UNIVERSITY OF WISCONSIN–MADISON

Center for Limnology
University of Wisconsin–Madison

# A Short Introduction into Data Science with R

Robert Ladwig
University of Wisconsin-Madison

rladwig2@wisc.edu

@hydrobert

# Disclaimer

- **material is based on the AWESOME workshop by Rachel Pilla (now postdoc at Oak Ridge National Lab, @rmpilla) "Introduction to R" (*Computer Science in Modern Biology* at Miami University, Ohio)**
- **with help by Alva Strand and Nicole Berry**
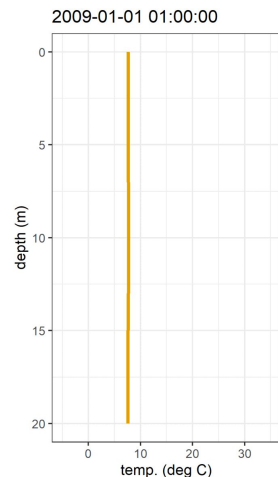
# About myself

- **Physical limnologist**
    - use mathematical models to explore mixing dynamics and water quality in lakes


- developing and applying open-source and open-access software in **R, Python and Matlab**
- background in civil engineering and geology

**Aquatic ecosystem modeling**

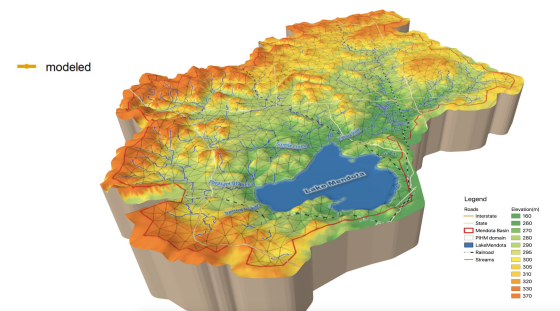mechanistic understanding of cause and effects,

2009-01-01 01:00:00



$$\frac{\partial T}{\partial t} = \frac{1}{A}\frac{\partial}{\partial z}\left(A(v_t^{\cdot} + v^{\cdot})\frac{\partial T}{\partial z}\right) + \frac{1}{\rho_0 c_p}\frac{\partial H_{sol}}{\partial z} + \frac{dA}{dz}\frac{H_{geo}}{A\rho_0 c_p}$$

$$\frac{\partial [O_2]}{\partial t} = NEP = GPP - ER \pm D$$

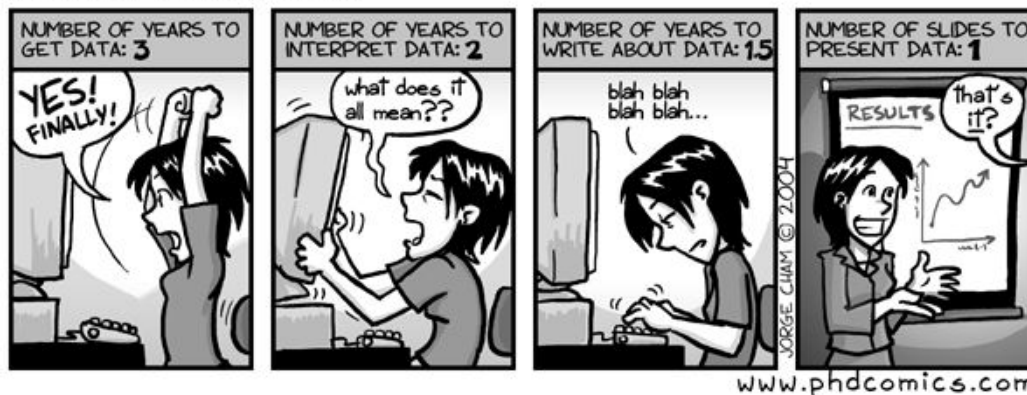$$\frac{\partial P}{\partial t} = D\frac{\partial^2 P}{\partial z^2} + \lambda P(1 - \frac{P}{K})$$
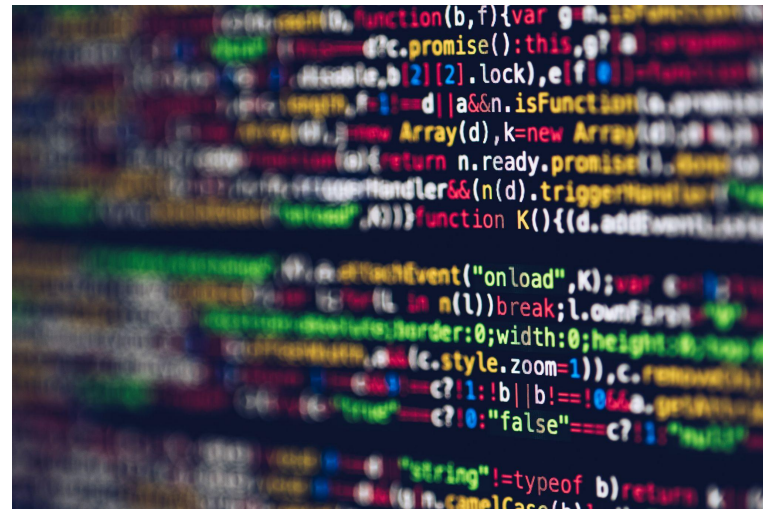
# What's data science?

