

Curating Code

Talya Cooper
New York University

Greg Janée
University of California, Santa Barbara

Kay P Maye
Tulane University

Nick Ruhs
Florida State University

Seth Erickson
University of California, Santa Barbara

Case Study: Willoughby-Hoye Scripts

“While preparing a manuscript, to our surprise, attempts by team members to replicate these results produced different [results] despite using the same [source code] outlined by Willoughby et al.” ([Neupane et al., 2019](#))

Case Study: “Willoughby-Hoye Scripts”

```
def read_gaussian_outputfiles():
    list_of_files = []
    for file in glob.glob('*.*out'):
        list_of_files.append(file)
    return list_of_files
```

Excerpt of original python code from (Willoughby et al., 2014)

```
def read_gaussian_outputfiles():
    list_of_files = []
    for file in glob.glob('*.*out'):
        list_of_files.append(file)
    # ...
    list_of_files.sort()
    return list_of_files
```

Excerpt of revised code (Yuhen Luo and Rui Sun, 2019)

Is it realistic to expect a curator to catch this sort of mistake?

Case Study: “Willoughby-Hoye Scripts”

“This simple glitch in the original script calls into question the conclusions of a significant number of papers on a wide range of topics in a way that cannot be easily resolved from published information **because the operating system is rarely mentioned.**” (Neupane et al., 2019)

Learning Outcomes (pt.1)

By the end of this workshops, you should be able to ...

- Identify and explain real-world examples of how differences between **computing environments** lead to **reproducibility problems**.
- Navigate significant **differences between major operating systems** and understand how they can affect the reusability of software
- Identify several **common programming languages** used in research computing and the file types and tools associated with each.

Learning Outcomes (pt .2)

- Identify and document **dependencies** used in research software for common programming languages.
- Identify and improve different types of **documentation** associated with source code and research software.
- Describe differences between **software licenses** and non-software licenses, and differences between different kinds of open source software licenses.
- Implement commonly used **project organization strategies** used in the research community.

Key Terms

- **Source Code** - *plain text*, written in a programming language, that can be run on a computer using an interpreter or compiled into runnable machine code by a compiler.
 - Example: contents of a python/R “script”
- **Software** - encompasses source code, compiled programs, operating systems, libraries, dependencies, packages, scripts, etc..
 - Example: Linux or Windows operating systems; the NumPy package.
- **Source Code Repository** - a centralized location where code is stored and shared, usually, often conjunction with a version control system. (Not to be confused with a data repository!).
 - Example: GitHub and GitLab

Research Software

“Research Software includes **source code files, algorithms, scripts, computational workflows and executables that were created during the research process or for a research purpose**. Software components (e.g., operating systems, libraries, dependencies, packages, scripts, etc.) that are used for research but were not created during or with a clear research intent should be considered software in research and not Research Software. This differentiation may vary between disciplines. The minimal requirement for achieving computational reproducibility is that all the computational components (Research Software, software used in research, documentation and hardware) used during the research are identified, described, and made accessible to the extent that is possible.”
[\(Gruenpeter et al. 2021\)](#)

Overview

1. Introduction
2. Computing Platforms
3. Programming Languages
4. Dependencies
5. Documentation
6. Licensing *
7. Project Structure *
8. Exercise

* - not covered due to time constraints

Example Dataset

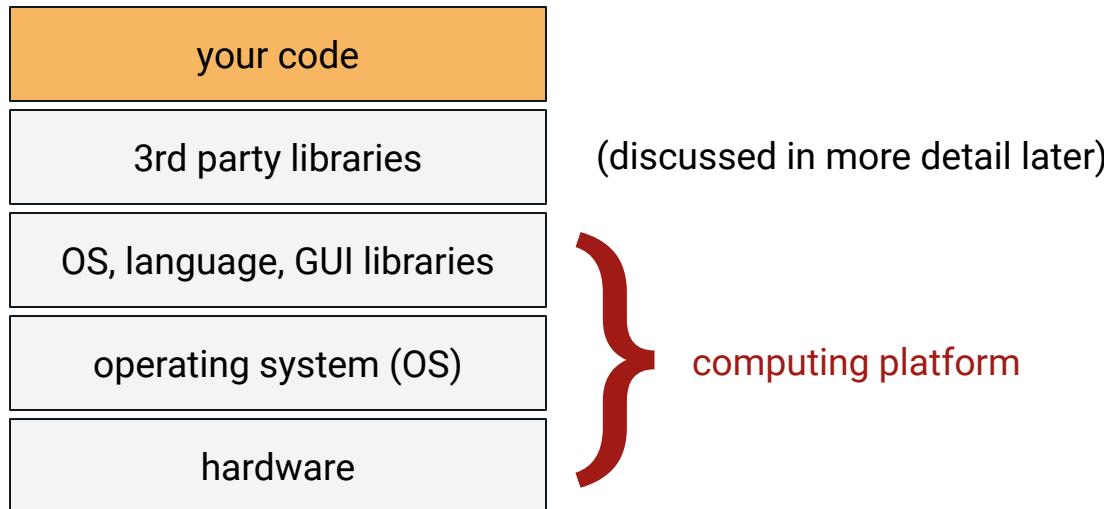
- *Code, data, and metadata document for the manuscript:*
“Density-dependence in wolf resource selection study designs”
- The submission includes data and analysis code from study of wolf habitat preference.
- We'll return to this example in each section of the workshop.
- See files in course materials.
- This is not an authentic dataset! (though it is derived from one)

2. Computing Platforms

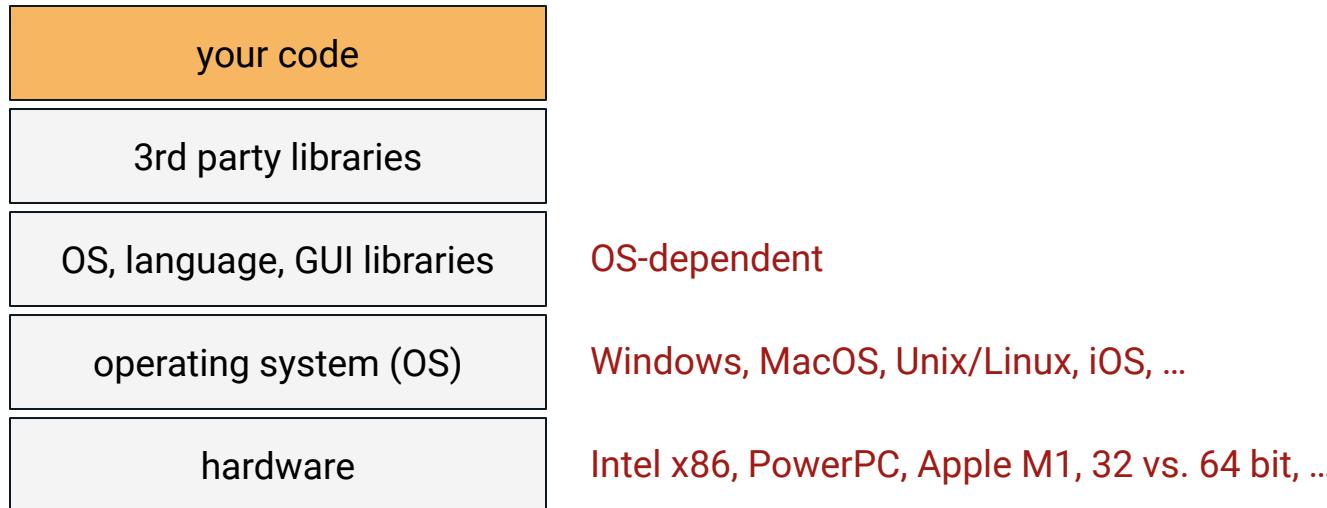
Learning Outcomes

- Understand elements of computing platforms
- Understand differences between compiled and interpreted languages
- Understand obstacles to reproducibility
- Understand implications for curation

Computing platforms



Central obstacle to reproducibility



Compiled vs interpreted languages

- Portability of code depends greatly on whether the code is written in a **compiled** language or in an **interpreted** language

Compiled languages

- Must be “compiled” into a “binary/executable/program/app”
 - Once compiled, the binary is entirely platform-specific and it is directly executable; the source code is not needed
- Compiled examples:
 - Languages: C, C++, FORTRAN, Go, Rust
 - Programs: most every program/app on your laptop or phone
- Code in compiled languages is often platform-dependent
- Compilation process is difficult

Interpreted languages

- Require a program called an “interpreter” to run
 - Code is directly executed from source code every time
- Examples: Python, R, Stata, SAS, MATLAB, Java, Perl, Mathematica
- Typically excellent portability

Language type comparison

	Compiled	Interpreted
Distribution form	Source code	Source code
Your responsibility	Compile code and 3rd party libraries	Load 3rd party libraries from central repository
OS provides (intrinsically or by add-on)	Compiler(s), package manager	Interpreter

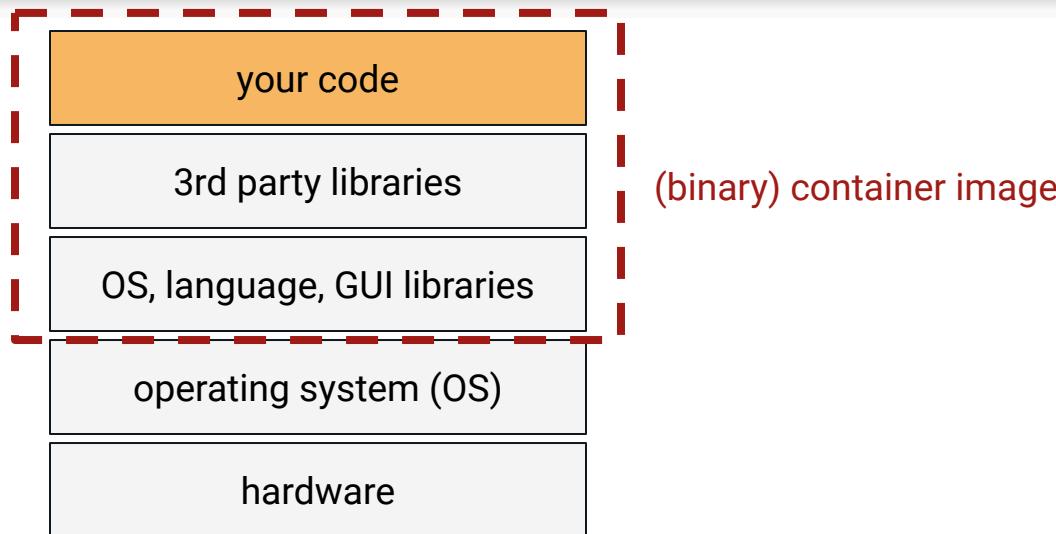
Where platform incompatibilities get addressed

	Compiled	Interpreted
Distribution form	Source code	Source code
Your responsibility	Compile code and 3rd party libraries	Load 3rd party libraries from central repository
OS provides (intrinsically or by add-on)	Compiler(s), package manager	Interpreter

Containerization

- Relatively new technology in libraries & curation
- A “container image” bundles a program with all its library dependencies, down to the operating system libraries
 - and possibly data as well

Computing platforms



Containerization features

- Images are run from/within a **container engine** (e.g., Docker)
- Images can be stored in repositories
 - Can serve as building blocks, have their own dependencies
- Images are opaque
- Images retain dependencies on OS, machine architecture

Language type comparison

	Compiled	Interpreted	Containerized
Distribution form	Source code	Source code	Container image
Your responsibility	Compile code and 3rd party libraries	Load 3rd party libraries from central repository	None
OS provides (intrinsically or by add-on)	Compiler(s), package manager	Interpreter	Container engine (that supports container's OS)

Implications for curation

- Interpreted code is highly portable
 - Document language and dependency versions
- Compiled code is likely **not** runnable without significant work
 - Document the environment it was developed in
- Container images: reproducible, but opaque
 - Capture build file

Four additional incompatibility sources

- Pathnames
- Character encoding
- Text file format
- Nondeterminism

Pathnames

- All platforms support similar concepts of **files** (named unit of information storage) and **directories** (named collection of files), organized hierarchically
- **Pathnames** are used for identification
 - /Users/alice/my_project/output_data/final.csv

Pathname incompatibilities

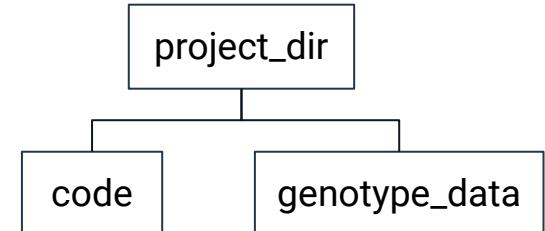
- Repertoire of allowable characters differs across platforms
- Case-sensitivity differs
 - README.txt, Readme.txt, REadME.TXT
- Pathname syntax
 - Unix: /Users/alice/my_project/output_data/final.csv
 - Windows: C:\Users\alice\my_project\output_data\final.csv

Absolute vs. relative pathnames

- Absolute pathname
 - Identifies starting from filesystem root
 - /Users/alice/my_project/output_data/final.csv
- Relative pathname
 - Identifies relative to another (contextual) directory
 - output_data/final.csv

Quiz: Which are good path/filenames?

