

MODULE 03

DATA SCIENCE BOOTCAMP

# Numpy and Pandas

# Reminder

1

Please turn on Zoom camera the whole duration of classes.

2

At the start of all classes, please rename yourselves to: Name + Last 3 digits and letter of your NRIC. Example: John Tan (123A)

# Agenda

- Python Modules
- Numpy
- Pandas
- Types of Joins





# Attendance Photo Taking

MODULE 3: NUMPY AND PANDAS

# NumPy



# Python Modules

A python module allows you to logically organize your Python Code.

Module contents available to the caller with the import statement.

Modules include:

- Numpy
- SciPy
- Matplotlib
- Pandas

Python

```
import <module_name>
```



# NumPy

NumPy or Numerical Python is a general-purpose array processing python package for scientific computing. It consists of numerous powerful features inclusive of:

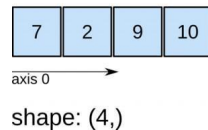
- A robust multi-dimension array object with many useful functions
- An enormous number of routines, including shape manipulation, logical, mathematical &
- many more to operate on [NumPy Array objects](#)
- NumPy is also utilized as a generic multi-dimension data container
- A wide set of databases can also be integrated with NumPy



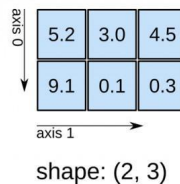
# Numpy Array

- A numpy array is a grid of values, all of the same type, and is indexed by a tuple of nonnegative integers.
- Basically, it is a multidimensional or n-dimensional array of fixed size with homogeneous elements( i.e., the data type of all the elements in the array is the same)

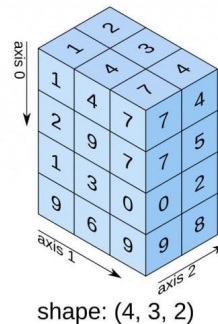
1D array



2D array



3D array





# Difference Between List And Numpy Array

- A user can treat [lists](#) as arrays.
- However, user cannot constraint the type of elements stored in a list.
- If you create arrays using the array module, all elements of the array must be of the same type



# Benefits of Numpy Array

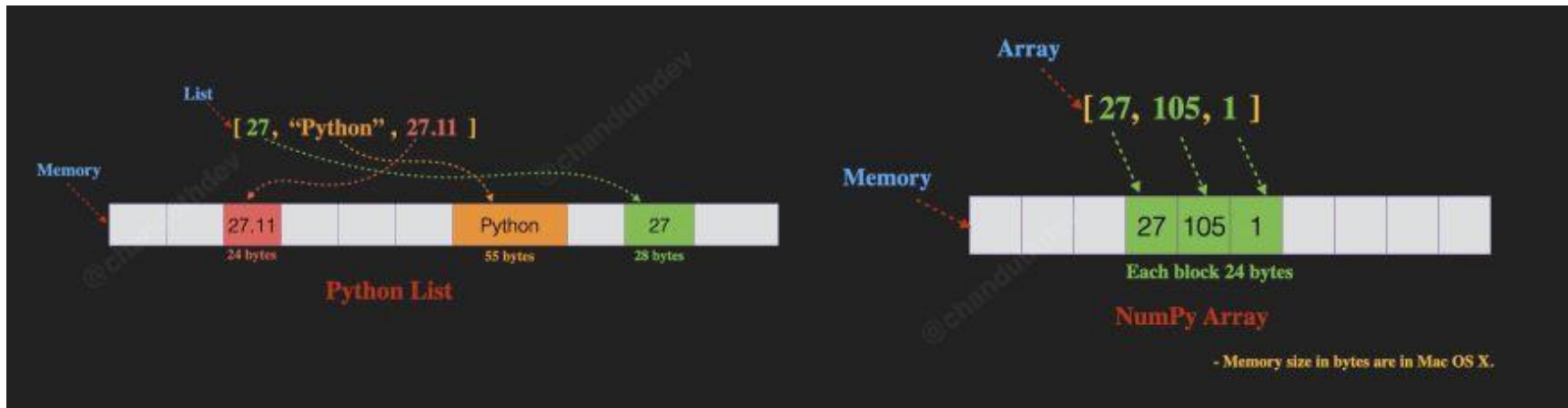
Numpy data structures perform better in:

- **Size** - Numpy data structures take up less space
- **Performance** - they have a need for speed and are faster than lists
- **Functionality** - SciPy and NumPy have optimized functions such as linear algebra operations built in.