- combining math & statistics, programming and machine learning to explore any kind of data
- data is often noisy or unstructured → cleaning is big part of the job
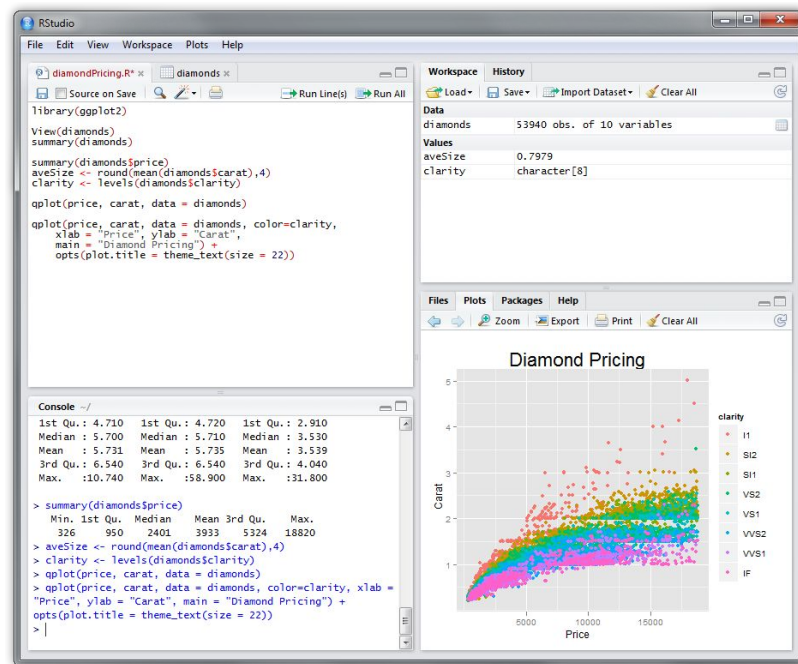- buzzwords like big data

# What's the plan for today?

- short introduction into the R scripting language
  - all material on github: https://github.com/robertladwig/introdatasciencer
  - GitHub: uses Git (version control), fantastic for coding and collaborations
- hands-on coding to introduce you to data analysis
- **at the end**: able to know basic R commands, load data, analyse data, plot data



Markus Spiske

# Install R and RStudio

- Install R from here (language):
  https://www.r-project.org/
- Install RStudio (editor and GUI):
  https://rstudio.com/products/rstudio/download/#download

# What's R?

- FREE and OPEN SOURCE statistical and computational software

Most Popular Programming Languages for Data Science (KDnuggets Software Poll)

