# Data Manipulation Seminar 4

**ICT233 Data Programming**

**Veronica Hu**          **huhe001@suss.edu.sg**

# RECAP

## S3 Key Learning Objectives

1) Appreciate the features and usage possibilities of the Pandas library as a data analytics package.

2) Understand the basic usage of the Python Pandas library, including loading files, counting data, and determining item structure and types in the data.

3) Learn the basic manipulation of data using Pandas, such as row and column selection, and itemized or vector operations on Pandas DataFrame.

4) Conduct operations on Pandas DataFrame, including subsetting, slicing, and indexing.

5) Present and visualize data in Pandas DataFrame using charting and plotting libraries like Matplotlib and Seaborn.

# SEMINAR OVERVIEW

## Data Manipulation – LEARNING OBJECTIVES

1) Assemble datasets together for analysis using Pandas

2) Understand the needs of concatenating datasets and performing the operations on them

3) Understand the needs of merging datasets and performing the operations on them

4) Learn what missing data are and how they are created

5) Work with data issues such as missing and incomplete data during analysis

6) Learn how to use pivot, melt, and normalization operations on datasets

# ETL Process

**Seminar 4**

### Extract

- One or more source systems containing customer, financial, or product data (CRM, Accounting system, Warehouse, MES)
- Files types - Flat files, XML, Oracle, IBM DB2, SQL Server,, IBM Websphere MQ, ODBC, JDBC, Hadoop Distributed File System (HDFS), Hive/HCatalog, JSON, Mainframe (IBM z/OS), Salesforce.com, SAP/R3
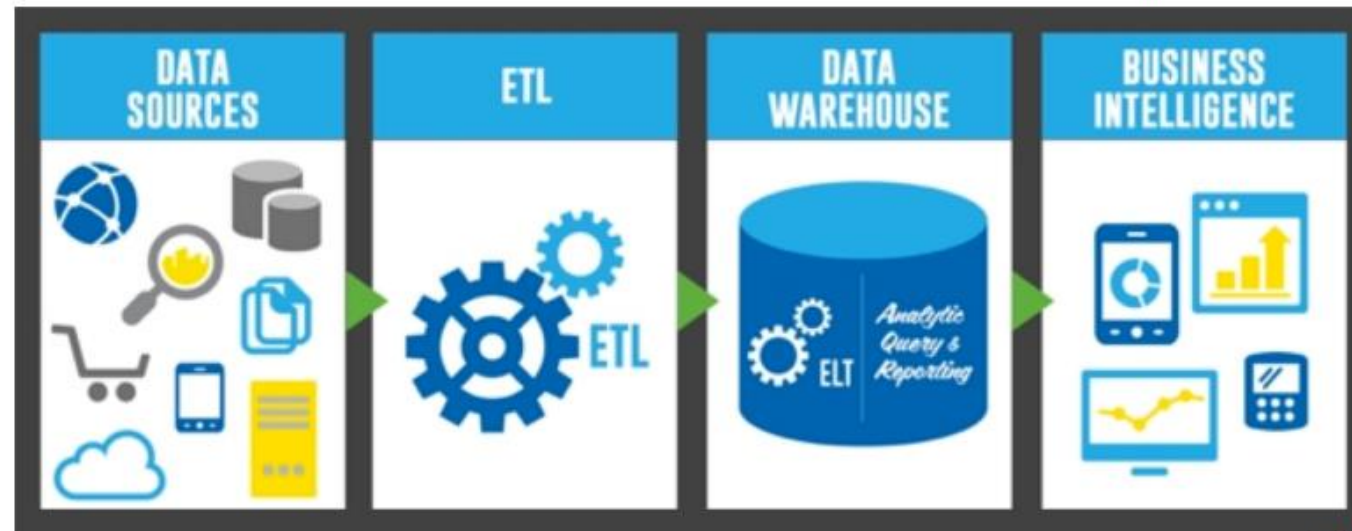
### Transform

- Applying business rules, cleansing, and validating the data.
- Aggregation, Copy, Join, Sort, Merge, Partition, Filter, Reformat, Lookup
- Mathematical: +, -, x, /, Abs, IsValidNumber, Mod, Pow, Rand, Round, Sqrt, ToNumber, Truncate, Average, Min, Max
- Logical: And, Or, Not, IfThenElse, RegEx, Variables
- Text: Concatenate, CharacterLengthOf, LengthOf, Pad, Replace, ToLower, ToText, ToUpper, Translate, Trim, Hash
- Date: DateAdd, DateDiff, DateLastDay, DatePart, IsValidDate
- Format: ASCII, EBCDIC, Unicode

