

Data Science for Business Analytics

Introduction

HEC Lausanne

Aleksandr Shemendyuk [aleksandr.shemendyuk@unil.ch]

Agenda

1 Introduction

2 Organisation

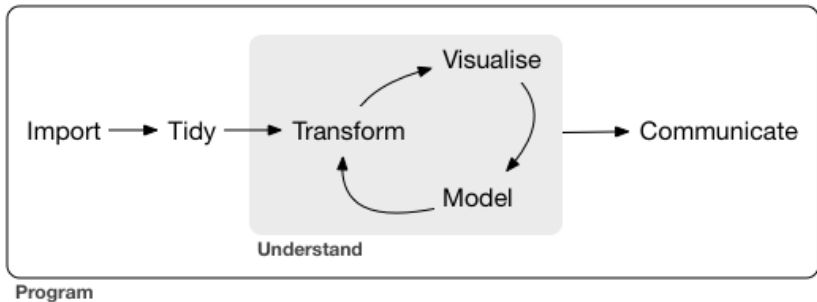
3 R

A little about me

- Born and raised in Riga, Latvia
- **Education:**
 - ▶ B.Sc. in Applied Mathematics 2017 (HSE University, Moscow)
 - ▶ M.Sc. in Statistical Modelling 2019 (HSE University, Moscow)
 - ▶ Ph.D. in Actuarial Science 2024 (HEC Lausanne, UNIL)
- **Experience:**
 - ▶ Teaching assistant at HSE University and HEC Lausanne
 - ▶ Research using real data
 - ▶ Telegram Bot development in Python
 - ▶ AI in research and personal projects
 - ▶ AI automation
- **Hobbies:**
 - ▶ Chess
 - ▶ Guitar
 - ▶ Finance

What you will learn

- Workflow of the 'modern' Data Scientist



Statistical computing & data science

- What's the difference between data science and statistics?

“A data scientist is just a sexier word for a statistician.”

— Nate Silver (outdated)

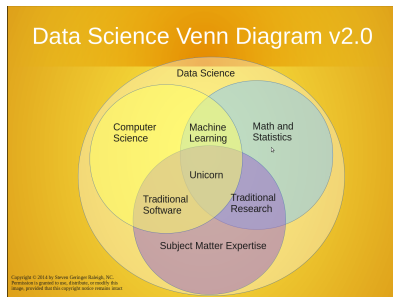
“A data scientist is a better computer scientist than a statistician and is a better statistician than a computer scientist.”

— Unknown (still accurate)

- What does a data scientist do?
 - ▶ *Jack of all trades.*
 - ▶ Transform data into valuable information!
 - ▶ A data scientist spends a significant portion of time processing data and less time modelling data.

What is Data Science?

- **Wikipedia:** “the extraction of knowledge from data”
- The precise definition is a bit unclear and controversial...
- Practitioners “agree” on the components of data science:
 - ▶ database management
 - ▶ gathering and cleaning
 - ▶ exploratory analysis
 - ▶ predictive modelling
 - ▶ data summary and visualization



Applications



Figure 1: Some of the hiring partners of *The Data Incubator*

- E-marketing
- Recommender systems
- Sport analytics
- Biotechnology
- Image or speech recognition
- Fraud and risk detection
- Social media
- redbfit scoring
- E-commerce
- Government analysis
- Gaming
- Price comparisons
- Airline routes planning
- Delivery logistics

Agenda

1 Introduction

2 Organisation

3 R

Course Description

- **Lectures** (10:15-12:00):
 - ▶ Focus on introducing the concepts.
 - ▶ Details are sometimes left for you to read up on
- **Exercise sessions** (12:15-14:00):
 - ▶ Work on assignments and project.
- **Instructor:** Professor Alex (Aleksandr Shemendyuk).
- **Captains** (Teaching assistants):
 - ▶ Ilia Azizi
 - ▶ Léo Wenger
 - ▶ Elwin Freudiger

Grading

- The **Final grade** is composed based on two equally weighted components:
 - ▶ Online Exam
 - ▶ Group Project
- **Exam:**
 - ▶ Online closed/**open**-book exam (50%)
- **Group Project:**
 - ▶ Project proposal (5%)
 - ▶ Project update (5%)
 - ▶ Video presentation (15%)
 - ▶ Final report (25%)

For the group project:

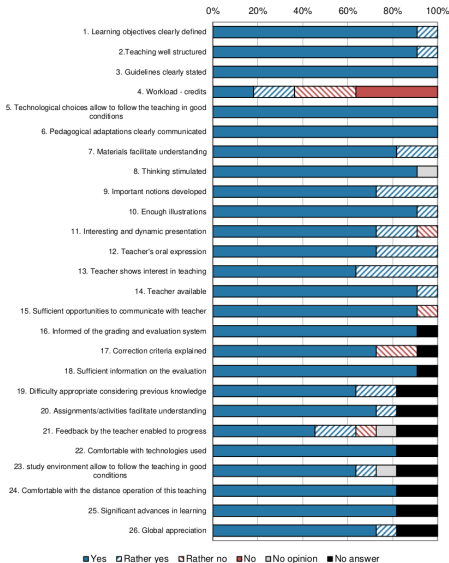
- Groups of 3 (or 4) members
- More details to come...

Milestones

Date	Assignment
October 14	Project proposal
November 11	Project update
December 16	Video presentation
December 23	Final Report

Course Evaluations in Fall 2022

- **Workload - credits:** ~33% of students found it too high!



Schedule outline

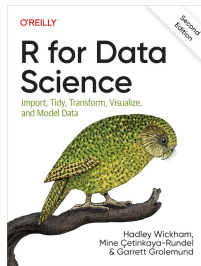
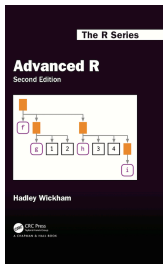
Date	Topic	Reading
September 23	Introduction	
September 30	Data Structures and Subsetting	ADVR 3+4
October 7	Control Flows and Functions	ADVR 5, 6, 9, 11
October 14	Data Wrangling	R4DS 4, 6, 8, 13-20
October 21	Visualisation	R4DS 2, 10-12
October 28	R Projects and Quarto	R4DS 3, 5, 7, 9, 29, 30
November 4	Project Coaching	
November 11	Data Wrangling	R4DS 4, 6, 8, 13-20
November 18	Visualisation	R4DS 2, 10-12
November 25	Project Coaching	
December 2	Presentations / Dashboards / Interactivity	R4DS 29+30, htmlwidgets
December 9	Projects Coaching	
December 16	Projects Presentations	

(numbers in the third column are book chapters)

Course website

- Course website:
 - ▶ <https://moodle.unil.ch/course/view.php?id=32207>
 - ▶ Registration key: **dsfba2024**
 - ▶ **A discussion forum**
 - ▶ PDFs of the lectures
 - ▶ Additional resources

Additional resources



- Books:
 - ▶ Advanced R
 - ▶ R for data science
- Additionally:
 - ▶ Rstudio cheat sheets
 - ▶ The CRAN website

Best place to look for answers?



... and our forum on Moodle.

Agenda

1 Introduction

2 Organisation

3 R

S and R

- S
 - ▶ A statistical programming language
 - ▶ First appeared in 1976
 - ▶ Developed by John Chambers and (in earlier versions) Rick Becker and Allan Wilks of Bell Labs
 - ▶ John Chambers, *[the aim is] to turn ideas into software, quickly and faithfully*
- R
 - ▶ Modern implementation of S
 - ▶ First appeared in 1993
 - ▶ Created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand
 - ▶ Currently developed by the *R Development Core Team*

Some “technical” details about R

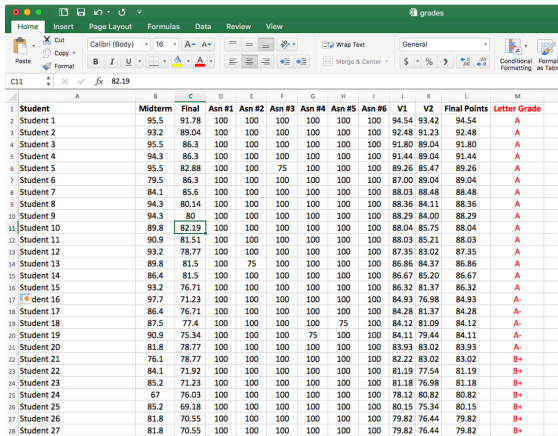
- Technical features:
 - ▶ **Available for Windows, macOS, and Linux**
 - ▶ Source code written primarily in C, Fortran, and R
 - ▶ Multi-paradigm: object-oriented, functional, procedural
 - ▶ Dynamically typed
 - ▶ Scripting language (interpreted)
 - ▶ **Wide variety of statistical and graphical techniques**
 - ▶ **Easily extensible through functions and packages**
 - ▶ **Read/write from/to various data sources**

What about Excel?



source: fantasyfootballanalytics.net

Excel is great for certain things...