- AC/DC Voltages.dat
- Buñuel Luis.xml
- /raid/data/spot/554/228/55422810503291849151J-9.tiff
- ../genotype_data/hap_G100.A1.O_X.R
- BT_Palmyra-2015.csv (look close)



Implications for curation

- Watch out for absolute pathnames
 - Actual example:
 - `setwd("C:\\Users\\soneil\\Documents") # change me`
- Encourage use of simple characters that transport well
 - Yes: A-Z a-z 0-9 . - _
 - No: é, – (en dash), (space), 😊

Other platform incompatibilities

- Character mis-encoding
 - “BenoÃÆt LalibertÃÂ©” instead of “Benoît Laliberté”
- Subtle differences in plain text file formats

Character encoding

- Text characters are internally represented as numbers; the assignment of characters to numbers is called an **encoding**
 - E.g., 65 = "A"
- Unicode UTF-8 is the *de facto* standard encoding today, but older, platform-specific encodings are still encountered

Encoding misinterpretation example

- Given “Benoît Laliberté”
- Encoded as UTF-8:

- 66 101 110 111 195 174 116 32 76 97 108 105 98 101 114 116 195 169
 - B e n o î t L a l i b e r t é

- Decoded as Windows-1252:
- Decoded as MacRoman:

- B e n o ÿ t L a l i b e r t é

- B e n o √ E t L a l i b e r t é

Implication for curation

- Watch out for character encoding problems
- Researchers don't always recognize, or know how to fix

Text file format

- Text files differ in how lines are delimited internally
 - Unix:
 - This is a line LF
 - So is this LF
 - MacOS
 - This is a line CR
 - So is this CR
 - Windows
 - This is a line CR LF
 - So is this CR LF

Implication for curation

- In most cases, interpreted languages (Python, R) automatically detect text file format and adapt automatically
- Incompatibilities still surface, though, and may cause errors

Nondeterminism

- Code is generally deterministic: same input produces same output every time
- Nondeterministic code: results can vary
 - Intended
 - Artifact of internal choices, ambiguity
 - A bug

Nondeterministic: intended

- Monte Carlo simulation



Image source: [NOAA, National Hurricane Center](https://www.nhc.noaa.gov/)

Nondeterministic: artifact

- Consider sorting these students in descending order by grade:
 - | Student | | Grade |
|---------|--|-------|
| Alice | | 90 |
| Bob | | 75 |
| Cindy | | 90 |
- “Alice,Cindy,Bob” or “Cindy,Alice,Bob”?

Nondeterministic: bug

- Code may depend on behavior that is not guaranteed
- Ex: documentation on Python's "os.listdir" function:
 - "Return a list containing the names of the entries in the directory given by path. **The list is in arbitrary order**, and does not include the special entries '.' and '..' even if they are present in the directory."

Implication for curation

- Code that produces (slightly) different results may be intentional or unintentional, may or may not be significant, may or may not reflect a bug
- Requires further investigation

Discussion: deposit example (1/2)

- Platform: R
 - Interpreted language, portability should be good
 - But documentation of platform could be better: version of R is documented in README, but library versions are missing
 - Consider using `renv` to automatically capture all version information
- Files and pathnames
 - Reads and writes CSV files, no issues there
 - But uses absolute pathname to set working directory
 - Line 48: `setwd ("C:\\Users\\soneil\\Documents")`
 - In this case, not needed, but generally, replace with a relative pathname

Discussion: deposit example (2/2)

- Repeatability
 - Sets random number seed, but that line is buried in the code
 - Line 81: `set.seed(0123)`
 - Better to set at top of file, or read in from input file
- Hardware requirements are specified in the README.

3. Programming Languages

Learning Outcomes

- Describe common programming languages used in academic research and their applications
- Identify key file types and extensions associated with MATLAB, Python, and R
- Introduce common Integrated Development Environments (IDE) and Notebooks for these programming languages

What is a Programming Language?

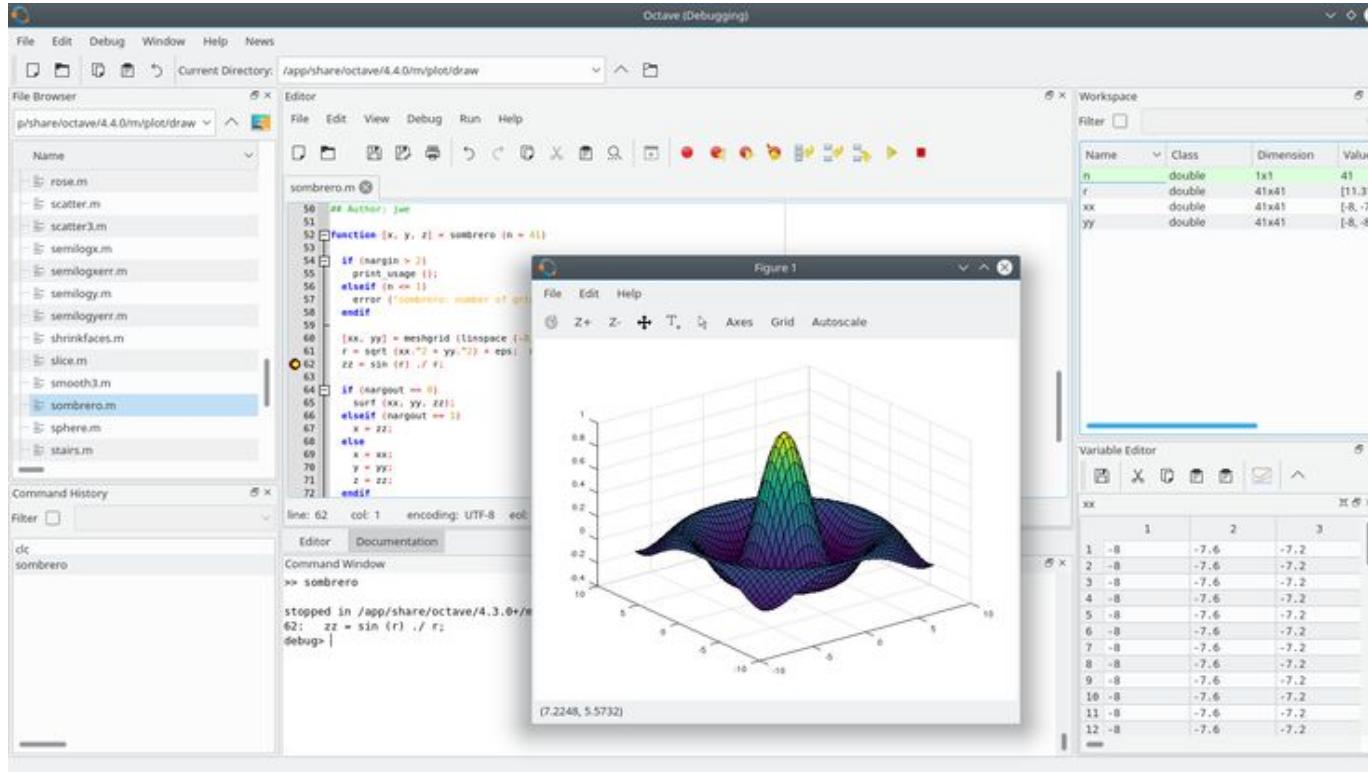
- A computer language used to develop software programs, scripts, or other sets of instructions for computers or computing devices to perform specific tasks
- During today's workshop, we will be focusing on MATLAB, Python and R
- Focusing on these languages as they are the most commonly encountered by DCN curators
- All three languages (and many encountered by curators) are interpreted languages

MATLAB Overview

- Numerical computing environment built around the MATLAB scripting language
- Licensed through Mathworks as a proprietary, paid product
- Several open-source alternatives with similar (but not always exact) syntax exist, with GNU Octave being the most common
- Many external toolboxes have been developed for MATLAB, most of which require a separate subscription



GNU Octave: “Open Source MATLAB”



Python Overview

- High-level, general purpose programming language with applicability cross a variety of field and research areas
- Easily extendable to other programming languages (i.e. C++)
- Applicability is enhanced by a large number of externally developed libraries and packages.
- Open-source software (Libraries are also open-source)



R Overview

- R is an open-source programming language that is specifically designed for complex data and statistical analysis
- The programming language itself (and its default console) is often referred to as Base R
- Community-developed “packages” have greatly enhanced the capabilities of R



Common Research Applications

MATLAB	Python	R
	Data Analysis	
	Data Visualization	
	Machine Learning	
Data Modeling	Web Development	Mapping
Simulation	Automation or Scripting	Spatial data analysis
Linear Algebra	Text Mining	Text mining
	Financial Analysis	

Common Research Applications

What are some research applications that you have seen for these languages that are not listed previously?

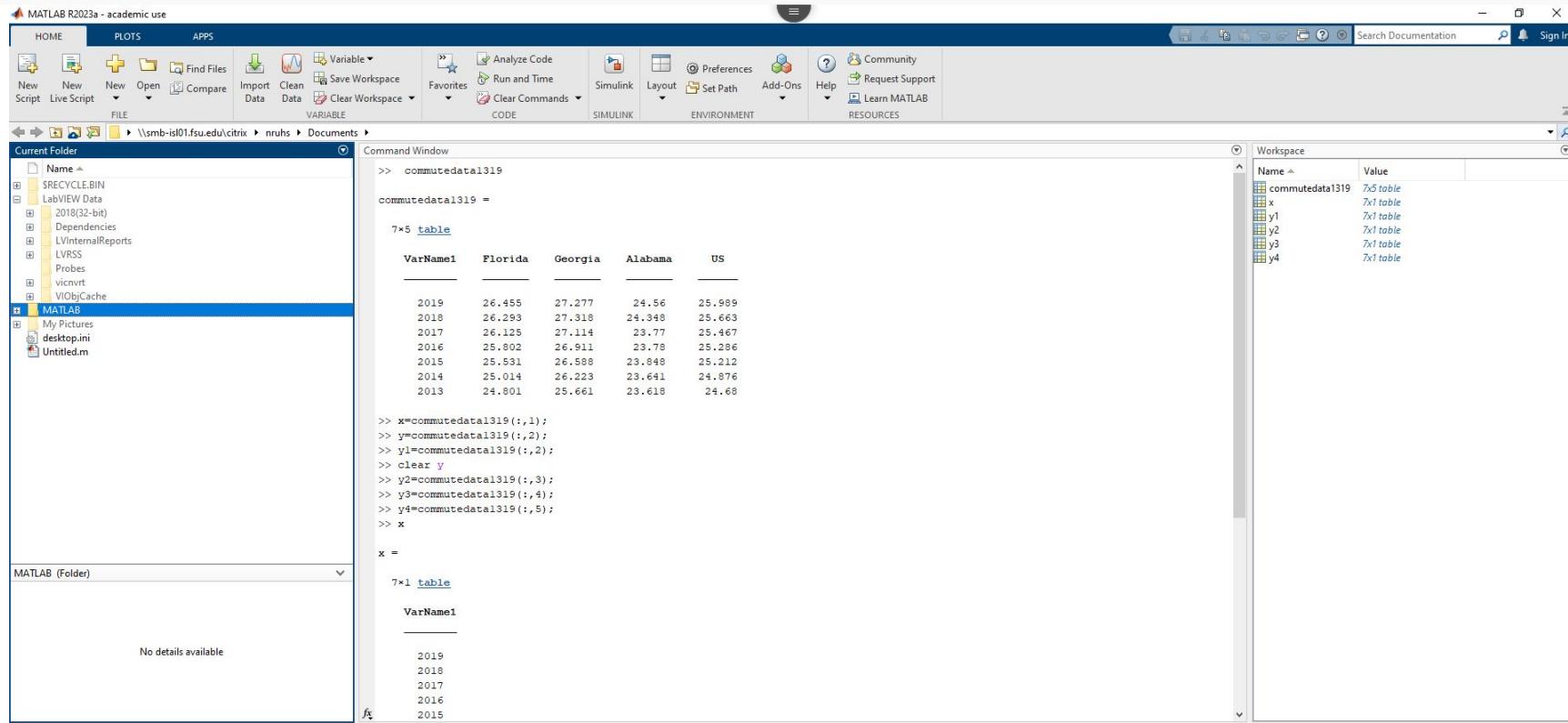
File Formats

- Each programming language has a specific file format for programming files or scripts.
 - MATLAB - .m
 - Python - .py
 - R - .r
- The scripts can be read as plain text files (.txt) by a text editor, but require the specific software to run
- MATLAB also uses .mat files to store MATLAB formatted data

Integrated Development Environments

- Software applications that provide a suite of tools for code and software development
- IDEs combine several programming tasks into one application
- IDEs allow you to develop, write, and check code written in a particular programming language all in one application
- Some IDEs are specific to certain programming languages, while others are compatible with multiple languages

MATLAB



R Studio

The screenshot shows the R Studio interface with the following components:

- File Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Toolbar:** Includes icons for file operations like Open, Save, New, Print, and Go to file/function.
- Script Editor:** An R script titled "example_script.R" containing the following code:

```
1 # Loading two packages into your library
2 # tidyverse and gapminder
3 library(tidyverse)
4 library(gapminder)
5
6 # Modify data
7 gapminder2007 = gapminder %>%
8   filter(year == 2007)
9
10 # Plot data
11 gapminder2007 %>%
12   ggplot(aes(x = gdpPercap, y = lifeExp)) +
13   geom_point()
14
15 # Statistical test
16 t.test(lifeExp ~ gdpPercap > 20000, data = gapminder2007)
```
- Console:** Displays the output of the t-test command:

```
> t.test(lifeExp ~ gdpPercap > 20000, data = gapminder2007)

Welch Two Sample t-test

data: lifeExp by gdpPercap > 20000
t = -14.014, df = 127.85, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-18.30652 -13.77657
sample estimates:
mean in group FALSE mean in group TRUE
63.27946    79.32100
```
- Environment Viewer:** Shows the Global Environment with "gapminder2007" loaded, containing 142 observations of 6 variables.
- Plots:** A scatter plot in the bottom right corner showing "lifeExp" on the y-axis versus "gdpPercap" on the x-axis. The x-axis ranges from 0 to 50,000 with major ticks every 10,000. The y-axis ranges from 40 to 80 with major ticks every 10 units. The plot shows a strong positive correlation, with data points clustered between gdpPercap values of 0 and 50,000 and lifeExp values of 40 and 80.