### Load

- Load the results into one or more target systems such as a data warehouse, datamart, or business intelligence reporting system.
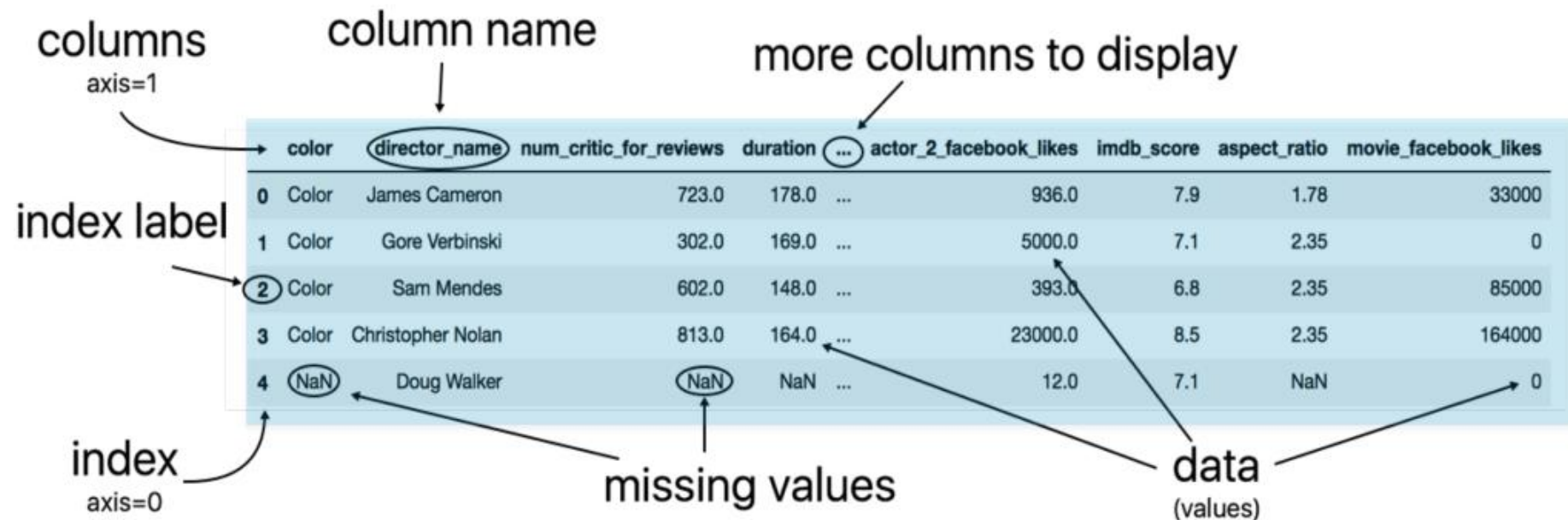- Output: Flat files, XML, Oracle, IBM DB2, SQL Server, Teradata, Sybase, Vertica, Netezza, Greenplum, ODBC, JDBC, Hadoop Distributed File System (HDFS), Hive/HCatalog, Mainframe (IBM z/OS), Salesforce.com, Tableau, QlikView

# Chapter 1: Assembling Data

## 1.1 Combining DataFrames

- Tidy data can be seen to meet the following criteria:

    - Each row in an observation

    - Each column is a variable

    - Each type of observational unit forms a table

# Chapter 1: Assembling Data

## 1.1 Combining DataFrames

Data cleaning/tidy processes

- Handling of missing values

    - Replacing nulls with values

        - Using another value

        - Using existing data

    - Drop the data from our data set

- Outliers

# Chapter 1: Assembling Data

## 1.1 Combining DataFrames

- The need to combine Data / DataFrame

  - Finding the data you need

  - e.g. find stock prices within the tech industry

- When splitting the data into separate tables

  - Advantage

    - Reduce redundant information

  - Drawback

    - The need to combine relevant data to answer questions



Company information



Stock information

# Chapter 1: Assembling Data

## 1.1 Combining DataFrames - Concatenation

- Horizontally (axis=1 )

- Vertically (axis=0 )

