

Course Name and Number: DATA 607 – Data Acquisition and Management

Credits: 3 cr.

Prerequisite(s): none

How is this course relevant for data analytics professionals?

Most data analytics professionals spend *most* of their time getting data and preparing it for analysis. This is the course that teaches these key skills, as we work with both structured and unstructured data.

Course Description:

In this course students will learn about core concepts of contemporary data collection and its management. Topics will include systems for collecting data (real time, sensors, open data sets, etc.) and implications for practice; types of data (textual, quantitative, qualitative, GIS, etc.) and sources; an overview of the use of data, including what and how much should be collected and the distinction between data, information, and knowledge from a data-centric point of view; provenance; managing data with and without databases; computer and data security; data cleaning, fusing, and processing techniques; combining data from different sources; storage techniques including very large data sets; and storing data keeping in mind privacy and security issues.

Students will be required to create a working system for a large volume of data using publicly available data sets.

Course Learning Outcomes:

By the end of the course, students should be able to:

- Load data into R from various data sources, including CSV files, Excel spreadsheets, relational databases, APIs, and web pages.
- Perform various data cleansing and transformation work, including splitting, combining; resampling; variable creation; data aggregation; sorting and filtering data; strategies for working with outliers and missing data; data visualization and analysis in support of data cleansing activities.
- Understand different information architectures, data types, and data structures.
- Understand relational and non-relational database design and querying.
- Provide context for data science

Program Learning Outcomes addressed by the course:

- Business Understanding. Apply frameworks and processes to build out data analytics solutions from understanding of business goals.
- Data Culture. Embody and champion the highest standards for the ethical and moral use of data; understand issues related to data privacy and data security.
- Solid foundational data programming skills, using industry standard tools, essential algorithms, and design patterns for working with structured data, unstructured data and big data.
- Data understanding. Collect, describe, model, explore and verify data.
- Data preparation. Selecting, cleaning, constructing, integrating, and formatting data.

Assignments and Grading:

Assignments (6 x 50)	30%
Projects (3 x 90)	27%
Final Project Proposal (1 x 20)	2%
Final Project (1 x 150)	15%
Final Project Presentation (1 x 30)	3%
Discussion Participation (14 x 10)	14%
Data Science in Context Presentation (1 x 50)	5%
TidyVerse recipes	4%
TOTAL	100%

Quality of Performance	Letter Grade	Range %
Excellent - work is of exceptional quality	A	93 - 100
	A-	90 - 92.9
Good - work is above average	B+	87 - 89.9
Satisfactory	B	83 - 86.9
Below Average	B-	80 - 82.9
Poor	C+	77 - 79.9
	C	70 - 76.9
Failure	F	< 70

Notes

- All discussions, projects, and assignments--unless otherwise noted--are due end of day on Sundays.

Late projects are not accepted. However, there are eight assignments and four projects assigned, and your final grade is based on your six highest-scoring assignments and your three highest-scoring projects.

- Each course week will be available on the previous Friday at 6:00 a.m. ET.
- **Course Completion Requirements.** To pass this course, you must complete at least six assignments, three projects, the final, and make the final presentation. If you cannot deliver your final presentation in our 5/17 Meetup, you'll need to make available a recorded version of your final presentation before 5/17. Final grades will be submitted on 5/18.
- There are some **short ungraded hands on labs** that will help you prepare for your weekly programming assignments and projects. You don't need to turn these in.
- **"Discussion", "Data Science in Context Presentations", and "TidyVerse Recipes"** While this material is important, please note that this work only makes up only 23% of your grade. Please do the readings and participate in the discussions and any discussion-related group assignments, make your Data Science in Context presentations, and participate in the creation and editing of TidyVerse recipes on the shared GitHub site. At the same time, if you have limited time for the course, please remember to invest most of your efforts in completing the projects and assignments. The assignments merit close attention because they will help you to be successful on the projects.
- **Reproducibility Requirement, Testing Requirement, But Not Perfection!** Students are responsible for providing all code and data so that I can test your work. If you turn in code that does not run, you will not receive credit, unless you also include an explanatory note at the time of submission. At the same time, you don't need to turn in perfect code. Generous partial credit will be given for deliverables that are timely, tested, and reproducible. Cutting corners—as long as they are documented at the time of submission—is also acceptable.
- **Groupwork** is encouraged on most projects and assignments and required on Project 3. Effective virtual collaboration is highly valued in the data science marketplace; because of its interdisciplinary nature, much of the work that needs to be done requires more than one person, and increasingly often at multiple locations.
- **Earning a Grade of A.** If you complete the course work correctly and on time, you'll comfortably pass the course. A grades will be reserved for students that go above and beyond, such as consistently taking on challenge assignments.