# Difference Between Numpy Array & Python Lists

Numpy array tend to be more efficient in memory and storage.



MODULE 3: NUMPY AND PANDAS

# Pandas



# What is Pandas?

- Pandas is an open source Python package that is most widely used for data science/ data analysis and machine learning tasks.
- It offers powerful, expressive and flexible data structures that make data manipulation and analysis easy, among many other things. The DataFrame is one of these structures.



# DataFrame

DataFrame is a 2-dimensional labeled data structure with columns of potentially different types.

Columns				
	Name	Score	Attempts	Qualify
0	Anastasia	12.5	1	yes
1	Dima	9.0	3	no
2	Katherine	16.5	2	yes
3	James	NaN	3	no
4	Emily	9.0	2	no

Rows

Data

Pandas DataFrame

© w3resource.com



# What You Can Do With DataFrames Using Pandas

Pandas makes it simple to do many of the time consuming, repetitive tasks associated with working with data, including:

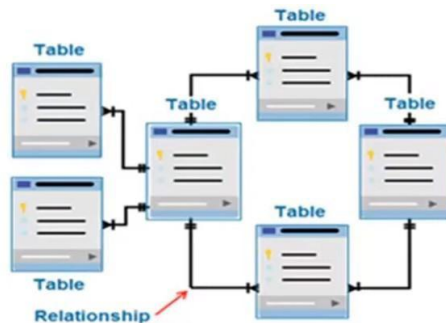
- Data cleansing
- Data fill
- Data normalization
- Merges and joins
- Data visualization
- Statistical analysis
- Data inspection
- Loading and saving data
- And much more



# Real World Use Case With Pandas

In general, the data in financial institutions are stored in relational data format. Most of the raw data are not in their best quality and are stored in different databases.

Hence, we utilized the pandas functions to efficiently clean and join data among the different data tables so as to perform data analytics of quality. Python programs are often utilized pandas to automate the data cleaning process for their data analytics use.

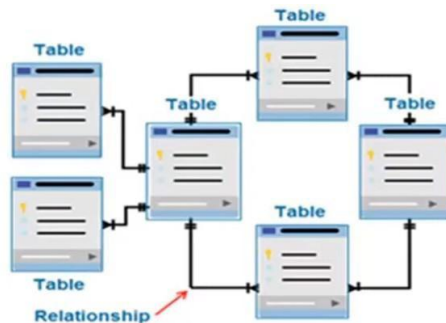




# Real World Use Case With Pandas

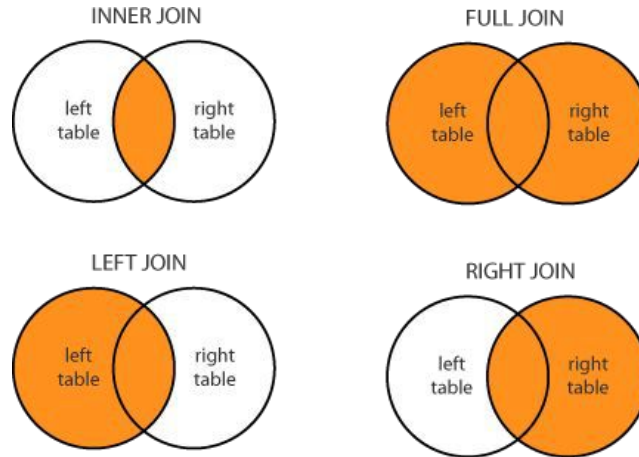
In general, the data in financial institutions are stored in relational data format. Most of the raw data are not in their best quality and are stored in different databases.

