**NIDA**
WISDOM *for* Change

DADS 5001

# DATA ANALYTICS AND DATA SCIENCE TOOLS AND PROGRAMMING

INTRODUCTION &
PANDAS 1

Asst.Prof.Thitirat Siriborvornratanakul, Ph.D. (thitirat@as.nida.ac.th)
Graduate School of Applied Statistics, National Institute of Development Administration

1

## OUTLINE

- × Course's Introduction
- × Mini-project
- × Package installation

- × Exploratory Data Analysis
- × List vs. Pandas vs. NumPy

- × Pandas

2

2

# COURSE'S INTRODUCTION

3

3
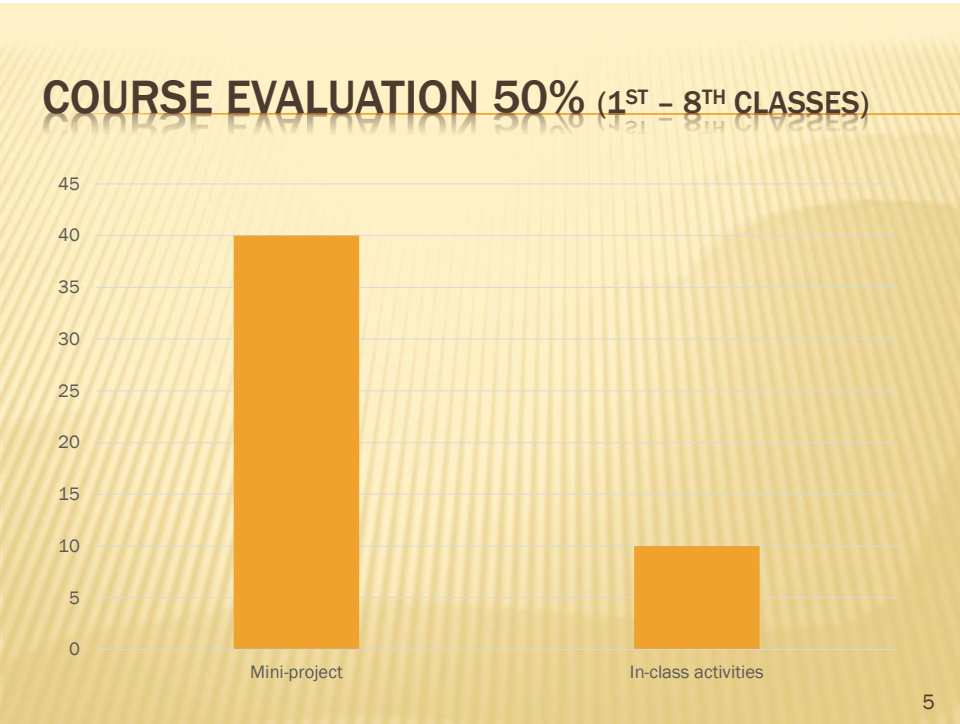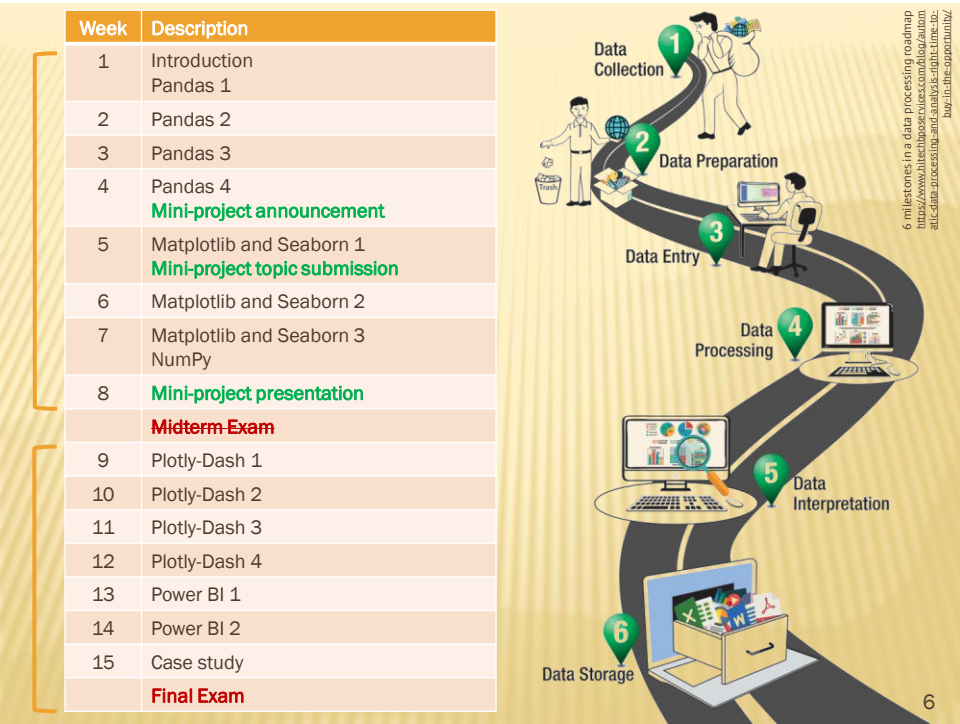
# MICROSOFT TEAMS