Policy on Sharing and "Stealing" Code. In this course, you may collaborate, and you may take base code from whatever sources you wish. But you must document what you started with, and what you added, so you are graded only on your own contributed work!

Course Learning Materials



Required Texts:

- *R for Data Science* (2e) by Hadley Wickham and Garrett Golemund. This is the primary text for the course. Freely readable here: <https://r4ds.hadley.nz/>. The first edition is also available in print.
- *Text Mining with R: A Tidy Approach*, Julia Silge and David Robinson. O'Reilly, 2017. Freely readable
- Max Kuhn and Kjell Johnson, *Feature Engineering and Selection: A Practical Approach for Predictive Models* (Chapman & Hall/CRC Data Science Series) 1st Edition, 2019. Freely readable at <https://bookdown.org/max/FES/intro-intro.html>

Print copies of each of these texts is available for download.

Recommended Texts:

- Any book on SQL, such as [*The Language of SQL, 3rd Edition*](#) by Larry Rockoff. ISBN: 978-0137632695. My favorite text for PostgreSQL is [*Practical SQL, 2nd Edition*](#), by Anthony deBarros. Alternatively, there are many excellent on-line resources, such as <http://sqlzoo.net>.

Relevant Software, Hardware, or Other Tools:

We will make use of the R programming environment and the RStudio IDE. We will use other open source software, including PostgreSQL and MongoDB. Details for obtaining and installing the appropriate software will be provided in the course materials. All of the software will work on (or from) both PCs and Macs.

Contact Information:

Andy Catlin
andrew.catlin@sps.cuny.edu
616-638-8344

How This Course Works:

Meetups take place every week on Wednesdays from 6:45 p.m. to 7:45 p.m. ET. Please see course site for specific dates. You are strongly encouraged to attend; all meetups will be recorded. You are not required to attend the meetups, but you are responsible for watching the recording if you're not able to attend.

Occasional Weekend Office Hours A few times during the semester, we'll have optional additional office hours on topics of interest, especially around data engineering.

Regular Office Hours can be scheduled by e-mail appointment. If you need extra help and are willing to invest the time and effort to be successful, I'll make the time to help you. But...you should not be asking for extra help on a project the day before it's due, since this indicates that you're not investing the time and effort to be successful.

You are encouraged to ask questions on the "Ask Your Instructor" forum on the course discussion board where other students will be able to benefit from your inquiries. I can set up a Zoom session for screen sharing. For the most part, you can expect me to respond to questions by email within one business day.

In addition to the graded discussion items on the Blackboard Learning Management System course site, there is a Slack channel for support and general questions here:

https://join.slack.com/t/sps-16h7212/shared_invite/zt-1ni7pxv75-1WmKD8xJ~udY65c3k_JPKA

<https://data607spring2023.slack.com>

Here is the link to sign up for your Data Science in Context 5-minute presentation:

<https://doodle.com/poll/inakz9m4vams7wqd>

Here is our weekly meetup link:

<https://zoom.us/j/99498039934?pwd=Q0R1MVNMbi9oYjRGaEhmVVJCMFEzZz09>

Meeting ID: 994 9803 9934, Passcode: 139406,

Use computer audio or call +1 646 876 9923 US (New York)