Hence, we utilized the pandas functions to efficiently clean and join data among the different data tables so as to perform data analytics of quality. Python programs are often utilized pandas to automate the data cleaning process for their data analytics use.



# Type of Joins

A Join clause is used to combine rows from two or more tables, based on a related column between them.



# Type of Joins

Join Clause	Description
(INNER) JOIN	Returns records that have matching values in both tables
LEFT (OUTER) JOIN	Returns all records from the left table, and the matched records from the right table
RIGHT (OUTER) JOIN	Returns all records from the right table, and the matched records from the left table
FULL (OUTER) JOIN	Returns all records when there is a match in either left or right table



# Pd.merge()

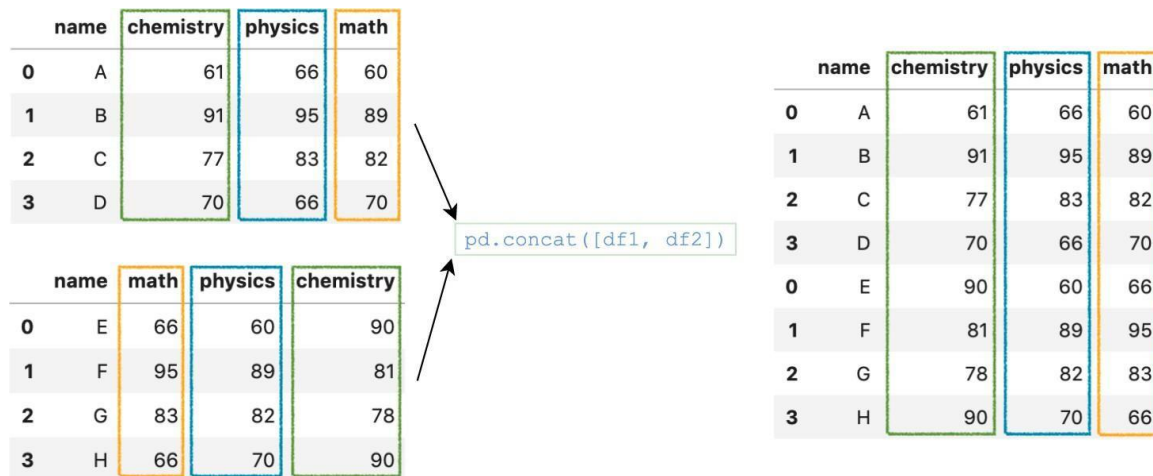
When you want to combine data objects based on one or more keys in a similar way to a relational database, merge() is the tool you need.

left					right					Result						
	key1	key2	A	B		key1	key2	C	D		key1	key2	A	B	C	D
0	K0	K0	A0	B0	0	K0	K0	C0	D0	0	K0	K0	A0	B0	C0	D0
1	K0	K1	A1	B1	1	K1	K0	C1	D1	1	K0	K1	A1	B1	NaN	NaN
2	K1	K0	A2	B2	2	K1	K0	C2	D2	2	K1	K0	A2	B2	C1	D1
3	K2	K1	A3	B3	3	K2	K0	C3	D3	3	K1	K0	A2	B2	C2	D2
										4	K2	K1	A3	B3	NaN	NaN
										5	K2	K0	NaN	NaN	C3	D3



# Pd.concat()

The **concat** function does all of the heavy lifting of performing concatenation operations along an axis. It is similar to appending the data rows of different dataset.



# Pandas Operations

In the jupyter notebook, we will cover a lot of the pandas operations that are essential to the data analytics process such as:

- Slice and dice of datasets by different selection methods
- Assignment of new values to dataframe
- Filtering with multiple conditions
- Mathematical operations
- Join and appending datasets



# Recap time!

**What are your favorite  
takeaways on NumPy &  
Pandas today?**

*Let's share with each other!*

# Some things to take note...

**Link and resource could be accessed in the Learning Portal.**

[https://elearn.verticalinstitute.com/users/sign\\_in](https://elearn.verticalinstitute.com/users/sign_in)







# Attendance Photo Taking

MODULE 03

DATA SCIENCE BOOTCAMP

# Thank you!