Notebooks

- Notebooks are “executable papers” combining richly formatted text, runnable code, and data visualization in a single file.
- Support different programming languages, including Python and R.
- Text is formatted using markdown
- Notebook files use special file extensions (e.g., “.ipynb”, “.Rmd”)
- Not simple text files: need to be “rendered” with the notebook application: Jupyter Notebooks, RStudio Notebooks, Quarto.

Jupyter

File Edit View Run Kernel Tabs Settings Help

Launcher README.md Lorenz.ipynb Terminal 1 Console 1 Data.ipynb Python 3 (ipykernel)

The Lorenz Differential Equations

Before we start, we import some preliminary libraries. We will also import (below) the accompanying `lorenz.py` file, which contains the actual solver and plotting routine.

```
[1]: %matplotlib inline
from ipywidgets import interactive, fixed
```

We explore the Lorenz system of differential equations:

$$\dot{x} = \sigma(y - x)$$

Output View

Parameter	Value
sigma	10.00
beta	2.67
rho	28.00



```
lorenz.py
```

```
1 from matplotlib import pyplot as plt
2 from mpl_toolkits.mplot3d import Axes3D
3 import numpy as np
4 from scipy import integrate
5
6 def solve_lorenz(sigma=10.0, beta=8./3, rho=28.0):
7     """Plot a solution to the Lorenz differential
8     equations.
9
10    max_time = 4.0
11    N = 30
12
13    fig = plt.figure()
14    ax = fig.add_axes([0, 0, 1, 1], projection='3d')
15    ax.axis('off')
```

Simple 1 2 3 Python

Ln 1, Col 1 Spaces: 4 lorenz.py

Curation Questions

- Are you able to identify the **correct programming language** for the code files?
- Are source code files named appropriately, with a **file extension** that matches the programming language?
- Are you able to **open** the code files and **view** them using the correct programming language?

Discussion: Deposit Example

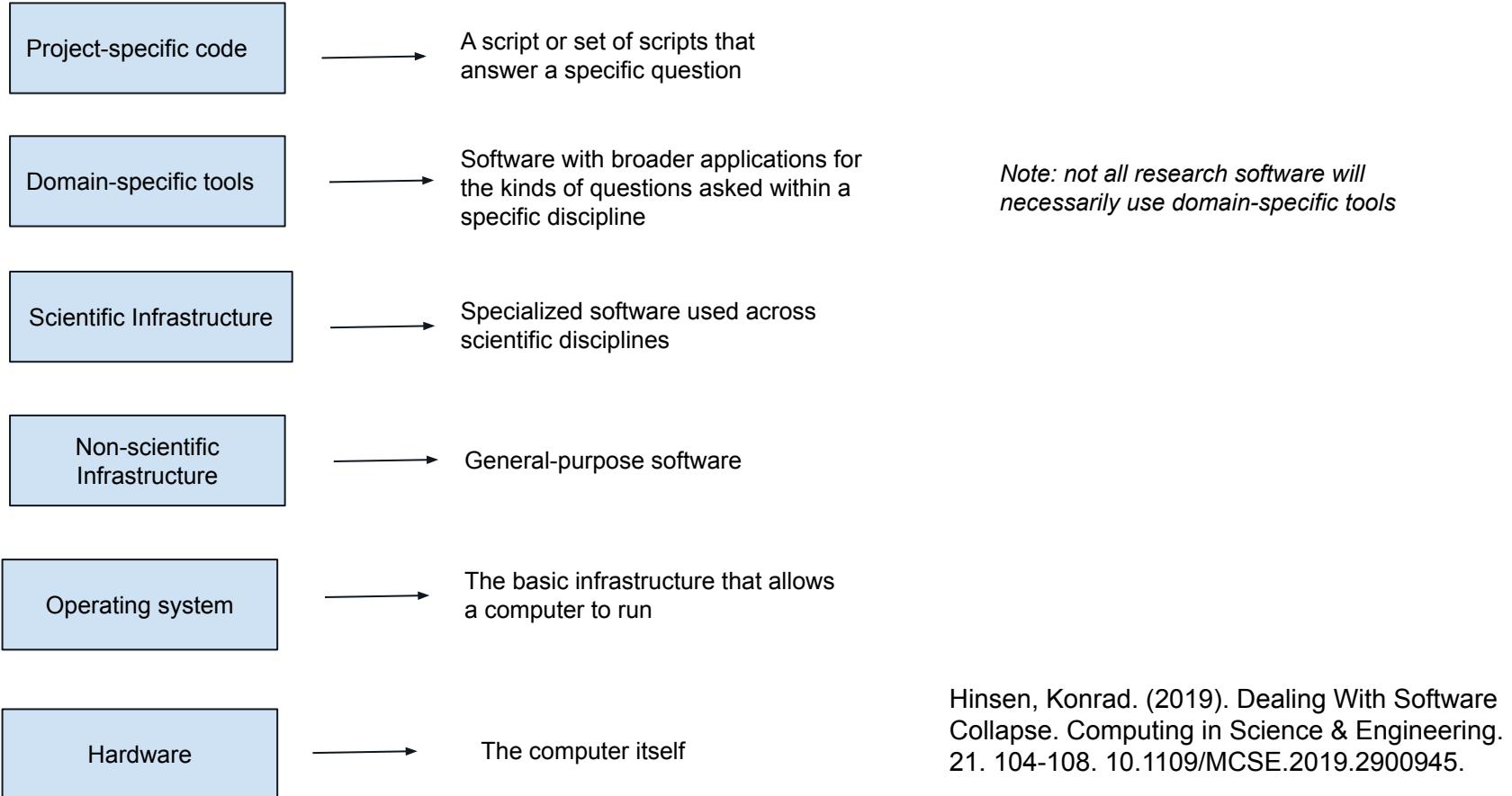
Let's revisit our deposit example (Code, data, and metadata document for the manuscript: **Density-dependence in wolf resource selection study designs**)

- File extension
 - Script uses a .txt extension in the deposit
 - While the code can be read in a text editor, it needs to be run using the programming language it was written in.
 - File extension should reflect the programming language
 - A better name for the script would be `code.R`

4. Dependencies

Learning Outcomes

- Define “dependency” and learn why dependencies are used
- Identify when dependencies are being used in code
- Gain familiarity with package managers and versioning
- Learn several ways to record dependencies—via README, programmatically, and through the use of containerization



Hinsen, Konrad. (2019). Dealing With Software Collapse. Computing in Science & Engineering. 21. 104-108. 10.1109/MCSE.2019.2900945.

Most code we curate
will be at these
levels

Project-specific code

Domain-specific tools

Scientific Infrastructure

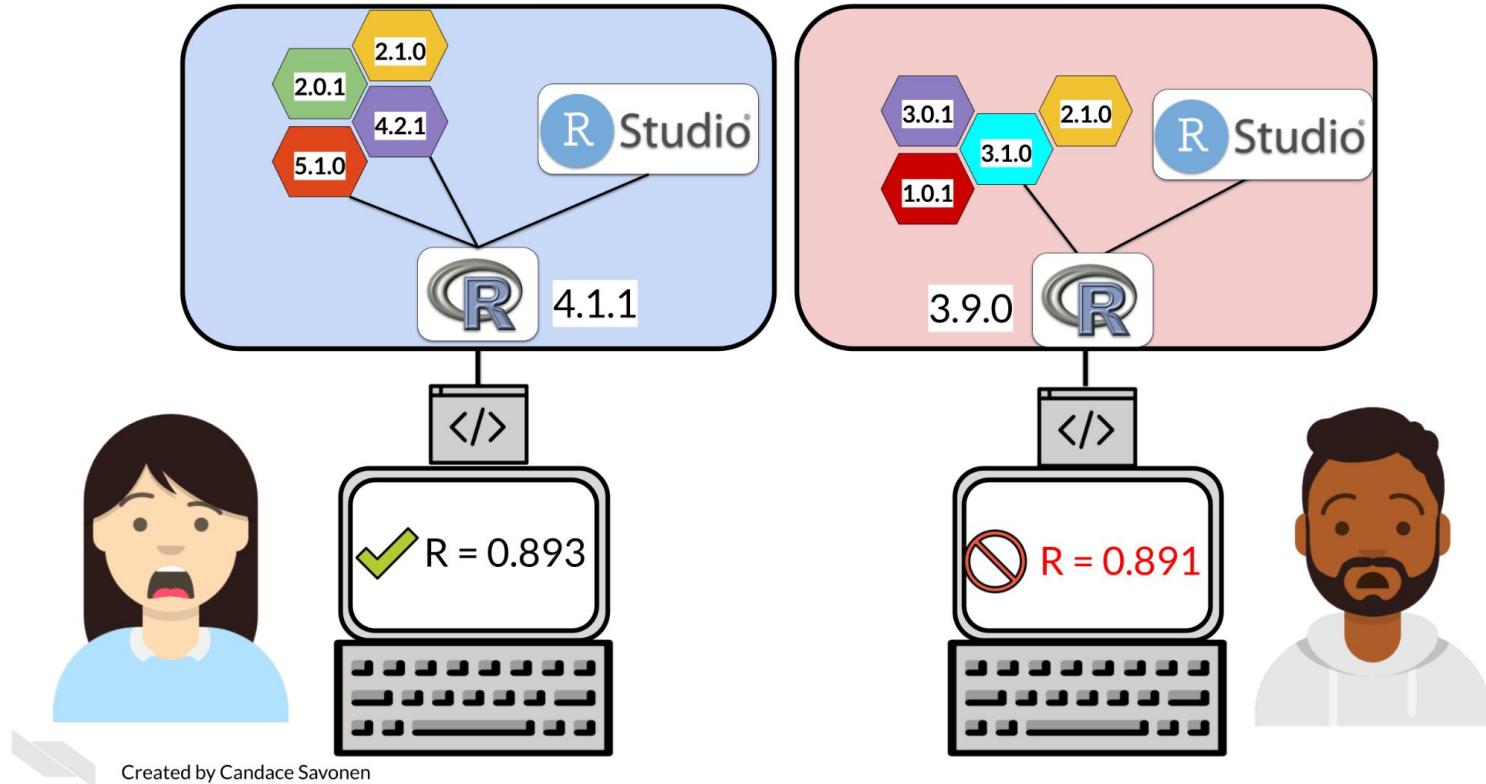
Non-scientific
Infrastructure

Operating system

Hardware

Dependencies

Ruby's local computing environment Avi's local computing environment

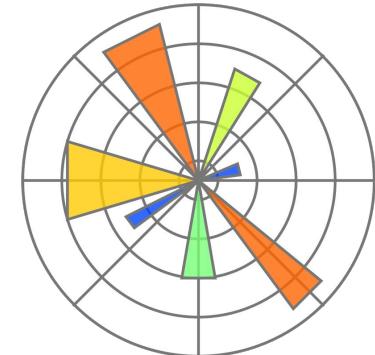


Activity - Which Python?

- Open a Terminal on a Mac, or a command prompt on a Windows PC
- At the prompt, type

```
python --version
```

Why use dependencies if they cause so many problems?



Ethical. Reusable. Better.

