**Rotman**

# INTRO TO R

R Workshop – Part 1 Overview & Basics / 1

Rotman School of Management
UNIVERSITY OF TORONTO

# Plan for the 4-Session Workshop

- Part 1: Overview & Basics (Session 1, 2)

- Part 2: Data Manipulation (Session 2, 3)

- Part 3: Data Visualization (Session 3)

- Part 4 - 1: Modeling Workflow (Session 4)
- Part 4 - 2: Time Series & Some Finance Applications (Session 4)

# Plan for Part 1

- Intro
  - What is R and what can R do?
  - Setup R
  - Motivation examples

- R programing and Data Science
  - Basics of R programming
  - Data science with R

- Learning Resources and Road Map
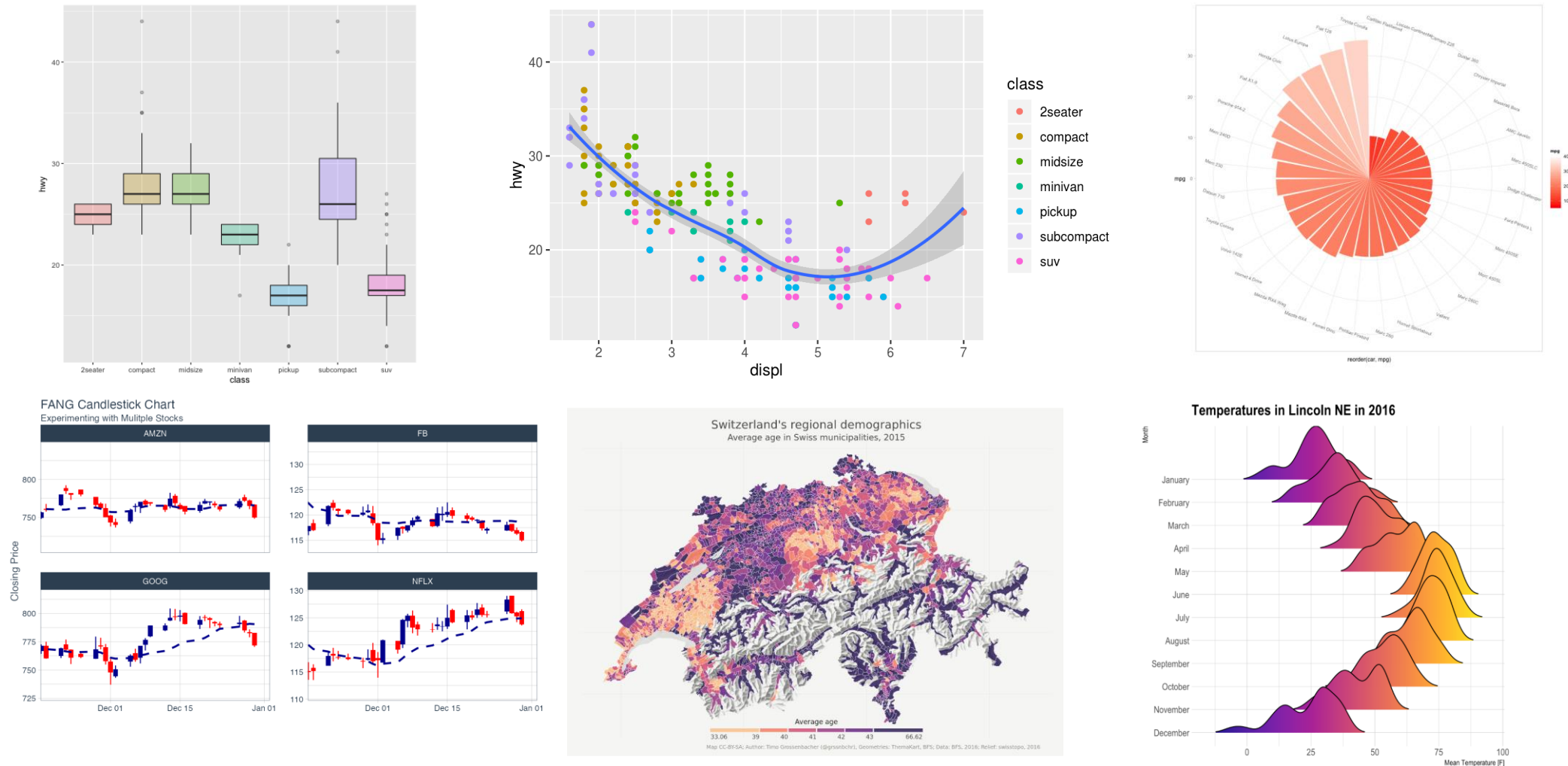
# What's R?