Student	Midterm	Final	Asn #1	Asn #2	Asn #3	Asn #4	Asn #5	Asn #6	V1	V2	Final Points	Letter Grade
Student 1	95.5	91.78	100	100	100	100	100	100	94.54	93.42	94.54	A
Student 2	93.2	89.04	100	100	100	100	100	100	92.48	91.23	92.48	A
Student 3	95.5	86.3	100	100	100	100	100	100	91.80	89.04	91.80	A
Student 4	94.3	86.3	100	100	100	100	100	100	91.44	89.04	91.44	A
Student 5	95.5	82.88	100	100	75	100	100	100	89.26	85.47	89.26	A
Student 6	79.5	86.3	100	100	100	100	100	100	87.00	89.04	89.04	A
Student 7	84.1	85.6	100	100	100	100	100	100	88.03	88.48	88.48	A
Student 8	94.3	80.14	100	100	100	100	100	100	88.36	84.11	88.36	A
Student 9	94.3	80	100	100	100	100	100	100	88.29	84.00	88.29	A
Student 10	89.8	82.19	100	100	100	100	100	100	88.04	85.75	88.04	A
Student 11	90.9	81.51	100	100	100	100	100	100	88.03	85.21	88.03	A
Student 12	93.2	78.77	100	100	100	100	100	100	87.35	83.02	87.35	A
Student 13	89.8	81.5	100	75	100	100	100	100	86.86	84.37	86.86	A
Student 14	86.4	81.5	100	100	100	100	100	100	86.67	85.20	86.67	A
Student 15	93.2	76.71	100	100	100	100	100	100	86.32	81.37	86.32	A
Student 16	97.7	71.23	100	100	100	100	100	100	84.93	76.98	84.93	A-
Student 17	86.4	76.71	100	100	100	100	100	100	84.28	81.37	84.28	A-
Student 18	87.5	77.4	100	100	100	100	75	100	84.12	81.09	84.12	A-
Student 19	90.9	75.34	100	100	100	100	75	100	84.11	79.44	84.11	A-
Student 20	81.8	78.77	100	100	100	100	100	100	83.93	83.02	83.93	A-
Student 21	76.1	78.77	100	100	100	100	100	100	82.22	83.02	83.02	B+
Student 22	84.1	71.92	100	100	100	100	100	100	81.19	77.54	81.19	B+
Student 23	85.2	71.23	100	100	100	100	100	100	81.18	76.98	81.18	B+
Student 24	67	76.03	100	100	100	100	100	100	78.12	80.82	80.82	B+
Student 25	85.2	69.18	100	100	100	100	100	100	80.15	75.34	80.15	B+
Student 26	81.8	70.55	100	100	100	100	100	100	79.82	76.44	79.82	B+
Student 27	81.8	70.55	100	100	100	100	100	100	79.82	76.44	79.82	B+

source: github.com/jdwilson4

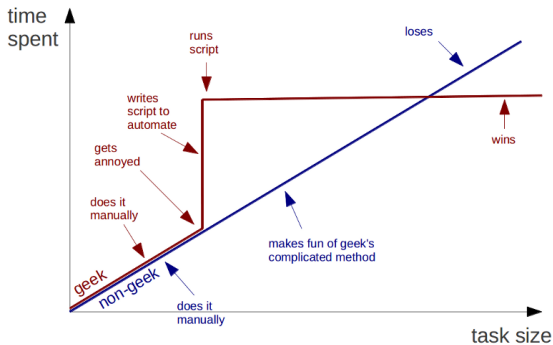
...but not everything

R's advantages:

- **Easier automation**
- **Better reproducibility**
- Faster computation
- Supports larger data sets
- Reads any type of data
- More powerful data manipulation capabilities
- Easier project organisation
- Easier to find and fix errors
- Free & open source
- Advanced statistics capabilities
- State-of-the-art graphics
- Runs on many platforms
- Anyone can contribute packages to improve its functionality