DATA CURATION NETWORK

datacurationnetwork.org

Loading dependencies in code

Dependencies in R

```
# Load libraries
library(RMark)
library(knitr)
library(lme4)
library(ggplot2)
library(doby)
```

Dependencies in MATLAB

```
MATLAB Version: 9.13.0.2080170 (R2022b) Update 1
MATLAB License Number: [REDACTED]
Operating System: macOS Version: 12.6 Build: 21G115
Java Version: Java 1.8.0_202-b08 with Oracle Corporation Java HotSpot(TM) 64-Bit Server VM mixed mode
-----
MATLAB Version 9.13 (R2022b)
Control System Toolbox Version 10.12 (R2022b)
Deep Learning Toolbox Version 14.5 (R2022b)
Signal Processing Toolbox Version 9.1 (R2022b)
Statistics and Machine Learning Toolbox Version 12.4 (R2022b)
>> untitled
```

Dependencies in Python

```
: from math import floor
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

Data and analysis code for "Capturing High Resolution Plant Movement in the Field"

Heuschele, Deborah; Furuta, Daniel; Smith, Kevin; Marchetto, Peter (2022)



Persistent link to this item

<https://doi.org/10.13020/7e6x-8f36>

<https://hdl.handle.net/11299/226273>

Services

[Full Metadata \(xml\)](#)

[View Usage Statistics](#)

Keywords

agronomy

Collection period

2019 to 2020

Date completed

2022

Title

Data and analysis code for "Capturing High Resolution Plant Movement in the Field"

Published Date

2022-02-07

Authors

Heuschele, Deborah
Furuta, Daniel
Smith, Kevin
Marchetto, Peter

Group

USDA Agricultural Research Service

Author Contact

Heuschele, Deborah (Jo.Heuschele@usda.gov)

Type

Dataset
Field Study Data
Programming Software Code

Abstract

This is the dataset and analysis code for the paper "Capturing High Resolution Plant Movement in the Field", in which we demonstrate a system to capture high resolution plant motion with small grains in natural conditions.

Description

Plant motion data, wind speed and direction data, and analysis code used for the paper.

Funding information

Sponsorship: MNDRIVE Robotics Sensing and Advanced Manufacturing initiative; Minnesota Department of Agriculture Grant No. 122130; University of Minnesota Rapid Agricultural Response Fund Grant No. AES00RR234

License

[CC0 1.0 Universal](#)

Suggested Citation

```
1 The scripts were developed with Python 3.7.13 and the following non-core package versions:  
2 Python 3.7.13  
3 numpy 1.21.5  
4 dateutil 2.8.2  
5 matplotlib 3.5.1  
6 pandas 1.3.5  
7 seaborn 0.11.2  
8 windrose 1.6.8  
9
```

Users > talyacooper > Downloads > analysis_code > visualization >  plot_windroses.py > ...

```
1 #!/usr/bin/env python3  
2 # -*- coding: utf-8 -*-  
3 """  
4 Created on Wed Jan  5 10:29:02 2022  
5  
6 This script creates windroses showing wind trends over the collected data.  
7 To adapt to your needs, replace the file locations, date ranges,  
8 and wind ranges to your situation.  
9 """  
10  
11 import windrose  
12 from matplotlib import pyplot as plt  
13 import pandas as pd  
14 import seaborn as sn  
15 import numpy as np  
16
```

Project-specific code

A script or set of scripts that answer a specific question

Domain-specific tools

Software with broader applications for the kinds of questions asked within a specific discipline

Scientific Infrastructure

Specialized software used across scientific disciplines

Non-scientific Infrastructure

General-purpose software

Operating system

The basic infrastructure that allows a computer to run

Hardware

The computer itself

`plot_windroses.py`

The deposited script

`windrose 1.6.8`

Specialized meteorological package for Python

`numpy 1.21.5`

Broadly used package for high-level mathematical functions

`Python 3.7.13`

The Python language

`macOS 10.9`

The version of macOS needed to run Python 3

`64-bit processor`

The computer chip needed to run macOS 10.9

```
: from math import floor
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
dependencies:
- absl-py@0.9.0=py36_0
- appnope@0.1.0=py36h9f0ad1d_1001
- astor@0.7.1=py_0
- attrs@19.3.0=py_0
- backcall@0.1.0=py_0
- bleach@3.1.3=pyhc360c0e_0
- blosc@1.19.0=h4a8c4bd_0
- bzip2@1.0.8=h0b31af3_2
- c-ares@1.15.0=h0d97ff_1001
- ca-certificates@2020.6.20-hecdca079_0
- cctools@927.0.2=h5ba7a2e_4
- certifi@2020.6.20=py36h9f0ad1d_0
- cycler@0.10.0=py_2
- dbus@1.13.6=h2f22bb5_0
- decorator@4.4.2=py_0
- defusedxml@0.6.0=py_0
- entrypoints@0.3=py36h9f0ad1d_1001
- expat@2.2.9=h4a8c4bd_2
- freetype@2.10.1=h8da9a1a_0
- gast@0.3.3=py_0
- gettext@0.19.8.1=h46ab8bc_1002
- glib@2.58.3=py36hbcc7cff_1003
- grpcio@1.23.0=py36h7c1f37e_1
- h5py@2.10.0=ompi_py36h106b333_102
- hdf5@1.10.5=ompi_h3e39495_1184
- icu@64.2=h6de7cb9_1
- importlib-metadata@1.5.0=py36h9f0ad1d_1
- importlib_metadata@1.5.0=1
- ipykernel@5.1.4=py36h5ca1d4c_0
- ipython@7.13.0=py36h95af2a2_1
- ipython_genutils@0.2.0=py_1
- ipywidgets@7.5.1=py_0
- jedi@0.16.0=py36h9f0ad1d_1
```

```
39   - jinja2@2.11.1=py_0
40   - jpeg@9c=h1de35cc_1001
41   - jsonschema@3.2.0=py36h9f0ad1d_1
42   - jupyter@1.0.0=py_2
43   - jupyter_client@6.1.0=py_0
44   - jupyter_console@6.1.0=py_1
45   - jupyter_core@4.6.3=py36h9f0ad1d_
46   - keras@2.3.1=py36_0
47   - keras-applications@1.0.8=py_1
48   - keras-preprocessing@1.1.0=py_0
49   - kiwisolver@1.1.0=py36h863e41a_1
50   - krb5@1.16.4=h1752a42_0
51   - ld64@450.3=h3c32e8a_4
52   - libblas@3.8.0=16_openblas
53   - libblas@3.8.0=16__openblas
54   - libclang@9.0.1=default_hf57f61e_0
55   - libcxx@9.0.1=1
56   - libedit@3.1.20170329=hcfe32e1_100
57   - libffi@3.2.1=h4a8c4bd_1007
58   - libgfortran@4.0.0=2
59   - libgpuarray@0.7.6=h1de35cc_1003
60   - libiconv@1.15=h0b31af3_1006
61   - liblapack@3.8.0=16_openblas
62   - libllvm@9.0.1=h1b3eb9_0
63   - libopenblas@0.3.9=h3d69b6c_0
64   - libpng@1.6.37=hbbe82c9_1
65   - libpq@12.2=h554dc5a_0
66   - libprotobuf@3.11.4=h0d174df1_0
67   - libsodium@1.0.17=h01d97ff_0
68   - llvm-openmp@9.0.1=h28b9765_2
69   - mako@1.1.0=py_0
70   - markdown@3.2.1=py_0
71   - markupsafe@1.1.1=py36h37b9a7d_1
72   - matplotlib-base@3.2.1=py36h83d3ec1_0
73   - mistune@0.8.4=py36h0b31af3_1000
74   - mock@4.0.2=py36h9f0ad1d_0
75   - nbconvert@5.6.1=py36_0
76   - nbformat@5.0.4=py_0
77   - ncurses@6.1=h0a44026_1002
78   - notebook@6.0.3=py36_0
79   - nspr@4.29=h0a44026_1000
80   - nss@3.47=h080d9_0
81   - numexpr@2.7.1=py36hcc1bba6_1
82   - numpy@1.18.1=py36hde6bac1_0
83   - openssl@1.1.1g=h0b31af3_0
84   - palettable@3.3.0=py_0
85   - pandas@1.0.3=py36hcc1bba6_0
86   - pandoc@2.9.2=0
87   - pandocfilters@1.4.2=py_1
88   - parso@0.6.2=py_0
89   - patsy@0.5.1=py_0
90   - pcrc@8.44=h4a8c4bd_0
91   - pexpect@4.8.0=py36h9f0ad1d_1
92   - pickleshare@0.7.5=py36h9f0ad1d_1001
93   - pip@20.0.2=py_2
94   - prometheus_client@0.7.1=py_0
95   - prompt_toolkit@3.0.4=py_0
96   - prompt_toolkit@3.0.4=py_0
97   - protobuf@3.11.4=py36h4a8c4bd_0
98   - ptyprocess@0.6.0=py_1001
99   - pygments@2.6.1=py_0
```

Swiger, B. M., Liemohn, M. W., & Ganushkina, N. Y. (2020). Improvement of Plasma Sheet Neural Network Accuracy With Inclusion of Physical Information. *Frontiers in Astronomy and Space Sciences*, 7. <https://doi.org/10.3389/fspas.2020.00042>

Packages and package managers

Package

A collection of code, data, and documentation that can be distributed and re-used. Also referred to in some languages as a **library** or **module**.

Package manager

A program that does its best to keep track of the different software installed on a computer and their dependencies on one another.

If you download a dependency as a **package** from a **package manager**, the **package manager** should ensure that you also install that package's **direct AND indirect dependencies** (and those dependencies' dependencies)



The Comprehensive R Archive Network

Documenting Dependencies - README

Hints for the Installation

`librosa` uses `soundfile` and `audioread` to load audio files.

 Note that older releases of `soundfile` (prior to 0.11) do not support MP3, which will cause `librosa` to fall back on the `audioread` library.

`soundfile`

If you're using `conda` to install `librosa`, then audio encoding dependencies will be handled automatically.

If you're using `pip` on a Linux environment, you may need to install `libsndfile` manually. Please refer to the [SoundFile installation documentation](#) for details.

`audioread` and MP3 support

To fuel `audioread` with more audio-decoding power (e.g., for reading MP3 files), you may need to install either `ffmpeg` or `GStreamer`.

 Note that on some platforms, `audioread` needs at least one of the programs to work properly.

If you are using Anaconda, install `ffmpeg` by calling

```
conda install -c conda-forge ffmpeg
```

<https://librosa.org/>

When to document dependencies in a README

- It never hurts!
- Always specify version information
- Document API and data dependencies in the README
- Documenting dependencies in the README may be sufficient for a simple script that calls a small number of external packages or libraries
- README is the place to explain specific issues, challenges, or instructions

Documenting Dependencies—Programmatic

Python

main / tile2net / requirements-dev.txt

Mary-h86 initial commit

1 contributor

32 lines (31 sloc) | 271 Bytes

```
1 centerline
2 jupyter
3 geopandas
4 geopy
5 matplotlib
6 momepy
7 numba
8 opencv-python
9 osmnx
10 pillow
11 pyarrow
12 rasterio
13 scikit-image
14 scikit-learn
15 shapely
16 tqdm
17 torch
18 torchvision
```

R

Set up package dependencies for compatibility with `renv`

Source: [R/tar_renv.R](#)

Write package dependencies to a script file (by default, named `_targets_packages.R` in the root project directory). Each package is written to a separate line as a standard `library()` call (e.g. `library(package)`) so `renv` can identify them automatically.

Usage