# What's R?

- FREE and OPEN SOURCE statistical and computational software
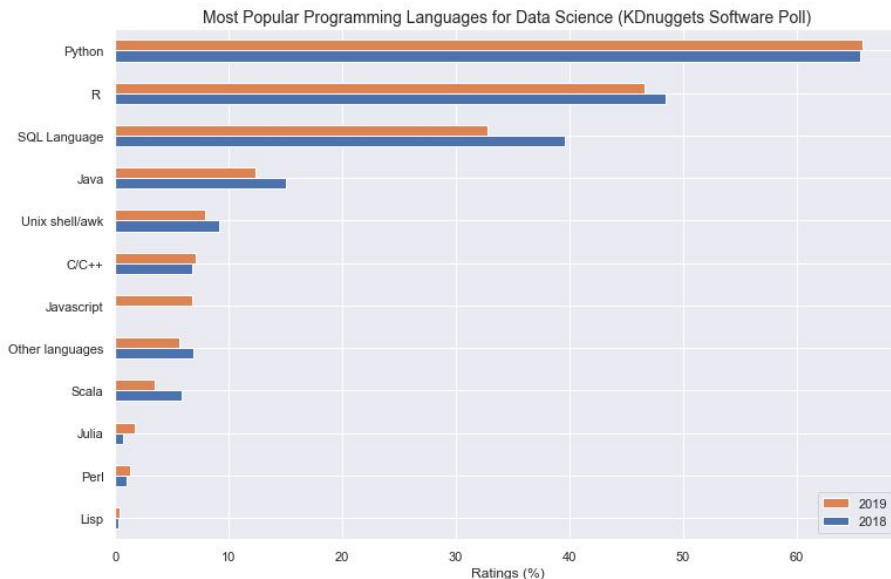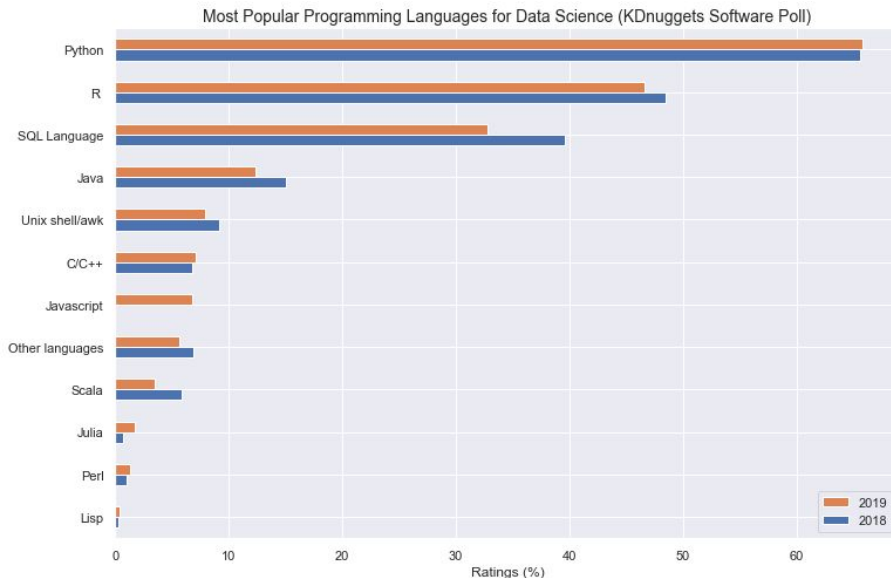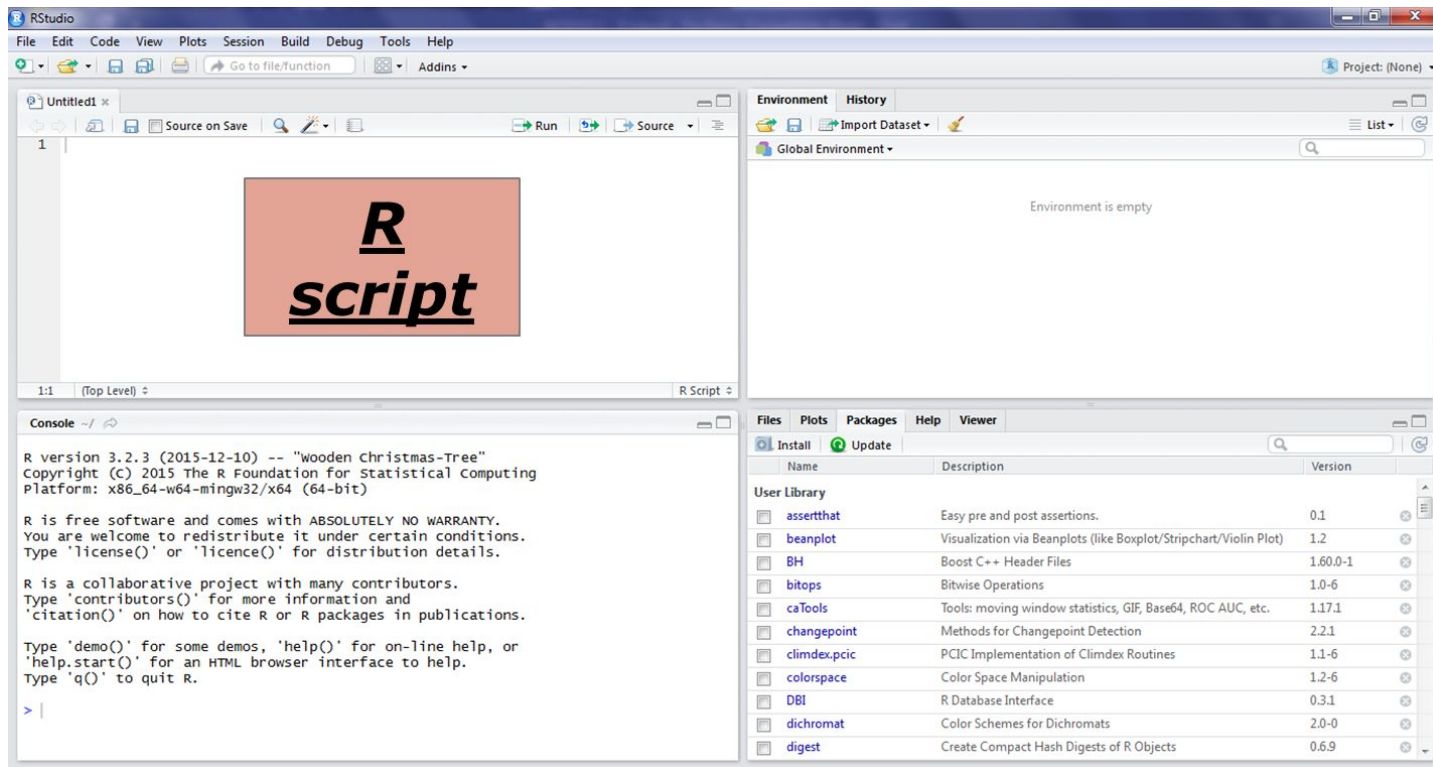- big community: easy to find solutions and troubleshooting online