× (1/2567) DADS5001

+ **Official announcement:** General channel

+ **Class material (all in English):** General channel > Files tab > เอกสารประกอบของคลาส (Class Materials) folder

+ **Syllabus:** General channel > Files tab > เอกสารประกอบของคลาส (Class Materials) folder > COURSE_SYLLABUS.pdf

+ **Latest class schedule:** General channel > Files tab > เอกสารประกอบของคลาส (Class Materials) folder > CLASS_SCHEDULE.docx
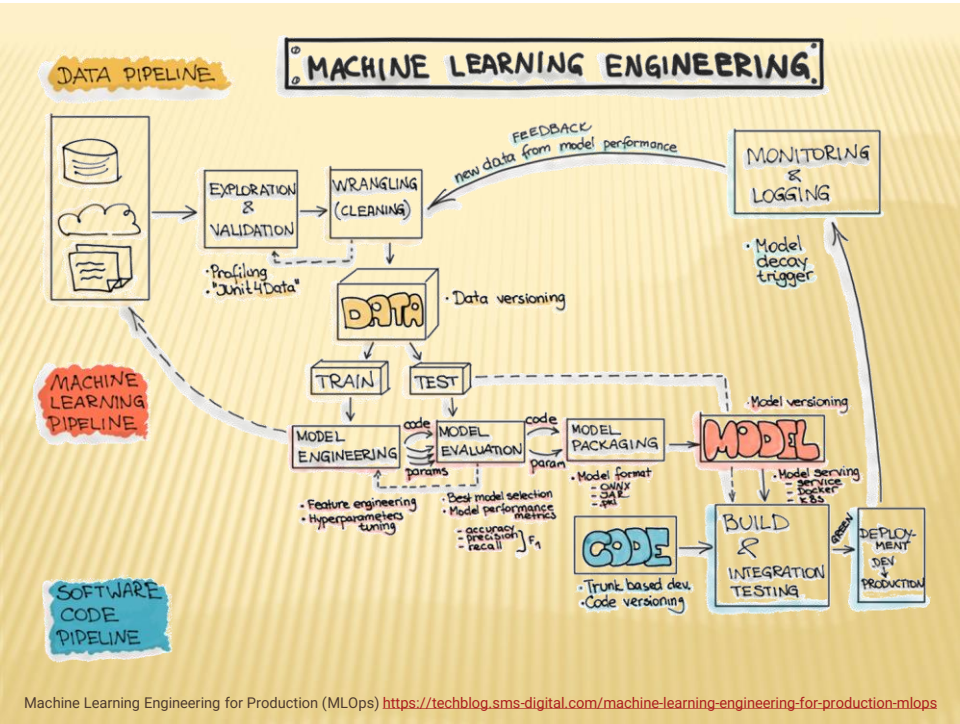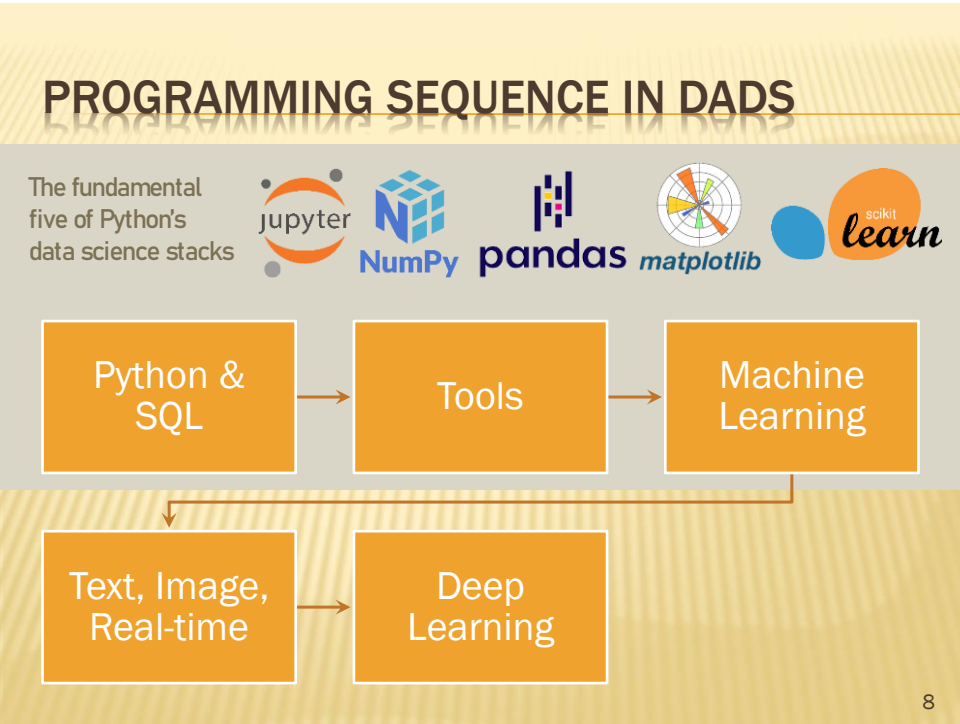
4