```
tar_renv(
  extras = c("bs4Dash", "clustermq", "future", "gt", "markdown", "pingr", "rstudio",
            "shiny", "shinybusy", "shinyWidgets", "visNetwork"),
  path = "_targets_packages.R",
  callr_function = callr::r,
  callr_arguments = targets::tar_callr_args_default(callr_function),
  envir = parent.frame(),
  script = targets::tar_config_get("script")
)
```



When to create programmatic documentation of dependencies

- It never hurts!
- Particularly helpful for more complex code that calls a large number of external packages and libraries
- Required for certain kinds of virtual environments, like Binder, if a depositor will want their code to function in that setting

Capturing dependencies: Containerizing

- Containerization tools like Docker and Singularity create a complete image of a piece of software's runtime environment—including all of the correct versions of the software's dependencies
- ReproZip is another tool designed for packaging software applications for reproducibility; it might be overkill for a simple script but is excellent for the preservation of more extensive projects that will not be updated.

When to capture dependencies with a container

- This is often done when software is still in development, to ensure developers are making changes in the exact same environment with the same versions of dependencies
- Don't forget to document how to use/run the container!
- If the software is no longer being developed, is complex and highly environment-specific, and you want to be able to run it in the future, preserving it in a ReproZip package is one of the best ways to ensure future reproducibility
 - Almost a “special collections” level of research software preservation

Discussion: deposit example

- README indicates a now-obsolete version of R, which may make it harder to run the code.
- No version information for dependencies
 - Because the version of R is older, someone who wanted to reuse this code would need to do some research to find the correct version of the dependencies that would work with that version of R – could lead to some issues with indirect dependencies, too!

Dependency issues to watch out for

What are some dependency-related issues you might look out for while curating?

What are some dependency-related issues we notice in our code example?

Dependency issues to watch out for

- Dependencies are not described anywhere in **documentation**
- Dependencies are called in the code but are **not actually used** when the code is run
- **No version information** is provided about dependencies
- When the code is run, **dependencies are no longer available** or an error message pops up when it is installed
- Dependencies are **proprietary software** and cannot be accessed
- Code requires **specialized hardware or firmware** (eg. a GPU, an AR/VR device.)
- Code calls resources from elsewhere and requires an **API or an internet connection**
- Code contains **private information**, particularly an API key

Code Curation Questions

- Are hardware and software **dependencies** adequately documented?
- Are **files referenced** in the source code included in the submission?

5. Reviewing Documentation

Learning Outcomes

By the end of this module, participants will be able to:

- Identify types of research software documentation and how they can be useful during the curation process
- Make suggestions for research software documentation (README, Commenting, Markdown, etc.)
- Review and offer edits to README files/outlines for research software

Reviewing Documentation as a Curator

Each step of the curation process involves reviewing provided documentation to help understand the structure and function of code deposit.

Check files and read documentation (risk mitigation, file inventory, appraisal/selection)

Understand the data (or try to), if not... (run files/environment, QA/QC issues, readme)

Request missing information or changes (tracking provenance of any changes and why)

Augment metadata for findability (DOIs, metadata standards, discoverability)

Transform file formats for reuse (data preservation, conversion tools, data viz)

Evaluate for FAIRness (licenses, responsibility standards, metrics for tracking use)

Document your curation activities (Curator Log, correspondence)

What is documentation?

Documentation is the use of text to describe the software's original environment, intended use, and outputs.

Documentation occurs throughout the research software lifecycle, so it is important to help researchers understand the value of documenting their processes and procedures before curation as this information is often referenced during the curation process.

Internal Documentation

Comments and notes made within the code/software

External Documentation

Comments and notes made outside of the code/software

External Documentation



README Files

README files are used to describe code included in a code repository. The main goal of a README file is provide information on the following:

- Description of Project
- Contributors
- Technical Requirements (System and Language Dependencies & Working Directory)
- Any Known Bugs/Issues
- Output Files Created and Storage Location
- Naming Conventions
- License
- Acknowledgements
- Contact Information
- Any information significant to the structure, usability, and context of the code file(s)

The goal is to facilitate reuse.

README Exemplar

See file in workshop materials

Whole Group Activity



README File Investigation

As a group, we will review the README file below and consider –

- The usefulness of the file
 - Potential edits to the file
-

****Attitude**** is a Python module for fitting the orientation of planes! Its core mission is to compute planar fits and transform them (*along with meaningful error distributions*) into a form that can be used for geology. It is designed to support the collection of structural measurements from remote sensing data.

The documentation for the module can be found at

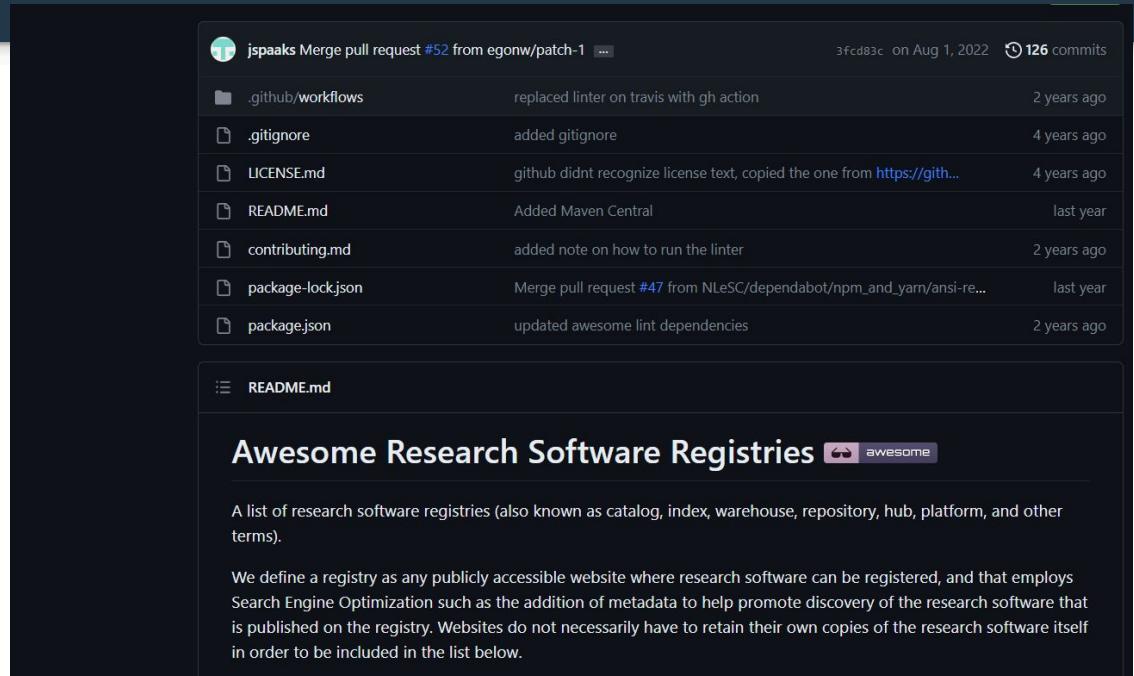
<https://davenquinn.github.io/Attitude>, and a preprint publication describing the underlying statistical method is on [EarthArXiV](<https://eartharxiv.org/4enzu/>).

Reviewing GitHub Documentation



READMEs on GitHub

GitHub uses the README file as the front page of a code repository.



The screenshot shows a GitHub repository's commit history and its README file content. The commit history lists several changes made by user 'jspaaks' over the past two years, including updates to workflows, .gitignore, LICENSE.md, README.md, contributing.md, package-lock.json, and package.json. The README file itself is titled 'Awesome Research Software Registries' and defines what a registry is, noting it's a publicly accessible website for registering research software.

File / Action	Description	Time Ago
.github/workflows	replaced linter on travis with gh action	2 years ago
.gitignore	added gitignore	4 years ago
LICENSE.md	github didnt recognize license text, copied the one from https://github.com	4 years ago
README.md	Added Maven Central	last year
contributing.md	added note on how to run the linter	2 years ago
package-lock.json	Merge pull request #47 from NLeSC/dependabot/npm_and_yarn/ansi-re...	last year
package.json	updated awesome lint dependencies	2 years ago

README.md

Awesome Research Software Registries 

A list of research software registries (also known as catalog, index, warehouse, repository, hub, platform, and other terms).

We define a registry as any publicly accessible website where research software can be registered, and that employs Search Engine Optimization such as the addition of metadata to help promote discovery of the research software that is published on the registry. Websites do not necessarily have to retain their own copies of the research software itself in order to be included in the list below.

The screenshot shows the GitHub interface for the repository 'NLeSC / awesome-research-software-registries'. The 'Issues' tab is selected, displaying 12 open issues. The search bar at the top contains the query 'is:issue is:open'. Below the search bar are buttons for 'Labels' (10), 'Milestones' (0), and a green 'New issue' button. The issues listed are:

- ① Add research software directory instances (#55 opened on Nov 2, 2022 by jspaaks)
- ① review agu list of registries (#54 opened on Oct 5, 2022 by jspaaks)
- ① clariah dev tools (#51 opened on May 23, 2022 by jspaaks)
- ① Awesome open source projects from Utrecht University (#50 opened on May 20, 2022 by kequach)

At the bottom of the page, there are dropdown menus for 'Author', 'Label', 'Projects', 'Milestones', and 'Assignee', followed by a 'Sort' dropdown.

The researcher may use the ISSUES menu to add further documentation about a software.

Issues can be created and addressed by all members in a repository.

GitHub Wiki

<> Code Issues Pull requests Actions Projects Wiki Security Insights Settings

Researchers may use GitHub wikis to provide software documentation

About wikis

You can host documentation for your repository in a wiki, so that others can use and contribute to your project.

Wikis are available in public repositories with GitHub Free and GitHub Free for organizations, and in public and private repositories with GitHub Pro, GitHub Team, GitHub Enterprise Cloud and GitHub Enterprise Server. For more information, see "[GitHub's plans](#)".

Every repository on GitHub.com comes equipped with a section for hosting documentation, called a wiki. You can use your repository's wiki to share long-form content about your project, such as how to use it, how you designed it, or its core principles. A README file quickly tells what your project can do, while you can use a wiki to provide additional documentation. For more information, see "[About READMEs](#)".

External Documentation: File Names

File names can be consistently arranged to ensure discoverability and reusability of research software. File naming conventions should be documented in the README file (more to come!.. Potential metadata include:

- Date (following the ISO 8601 standards: YYYYMMDD or YYYY-MM-DD)
- Time
- Sample/Subject/Image ID Numbers
- File Description
- Versioning Information (mostly for files that do not undergo automatic version control)
- Output
- Data type
- Project/Project Component Title
- Researcher/Assistant Name

Abbreviations will come in handy as you want your file names to be as short as possible.

Avoid using spaces; instead, **use underscores(_)** or **hyphens(-)**.

Use one case when possible to help in retrieval.

Sample File Names

- **YYYY-MM-DD_ProjectTitle_SubjectID_Version.FileType**
- **2012-12-20_birdmiganalysis_sue_v1.csv**
- **YYYYMMDD_ProjectTitle.FileType**
- **20221015bookclustering_analysis.py**
 - *Versioning information for code files is typically included in version control protocols.*

If you notice that a deposit has a standard file naming structure, suggest that this be structure be outlined in the README file.

File Names

WARNING!

Be sure to remind researchers to avoid absolute paths! Relative paths are preferred.

Absolute Path: "C:\\Users\\rholmes\\Documents.csv")
Relative Path:(“..\\Documents.csv”)

Internal Documentation Strategies



Internal Documentation Strategies: Commenting

Comments are natural language text strings that describe software code. **Commenting is useful for annotating software code as they help readers understand the complexities of the software's language.** Comments are usually made when—

- An action will occur that results in a new or a modified variable, including adding columns to data frames or other similar edits
- An action will occur that results in an output, both saved to your system or printed in the console
- User input will be required
- Conditions needs must be met before running a line/block of code
- An output is to be expected
- An explanation of a process needs to be provided (why was a certain block of code run?)
- Users need to change alter code for their purposes



Identifying Comments in Python

Comments in Python are created begin with a hashtag (#).

Additionally, you can make block comments by surrounding text with triple double quotes ("").

```
In [1]: import pandas as pd  
In [2]: #import data file  
df = pd.read_excel("C:/Users/MayeKapounyers/Box/Librarianship/faculty_research_interests.xlsx")
```

```
In [ ]: """  
Creation Date: 10/15/2022  
Modification Date: 6/12/2023  
  
Dependencies and Requirements  
- Python 3.7 or higher  
- Pandas 1.0.3 or higher  
- IDE: Jupyter  
  
Other required files:  
- faculty_research_interests.xlsx  
  
Summary: This code was created to restructure faculty research interest data.  
"""
```

Identifying Comments in R



To create comments in R, begin the comment text with a hashtag (#).

A screenshot of the RStudio IDE interface. The top panel shows a script editor with two functions defined:

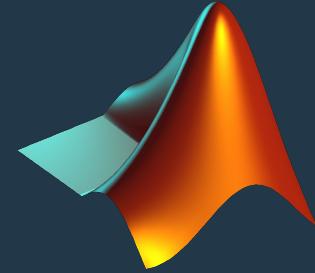
```
1 # Function to calculate the circumference of a circle
2 circumference <- function(radius) {
3   2 * pi * radius
4 }
5
6 # Function to calculate the area of a circle
7 area <- function(radius) {
8   pi * radius^2
9 }
```

The bottom panel shows a text editor with a block comment:

```
1 """
2 Author: Jayden Eaves
3 Contact: jeaves@ggolden.com
4 Affiliation: Gee's Golden Code Academy
5
6 Creation Date: 6/7/2023
7 License Information: None
8 Dependencies: R (version 2.3.1 or higher)
9
10 Other Required Files:
11 - data.csv (input file)
12 - output.png (output visualization)
13
14 Summary: This code calculates the dimensions of a circle based on its circumference
15 """
```

Additionally, you can make block comments by surrounding text with triple double quotes (""""").

Identifying Comments in MATLAB



To create comments in MATLAB, begin the comment text with a percent symbol (%).

```
% Calculate the circumference  
circumference = 2 * pi * radius;
```

```
% Calculate the area  
area = pi * radius^2;
```

Additionally, you can make block comments by surrounding text in the %{ %}

construction.