- Concatenation with different indices

    - Missing data is introduced

```
# axis default as 0
verticall_stacked = pd.concat([df1,df2,df3], axis=0)
display(verticall_stacked)
```

|   | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | 1.0 | 1.0 | 1.0 | 1.0 |
| 1 | 1.0 | 1.0 | 1.0 | 1.0 |
| 2 | 1.0 | 1.0 | 1.0 | 1.0 |
| 3 | 1.0 | 1.0 | 1.0 | 1.0 |
| 0 | 2.0 | 2.0 | 2.0 | 2.0 |
| 1 | 2.0 | 2.0 | 2.0 | 2.0 |
| 2 | 2.0 | 2.0 | 2.0 | 2.0 |
| 3 | 2.0 | 2.0 | 2.0 | 2.0 |
| 0 | 3.0 | 3.0 | 3.0 | 3.0 |
| 1 | 3.0 | 3.0 | 3.0 | 3.0 |

```
horizontal_stacked = pd.concat([df1,df2,df3], axis=1)
display(horizontal_stacked)
```

|   | 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 3.0 | 3.0 | 3.0 | 3.0 |
| 1 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 3.0 | 3.0 | 3.0 | 3.0 |
| 2 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 3.0 | 3.0 | 3.0 | 3.0 |
| 3 | 1.0 | 1.0 | 1.0 | 1.0 | 2.0 | 2.0 | 2.0 | 2.0 | 3.0 | 3.0 | 3.0 | 3.0 |

# Chapter 1: Assembling Data

## 1.1 Combining DataFrames - Concatenation

Reflections

- When do we need to do concatenation operations on data?

- How do we calculate daily average sales for a retail shop, for the past 3 months?

    - When each file contains daily transaction details

# Chapter 1: Assembling Data

## 1.2 Merging DataFrames

Merging Multiple Data Sets

| Pandas | SQL | Description |
| --- | --- | --- |
| left | left outer | Keep all the keys from the left |
| right | right outer | Keep all the keys from the right |
| outer | full outer | Keep all the keys from left & right |
| inner | inner | Keep only the keys that exist in both left and right |



**LEFT JOIN**

**FULL OUTER JOIN**

**INNER JOIN**

**RIGHT JOIN**

```
what_sites_were_visited = site.merge(visit, left_on='name', right_on='site', how='inner')
display(what_sites_were_visited)
```

| | name | location | id | site | date |
| --- | --- | --- | --- | --- | --- |
| 0 | DR-1 | location 1 | 1 | DR-1 | 2018-02-20 |
| 1 | DR-1 | location 1 | 2 | DR-1 | 2018-02-22 |
| 2 | DR-2 | location 2 | 3 | DR-2 | 2018-02-25 |

# Chapter 1: Assembling Data

## 1.2 Merging DataFrames

Reflection

- When do we need to do "inner" or "outer" join or merge operations on datasets ?

  - OUTER JOIN (LEFT or RIGHT)

    - Use when you want all rows from one table.

    - Retrieves only matching rows from the other table.

  - FULL OUTER JOIN:

    - Use when you want to get all rows from both tables.

# Chapter 2:  Handling Missing Data and Tidying up Information

## 2.1 Missing Information

Introduction

- Pandas displays missing values as NaN.

- NaN is not be equivalent to 0 or an empty string, ''

- Test for missing values

```
# import missing value defined in numpy library
from numpy import NaN, NAN, nan

import pandas as pd

print(pd.isnull(NaN))
```

True

```
print(pd.notnull(NaN), pd.notnull(888), pd.notnull('test'))
```

False True True

# Chapter 2: Handling Missing Data and Tidying up Information

## 2.1 Missing Information

- Cleaning Missing Data

  - Testing for null values

  - Replacing nulls with values

    - Fill in with another value

    - Fill in using existing data

      - Fill Forward/Backward

      - Interpolate

  - Drop the data from our data set

- Calculations with Missing Data

# Chapter 2: Handling Missing Data and Tidying up Information

## 2.2 Tidying and Organizing Information

- Multiple Columns with Same Variable

  - 'wide' data

  - melt()

| | religion | <$10k | $10-20k | $20-30k | $30-40k | $40-50k | $50-75k | $75-100k | $100-150k | >150k | Don't know/refused |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Agnostic | 27 | 34 | 60 | 81 | 76 | 137 | 122 | 109 | 84 | 96 |
| 1 | Atheist | 12 | 27 | 37 | 52 | 35 | 70 | 73 | 59 | 74 | 76 |
| 2 | Buddhist | 27 | 21 | 30 | 34 | 33 | 58 | 62 | 39 | 53 | 54 |
| 3 | Catholic | 418 | 617 | 732 | 670 | 638 | 1116 | 949 | 792 | 633 | 1489 |
| 4 | Don't know/refused | 15 | 14 | 15 | 11 | 10 | 35 | 21 | 17 | 18 | 116 |

