

DS002 Intro to Data Science

Class DS002: Intro to Data Science

Instructor Douglas Goodwin
Visiting Assistant Professor in Media Studies, Scripps College
Lang 228

Contact dgoodwin@scrippscollege.edu

Class hours MW 11:00–12:15

Office Hours MW 1:00–3:00pm

Discord <https://discord.gg/gtMVkAMV>

Textbook **Data Science from Scratch, 2nd Edition**
Joel Grus.
ISBN: 9781492041139

Classroom Zoom, Steele 229

Description This course is the second part of a two-semester introduction to computer programming and data science. Students will explore, using Python and other tools the nuances of gathering, visualizing and analyzing data to gain insight and intuition with data. Students will be introduced to various data manipulation/analysis, and statistical methods by building their own code from scratch. They will also consider the ethical implications and limitations of creating models of data.

Course Goals

Students completing this course will be able to:

Define and explain key concepts in the data science

Learn how to analyze real-world data.

Gain fluency in basic programming skills in Python with a focus on statistical modeling and machine learning.

Develop and use essential python programming skills in data analysis, data visualization, and machine learning

We begin by reviewing Python and establishing repositories on GitHub. After review and setup you will read two chapters from “Data Science from Scratch” (DSfs) each week and reproduce the code in the book to build your own code libraries. You are encouraged to write your own code (and tests) as long as you use the same function names. Your code won’t be efficient, but it will work and your APIs will match dedicated libraries such as *NumPy* (Numerical Python), *pandas* (Python Data Analysis Library), and *scikit-learn* (machine learning in Python).

Please complete the readings before and start your code before class on Monday. I will lecture on the chapters then we will go through your code together during Monday class.

Wednesdays are lab days: work alone or in small groups to use your code to complete the lab assignments.

Points

| Activity | Points |
|-----------------------|------------|
| 10 Weekly Assignments | 30 points |
| 1 Presentation | 20 points |
| 10 in-class exercises | 30 points |
| Project | 30 points |
| TOTAL | 110 points |

I give you a little extra room in case you miss an exercise or two.

Points and letter grades

This class will be easy if you keep up with the readings and come to class.

| Point range | Grade |
|-------------|-------|
| 90–110 | A |
| 80–90 | B |
| 70–80 | C |
| 60–70 | D |
| < 60 | ... |

Weekly assignments

You will implement the code in the chapters from Data Science from Scratch (DSfS) to build a library of code to use in Deepnote. Push your assignments to your GitHub repository by Monday morning before class. We will use import your code and use it to complete in-class exercises with Deepnote.

Please DON'T copy and paste code! Typing it will help you get familiar with and synthesize the code.

Note the liberal use of the Python's `assert` statement in the sources files. A clean import will give you some assurance that your code is usable. There are more involved ways to test code, even a style of programming called TDD (Test-Driven Development). `assert` statements sprinkled throughout your source code give you TDD-lite!

In-class exercises

Use your code library to solve Data Science problems related to each week's theme. I am creating these exercises now—please ask if you have an idea for an exercise!

Presentation

DSfS uses a *Crow icon* to indicate relevant subjects not covered adequately in the book. Each of you will select one of the “Crows” to present to the class on a Monday. 15 minutes, you may work alone or in small groups.

Projects

Complete a small, real-world project in the final week of the class. You can start with one of the in-class exercises or dream up a project of your own. Projects should use external data and be executed in Deepnote or on Google Colab.

Accessibility

If you have a documented disability (physical or cognitive) that may impair your ability to complete assignments or otherwise participate in the course and satisfy course criteria, please meet with us at your earliest convenience to identify, discuss, and document any feasible instructional modifications or accommodations. You should also contact the Accessible Education Office to request an official letter outlining authorized accommodations.

Credits

I will lean heavily on the content in DSfS. Some of the material in this course is based on other classes. We have also heavily drawn on materials and examples found online and tried our best to give credit by linking to the original source. Please contact us if you find materials where the credit is missing or that you would rather have removed.

Weeks

1. W01 01/17
2. NO CLASS
3. CH01 Introduction
4. W02 01/24
5. CH02 A Crash Course in Python, pt1
6. CH02 A Crash Course in Python, pt2
 1. Class Exercise
7. W03 01/31
8. Python Assignment DUE BEFORE CLASS
9. CH03 Visualizing Data
10. CH04 Linear Algebra

1. Exercise

11. W04 02/07

12. Viz Assignment DUE BEFORE CLASS

13. CH05 Statistics

14. CH06 Probability

1. Coinflip Class Exercise, Monty Hall Problem

15. W05 02/14

16. Pandas Assignment DUE BEFORE CLASS

17. CH07 Hypothesis and Inference

18. CH08 Gradient Descent

1. Class Exercise

19. W06 02/21

20. Assignment DUE BEFORE CLASS

21. CH09 Getting Data

22. CH10 Working With Data

1. Class Exercise

23. W07 02/28

24. Getting data Assignment DUE BEFORE CLASS

25. CH11 Machine Learning

26. CH12 k-Nearest Neighbors

1. Class Exercise

27. W08 03/07

28. KNN Assignment DUE BEFORE CLASS

29. CH13 Naive Bayes

30. CH14 Simple Linear Regression

1. Class Exercise

31. W09 03/14 SPRING BREAK

32. W10 03/21

33. Linear Regression Assignment DUE BEFORE CLASS

34. CH15 Multiple Regression

35. CH16 Logistic Regression

1. Class Exercise

36. W11 03/28

1. Regression Assignment DUE BEFORE CLASS

2. CH17 Decision Trees

3. CH18 Neural Networks

4. Class Exercise

37. W12 04/04

1. Neural Networks Assignment DUE BEFORE CLASS

2. CH19 [Deep Learning]

3. CH20 Clustering

4. Class Exercise

38. W13 04/11

1. Clustering Assignment DUE BEFORE CLASS
2. CH21 Natural Language Processing
3. CH22 Network Analysis
4. Class Exercise

39. W14 04/18

1. NLP Assignment DUE BEFORE CLASS
2. CH23 Recommender Systems
3. CH24 Databases and SQL
4. Class Exercise

40. W15 04/25

1. SQL Assignment DUE BEFORE CLASS
2. CH25 MapReduce
3. CH26 Data Ethics
4. Class Exercise

41. W16 04/02

1. PROJECTS
2. PROJECTS
3. Share on Discord