```
%{  
this code  
will generate a new set  
of dimensions for your research  
%}
```

Internal Documentation: Comments

```
# Function to calculate the factorial of a number
def factorial(n):
    if n == 0:
        return 1
    else:
        return n * factorial(n - 1)

# input
num = int(input("Enter a number: "))

# Check the input
if num < 0:
    print("Factorial is undefined for negative numbers.")
elif num == 0:
    print("The factorial of 0 is 1.")
else:
    result = factorial(num)
    print(f"The factorial of {num} is {result}.")
```

```
# Function to calculate the factorial of a number
def factorial(n):
    if n == 0:
        return 1
    else:
        return n * factorial(n - 1)

# Get input from user
num = int(input("Enter a number: "))

# Check whether the input is non-negative
if num < 0:
    print("Factorial is undefined for negative numbers.")
elif num == 0:
    print("The factorial of 0 is 1.")
else:
    result = factorial(num)
    print(f"The factorial of {num} is {result}.")
```

Internal Documentation: Comments

```
1 # Create a list of numbers
2 numbers <- c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
3
4
5 even_sum <- 0
6
7
8 for (number in numbers) {
9   # Check the number
10  if (number %% 2 == 0) {
11    # Check evenness
12    even_sum <- even_sum + number
13  }
14}
15
16 cat("The sum of even numbers in the list is:", even_sum, "\n")|
```

Internal Documentation: Comments

```
1 # Create a list of numbers
2 numbers <- c(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)
3
4 # Initialize a variable to store the sum of even numbers
5 even_sum <- 0
6
7 # Iterate through the list of numbers
8 for (number in numbers) {
9   # Check if the number is even
10  if (number %% 2 == 0) {
11    # If it's even, add it to the sum
12    even_sum <- even_sum + number
13  }
14}
15
16 # Print the sum of even numbers
17 cat("The sum of even numbers in the list is:", even_sum, "\n")
```

Internal Documentation: Header Text

When reviewing a code deposit, the header can provide useful information about the content of a code file. This text usually includes –

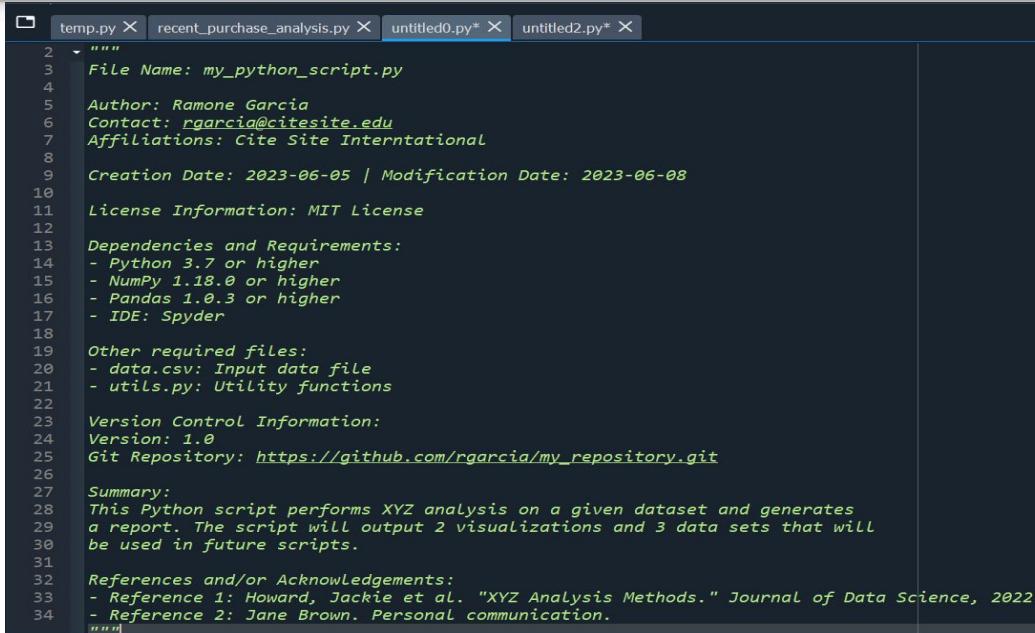
- File Name
- Author Information (name, contact, affiliations)
- Creation Date and Modification Dates
- License Information
- Dependencies and Requirements
- Other Required Files
- Version Control Information
- Summary
- References and/or Acknowledges

Header Texts are often condensed/individualized versions of the README file.

R Header Text Example

```
1 """
2 Author: Jayden Eaves
3 Contact: jeaves@ggolden.com
4 Affiliation: Gee's Golden Code Academy
5
6 Creation Date: 6/7/2023
7 License Information: None
8 Dependencies: R (version 2.3.1 or higher)
9
10 Other Required Files:
11 - data.csv (input file)
12 - output.png (output visualization)
13
14 Summary: This code calculates the dimensions of a circle based on its circumference
15 """
```

Python Header Text Example



```
temp.py X recent_purchase_analysis.py X untitled0.py* X untitled2.py* X
2 """
3 File Name: my_python_script.py
4
5 Author: Ramone Garcia
6 Contact: rgarcia@citesite.edu
7 Affiliations: Cite Site Interntational
8
9 Creation Date: 2023-06-05 | Modification Date: 2023-06-08
10
11 License Information: MIT License
12
13 Dependencies and Requirements:
14 - Python 3.7 or higher
15 - NumPy 1.18.0 or higher
16 - Pandas 1.0.3 or higher
17 - IDE: Spyder
18
19 Other required files:
20 - data.csv: Input data file
21 - utils.py: Utility functions
22
23 Version Control Information:
24 Version: 1.0
25 Git Repository: https://github.com/rgarcia/my\_repository.git
26
27 Summary:
28 This Python script performs XYZ analysis on a given dataset and generates
29 a report. The script will output 2 visualizations and 3 data sets that will
30 be used in future scripts.
31
32 References and/or Acknowledgements:
33 - Reference 1: Howard, Jackie et al. "XYZ Analysis Methods." Journal of Data Science, 2022.
34 - Reference 2: Jane Brown. Personal communication.
"""

```

MATLAB Header Text Example

```
1 % FUNCTION NAME:  
2 %     avg  
3 %  
4 % DESCRIPTION:  
5 %     Compute the average of three inputs  
6 %  
7 % INPUT:  
8 %     in1 - (double) Input one  
9 %     in2 - (double) Input two  
10 %    in3 - (double) Input three  
11 %  
12 % OUTPUT:  
13 %    out - (double) Calculated average of the three inputs  
14 %  
15 % ASSUMPTIONS AND LIMITATIONS:  
16 %    None  
17 %  
18 % REVISION HISTORY:  
19 %    05/02/2018 - mmyers  
20 %        * Initial implementation  
21 %
```

Command Window

MATLAB Command Window

Best Practice Resources

[Google Style Guides](#)

[PEP8](#) (Python)

[MatLab Style Guide](#)

[ISO 8601](#)

Small Group Activity



Comment Edits

The following code snippet creates a new column based on values in another column. What internal documentation edits might you suggest for the author?

```
1 # Load the required libraries
2 library(dplyr)
3
4 # Read the data file
5 data <- read.csv("input_data.csv")
6
7 # Create a new column
8 data$new_column <- ifelse(data$old_column > 10, "Greater than 10", "Less than or equal to 10")
9
10 # Remove a column
11 data <- data[, -ncol(data)]
12
13 # Display the head of the data frame
14 head(data)
15 write.csv(data, "output_data.csv", row.names = FALSE)
16 |
```

Potential Curation Actions

Create meaningful header text.	Ask depositor, “Do you anticipate any of these files being separated?”
Ensure interoperable metadata	If metadata standards are used, link to documentation
Promote related literature	Ensure that literature is referenced in documentation (README)
Ensure human-approachable code	Provide depositors with access to best practices resources like PEP8
Ensure system interoperability	Suggest system and directory information be included in documentation (README)

Discussion: Deposit Example

Visit the code example included in the workshop materials (*Code, data, and metadata document for the manuscript: Density-dependence in wolf resource selection study designs*) and consider the following:

1. How can the README file be edited to give a clear explanation of what the code does and its expected output?
2. What are the benefits and disadvantages of installing libraries at the end of the code?
3. What edit(s) can be made to avoid including the researcher's working directory?

Questions, Concerns, Comments

Thanks for completing this
module with us!

Final Exercise: Following a Checklist

- Are **file paths in the source code portable?**
- Are source code files named appropriately, with a **file extension** that matches the programming language?
- Are hardware and software **dependencies** adequately documented?
- Are **files referenced** in the source code included in the submission?
- Are important functions and components **documented with comments** in the source code?
- Does the code include **README**?
- Does the code include a **license** and is it appropriate?
- Are files in the submission **organized in a consistent and logical manner?**

Additional Slides

The following slides were not covered during the pilot workshop. They are included for reference

- Draft Software Licenses Section
- Draft Project Structure Section
- Markdown/Markup

Software Licenses

Note: This section was not covered in the pilot workshop due to time constraints. Slides are included for reference and to facilitate future work.

What is a software license?

- *A license is legal permission for others to copy, modify, and share your work*

What is a software license?

- The Open Source Definition: <https://opensource.org/osd>
 - Allows free distribution
 - Applies to software source code
 - Allows derived works
 - (...)
- Common Open Source Licenses in Science:
 - [BSD](#) (Astropy, NumPy, Jupyter Notebook, scikit-learn)
 - [MIT](#) (Julia, GPT-2, GitLab)
 - [Apache](#) (TensorFlow, Natural Language Toolkit)
 - [GPL](#) (R, FFTW, SageMath, Linux)

Copyleft vs Permissive Licenses

- Copyleft licenses: copyright licenses that include “share alike” restrictions
 - GPL: derived works must be distributed under the GPL
 - “Lesser” GPL: provides an exception for linked libraries
 - Affero GPL: covers distribution via web-based applications
- Permissive licenses place few restrictions on use
 - BSD (3-clause): derived works must include original copyright and license text
 - MIT: explicitly permits sublicensing

tl;dr Legal (<https://tldrlegal.com>) has useful, short summaries of these licenses

Example: BSD 3-Clause

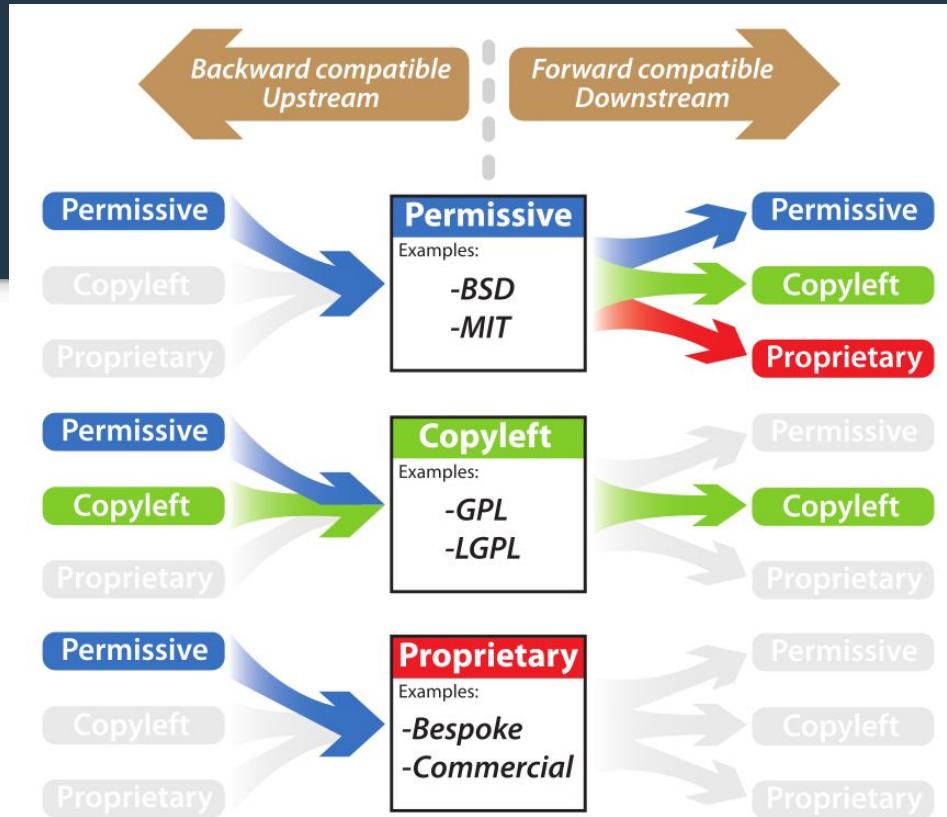
Copyright <YEAR> <COPYRIGHT HOLDER>

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met.

1. **Redistributions of source code** must retain the above copyright notice, this list of conditions and the following disclaimer.
2. **Redistributions in binary form** must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
3. Neither the name of the copyright holder nor the names of its contributors may be used to **endorse or promote products** derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR...

License compatibility



Project Structure

Note: This section was not covered in the pilot workshop due to time constraints. Slides are included for reference and to facilitate future work.

Objectives

- Suggest some principles to follow for organizing projects
- Provide guiding questions to ask yourself and a depositor about project structure as you curate

Three rules of project organization

- Separate directories by function.
- Separate files into inputs and outputs.
- Make directories portable

Gentzkow M, Shapiro JM. Code and Data for the Social Sciences: A Practitioner's Guide; 2014. Available from:
<http://web.stanford.edu/~gentzkow/research/CodeAndData.pdf>.

Data = “input”



1	Data
2	version1_epoch_stats.hdf5
3	version1_epoch_stats.pkl
4	version1_observed.hdf5
5	version1_observed.pkl
6	version1_predicted.hdf5
7	version1_predicted.pkl
8	version1_traindata.hdf5
9	version1_traindata.pkl
10	version2_epoch_stats.hdf5
11	version2_epoch_stats.pkl
12	version2_observations.hdf5
13	version2_observations.pkl
14	version2_predictions.hdf5
15	version2_predictions.pkl
16	version2_train.hdf5
17	version2_train.pkl
18	Figures
19	calculated_v1_loss.pdf
20	calculated_v1_skill.pdf
21	calculated_v2_loss.pdf
22	calculated_v2_skill.pdf
23	density_scatter_v1.png
24	density_scatter_v2.png
25	PlottingCode
26	Calculating Metrics.ipynb
27	Calculating Metrics.pdf
28	LICENSE.txt
29	conda_environment.yml
30	datamodel_comparisons.py
31	plot_scatter_energy.py
32	README.txt

Figures = “output”

Plotting code = scripts

Swiger, B. M., Liemohn, M. W., & Ganushkina, N. Y. (2020). Data From Improvement of Plasma Sheet Neural Network Accuracy With Inclusion of Physical Information. DOI: 10.7302/559r-t639

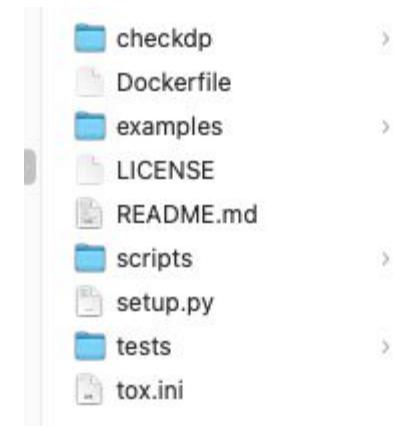
Other kinds of project organization

Common file organization for academic papers (example)



Common organizational structure for a software package (from CheckDP)

Wang, Y., Z. Ding, D. Kifer, and Danfeng Zhang.
2020. CheckDP: An Automated and Integrated Approach for Proving Differential Privacy or Finding Precise Counterexamples. DOI:
<https://doi.org/10.1145/3372297.3417282>



Key questions

- Does the project structure make sense for the type of project you are curating?
- Does the project structure conform to similar types of projects (code and data/software libraries/etc.)?
- Can you clearly describe how the files are organized/why they are organized in a particular way?
- Is the project structure described in the README?

Be CAREFUL when you reorganize a project!

- Altering the structure of a project will cause it to BREAK if you do not make corresponding changes to the code itself!