Most Popular Programming Languages for Data Science (KDnuggets Software Poll)

# What's R?

- FREE and OPEN SOURCE statistical and computational software
- big community: easy to find solutions and troubleshooting online
- widely used in sciences (esp. statistics and visualization) and rapidly growing in popularity

Most Popular Programming Languages for Data Science (KDnuggets Software Poll)

# What's R?

- FREE and OPEN SOURCE statistical and computational software
- big community: easy to find solutions and troubleshooting online
- widely used in sciences (esp. statistics and visualization) and rapidly growing in popularity
- can handle more advanced computations, statistical analyses and bigger data files than Excel

Most Popular Programming Languages for Data Science (KDnuggets Software Poll)

# What's R?

- FREE and OPEN SOURCE statistical and computational software
- big community: easy to find solutions and troubleshooting online
- widely used in sciences (esp. statistics and visualization) and rapidly growing in popularity
- can handle more advanced computations, statistical analyses and bigger data files than Excel
- lots of styles for coding

Most Popular Programming Languages for Data Science (KDnuggets Software Poll)

| Language | |
|---|---|
| Python | |
| R | |
| SQL Language | |
| Java | |
| Unix shell/awk | |
| C/C++ | |
| Javascript | |
| Other languages | |
| Scala | |
| Julia | |
| Perl | |
| Lisp | |