4

## COURSE EVALUATION 50% (1ST – 8TH CLASSES)

A bar chart showing two categories:
- Mini-project: 40
- In-class activities: 10

Y-axis scale: 0, 5, 10, 15, 20, 25, 30, 35, 40, 45

5

5

| Week | Description |
|------|-------------|
| 1 | Introduction<br>Pandas 1 |
| 2 | Pandas 2 |
| 3 | Pandas 3 |
| 4 | Pandas 4<br>**Mini-project announcement** |
| 5 | Matplotlib and Seaborn 1<br>**Mini-project topic submission** |
| 6 | Matplotlib and Seaborn 2 |
| 7 | Matplotlib and Seaborn 3<br>NumPy |
| 8 | **Mini-project presentation** |
|  | **Midterm Exam** |
| 9 | Plotly-Dash 1 |
| 10 | Plotly-Dash 2 |
| 11 | Plotly-Dash 3 |
| 12 | Plotly-Dash 4 |
| 13 | Power BI 1 |
| 14 | Power BI 2 |
| 15 | Case study |
|  | **Final Exam** |

Data Collection 1
Data Preparation 2
Data Entry 3
Data Processing 4
Data Interpretation 5
Data Storage 6

6 milestones in a data processing roadmap
https://www.hitechbposervices.com/blog/autom
atic-data-processing-and-analysis-right-time-to-
buy-in-the-opportunity/

6

6

Machine Learning Engineering for Production (MLOps) https://techblog.sms-digital.com/machine-learning-engineering-for-production-mlops

7

# PROGRAMMING SEQUENCE IN DADS



The fundamental five of Python's data science stacks

8

8

9



10

## MINI-PROJECT: INTRO (2/3)

"Most important, to squeeze insights out of Big Data, you have to ask the right questions."



11

## MINI-PROJECT: INTRO (3/3)



https://www.facebook.com/100057174082208/posts/pfbid022YgWGJAuUTXkPXN9gho825NkL9NVc7qNnd5dhTjZg8QnBiT6ujy1K7xzNZ82bEHnI/?mibextid=cr9u03 , https://www.facebook.com/d4biz/photos/a.3233628596865997/3229272480634942

12

# MINI-PROJECT: HOW TO

1. Having a good question is a good start, so as an initial assumption.
2. Check data availability/accessibility (realistic and unprocessed data, not the already-prepped neat dataset with abundant tutorials/examples)
3. Prove or disprove your assumption with data analytics and data science
   + Programming tools and (static) visualization must be included significantly.
   + The data source may be shared among some students but what comes afterward (questioning, cleaning, processing, EDA, visualizing) must be done individually.
4. Submit as a public GitHub link (one link per one group), including codes and development journey



13

# MINI-PROJECT: EXAMPLES (1/2)

* ส่อง Mega trend "EV Car" https://github.com/ssorawits/Project.git
* ผลกระทบจากเศรษฐกิจทำให้อัตราเด็กเกิดใหม่ลดลงหรือไม่ https://github.com/pnithida/6420422004_MiniProject
* PM2.5 in Thailand(Y2019-Y2021) https://github.com/Porrakij/PM2.5-Data-Analysis
* สำรวจกระเป๋าตังคนไทย รายได้ และ รายจ่าย เฉลี่ยของครัวเรือน https://github.com/koraweep/DADS5001_Mini-Project
* จังหวัดเฝ้าระวัง เพื่อป้องกันและลดจำนวนผู้เสียชีวิตจากอุบัติเหตุทางถนน https://github.com/ploychenya/DADS5001_mini_project
* Big Data กับโลกอสังหาริมทรัพย์ https://github.com/Hakulani/miniprojectDADS5001
* ข้อมูลมูลค่าการนำเข้าสินค้าของประเทศไทยจากทั่วโลกในช่วงปี 2002 – 2022 https://github.com/nacknatthawit/DADS5001_6420412006
* Thailand_Accident_Jan-Jun2022 https://github.com/waewma/Thailand_Accident_Jan-Jun2022
* ยอดการค้นหา Google trend ด้วย keyword "Bitkub" จะส่งผลดีต่อบริการของตัวบริษัทหรือไม่ https://github.com/kimteespk/DADS5001_Miniproject_bitkub_googletrend_6510412011
* Vending machine analysis https://github.com/kikkalo/6420412010
* Thailand Labour Demand Trend Y2015-2020 https://github.com/Sujitra17/Mini-Project

