

CMP 464 DATA HANDLING AND ANALYSIS

Liang Zhao
Fall 2021

COURSE INFORMATION

Instructor: Liang Zhao

Office: GI 101A

Office Hours: M/W 1:00 - 3:00PM

Email: Liang.Zhao1@lehman.cuny.edu

Prerequisites: Programming Methods II, Discrete Mathematics

Textbook: Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython (2nd Edition) [Amazon](#) [O'Reilly](#)

DATA HANDLING

1. Import and export data in various formats
2. Modify and transform data
3. Handle outliers and missing values
4. Split, group, or combine data from multiple sources
5. Sorting, merging, ...

DATA ANALYSIS

1. Identify and extract useful data features
2. Visualize and summarize data distributions and statistics
3. Discover relationships among features
4. Provide estimates or predictions

MAIN TOPICS

1. A fast-paced introduction to Python programming
2. Arrays and vectorized computations with NumPy and SciPy
3. Data cleaning and processing with Pandas
4. Plotting and visualization with Matplotlib
5. Interacting with Data from Web or Databases
6. Managing time series data, images, texts
7. Introduction to Python modeling libraries
8. Data visualization with Tableau
9. Big data analytics with Apache Spark

GRADING POLICY

Homework – 20%

Tests – 20%

Midterm Project – 30%

Final Project – 30%

- Homework are submitted on Blackboard
- Late homework submissions will receive a 40% penalty
- Tests are held during class time. They are in-person tests.

MISC.

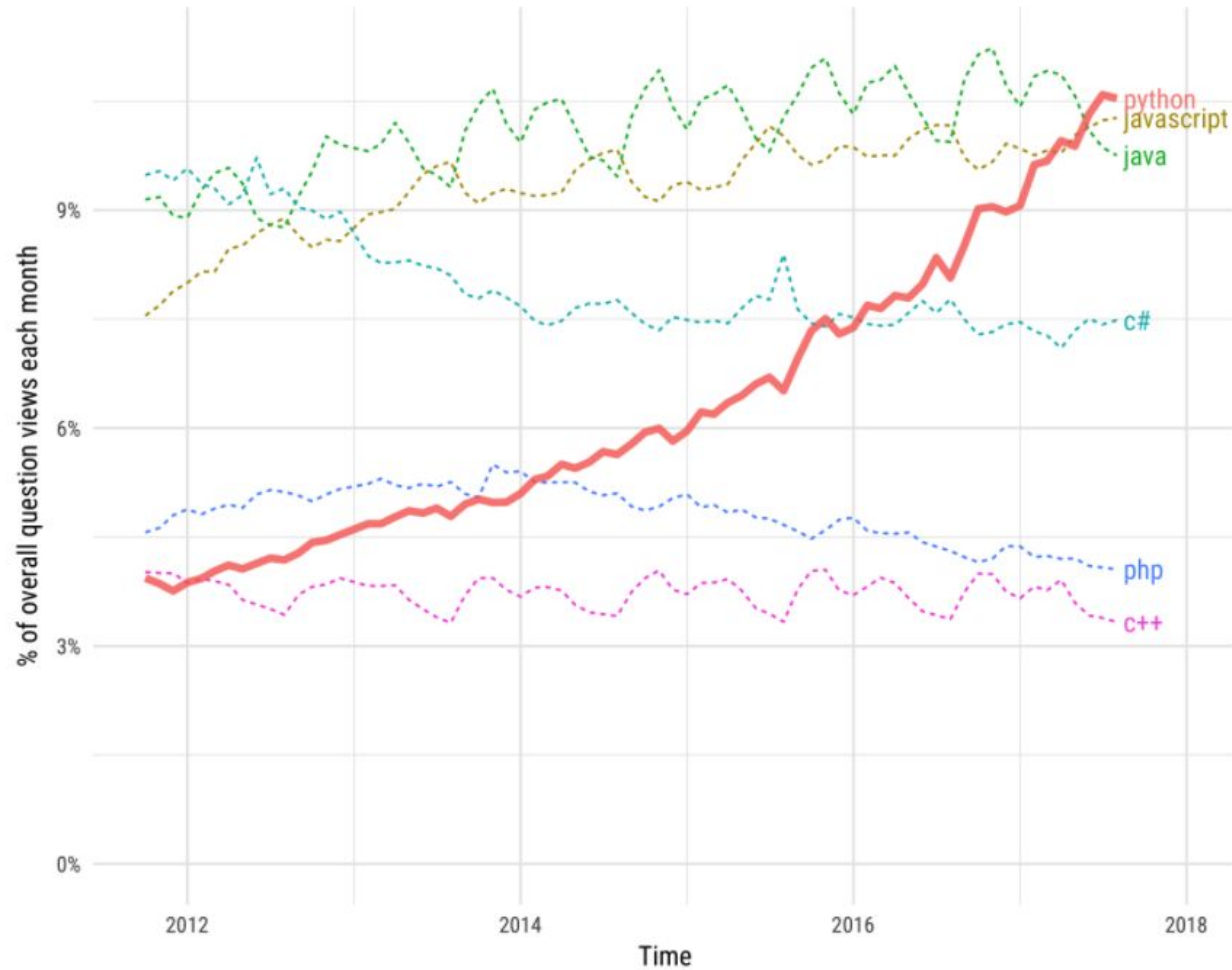
- Expectations
- Honor code
- Email
- Accommodating disabilities
- Consent to Zoom recording
- Additional Resources: Code and Data used in textbook: <https://github.com/wesm/pydata-book>
- Blackboard
- Github: to be announced
- Course videos: to be announced

