# Python for Economics

Zhenhua He
Afternoon session, 9/17/2021

# Table of Contents

This course is divided into numbered lessons

# Python Libraries Covered

**matplotlib** — Plotting data

**pandas** — Analyzing, cleaning, and manipulating data

# Lesson 13
# Data visualization with Matplotlib

Use Python Matplotlib library for data visualization

# Learning Objectives

After this lesson, you will know how to make:

- Scatter plot and Line plot

- Color map

- Contour figures

- 3D figures

  ▪ Surface plots

  ▪ Wire-frame plot

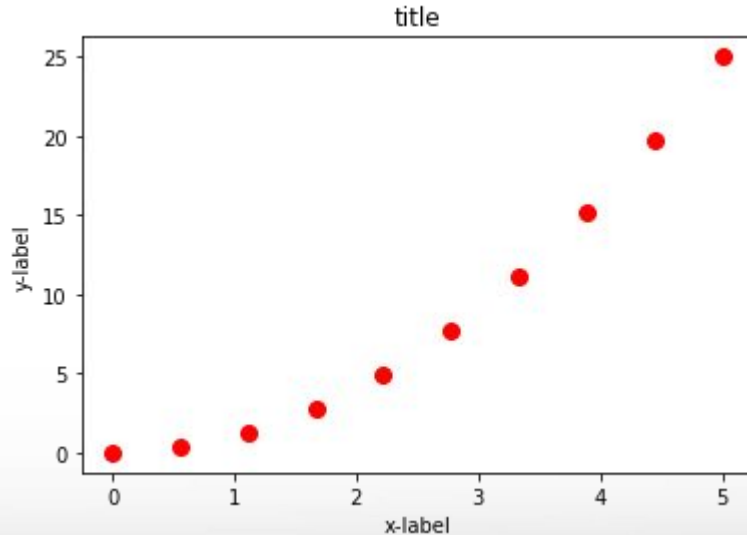  ▪ Contour plots with projections

# Anatomy of a Scatter Plot

**Marker**

- style

- size

- color

**Figure**

- title

- xlabel

- ylabel

# Scatter plot - Marker symbols

| marker | symbol | description |
|--------|--------|-------------|
| "." | • | point |
| "," | · | pixel |
| "o" | ● | circle |
| "v" | ▼ | triangle_down |
| "^" | ▲ | triangle_up |
| "<" | ◄ | triangle_left |
| ">" | ► | triangle_right |
| "1" | Y | tri_down |
| "2" | ⅄ | tri_up |
| "3" | ⊰ | tri_left |
| "4" | ⊱ | tri_right |
| "8" | ● | octagon |
| "s" | ■ | square |
| "p" | ⬟ | pentagon |
| "P" | ✚ | plus (filled) |
| "*" | ★ | star |

# Hot Tip!

Give a module a nickname with `as`

```
import matplotlib.pyplot as plt

import numpy as np

import pandas as pd
```

# Examples and Exercises

Go to Google Classroom assignment "Scatter Plot"

Tasks

- Follow instructions for the examples
- Work on the exercises (**due** by 9/17 6:00 PM )

# Line plot

Simple line styles can be defined using the strings "solid", "dotted", "dashed" or "dashdot".

# Examples and Exercises

Go to Google Classroom assignment "Line Plot"

Tasks

- Follow instructions for the examples
- Work on the exercises (**due** by 9/17 6:00 PM)

# Subplots

# Exercises and Homework

Go to Google Classroom assignment "Subplots"

Tasks

- Follow instructions for the examples
- Work on the exercises (**due** by 9/17 6:00 PM)
- Work on the homework (**due** by 9/23 11:59 PM)

# Color map + savefig()



color map

- pcolor

- imshow

savefig()