14

14

# MINI-PROJECT: EXAMPLES (2/2)

✖ Are Thai healthcare related stock still captivating? https://github.com/mmaiip/project-5001

✖ Global warming https://github.com/Hellper1/DADS_5001_miniproject

✖ Video Game Sales https://github.com/MeenWhile/Tools-Mini-Project

✖ World-Population Analysis https://github.com/o-joe-v/World-Population

✖ Thailand and World Happiness
https://github.com/HikariJadeEmpire/TH_WLRD_Happiness_Project

✖ สถานการณ์อุตสาหกรรมรถยนต์ไทย https://github.com/crispyporkwithholybasil/5001-DADS-miniproject

✖ "Global Cost Of Living" ถ้าต้องย้ายประเทศเราจะไปไหนดีล่ะ?
https://github.com/Bonita1996/Project

✖ Hollywood Insight https://github.com/y-lims/DADS5001_Hollywood_Insight

✖ ค่าไฟแพงมาจากอะไร? สถานการณ์พลังงานไฟฟ้าในประเทศไทย
https://github.com/kkengg/Pikachu-Project

15

15

# PACKAGE INSTALLATION

16

16

## PIP INSTALL VS. CONDA INSTALL



17

## REQUIRED LIBRARIES

```
conda install jupyter

conda install pandas
conda install numpy

conda install matplotlib
conda install seaborn
```

18

# EXPLORATORY DATA ANALYSIS

19

19

# EXPLORATORY DATA ANALYSIS (1/4)

× Originally developed by American mathematician John Tukey in the 1970s, EDA techniques continue to be a widely used method in the data discovery process today.
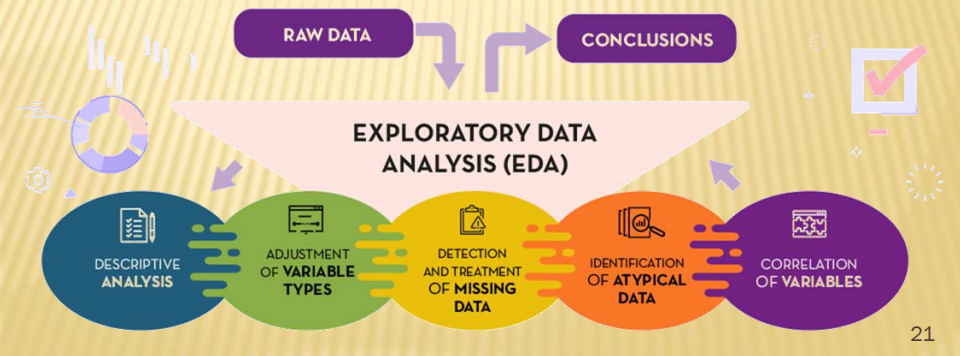


20

20

# EXPLORATORY DATA ANALYSIS (2/4)

× EDA is used by data analysts/scientists to <u>analyze</u> and <u>investigate</u> data sets and <u>summarize</u> their main characteristics, often employing <u>data visualization methods</u>.

