

# Principles of Data Science

PSYC 259

Winter 2022

E-mail: [franchak@ucr.edu](mailto:franchak@ucr.edu)

Office Hours: By appointment

Office: <https://ucr.zoom.us/my/johnfranchak>

Web: [padlab.ucr.edu](http://padlab.ucr.edu)

Class Hours: Mondays 1:00pm-3:50pm

Class Room: Psychology 1205

---

## TA Information

- Jacob Elder
- Office Hours: TBD
- E-mail: [jacob.elder@email.ucr.edu](mailto:jacob.elder@email.ucr.edu)
- Zoom: <https://ucr.zoom.us/my/jelder>

## Course Description

Most quantitative courses (importantly) focus on the final steps of data analysis—conducting and understanding statistical tests. However, much of the work in data science is taking raw data, often from multiple, incompatible sources, and processing those data into a usable form. This course will emphasize the importance of robust, documented, and automated workflows for processing data to save time, reduce errors, improve reproducibility, and facilitate collaboration among multiple researchers. We will also spend time on data visualization and communication—an important part of creating, checking, and collaborating on data workflows. We will use the R programming language, Github, and Rmarkdown to work through examples, but the focus is on concepts/best practices that can be applied to any software or programming language. The course is open to students who have little programming experience or experience with R. The goal is for students at all levels of programming experience to set goals to improve their data science skills.

## Course Objectives

The goals of this course are for you to critically analyze and improve your data analysis workflows. Implementing robust, automated procedures for handling data will allow you to:

- Foster open science through increased transparency, reproducibility, and easier data sharing
- Increase the fidelity of your data and analyses by detecting and preventing errors
- Better understand and communicate about your data
- Save time by preventing errors, automating tasks, and reusing code
- Facilitate collaboration with organized and documented workflows

## Course Materials

- The course [Github page](#) has links to project files, readings, and the syllabus
- Readings from *R for Data Science* are available [online](#).
- PDFs of other course readings referenced below are available on Github in the syllabus folder.

## Other helpful resources

- RStudio [cheatsheets](#) for base R, data import, ggplot, R Markdown, and other packages.
- [STAT 545](#), another good, free textbook on data wrangling and exploration
- [Happy Git with R](#), for all of your Github-Git-R-RStudio needs
- [R Markdown: The Definitive Guide](#)

## Course Policies

Please follow the course [Github page](#) for updates to the syllabus and materials. I will communicate changes through email and/or Canvas announcements.

## Assignments and Grading

Grading for the class is S/NC. Your grade will be based on the following:

Component	Weight
Participation	20%
Weekly assignments	20%
Workflow critique	20%
Final project	40%

### *Participation*

You are expected to attend each class and participate in class discussions. Readings should be completed prior to class so that you can contribute to discussions.

### *Weekly Assignments*

I will assign short practical assignments after most weeks of class. These exercises will develop your skills in R programming and other concepts that we discuss. You are allowed to work in groups of 2-3 students to complete these exercises if you would like, but each student should turn in their own assignment. Please indicate in the top of your assignment the names of the students you worked with. Some class time will be given so that you can start work on the assignment with the help of the instructor/TA. Office hours are also available for help on assignments.

### *Workflow Self-Critique*

You should choose a current or past data analysis project that you have worked with (or one from your lab if you are a newer student). In a short paper (3-4 pages), you will first describe the end-to-end workflow of your data. What are sources of raw data? How are those sources combined and/or processed? What research personnel, computing resources, software, and hardware devices are involved? What is the end product needed for statistical analysis? Next, you will take a

critical eye to your workflow and identify 1) Where are errors most likely to occur, 2) What time-consuming steps could be automated, and 3) how your workflow could be made more transparent and reproducible. You will briefly present your workflow critique in class. Each student will work individually on this assignment.

### *Final Project*

In your final project, you will improve the data workflow that you chose using skills learned in this class. Your final project should be shared with the instructor and TA through an online repository (such as RStudio Online, Github, OSF, or Code Ocean) and allow your end-to-end data workflow to be reproduced (e.g., include raw data files, functions that implement processing steps, etc.). You can use whatever programming languages are necessary (it doesn't need to just be in R), but you should consult with the instructor if R will not be used to ensure that the instructor can run your code (or alternatively, that you demonstrate your workflow to the instructor).