- save the current figure

# Examples and Exercises
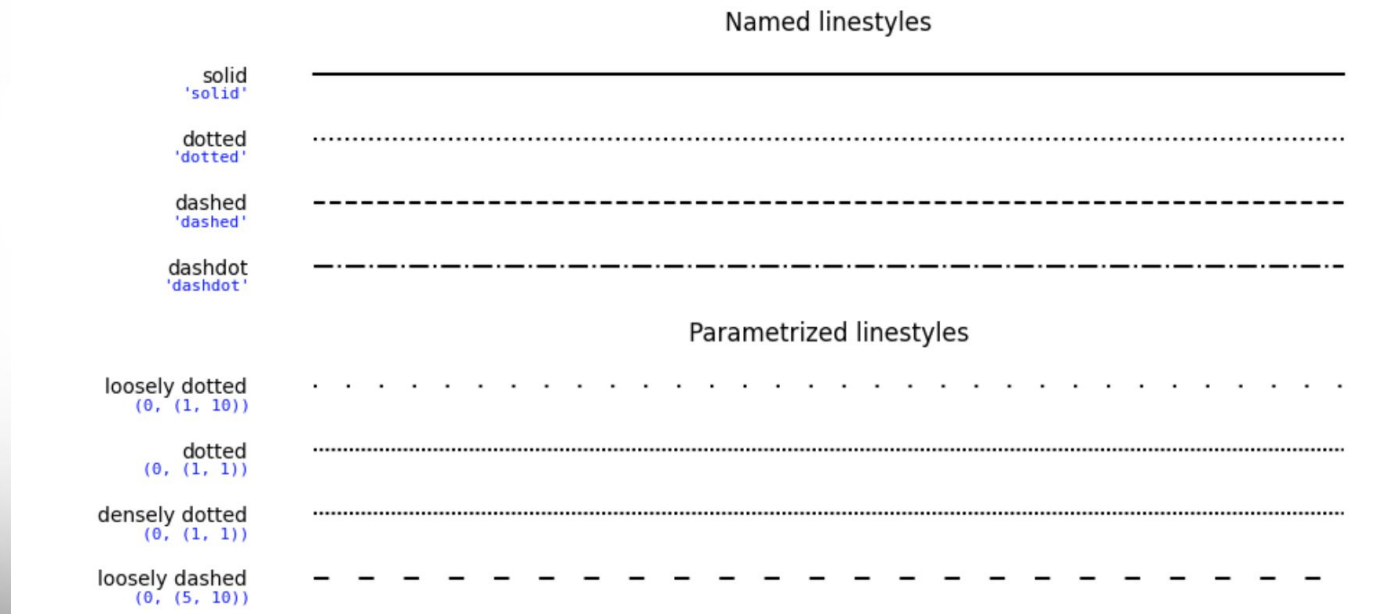
Go to Google Classroom assignment "Color Plot"

Tasks

- Follow instructions for the examples

- Work on the exercises (**due** by 9/17 6:00 PM )

# Break Time Reminder Slide

10 minutes break

# Lesson 14
# Pandas

Use Python Pandas library to manipulate data

# Learning Objectives

After this lesson, you should know how to:

- Create a DataFrame

- Drop Entries

- Index, Select, and Filter data

- Sort data

- Handle missing and duplicate data

- Input and Output

# Pandas VS NumPy

| NumPy | Pandas |
|---|---|
| Faster mathematical operations ✅ | Slower mathematical operations |
| Only supports integer index | Customized index ✅ |
| must use structured arrays | Easily handles different data types ✅ |
| better performance when number of rows is 50K or less | better performance when number of rows is 500K or more ✅ |
| more complicated to read and write files | simpler to read and write more file formats ✅ |

# Series

- One-dimensional labeled array
- Capable of holding any data type (integers, strings, floating point numbers, etc.)
- Example: time-series stock price data



| Index | | Value |
|:---:|:---:|:---:|
| A | → | 0 |
| B | → | 1 |
| C | → | 2 |
| D | → | 3 |
| E | → | 4 |

# Array refresher -> Series

- index
- values
- get a value
- get a set of values
- filtering

# Examples and Exercises

Go to Google Classroom assignment "Series"

Tasks

- Follow instructions for the examples
- Work on the exercises (**due** by 9/17 6:00 PM )

  Create a series -

  > index: datetime;

  > values: randomly generated stock price.

# DataFrame

- Primary Pandas data structure
- A dict-like container for Series objects
- Two-dimensional size-mutable
- Heterogeneous tabular data structure

| | C1 | C2 | C3 | C4 |
|---|---|---|---|---|
| A | 0 | x | 0.1 | True |
| B | 1 | y | 2.4 | False |
| C | 2 | z | 1.9 | True |
| D | NA | w | 8.3 | False |
| E | 9 | a | 6.8 | False |

# DataFrame Example

house sale data

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|
| id | date | price | bedrooms | bathrooms | sqft_living | sqft_lot | floors |
| 7129300520 | 20141013T0( | 221900 | 3 | 1 | 1180 | 5650 | 1 |
| 6414100192 | 20141209T0( | 538000 | 3 | 2.25 | 2570 | 7242 | 2 |
| 5631500400 | 20150225T0( | 180000 | 2 | 1 | 770 | 10000 | 1 |
| 2487200875 | 20141209T0( | 604000 | 4 | 3 | 1960 | 5000 | 1 |
| 1954400510 | 20150218T0( | 510000 | 3 | 2 | 1680 | 8080 | 1 |
| 7237550310 | 20140512T0( | 1.23E+0( | 4 | 4.5 | 5420 | 101930 | 1 |
| 1321400060 | 20140627T0( | 257500 | 3 | 2.25 | 1715 | 6819 | 2 |
| 2008000270 | 20150115T0( | 291850 | 3 | 1.5 | 1060 | 9711 | 1 |
| 2414600126 | 20150415T0( | 229500 | 3 | 1 | 1780 | 7470 | 1 |

# Creating a Data Frame

Ways to do so:

- from Dictionary

- from Numpy array

- Read file (read_csv, read_excel, read_stata, read_html, ...)