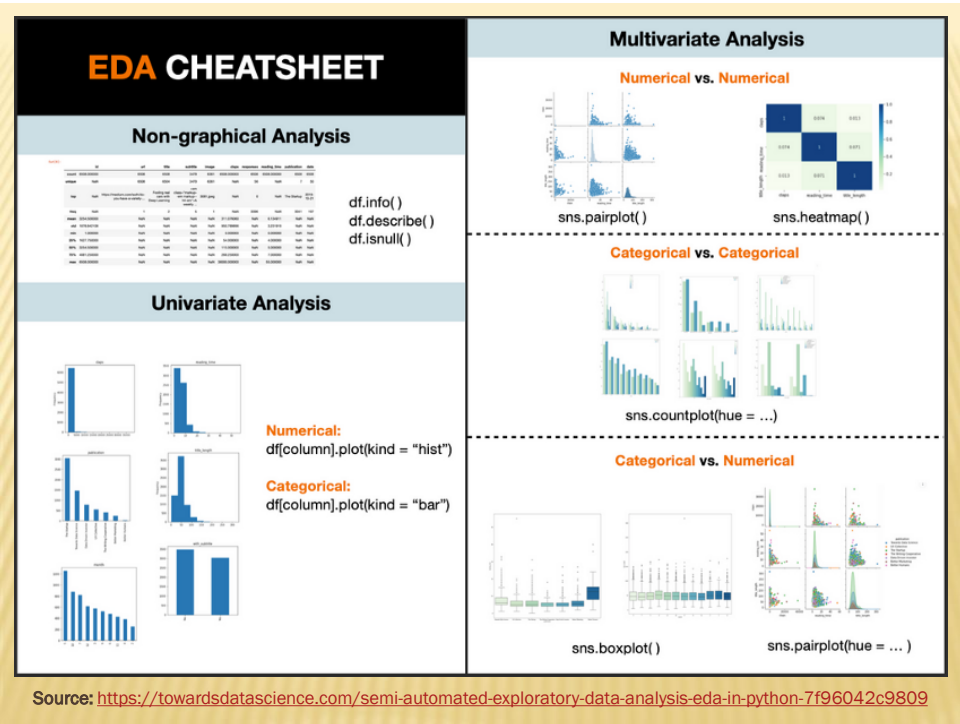× The main purpose of EDA is to help look at data before making any assumptions.



21

21

# EXPLORATORY DATA ANALYSIS (3/4)
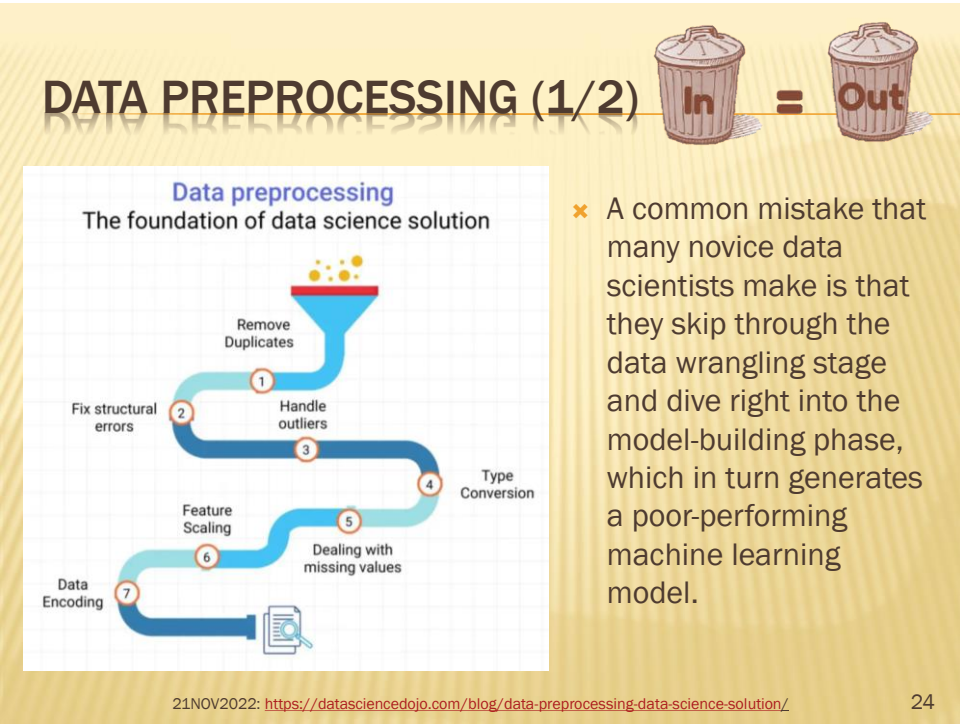
× EDA is an iterative process of



DATA TRANSFORMATION    DATA ANALYSIS    DATA VISUALIZATION

× What for?
+ Provide a better understanding of data set patterns, variables, and the relationships between them
+ Help identify obvious errors and detect outliers or anomalous events
+ See what data can reveal beyond the formal modeling or hypothesis testing task
+ Help determine if the statistical techniques we are considering for data analysis are appropriate
+ Ensure the results we produce are valid and applicable to any desired business outcomes and goals

× Once EDA is complete and insights are drawn, its features can then be used for more sophisticated data analysis or modeling, including machine learning and deep learning.

22

22

Source: https://towardsdatascience.com/semi-automated-exploratory-data-analysis-eda-in-python-7f96042c9809

23

# DATA PREPROCESSING (1/2)



Data preprocessing
The foundation of data science solution

× A common mistake that many novice data scientists make is that they skip through the data wrangling stage and dive right into the model-building phase, which in turn generates a poor-performing machine learning model.

