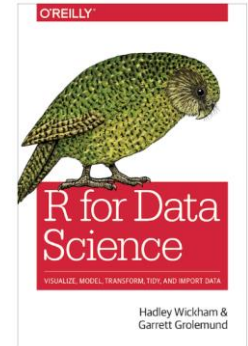
# It's always been data science



Source: Wikipedia

# The data science process



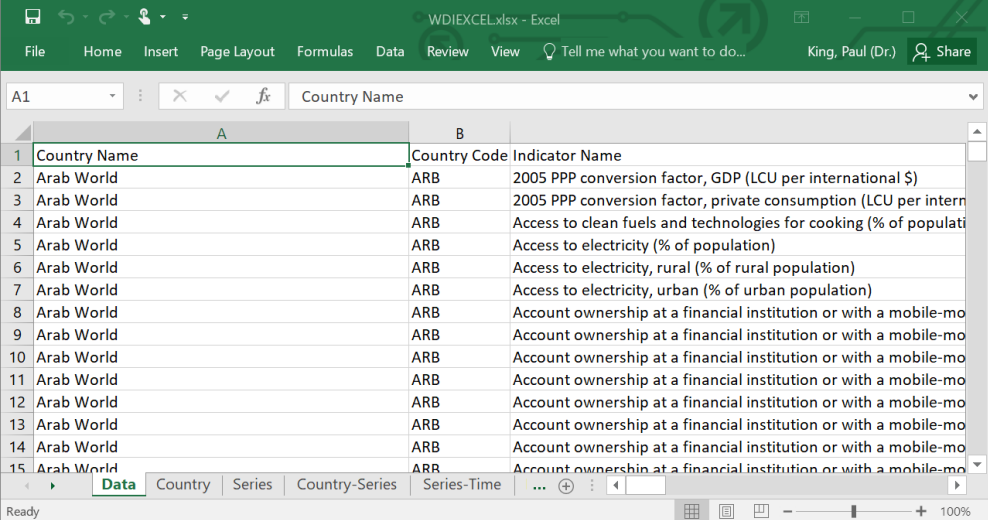Import → Tidy → Transform → Visualise / Model (Understand) → Communicate

Program

# Data science skills

# Case study: from this…

- Clean, structured data
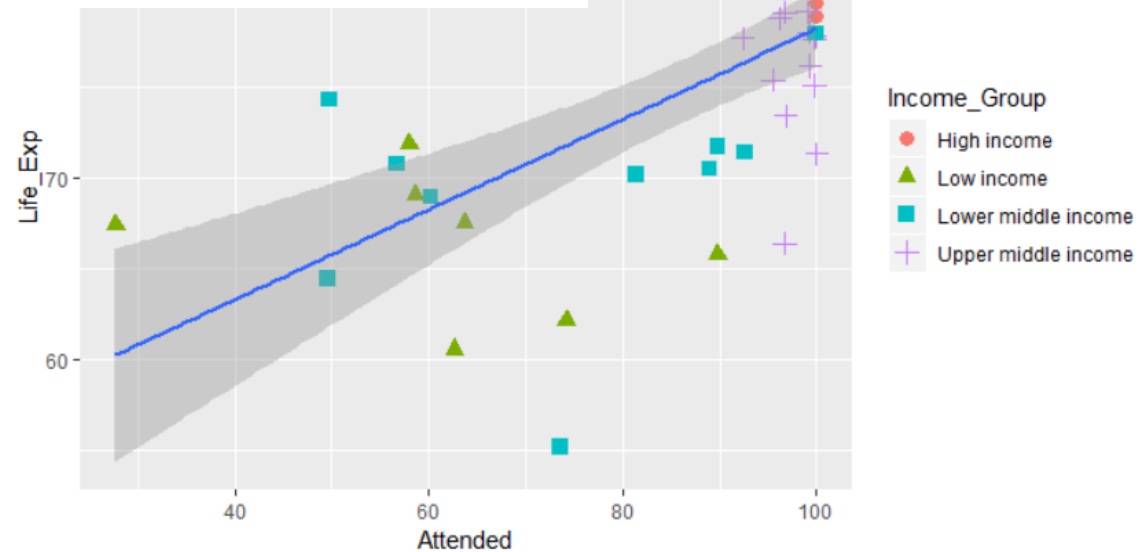- 5 sheets
- 264 country codes
- 1599 indicators
- 422,136 rows

# Case study: to this (reproducibly)…

| Income_Group<br><fctr> | Region<br><fctr> | Av_Fem_Life_Exp<br><dbl> |
|---|---|---|
| High income | East Asia & Pacific | 83.61193 |
| Upper middle income | East Asia & Pacific | 77.08483 |
| Low income | East Asia & Pacific | 75.07100 |
| Lower middle income | East Asia & Pacific | 71.70246 |
| High income | Europe & Central Asia | 83.34810 |
| Upper middle income | Europe & Central Asia | 77.75850 |
| Lower middle income | Europe & Central Asia | 75.48883 |
| Low income | Europe & Central Asia | 74.18700 |
| High income | Latin America & Caribbean | 80.21331 |
| Upper middle income | Latin America & Caribbean | 77.25928 |

# Case study: skills required

- Read in data

- Change data type

- Select and filter

- Search using regular expressions

- Reshape

- Plot and categorise

- Group and summarise

# Case study

# Module outline

- Introduction & infrastructure. Reproducible workflows and collaborative working.

- Reading tabular files (CSV, Excel). Data structures: data frames and vectors. Simple plots.

- Tidy data: wide vs tall tables: pivoting

- Calculations on tabular data

- Visualizing data

- Putting it together - a first complete project

- Checking data; data ethics, governance, and regulation.

- Non-tabular data (XML, JSON, text)

- Working with databases: SQL and relational databases; noSQL types & uses

- Big data tools

# Tools and infrastructure

# Tools

- R

- Rstudio

- Rmarkdown

- Shiny

- Leaflet

- GitHub

- Bookdown / Blogdown

# Bibliography

- R for Data Science
  https://r4ds.had.co.nz/index.html

- Efficient R programming
  https://bookdown.org/csgillespie/efficientR/

- R Markdown: The Definitive Guide
  https://bookdown.org/yihui/rmarkdown/

- Geocomputation with R
  https://geocompr.robinlovelace.net/

- See the Bookdown site
  https://bookdown.org/

UNIVERSITY OF
LEICESTER

# Data Science for Actuaries

Paul King, ATRC 28 June 2019