PYTHON

- Created by Guido van Rossum in 1991
- An interpreted, high-level, general-purpose programming language
- Open source and object oriented
- Has many data science libraries (NumPy, Pandas, Matplotlib, TensorFlow,...)
- Python 2.x (out-dated) vs. Python 3.x
- + Python shell
- + Python script
- + Jupyter notebook (Main programming environment for this class, installed through Anaconda)

Growth of major programming languages

Based on Stack Overflow question views in World Bank high-income countries



Python

Programming language

JavaScript

High-level programmin...

Java

Programming language

C++

Programming language

Go

Programming language

India

Past 12 months

All categories

Web Search

Interest over time

Average

Jan 6, 2019

May 12, 2019

Sep 15, 2019

The chart displays the relative interest in five programming languages over a 12-month period in India. The Y-axis represents interest levels from 0 to 100. Java (yellow) consistently shows the highest interest, peaking near 100. Python (blue) and JavaScript (red) follow, with Python generally higher. C++ (green) shows moderate, stable interest. Go (purple) has the lowest interest, remaining near zero. A bar chart on the left provides a summary of the average interest for each language.

Language	Average Interest
Python	~65
JavaScript	~55
Java	~90
C++	~20
Go	~2

JUPYTER NOTEBOOK

We will use Jupyter Notebook as our coding environment. It supports writing and executing python code in browser.

1. Install the Anaconda distribution of Python and start “Jupyter Notebook”.
2. Create a new notebook named “CMP464Week01”.
3. Type ``print(“Hello World!”)`` in the first cell.
4. Click the run button on the left to see execution results.

MORE ON JUPYTER NOTEBOOK

1. Click “+” button to add a new code cell.
2. Type ``print(1+2)`` and execute.
3. Add a new cell and convert its format from “code” to “markdown”
4. Click the up arrow to move the text cell to the top.
5. Type ``# Week 1 Python Tutorial`` and execute.
6. To print the notebook as a PDF file, click the “Print” option from the browser (not from the menu of jupyter notebook), and then choose “Print to PDF”. **All homework submissions for this class are required in PDF format.**

PYTHON IS SIMPLE

1. Print "Hello World"
2. Calculate the area of a square with length = 1.5
3. Add comments
4. Print sentence "The area is: ***"

VARIABLE AND TYPE

1. Calculate the Body Mass Index (BMI) of a person whose height is 1.73 and weight is 70.5.
 2. Calculate the BMI of a person with height = 1.79 and weight = 68.7
 3. What is the type of the variables used?
- Python types: int, float, bool, str
 - “+” acts differently on numbers and strings
 - Type conversion

METHODS

1. Write a method called “square” that calculates the square of a given number.
2. Write a method called “print_square()” that calculates and prints the square of a given number.

CLASSES

Define a class Book with attributes:

1. Name
2. Author
3. Price

Define the constructor method `__init__()`

Define “`print()`” method that prints the book information

LANGUAGE DESIGN

- Indentation, not braces
- Everything is an object (including functions)
- And, or, not
- Name conventions:
 - Class - CapWords,
 - Variable/method - snake_case
 - Constants - ALL_CAPS
- Name conventions for private and protected methods