21NOV2022: https://datasciencedojo.com/blog/data-preprocessing-data-science-solution/    24

24

## DATA PREPROCESSING (2/2)

1. Remove duplicates
   + Repeated entries will lead to the model overfitting.
   + Be careful when there are too many duplicate entries.
2. Fix structural errors
   + Structural errors in a dataset refer to the entries that either have typos or inconsistent spellings.
   + Check all unique values and their corresponding occurrence
3. Detect and handle outliers
   + Outlier is any value in a dataset that drastically deviates from the rest of the data points. It can mess up our machine-learning model if not taken care of.
   + Use the describe function on columns, visualize outliers with box plots, compute z-score on columns (99.7% of the data points within the range of -3 and +3 scores), etc.
4. Type conversion
   + Do type conversion when certain columns are not of valid data type (e.g., the object data type)
5. Dealing with missing values
   + Often, data set contains numerous missing values, which can be a problem. For example, it can play a role in the development of a biased estimator, or it can decrease the representativeness of the sample under consideration.
   + Drop rows with missing values, impute the missing values with central tendencies (e.g., mean, median, mode), forward/backward fill
6. Feature scaling
7. Data encoding

Credit (21NOV2022): https://datasciencedojo.com/blog/data-preprocessing-data-science-solution/

25

25

# PYTHON'S LIST VS. PANDAS VS. NUMPY

26

26

# PYTHON'S LIST

× https://wiki.python.org/moin/TimeComplexity

## list

The Average Case assumes parameters generated uniformly at random.

Internally, a list is represented as an array; the largest costs come from growing beyond the current allocation size (because everything must move), or from inserting or deleting somewhere near the beginning (because everything after that must move). If you need to add/remove at both ends, consider using a collections.deque instead.
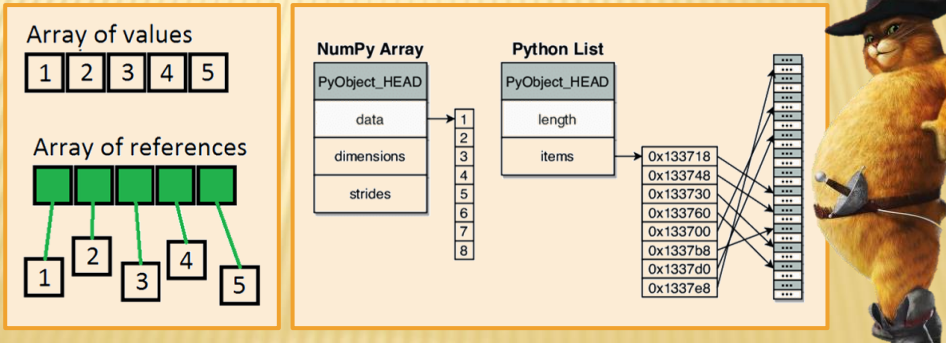
27

27

# ARRAY VS. PYTHON'S LIST

× In other programming languages, Array refers to the static memory allocation for holding a sequence of homogeneous data. Array is faster.

× Because Python's List stores references to data, a List can (seemingly) hold heterogeneous data. List is slower.

Array of values
1 2 3 4 5

Array of references

NumPy Array
PyObject_HEAD
data
dimensions
strides

Python List
PyObject_HEAD
length
items

0x133718
0x133748
0x133730
0x133760
0x133700
0x1337b8
0x1337d0
0x1337e8

28

29



30

## PANDAS (2008) VS. EXCEL (1985)

× Too much data slows down Excel.

× Excel is limited to 1,048,576 rows by 16,384 columns per worksheet. Exceeding the limit causes data loss.

× Pandas is one of the most popularly used Python modules for data analysis and manipulation. It's said to be a data analyst's best friend.

× No limits to the size of data, powerful data transformation, allow automation and extended file formats

× Pandas is good and intuitive for tabular data (e.g., csv, excel, sql) and it includes time series functionalities.

31

31

## PANDAS (2008) VS. NUMPY (2005)

× Pandas is built on top of two core Python libraries—matplotlib for data visualization and NumPy for mathematical operations.

× Pandas acts as a wrapper over these two libraries, allowing us to access many of matplotlib's and NumPy's methods with less code.

× Pandas for data analysts. NumPy for data scientists

× Pandas is popularly used for data analysis and manipulation. NumPy is primarily used for numerical calculations.

× Pandas is more natural for database-like data (e.g., csv, excel, sql). NumPy is more natural for numeric processing of data (e.g., signals, images).