Unit	Topic	Core Readings	Deliverables
Week 1 Jan 25 – Jan 29	Data Ethics; R: Data Types and Basic Operations	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapters 1, 2, 3, 7, 14, 15, 27, 28	Meetup on 01/25, 6:45 p.m. EST Week 1 Assignment
Week 2 Jan 30 – Feb 05	R and SQL	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapter 8, 9, 10, 23	Meetup on 02/01, 6:45 p.m. EST Week 2 Assignment
Week 3 Feb 06 – Feb 12	R: Character Manipulation and Date Processing	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapters 14, 15, 16, 17, 18, 19	Meetup on 02/08, 6:45 p.m. EST Week 3 Assignment
Week 4 Feb 13 – Feb 19	R: Exploratory Data Analysis; Data Imputation	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapters 11, 12, 13, 20	Meetup on 02/15, 6:45 p.m. EST Project 1
Week 5 Feb 20 – Feb 26	R: Working with Tidy Data	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapters 4, 5, 6	Meetup on 02/22, 6:45 p.m. EST Week 5 Assignment
Week 6 Feb 27 – Mar 05	R: Data Transformations; Feature Engineering	<i>Feature Engineering and Selection</i> , http://www.feateengineering/ chapter 1	Meetup on 03/01, 6:45 p.m. EST Project 2
Week 7 Mar 06 – Mar 12	Web Technologies; MongoDB	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapter 25	Meetup on 03/08, 6:45 p.m. EST Week 8 Assignment Project 3 Team Document
Week 8 Mar 13 – Mar 19	Scraping Web Pages	<i>R for Data Science (2e)</i> , https://r4ds.hadley.nz/ chapter 26	Meetup on 03/15, 6:45 p.m. EDT Project 3
Week 9 Mar 20 – Mar 26	Working with Web APIs	httr quickstart vignette, https://cran.r-project.org/web/packages/httr/vignettes/quickstart.html	Meetup on 03/22, 6:45 p.m. EDT; <i>Project 3 team presentations!</i> Week 9 Assignment Tidyverse CREATE due
Week 10 Mar 27 – Apr 02	Natural Language Processing	<i>Text Mining w/ R</i> , https://www.tidytextmining.com/ , ch 1-4	Meetup on 03/29, 6:45 p.m. EDT Week 10 Assignment
Spring Break Apr 03 – Apr 16	Spring Break	No Readings	No Meetups on 04/05 or 04/12
Week 11 Apr 17 – Apr 23	Recommender Systems	<i>Mining Massive Datasets</i> , http://www.mmds.org/ , ch 9	Meetup on 04/19, 6:45 p.m. EDT Week 11 Assignment
Week 12 Apr 24 – Apr 30	Graph Databases	Selected readings from web	Meetup on 04/26, 6:45 p.m. EDT Project 4; Final Project Proposals; TIDYVERSE EXTEND
Week 13 May 01 – May 07	Working with Data in the Cloud; Deployment	No Readings	Meetup on 05/03, 6:45 p.m. EDT Data Science in Context recordings due for students who did not present in class
Week 14 May 08 – May 14	Automated Machine Learning	No Readings	Meetup on 05/10, 6:45 p.m. EDT
Week 15 May 15 – May 17	Final Presentations	No Readings	Meetup on 05/17, 6:45 p.m. EDT Final Project Presentations

Accessibility and Accommodations

The CUNY School of Professional Studies is firmly committed to making higher education accessible to students with disabilities by removing architectural barriers and providing programs and support services necessary for them to benefit from the instruction and resources of the University. Early planning is essential for many of the resources and accommodations provided. Please see: http://sps.cuny.edu/student_services/disabilityservices.html

Online Etiquette and Anti-Harassment Policy

The University strictly prohibits the use of University online resources or facilities, including Blackboard, for the purpose of harassment of any individual or for the posting of any material that is scandalous, libelous, offensive or otherwise against the University's policies. Please see:

http://media.sps.cuny.edu/filestore/8/4/9_d018dae29d76f89/849_3c7d075b32c268e.pdf

ACADEMIC INTEGRITY

Academic dishonesty is unacceptable and will not be tolerated. Cheating, forgery, plagiarism and collusion in dishonest acts undermine the educational mission of the City University of New York and the students' personal and intellectual growth. Please see:

http://media.sps.cuny.edu/filestore/8/3/9_dea303d5822ab91/839_1753cee9c9d90e9.pdf

STUDENT SUPPORT SERVICES

If you need any additional help, please visit Student Support Services:

http://sps.cuny.edu/student_resources/