

Data Management and Visualization

BI 410/510 | Winter 2024 | Klamath 33 | MW 11 am - 12:50 pm | 4 credits

Overview

This course covers the non-statistical aspects of the data life cycle, including how to store, clean, visualize and communicate data (Figure 1). It is intended as a complement to statistics courses - we will cover how to get your data into shape for analysis, and how to communicate your findings visually. It is primarily a methods class and will be taught in R (but there is no expectation that students know R coming in).

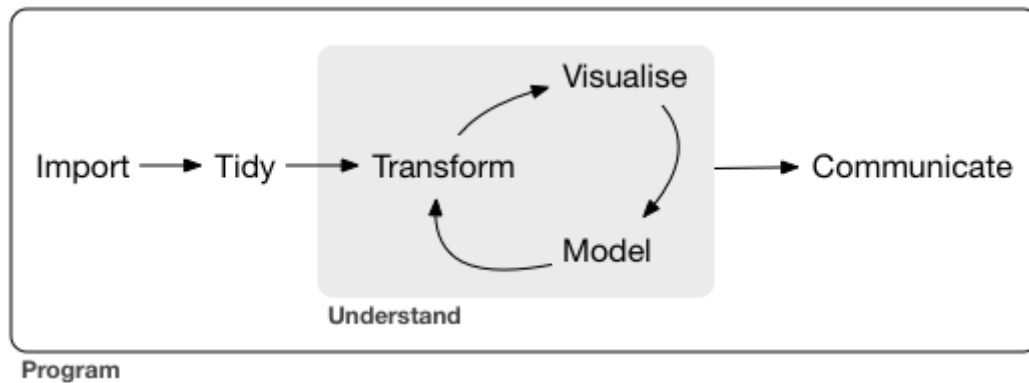


Figure 1: The data life-cycle (figure from Golemund and Wickham)

This course satisfies the “Analytical Approaches” requirement for Environmental Science majors and the “MAPS” requirement for Biology majors.

Instructor: Dr. Lauren Hallett (she/her) is an associate professor of plant community ecology specializing in ecological restoration. Data synthesis is fundamental to her work and she has a passion for developing and teaching open, reproducible science.

GE: Jasmin Albert (she/her) is a PhD student interested in how plant-pollinator interactions alter plant coexistence.

Email: hallett@uoregon.edu; pimsupaa@uoregon.edu

Office hours: 3:30-4:30 Wednesday (Hallett, Pacific 220)

Canvas site: Our website is accessible via the UO Canvas server, use your UO email and password to access the site. Pre-recorded videos and scripts will be available in modules for each week. Problem sets will be distributed and submitted via Canvas. <https://canvas.uoregon.edu/>

How we will contact you: Our class will communicate through our Canvas site. Announcements and emails are archived there and automatically forwarded to your UO email, and can even reach you by text. Check and adjust your settings under Account > Notifications.

Course material

Our primary course material will be **R for Data Science** by Garrett Golemund and Hadley Wickham, which is available for free online (<http://r4ds.had.co.nz/>); supplemental readings are uploaded to Canvas. Please make sure to read book sections and papers before the class in which they are assigned.

Objectives

By completing this course, students will be able to:

- 1) Interpret figures in scientific papers and popular media
- 2) Locate data relevant to biological and environmental questions
- 3) Understand the steps linking raw data to communicated findings
- 4) Create exploratory and publication-worthy graphs

Structure of the course

The course is broken into five approx. 2-week modules: Intro to R, Visualize, Transform, Wrangle, and Communicate. For each module the readings, scripts and assignments are posted in advance, along with short videos that overview key concepts and demonstrate their implementation in R.

Early in a module we will use class time to go over the concepts and scripts in depth, and later in a module we will use these times as “work sessions” for you to collaborate on the problem sets and other assignments. I encourage you to watch the posted videos before class. This makes for essentially a “flipped” classroom setting, allowing us to spend less class time on lectures and more time completing the scripts and making substantial progress on the problem sets, aided by having your classmates and us available.

There is a lot of value in working collaboratively, but there may be times when you need to miss class. If you have to miss a day just work through the scripts and videos. If you have to miss an extended period, *please* communicate with us about it so we can make a plan to keep you on track.

Participation will be based primarily on your completion of the scripts and secondarily by interactive engagement (either in class or via discussion boards). **The one day I expect you to prioritize attendance is peer review day.** Attendance during peer review is important to give and receive feedback with your fellow classmates. *If you cannot attend please communicate with us in advance.*

Class assignments and requirements

There are two main components required for successful completion of the course.

A. Problem sets Problem sets are designed to develop the skills you learn in class and to gain comfort in the R environment through practice. Problem sets will typically include designing and implementing code and interpreting code and figures. There will be four problem sets, to be submitted on Canvas by midnight the day they are due. Students are encouraged to collaborate on problem sets, and to come to “work day” classes with ideas and questions.