# Chapter 2: Handling Missing Data and Tidying up Information

## 2.2 Tidying and Organising Information

- Columns with Multiple Variables – e.g. Ebola dataset

  - a single column in a dataset may **represent** multiple variables

  - multiple steps to tidy the data

```
list(ebola)
```

```
['Date',
 'Day',
 'Cases_Guinea',
 'Cases_Liberia',
 'Cases_SierraLeone',
 'Cases_Nigeria',
 'Cases_Senegal',
 'Cases_UnitedStates',
 'Cases_Spain',
 'Cases_Mali',
 'Deaths_Guinea',
 'Deaths_Liberia',
 'Deaths_SierraLeone',
 'Deaths_Nigeria',
 'Deaths_Senegal',
 'Deaths_UnitedStates',
 'Deaths_Spain',
 'Deaths_Mali']
```

```
ebola.head()
```

|   | Date | Day | Cases_Guinea | Cases_Liberia | Cases_SierraLeone | Cases_Nigeria | Cases_Senegal |
|---|------|-----|--------------|---------------|-------------------|---------------|---------------|
| 0 | 1/5/2015 | 289 | 2776.0 | NaN | 10030.0 | NaN | NaN |
| 1 | 1/4/2015 | 288 | 2775.0 | NaN | 9780.0 | NaN | NaN |
| 2 | 1/3/2015 | 287 | 2769.0 | 8166.0 | 9722.0 | NaN | NaN |
| 3 | 1/2/2015 | 286 | NaN | 8157.0 | NaN | NaN | NaN |
| 4 | 12/31/2014 | 284 | 2730.0 | 8115.0 | 9633.0 | NaN | NaN |

# Chapter 2:  Handling Missing Data and Tidying up Information

## 2.2 Tidying and Organizing Information

- Variables in Both Rows and Columns

  - multiple steps to tidy the data

  - melt() / pivot_table()

```python
import pandas as pd
weather = pd.read_csv('weather.csv')
weather.head()
```

|   | id | year | month | element | d1 | d2 | d3 | d4 | d5 | d6 | ... |
|---|----|------|-------|---------|-----|------|------|-----|------|-----|-----|
| 0 | MX17004 | 2010 | 1 | tmax | NaN | NaN | NaN | NaN | NaN | NaN | ... |
| 1 | MX17004 | 2010 | 1 | tmin | NaN | NaN | NaN | NaN | NaN | NaN | ... |
| 2 | MX17004 | 2010 | 2 | tmax | NaN | 27.3 | 24.1 | NaN | NaN | NaN | ... |
| 3 | MX17004 | 2010 | 2 | tmin | NaN | 14.4 | 14.4 | NaN | NaN | NaN | ... |
| 4 | MX17004 | 2010 | 3 | tmax | NaN | NaN | NaN | NaN | 32.1 | NaN | ... |

# Chapter 2: Handling Missing Data and Tidying up Information

## 2.2 Tidying and Organizing Information

Normalization

- Definition in Database Design:

  - Process of organizing and structuring data to eliminate redundancy.

  - Improves data integrity by defining dependencies and relationships between data.

- Key Steps:

  - Create tables to group related data.

  - Define relationships between attributes (columns) of these tables.

# Chapter 2: Handling Missing Data and Tidying up Information

## 2.2 Tidying and Organising Information

Normalization

- Normalization in DataFrames

  - Starting Point:

    - Check if multiple observational units are represented in a single table.

  - Identify Redundancies:

    - Examine rows for cells or values that are repeated across rows.

  - Strategy:

    - Reorganize repeated information into separate tables.

    - Define relationships between these tables to maintain data integrity.

# SUMMARY

## Data Manipulation – LEARNING OBJECTIVES

1) Assemble Data sets together for analysis using Pandas

2) Understand the needs of concatenating data sets and performing the operations on them

3) Understand the needs of merging data sets and performing the operations on them

4) Learn what missing data are and how they are created

5) Work with data issues such as missing and incomplete data during analysis

6) Learn how to use pivot, melt and normalization operations on data sets