× Scikit-learn was originally developed to work well with NumPy array. It's recommended to use NumPy array with Scikit-learn due to mature data handling.
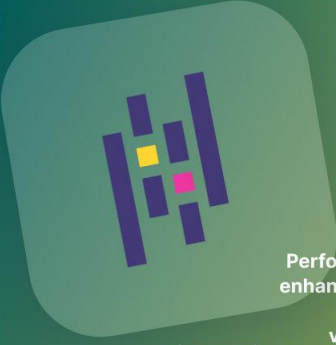
32

32

33



34

## PANDAS 2.0 (3/4)
What's new in 2.0.0 (April 3, 2023): https://pandas.pydata.org/docs/whatsnew/2.0.0.html

- × A new `dtype_backend` parameter for I/O operators that support creating Arrow-backed DataFrames.

```
pd.options.mode.dtype_backend = 'pyarrow'      # Set it globally
pd.read_csv(my_file, dtype_backend='pyarrow') # Set it on each DataFrame
```

- × Represent "missing values" and use better support for data types outside of numerical types
  - + NumPy has poor support for non-numerical data and a lack of missing values.
  - + Before pandas 2.0, there are different types of missing values; `np.nan` is for floating-point numbers; `None` and `np.nan` are for object types, and `pd.NaT` is for date-related types.
  - + In Pandas 1.0, `pd.NA` was introduced to avoid type conversion, but it needs to be specified manually by the user.
  - + The new `string[pyarrow]` column type is around 3.5 times more efficient.

```
pd.Series([1,2,3,4], dtype='int64[pyarrow]')
pd.Series(['foo', 'bar', 'foobar'], dtype='string[pyarrow]')
```

35

35

## PANDAS 2.0 (4/4)
What's new in 2.0.0 (April 3, 2023): https://pandas.pydata.org/docs/whatsnew/2.0.0.html

```
In [1]: df2 = pd.DataFrame({'a':[1,2,3, None]}, dtype='int64[pyarrow]')

In [2]: df2.dtypes
Out[2]:
a    int64[pyarrow]
dtype: object

In [3]: df2
Out[3]:
     a
0    1
1    2
2    3
3  <NA>
```

```
In [1]: import pandas as pd

In [2]: pd.__version__
Out[2]: '2.0.0'

In [3]: df = pd.read_csv('pd_test.csv')

In [4]: df.dtypes
Out[4]:
name       object
address    object
number      int64
dtype: object

In [5]: df.memory_usage(deep=True).sum()
Out[5]: 17898876

In [6]: df_arrow = pd.read_csv('pd_test.csv', dtype_backend="pyarrow", engine="pyarrow")

In [7]: df_arrow.dtypes
Out[7]:
name       string[pyarrow]
address    string[pyarrow]
number      int64[pyarrow]
dtype: object

In [8]: df_arrow.memory_usage(deep=True).sum()
Out[8]: 7298876
```
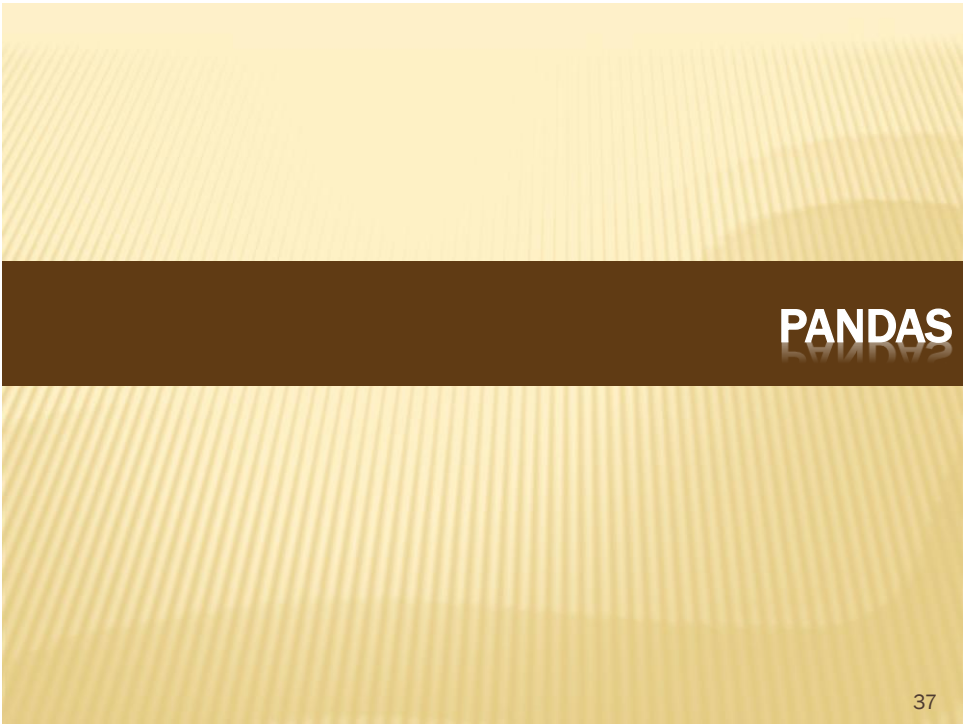
Credit:
https://www.reddit.com/r/Python/comments/12b7w3y/everything_you_need_to_know_about_pandas_200/

36

## PANDAS

37

37

## SERIES AND DATAFRAME

× `pandas.Series` is a 1D array holding data of any type.