2019
2018

Ratings (%)

# Open RStudio on your laptop now



**R script:**

write and edit code

color-coding for easier identification

automatically fills in parentheses and quotations

# Open RStudio on your laptop now



**R console:**

commands are run and results appear

> means 'ready to work'

# Open RStudio on your laptop now



**Environment:**

see which objects were created

type of objects (double, integers, etc.), size and dimensions

# Open RStudio on your laptop now



**Plots/Packages/ Help:**

plots will appear

see which packages are installed

help functions

15

# What are packages?

- bundles of tools and functions that others have developed to be used in R

# What are packages?

- bundles of tools and functions that others have developed to be used in R
- grouped to specific types of functions, analyses or datasets

# What are packages?

- bundles of tools and functions that others have developed to be used in R
- grouped to specific types of functions, analyses or datasets
- you can also create your own packages for specific tasks!



Cottonbro

# What are packages?

- bundles of tools and functions that others have developed to be used in R
- grouped to specific types of functions, analyses or datasets
- you can also create your own packages for specific tasks!
- currently, above 16,000 packages are available



Cottonbro

# Install one package



we need:

**tidyverse**

# Install one package



we need:

**tidyverse**

# Install one package



we need:

**tidyverse**

# Installing a package

- packages only need to be installed ONCE
- but, if you want to use it you'll need to load it:

  ```
  library(package name)
  ```

Hadley Wickham ✓
@hadleywickham

Replying to @ijlyttle

@ijlyttle a package is a like a book, a library is like a library; you use library() to check a package out of the library #rsats

8:34 AM · Dec 8, 2014 · Echofon

# Package information



under packages,
click on **dplyr**

# Package information



this lists all available functions

click on any function to see the help file

# Coding in R

- write commands in the script (upper left panel)
    - save it, edit it, revisit it later, etc.

# Coding in R

- write commands in the script (upper left panel)
  - save it, edit it, revisit it later, etc.
- code NOT automatically run when you hit Enter

# Coding in R

- write commands in the script (upper left panel)
    - save it, edit it, revisit it later, etc.
- code NOT automatically run when you hit Enter
- to run it:
    - "Run" button in upper right corner
    - "CTRL + Enter" (Windows)
    - "Command + Enter" (Mac)

# Coding in R



- add comments using #
  - additional information
  - will be ignored

# Coding in R



- add comments using #
    - additional information
    - will be ignored
- R is case sensitive
    - "Mean" ≠ "mean"

# Coding in R



- add comments using #
  - additional information
  - will be ignored
- R is case sensitive
  - "Mean" ≠ "mean"
- R doesn't care about spaces and tabs

# Coding in R



- add comments using #
  - additional information
  - will be ignored
- R is case sensitive
  - "Mean" ≠ "mean"
- R doesn't care about spaces and tabs
- need to close all parentheses and quotations

# Key components

- **functions**: allow you to manipulate data, apply calculations, run statistical analysis, much more!

# Key components

- **functions**: allow you to manipulate data, apply calculations, run statistical analysis, much more!
- **arguments**: defining information for functions, "customize" it

# Key components

- **functions**: allow you to manipulate data, apply calculations, run statistical analysis, much more!
- **arguments**: defining information for functions, "customize" it
- **objects**: pieces of data saved in R, can be called up, reused and manipulated

# Key components

- **functions**: allow you to manipulate data, apply calculations, run statistical analysis, much more!
- **arguments**: defining information for functions, "customize" it
- **objects**: pieces of data saved in R, can be called up, reused and manipulated

# Types of data

- **objects (data types):**
    - numeric (2.334, 3.14159)
    - integer (2, 2405, 54)
    - characters ("hello world")
    - logical (TRUE/FALSE)
    - complex (1+4$i$)

- **data structures:**
    - vector (1-D object with same data type)
    - matrix (2-D object with same data type
    - data frame (2-D object with different data types per column if needed, very useful!)

# Live Coding time!