- R = a language + an eco-system
  - A free and open-source programming language
  - An eco-system of many high-quality user-contributed libraries/packages

- In the past R is mostly known for its statistical analysis toolkits

- Nowadays R is capable of (and very good at) many other tasks
  - Tools that facilitate the whole data analysis workflow
  - Tools for web technology
  - Many more…

# What can R do – Statistics & related

- Statistics & Econometrics
  - Regressions
  - Time series analysis
  - Bayesian inference
  - Survival analysis
  - …

- Numerical Mathematics
  - Optimization
  - Solver
  - Differential equations
  - …

- Finance
  - Portfolio management
  - Risk management
  - Option pricing
  - …
- …

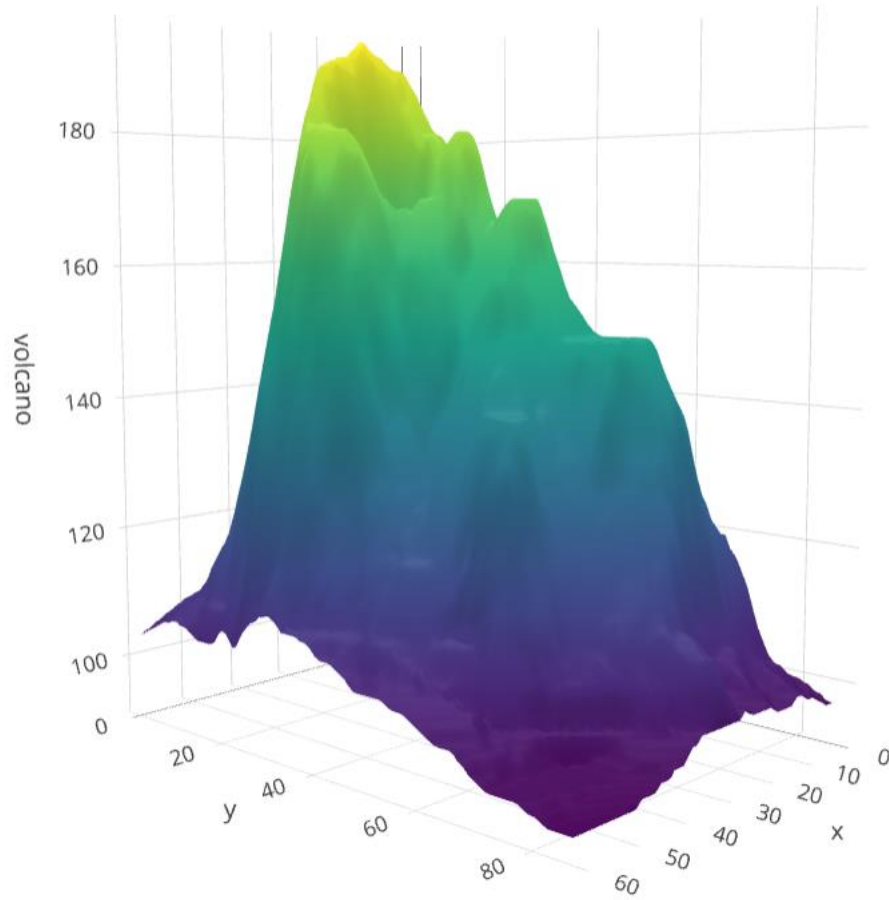See more R Packages on R Task View and more Empirical Finance Packages on R Task View - Finance
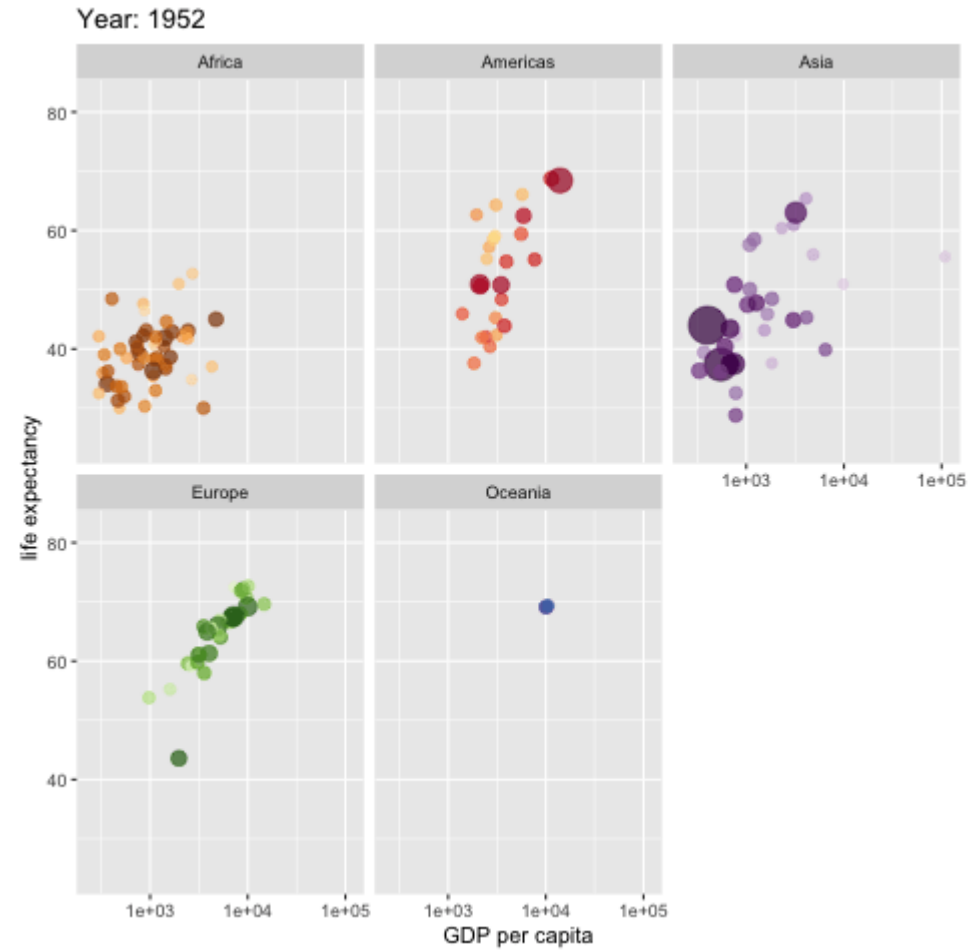
# What can R do – Graphics (static ones)



Ref. 1) https://www.r-graph-gallery.com/
2) https://timogrossenbacher.ch/bivariate-maps-with-ggplot2-and-sf/

# What can R do – Graphics (dynamic ones)



https://plot.ly/r/3d-surface-plots/;

https://gganimate.com/;

# What can R do – ML & NLP

- Machine learning
  - Statistical learning (clustering, decision tree, etc.)
    - [An Introduction to Statistical Learning (with Applications in R)](#)

  - Deep learning (neural networks)
    - [Tensorflow for R](#) (via [reticulate](#), an R to Python interface)
    - [R interface to Keras](#)
    - [Torch for R](#) (natively from R; similar as PyTorch in Python)