B. Final project The focal experience of the class will be to develop a research project that addresses a biological or environmental question with data. Students will be expected to identify a question, contextualize the question with a literature review, analyze data relevant to answering the question, and interpret and communicate that data with a workflow in R. In general, students will be using one of our pre-identified datasets, but if you have a data project of your own (particularly graduate students) we can discuss its suitability for this requirement.

Grade allocation

Grading will be based on a total of 200 points, where 90% of the points will earn an A, 80% a B, etc. Participation will reflect attendance and involvement in discussion and in-class exercises. The breakdown by assignments is as follows:

Assignment	Points
4 Problem sets (20 pts each)	80
Final project	
Part I: Annotated bibliography	10
Part I: Literature review and proposed workflow	25
Part II: Peer review	20
Part III: Final paper	35
Participation	
Completed scripts (2 pts each)	18
Interaction (in class or message board)	12

Policies

- 1) Please note that assignments are due on Canvas. If an assignment is late, we will deduct 10% of the total points allocated to that assignment, and we will deduct 10% for each additional late day.
- 2) Please do not come to class if you are sick! We expect all students to actively participate in exercises and contribute to discussion, but this can be either in class, on the discussion boards, or both. There are enough resources posted on Canvas that you can be successful even if you have to miss a few sessions. Please just communicate with us so we can make a plan for keeping you on track.
- 3) We will follow school policy of plagiarism and academic dishonesty. All students need to be familiar with the Student Conduct Code (<https://policies.uoregon.edu/vol-3-administration-student-affairs/ch-1-conduct/student-conduct-code>).
- 4) Use of artificial intelligence systems (e.g., ChatGPT, iA Writer, etc.) is allowed provided you note explicitly where in your work process you used AI (e.g. generating an outline or first draft) and which platform(s) you used. If you use code or text generated by an artificial intelligence system as part of an assignment submission you must attribute the text to the AI-based system that is its source. For example, if you include text generated by ChatGTP, you must cite the source as follows:

ChatGPT. (Year, Month, Day of query). "Text of your query/prompt." Generated using OpenAI. <https://chat.openai.com/>

Deadlines

Completed in-class scripts are due by midnight on:

F 1/12 Intro to R
F 1/26 ggplot2 and geoms
F 2/9 rearrange and pipelines
F 2/23 relational data and tidy data
F 3/1 better graphics

Problem sets are due by midnight on:

F 1/19 PS 1
F 2/2 PS 2
F 2/16 PS 3
F 3/1 PS 4

Final project milestones are due by midnight on:

M 1/29 Topic ID

M 2/12 Annotated bibliography M 2/26 Literature review and workflow plan

Su 3/3 Peer review

M 3/18 Final project

Please see the BI 410/510 calendar on Canvas (click on the Calendar icon in the dark green bar on the far left, then click on Bi 410/510 in the Calendars drop-down) for a calendar view of deadlines. You can also link this to your preferred calendar app to keep track of deadlines! To do so, click on “Calendar Feed” and copy the link to your calendar app.

Course topics and tentative schedule

The topics on the tentative outline are subject to change. This is a guess, but we will take as long as needed on each lesson. Topics and lessons generally correspond to the noted chapter numbers in the book, additional readings will be posted to Canvas and emailed the week prior to when they should be read. Please note that the reading information is also summarized within each module on Canvas.

Day	Date	Module	Lesson
1	M 1/8	Visualize	Plotting before analyzing
2	W 1/10	General	Overview of R, R Studio Chapter 1, 6.1-6.3
3	W 1/17	General	<i>Work day: swirl practice</i>
4	M 1/22	Visualize	ggplot2 : aesthetic mapping and facets 3.1-3.5
5	W 1/24	Visualize	ggplot2 : geometric objects, coordinate systems 3.6-3.10
6	M 1/29	Visualize	RMarkdown and importing data Chapter 11, 27.1-27.4
7	W 1/31	Visualize	<i>Work day</i>
8	M 2/5	Transform	Rearranging data and dplyr 5.1-5.5
9	W 2/7	Transform	Grouping, summarizing and piping with dplyr 5.6-5.7, 18.1-18.3
10	M 2/12	Transform	Workflows
11	W 2/14	Transform	<i>Work day</i>
12	M 2/19	Wrangle	Relational data and joins with dplyr 13.1-13.7
13	W 2/21	Wrangle	Tidy data and tidyr 12.1-12.7
14	M 2/26	Communicate	Beautiful graphs Chapter 28
15	W 2/28	Wrangle/Communicate	<i>Work day</i>
16	M 3/4	Communicate	Peer review
17	W 3/9	Communicate	<i>Work day</i>
18	M 3/11	Communicate	RMarkdown v2 27.5-27.6
19	W 3/13	Communicate	Odds and ends