```
# Load predictions and observations of test data.  
  
# Start with Version 1.  
observed_fname = \  
    '../Data/version1_observed.pkl'  
  
predicted_fname = \  
    '../Data/version1_predicted.pkl'
```

What would happen if I decided to make a subdirectory in “Data” called “version1” for all the version1 data?

Code Curation Questions

- Are files in the submission **organized in a consistent and logical manner?**

Markdown/Markup

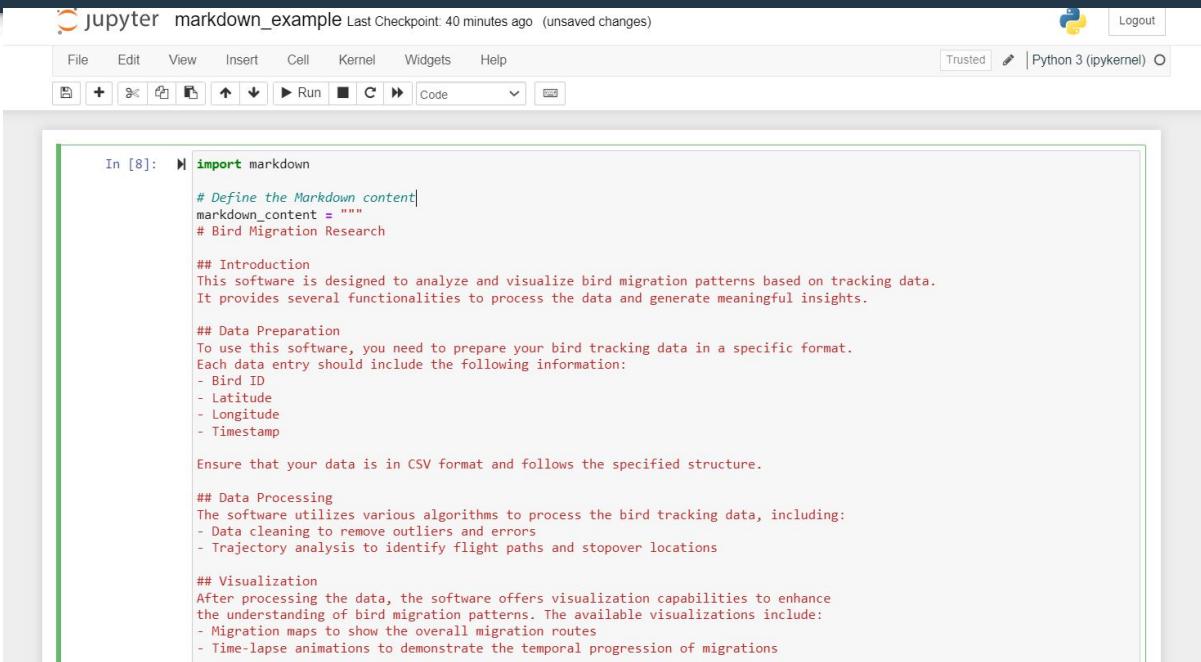
Note: This section was not covered in the pilot workshop due to time constraints. Slides are included for reference and to facilitate future work.

Internal Documentation Strategies: Markdown/Markup

Markdown is a language used to add visual features of text. The markdown process varies depending on the development environment. In Python and R, you can use packages to help create markdown text. However, you can also use certain development environments to render markdown specifications. MATLAB is similar as it natively supports markdown.

Markdown Library in Python

Install and import the markdown library to render text. Content can be stored in a variable.



The screenshot shows a Jupyter Notebook interface with the title "jupyter markdown_example". The notebook has a single cell labeled "In [8]". The cell contains the following Python code:

```
In [8]: M import markdown

# Define the Markdown content
markdown_content = """
# Bird Migration Research

## Introduction
This software is designed to analyze and visualize bird migration patterns based on tracking data. It provides several functionalities to process the data and generate meaningful insights.

## Data Preparation
To use this software, you need to prepare your bird tracking data in a specific format. Each data entry should include the following information:
- Bird ID
- Latitude
- Longitude
- Timestamp

Ensure that your data is in CSV format and follows the specified structure.

## Data Processing
The software utilizes various algorithms to process the bird tracking data, including:
- Data cleaning to remove outliers and errors
- Trajectory analysis to identify flight paths and stopover locations

## Visualization
After processing the data, the software offers visualization capabilities to enhance the understanding of bird migration patterns. The available visualizations include:
- Migration maps to show the overall migration routes
- Time-lapse animations to demonstrate the temporal progression of migrations
"""

print(markdown.markdown(markdown_content))
```

The rendered output of the Markdown content is displayed below the code cell. It includes a main heading "# Bird Migration Research", a section titled "Introduction" with a descriptive paragraph, a section titled "Data Preparation" with a list of required data fields, a note about CSV format, a section titled "Data Processing" with a list of algorithms used, and a section titled "Visualization" with a list of available visualizations.

Markdown Library in Python

Once text is stored, use the markdown library to output an HTML rendering.

```
## Getting Started
Follow the steps below to get started with the software:
1. Install the required Python packages: `pip install pandas matplotlib`
2. Prepare your bird tracking data in CSV format.
3. Configure the software with the data file path and desired parameters.

## Conclusion
This software provides a comprehensive framework for bird migration research.
It simplifies the process of data analysis and visualization, enabling researchers to gain valuable insights.

For detailed documentation and code examples,
refer to the project's GitHub repository: [Bird Migration Research Software](https://github.com/your/repository)

---
"""

# Convert Markdown to HTML
html_content = markdown.markdown(markdown_content)

# Save the HTML content to a file
with open("bird_migration_research.html", "w") as f:
    f.write(html_content)

print("Markdown content converted and saved as HTML.")
```

Markdown content converted and saved as HTML.

Markdown Library in Python

Bird Migration Research

Introduction

This software is designed to analyze and visualize bird migration patterns based on tracking data. It provides several functionalities to process the data and generate meaningful insights.

Data Preparation

To use this software, you need to prepare your bird tracking data in a specific format. Each data entry should include the following information: - Bird ID - Latitude - Longitude - Timestamp

Ensure that your data is in CSV format and follows the specified structure.

Data Processing

The software utilizes various algorithms to process the bird tracking data, including: - Data cleaning to remove outliers and errors - Trajectory analysis to identify flight paths and stopover locations - Statistical analysis to extract migration patterns and trends

Visualization

After processing the data, the software offers visualization capabilities to enhance the understanding of bird migration patterns. The available visualizations include: - Migration maps to show the overall migration routes - Time-lapse animations to demonstrate the temporal progression of migrations - Statistical charts to present migration trends and patterns

Getting Started

Follow the steps below to get started with the software: 1. Install the required Python packages: `pip install pandas matplotlib` 2. Prepare your bird tracking data in CSV format. 3. Configure the software with the data file path and desired parameters. 4. Run the software using the provided Python script.

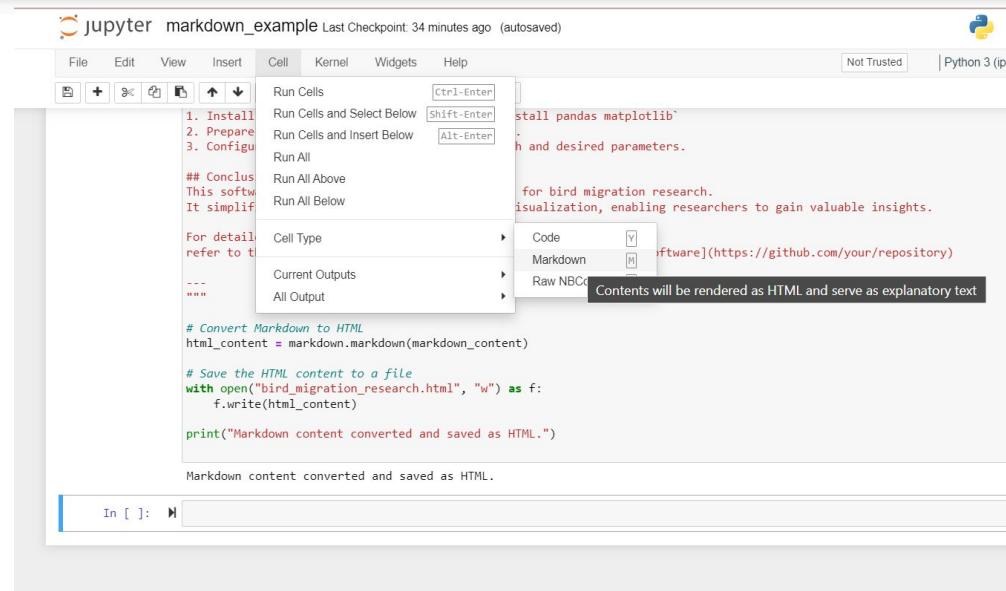
Conclusion

This software provides a comprehensive framework for bird migration research. It simplifies the process of data analysis and visualization, enabling researchers to gain valuable insights into bird behavior and migration patterns.

For detailed documentation and code examples, refer to the project's GitHub repository: [Bird Migration Research Software](#)

Markdown Cells in Python – Jupyter Notebook

Some Integrated Development Environments (IDEs) allow you to create markdown cells within your code.



The screenshot shows a Jupyter Notebook interface with the title "jupyter markdown_example Last Checkpoint: 34 minutes ago (autosaved)". The notebook has a "Cell" menu open, with the "Cell Type" option highlighted. A tooltip for "Cell Type" states: "Contents will be rendered as HTML and serve as explanatory text". The main area contains a cell with the following content:

```
#> %%conclusion
#> This software is designed to facilitate bird migration research by providing a user-friendly interface for data visualization, enabling researchers to gain valuable insights.
#>
#> For detailed information, refer to the accompanying documentation.
#>
#> ---
```

Below the cell, the output is shown as:

