# **Syllabus**

## **Principles of Data Science**

#### Instructor

- Prof. John Franchak
- **1** Psychology 2125
- 🗷 franchak@ucr.edu
- Schedule an appointment

#### Course details

- **ii** Mondays
- iii January 6 to March 10
- **Q** 9:00am-11:50am
- Psychology 1205

#### **TA** Information

- 🚨 TBD
- **\begin{aligned}
  \begin{aligned}
  \begin{alig**
- **■** SOMWHERE
- \tod@ucr.edu

# **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 Quarto 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 website has the schedule, lectures, and assignments.
- The Github page has all of the code repositories that we cover, and the source code for this website.
- Readings from R for Data Science are available online.
- PDFs of other course readings are linked from the schedule and lectures
- Other helpful resources are compiled for your future reference.

## **Course Policies**

Keep an eye on the course website 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:

Weight
20%
20%
20%
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. 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 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 Quarto and contains visualizations written in R.

#### **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!