Automation and reproducibility

Geeks and repetitive tasks



source: trendct.org

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

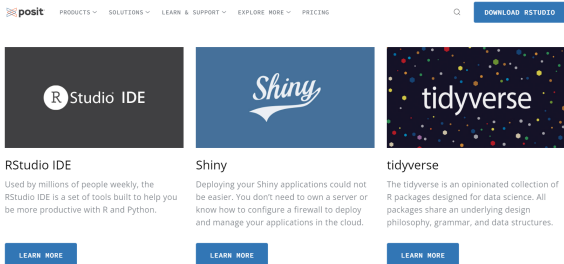
- The latest release (2017-11-30, Kite-Eating Tree) [R-3.4.3.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features and bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)

Questions About R

- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

source: cran.r-project.org

- An open-source integrated development environment (IDE)
- RStudio Desktop available for Windows, macOS, and Linux



The screenshot shows the Posit website header with navigation links: PRODUCTS, SOLUTIONS, LEARN & SUPPORT, EXPLORE MORE, and PRICING. A search icon and a "DOWNLOAD RSTUDIO" button are also visible. Below the header, three product cards are displayed:

- RStudio IDE**: A dark card with the RStudio logo. Description: "Used by millions of people weekly, the RStudio IDE is a set of tools built to help you be more productive with R and Python." A "LEARN MORE" button is at the bottom.
- Shiny**: A blue card with the Shiny logo. Description: "Deploying your Shiny applications could not be easier. You don't need to own a server or know how to configure a firewall to deploy and manage your applications in the cloud." A "LEARN MORE" button is at the bottom.
- tidyverse**: A dark card with the tidyverse logo. Description: "The tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures." A "LEARN MORE" button is at the bottom.

source: posit.co

Base R

- What is Base R?

“The package named ‘base’ is in a way the core of R and contains the basic functions of the language, particularly, for reading and manipulating data.”

— *R for Beginners*, Emmanuel Paradis

- Base R includes all default code for performing common data manipulation and statistical tasks.
- You might recognize some Base R functions:
 - ▶ `mean()`, `median()`, `lm()`, `summary()`, `sort()`
 - ▶ `data.frame()`, `read.csv()`, `cbind()`, `grep()`, `regexpr()`
 - ▶ Many many more...
- If you don't recognize any Base R functions, don't worry!

The tidyverse

- Common criticisms of Base R:
 - ▶ The code doesn't flow as well as other languages.
 - ▶ Function names/arguments are often inconsistent/confusing.
 - ▶ Base R functions sometimes don't return type-stable objects.
 - ▶ Base R functions are not refined to run as fast as possible.
 - ▶ Other complaints exist...
- So what is the **tidyverse**? A collection of R packages
 - ▶ designed for data science,
 - ▶ sharing an underlying design philosophy, grammar, and data structures.
- Often, it performs the same tasks as Base R, but:
 - ▶ Relies on a **pipe** operator to help with the flow of the code.
 - ▶ More descriptive function names and consistent inputs.
 - ▶ Type-stable.
 - ▶ Often faster than common Base R functions.

Core tidyverse packages

- `ggplot2`: declarative graphics, based on The Grammar of Graphics.
- `dplyr`: grammar of data manipulation.
- `tidyr`: functions that help you get to tidy data.
- `readr`: reading in rectangular data.
- `purrr`: enhancing R's functional programming (FP).
- `tibble`: a modern `data.frame`.
- `stringr`: makes working with strings as easy as possible.
- `forcats`: useful tools for working with factors.

See more on the tidyverse website.

Base R vs tidyverse

- Why ever use Base R?
 - ▶ Gets the job done!
 - ▶ To become an expert, you have to know Base R.
 - ▶ Some Base R functions are very common/useful, e.g., `mean()`.
- What should you learn first? Base R or tidyverse?
 - ▶ Some believe you should learn Base R first, others the tidyverse first.
 - ▶ Lately, more are shifting to tidyverse...

Install links

- R:
 - ▶ MacOS
 - ▶ Be mindful of the difference between Apple silicon (M1/M2) Macs and older Intel Macs
 - ▶ Windows
- RStudio Desktop:
 - ▶ MacOS 12+
 - ▶ Windows 10/11