```
Markdown content converted and saved as HTML.
```

Markdown Cells in Python – Jupyter Notebook

Use special characters to establish headings and other visual features.

Bird Migration Research

Introduction

This software is designed to analyze and visualize bird migration patterns based on tracking data. It provides several functionalities to process the data and generate meaningful insights.

Data Preparation

To use this software, you need to prepare your bird tracking data in a specific format. Each data entry should include the following information:

- Bird ID
- Latitude
- Longitude
- Timestamp

Ensure that your data is in CSV format and follows the specified structure.

Data Processing

The software utilizes various algorithms to process the bird tracking data, including:

- Data cleaning to remove outliers and errors
- Trajectory analysis to identify flight paths and stopover locations

Visualization

After processing the data, the software offers visualization capabilities to enhance the understanding of bird migration patterns. The available visualizations include:

- Migration maps to show the overall migration routes
- Time-lapse animations to demonstrate the temporal progression of migrations

Getting Started

Follow the steps below to get started with the software:

1. Install the required Python packages: `pip install pandas matplotlib`
2. Prepare your bird tracking data in CSV format.
3. Configure the software with the data file path and desired parameters.

Markdown Cells in Python – Jupyter Notebook

Bird Migration Research

Introduction

This software is designed to analyze and visualize bird migration patterns based on tracking data. It provides several functionalities to process the data and generate meaningful insights.

Data Preparation

To use this software, you need to prepare your bird tracking data in a specific format. Each data entry should include the following information:

- Bird ID
- Latitude
- Longitude
- Timestamp

Ensure that your data is in CSV format and follows the specified structure.

Data Processing

The software utilizes various algorithms to process the bird tracking data, including:

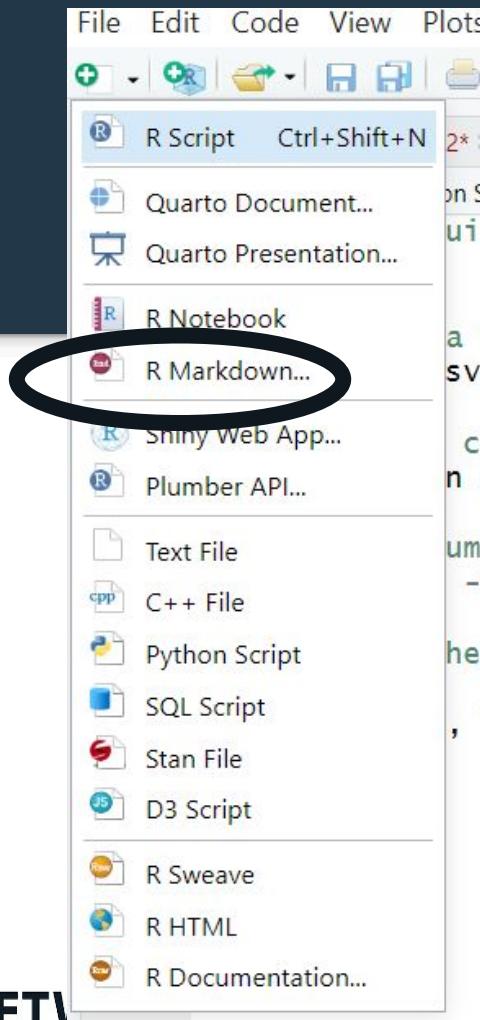
- Data cleaning to remove outliers and errors
- Trajectory analysis to identify flight paths and stopover locations

Visualization

After processing the data, the software offers visualization capabilities to enhance the understanding of bird migration patterns. The available visualizations include:

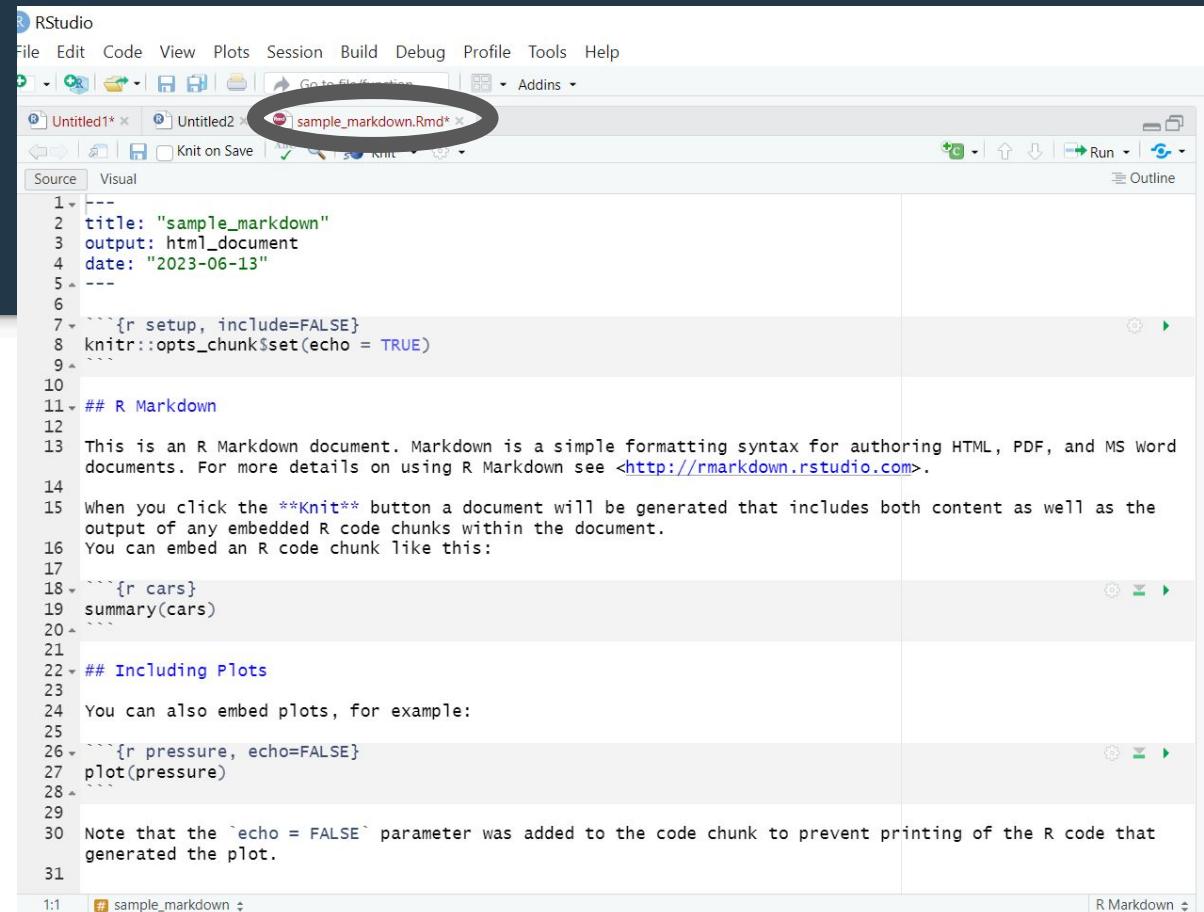
Markdown in R

R Studio allows you to create markup files that can include snippets of R code. (.rmd)



Markdown in R

R Studio allows you to create markup files that can include snippets of R code. (.rmd)



```
1 ---  
2 title: "sample_markdown"  
3 output: html_document  
4 date: "2023-06-13"  
5 ---  
6  
7 ```{r setup, include=FALSE}  
8 knitr::opts_chunk$set(echo = TRUE)  
9  
10 ## R Markdown  
11  
12 This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
13  
14 When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.  
15 You can embed an R code chunk like this:  
16  
17 ```{r cars}  
18 summary(cars)  
19 ```  
20  
21 ## Including Plots  
22  
23 You can also embed plots, for example:  
24  
25 ```{r pressure, echo=FALSE}  
26 plot(pressure)  
27 ```  
28  
29  
30 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.  
31
```

Markup in R

Select Knit to render your file in HTML or some other format.

sample_markdown

2023-06-13

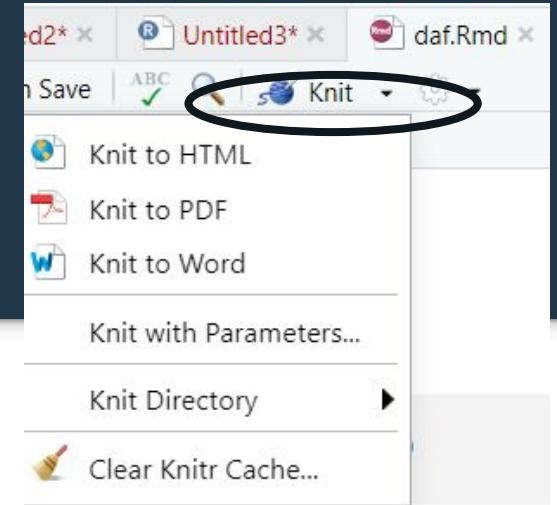
R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

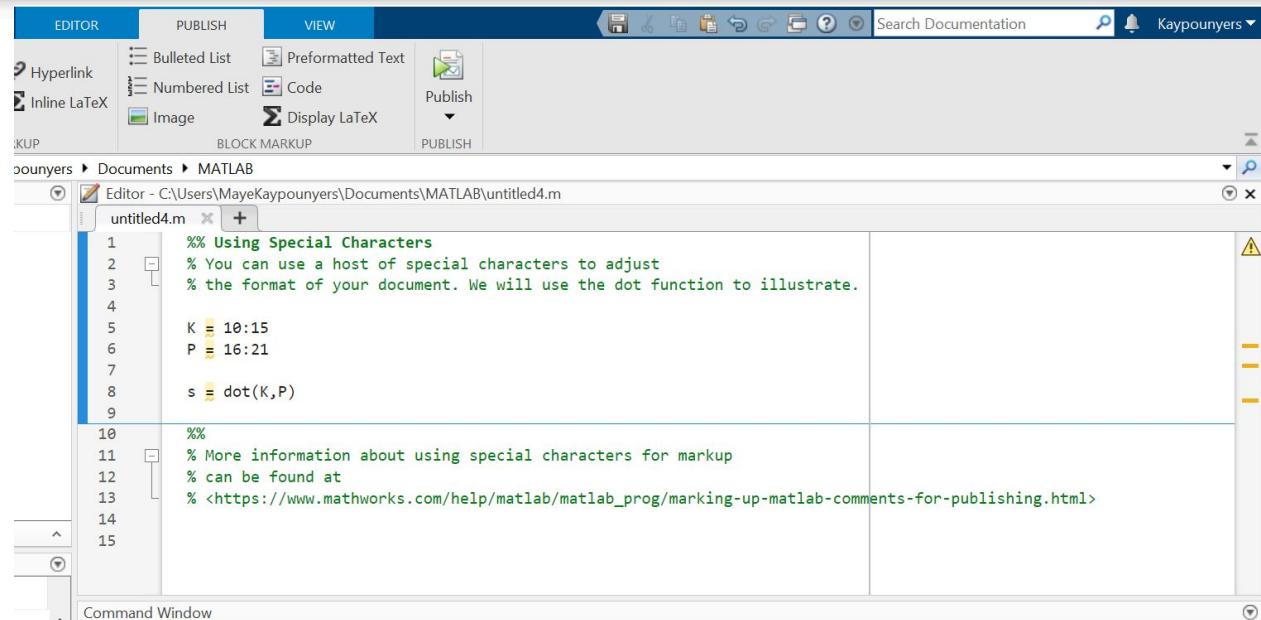
```
summary(cars)
```

```
##      speed         dist
## Min.   : 4.0   Min.   : 2.00
## 1st Qu.:12.0   1st Qu.: 26.00
## Median :15.0   Median : 36.00
## Mean   :15.4   Mean   : 42.98
## 3rd Qu.:19.0   3rd Qu.: 56.00
## Max.   :25.0   Max.   :120.00
```



Markdown in MATLAB

MATLAB natively supports markdown, so be sure to consult the use of markdown symbols for the system (ex. Levels of % symbol)



The screenshot shows the MATLAB Editor interface. The top menu bar has tabs for EDITOR, PUBLISH (which is selected), and VIEW. Below the menu are various toolbar icons. The main workspace shows a script named 'untitled4.m' with the following content:

```
%> Using Special Characters
% You can use a host of special characters to adjust
% the format of your document. We will use the dot function to illustrate.

K = 10:15
P = 16:21

s = dot(K,P)

%%
% More information about using special characters for markup
% can be found at
% <https://www.mathworks.com/help/matlab/matlab\_prog/markup-up-matlab-comments-for-publishing.html>
```

The code uses standard MATLAB syntax with additional Markdown-like comments (e.g., % Using Special Characters) and a special character (%) to denote levels of nesting.

Markdown in MATLAB

After creating your file,
select the publish
button to render the
text in HTML.

Using Special Characters

You can use a host of special characters to adjust the format of your document. We will use the dot function to illustrate.

```
K = 10:15  
P = 16:21  
  
s = dot(K,P)
```