Your project should be a report, either to demonstrate exploratory analyses or to communicate the results, that is written in R Markdown and contains visualizations written in R. In a brief presentation in the last week of class, you will include summarize the changes you implemented in the workflow, describe the problems you believe they have solved, and reflect on what future changes may need to be implemented.

### **Attendance**

In the event of absences due to illness or family obligations, please contact the instructor immediately so that we can make arrangements. Everyone is in a different situation, so I will always work with you to figure out a way to handle disruptions to learning. I want everyone in this class to succeed!

## Class Schedule

Readings should be completed prior to each class.

### Week 01, 01/03: Course Goals

- Goals of the class
- Logistics
- Readings and assignments

R4DS: [Introduction](#) and [Workflow Basics](#)

### Week 02, 01/10: Data Workflow

- File organization
- Version control
- Workflow example
- SKILLS: Basic I/O functions, getting help, using Github

R4DS: [Data import](#), [Tibbles](#), and [Projects](#)

OTHER READINGS: Minimizing mistakes article (Rouder et al., 2019), [Github beginner guide](#)

*GitHub Project Links:*

- [Instructions to set up R, RStudio, git, and Github](#)
- [Basic file import](#)
- [Multiple file import](#)
- [Import homework](#)

### Week 03, 01/17: Data Structure

- Data types (numbers, factors, strings, dates)
- Data organization (variables/observations)
- SKILLS: Factors, logic, data wrangling

R4DS: [Data transformation](#), [Tidy data](#), and [Factors](#)

OTHER READINGS: Tidy Data [[@Wickham2014](#)]

*GitHub Project Links:*

- [Data wrangling](#) (scripts 1 and 2)
- [Data wrangling homework](#)

### Week 04, 01/24: Automation

- Automating your analyses
- Writing more efficient code
- SKILLS: More data wrangling, iteration (map, for loops)

R4DS: [Functions](#), [Vectors](#) and [Iteration](#)

*GitHub Project Links:*

- [Data wrangling](#) (scripts 3 and 4)
- [Automation](#)
- [Tidying and automation homework](#)

## **Week 05, 01/31: Exploratory Data Analysis**

- Data validation
- Automating visualizations
- *SKILLS*: Basic visualizations with ggplot2

*R4DS*: [Data visualization](#) and [Exploratory data analysis](#)

*GitHub Project Links*:

- [Visualization and EDA](#)
- [Visualization and EDA group exercise](#)

## **Week 06, 02/07: Custom Functions Part 1**

- **WORKFLOW PRESENTATIONS GROUP 1**
- Finding new packages/APIs vs. writing your own functions
- Defining custom functions within a script
- *SKILLS*: Writing basic functions

## **Week 07, 02/14: Custom Functions Part 2**

- **WORKFLOW PRESENTATIONS GROUP 2**
- “Technical debt”, “design smells”, and code refactoring
- Sourcing functions
- Working with function arguments
- *SKILLS*: Writing advanced functions

*OTHER READING*: Technical debt [[@SuryanarayanaSamarthyam2014](#)]

## **Week 08, 02/21: Data Sharing and Reproducibility**

- Reuse-minded project management
- Reproducible reports
- Preserving programming environment and analyses
- *SKILLS*: R Markdown, package control

*R4DS*: [R Markdown](#)

*OTHER READINGS*: Transparency in psychological science [[@KleinHardwicke2018](#)] and Care and feeding of data [[@GoodmanPepe2014](#)]

## **Week 09, 02/28: Communication**

- Communicating through graphical styles
- Interactive plots for data exploration
- Manuscript preparation in R Markdown
- *SKILLS*: ggplot and extensions, papaja

*R4DS*: [Graphics for communication](#)

*OTHER READING*: Designing graphs for decision-makers [[@ZacksFranconeri2020](#)], [Chartjunk](#) from [[@Tufte1990](#); [@Tufte2001](#); [@Tufte2006](#)], and (optional) Graph construction [[@Witt2019](#)].

### **Week 10, 03/07: Project Presentations**

- How have you changed your workflow?
- What have you learned about your data?
- What problems are still unresolved?