- Natural language processing (including LLM)
  - Packages (e.g., [tidytext](#), [topicmodels](#), [ellmer](#))
  - Books (e.g., [Text Mining with R](#), [Supervised ML for Text Analysis in R](#))
  - Leveraging the Python Transformers library (e.g., [Transformers from R](#))

1. See more R Machine Learning Packages on [R Task View - ML & Statistical Learning](#)
2. See more R Natural Language Processing Packages on [R Task View - NLP](#)

# What can R do – Web & Reporting

- Web technology
  - Web scraping (e.g., rvest)
  - API wrapper (e.g., Bluesky: bskyr; Bigquery: bigrquery; Nasdaq Data Link: Quandl)
  - Shiny web app (https://shiny.posit.co/)


- Reporting
  - R Markdown (write reports, slides, blogs, books, etc. See a gallery here.)
  - Quarto (new authoring tool; multi-language and multi-engine;)


- … (see **R Task View** for more)

# R vs Excel and BI Tools vs Python



- Excel & Business Intelligence (BI) Tools (e.g., Tableau, Power BI, etc.)
  - 2-D tables as basic data structure
  - Good UI (User Interface) and minimum programming
  - Limited modeling tools
  - Not easy to reproduce an analysis (because it's hard to store UI clicks)
  - Not flexible enough for complicated analytics problems, i.e., problems with
    - Many data cleaning steps/pipelines
    - Many different models to try

- Python
  - Python is more general purpose; R is more specialized in statistical analysis
  - R is much easier to learn (in my opinion)
  - Recall that you use Python from R
    - R reticulate package provides the interoperability.

# Why learn R (What can R do **for You**)?

- Beyond Excel Data Analysis
  - I wish Excel could…

- Automate boring repeating tasks
  - e.g., daily data collection from different sources, weekly dashboard update

- Prototype ideas
  - e.g., a novel trading strategy, a new credit risk model

- Really, find anything that interests you and use R…

# Plan for Part 1

- <mark>Intro</mark>
  - What is R and what can R do?
  - <mark>Setup R</mark>
  - Motivation examples

- Overview of R programing and Data Science
  - Basics of R programming
  - Data science with R

- Learning Resources and Road Map

# Setup R (Install R & its Coding Environment)

| | R & RStudio | R & Notebook |
|---|---|---|
| Run locally (i.e., on your laptop) | • **Install R (https://www.r-project.org/)**<br><br>**Our Choice**<br><br>• **Install RStudio (https://posit.co/download/rstudio-desktop/)** | • Install R (https://www.r-project.org/)<br><br>• Install RStudio or Jupyter Notebook (https://jupyter.org/) |
| Run in the cloud | • Option 1: RStudio Cloud (https://posit.cloud/)<br><br>• Option 2: UofT JupyterHub RStudio (https://datatools.utoronto.ca/) | **Our Choice**<br>• **Option 1: Google Colab (https://colab.research.google.com/)**<br><br>• Option 2: UofT JupyterHub Notebook (https://datatools.utoronto.ca/) |

# What's RStudio? ⓡ Studio®

# RStudio Cloud

# RStudio at UofT Jupyterhub

# R Notebook in Google Colab

# R Notebook at UofT Jupyterhub

# Plan for Part 1

- <mark>Intro</mark>
    - What is R and what can R do?
    - Setup R
    - <mark>Motivation examples</mark>

- Overview of R programing and Data Science
    - Basics of R programming
    - Data science with R

- Learning Resources and Road Map

# A Few Examples

- Analyze portfolio performance

- Perform sentiment analysis on earning call transcripts
  - Web scraping
  - Sentiment dictionary vs Language model

- Make Web API calls to retrieve data
  - Get cryptocurrency trading data from CoinGecko

- Recognize handwritten digits
  - an example of deep learning

**PerformanceAnalytics Package**

# A Few Examples: What to Look For

- Focus on analysis workflow (by reading the code comments)
  - Import and manipulate data
  - Model data
  - Report and visualize results

- Don't focus on R syntax
  - By the end of the workshop, you will be able to understand the code

- Do notice everything is done in a sequential way
  - no conditional branching or looping