× `pandas.DataFrame` is a 2D array consisting of one or more `pandas.Series` (= column).

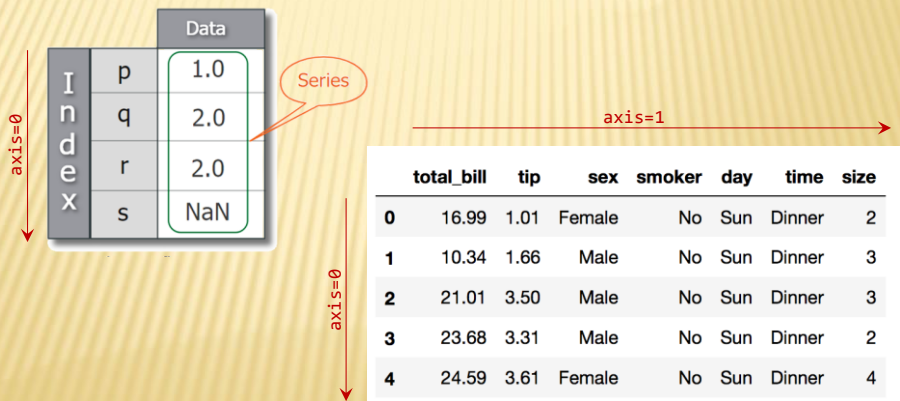| Series 1 | | | Series 2 | | | Series 3 | | | DataFrame | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Mango** | | | **Apple** | | | **Banana** | | | **Mango** **Apple** **Banana** | | |
| 0 | 4 | | 0 | 5 | | 0 | 2 | | 0 | 4 | 5 | 2 |
| 1 | 5 | + | 1 | 4 | + | 1 | 3 | = | 1 | 5 | 4 | 3 |
| 2 | 6 | | 2 | 3 | | 2 | 5 | | 2 | 6 | 3 | 5 |
| 3 | 3 | | 3 | 0 | | 3 | 2 | | 3 | 3 | 0 | 2 |
| 4 | 1 | | 4 | 2 | | 4 | 7 | | 4 | 1 | 2 | 7 |

38

# SERIES AND DATAFRAME



39

# AXIS

× `pandas.Series` is 1D so there is only one axis.
× `pandas.DataFrame` is 2D so there are two axes.



40

# LET'S CONTINUE IN colab

- Pandas1.ipynb
  - IPython.display: `display()`, `Markdown()`, `Latex()`, `Code()`, `HTML()`, `JSON()`, `clear_output()`, …
  - Data types: `pandas.Series`, `pandas.DataFrame`
  - Load, save, and render: `read_csv()`, `to_csv()`, `read_excel()`, `to_excel()`, `to_html()`, `to_latex()`, `to_json()`, …
  - Inspect, (unconditional) access, change, sort, and save data: `loc[]`, `iloc[]`, `at[]`, `iat[]`, `sort_index()`, `sort_values()`, …

- Pandas2.ipynb
  - Conditionally filter data with boolean indexing, `query()`, `filter()`, `select_dtypes()`
  - Handle missing data and duplicated data: `isna()`, `notna()`, `dropna()`, `fillna()`, `duplicated()`, `drop_duplicates()`

- Pandas3.ipynb
  - Aggregate data: `agg()`, `sum()`, `mean()`, `min()`, `max()`, `median()`, `mode()`, `std()`, `unique()`, `nunique()`, `count()`, `value_count()`, …
  - Transform data: `apply()`, `transform()`
  - Group data: `groupby()`, `filter()`

- Pandas4.ipynb
  - Combine data: `merge()`, `join()`, `compare()` , `concat()`
  - Reshape data: `pivot()`, `pivot_table()`, `melt()`
  - Table visualization: `pandas.DataFrame.style`, `style.format()`, `style.applymap()`, `style.apply()`, `style.bar()`, `style.highlight_max()`, `style.highlight_min()`, …

41

41

# END OF THIS CLASS

- News and announcement:
  - Microsoft Teams > General channel

- Class schedule:
  - Microsoft Teams > General channel > Files tab > เอกสารประกอบของคลาส folder > CLASS_SCHEDULE.docx

42

42