# Dictionary

For example, you have a car and its information is as below,

*   brand: Ford

*   model: Mustang

*   year: 1964

You can create a dictionary as below

```python
car_dict = {
 "brand": "Ford",
 "model": "Mustang",
 "year": 1964,
}
```

# Examples and Exercises

Go to Google Classroom assignment Pandas "DataFrame-1"

Tasks

- Follow instructions for the examples

- Work on the exercises (**due** by 9/17 6:00 PM )
  1. Create a nation_economics DataFrame - including columns of Country, Continent, GDP, Population, GDPPerCapita
  2. Data on the next slide

# Examples and Exercises

nation_economics data

| Country | Continent | GDP (Billion dollars) | Population (Millions) | GDPPerCapita (Thousand dollars) |
|---|---|---|---|---|
| United States | America | 18624.5 | 332.9 | 66.7 |
| China | Asia | 11218.3 | 1444.2 | 10.7 |
| Japan | Asia | 4936.2 | 126.1 | 43.6 |
| Germany | Europe | 3477.8 | 83.9 | 49.5 |
| India | Asia | 2259.6 | 1393.4 | 2.3 |
| United Kingdom | Europe | 2647.9 | 68.2 | 42.9 |
| France | Europe | 2465.5 | 65.4 | 44.0 |
| Italy | Europe | 1858.9 | 60.4 | 34.6 |
| Brazil | America | 1795.9 | 214.0 | 9.6 |
| Canada | America | 1529.8 | 38.1 | 48.1 |

# Break Time Reminder Slide

10 minutes break

# DataFrame: data retrieval

- Retrieve a column

- Retrieve multiple columns

- Retrieve a row

- Retrieve multiple rows

- Drop entries

# Examples and Exercises

Go to Google Classroom assignment Pandas "DataFrame-2"

Tasks

- Follow instructions for the examples

- Work on the exercises (**due** by 9/17 6:00 PM )

    From the nation_economics DataFrame,
    1. Retrieve the *GDPPerCapita* column
    2. Retrieve the *United Kingdom* row
    3. Drop the *Population* column
    4. Drop the *Canada* row

# DataFrame: operations/manipulation

- Selecting with slicing
- Filtering
- Sorting
    - sort by index
    - sort by values

# Examples and Exercises

Go to Google Classroom assignment Pandas "DataFrame-3"

Tasks

- Follow instructions for the examples

- Work on the exercises (**due** by 9/17 6:00 PM )

  From the national_economics DataFrame
  1. Select the last 5 rows
  2. Select the rows with the population greater than 100M
  3. Sort the DataFrame by GDPPerCapita in descending order

- Work on the homework (**due** by 9/23 11:59 PM)

# DataFrame: input and output

- Read/Write

- Different file formats

- describe()

# Capstone - Candlestick Chart

A financial chart to depict price movement.

Four data values per marker:

- High

- Low

- Open

- Close



Daily Candlestick Chart of NIFTY50

# Exercise and Homework

Go to Google Classroom assignment "Matplotlib-Candlestick"

Tasks

- Follow instructions for the examples
- Work on the exercises (**due** by 9/17 6:00 PM )
- Work on the homework (**due** by 9/23 11:59 PM)

# Day 2 wrap-up

almost time to go home

# Practice for next week

Most important skills to master

- List loops

- Filtering with conditionals

- Pandas DataFrame structure

Slides from today are available in Google Classroom

# Homework Assignments

- Lesson 9: "Lists and Strings"

- Lesson 10: "National Economic Data"

- Lesson 11: "Talking Cats"

- Lesson 12: "Array Quiz"

- Lesson 13: "Matplotlib - Subplots"

- Lesson 13: "Matplotlib - Candlestick chart"

- Lesson 14: "Pandas - DataFrame operations"

Please submit your homework assignments before 9/23 11:59 PM
Turn in your in-class exercises before 6:00 PM today

# Office Hours

Please come to our office hours for assistance
- `M` 10 - 11 am Blocker 219B
- `T` 10 - 11 am (on Zoom only)
- `W` 2 - 4:30 pm Blocker 219B
- `R` 2 - 3 pm Blocker 219B

Please join our slack channel for discussion
- Workspace `sweeterworkspace.slack.com`
- Channel `hprc-econ-fall-21` (private channel)

# New HPRC Help Resource

Bring Your Own Code (BYOC) sessions

These sessions are meant to help researchers overcome general Python programming hurdles in their research projects.

In person (Rooms 218A and 217B) or via zoom
Weekly on Wednesdays from 3-4:30pm through December 15.

Contact help@hprc.tamu.edu

# Pandas Cheat Sheet (continued learning)



https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf