

# Data Science Tools

---

Dr. Sethavidh Gertphol

# Outline

- \* IPython and Jupyter
- \* NumPy and SciPy
- \* Pandas
- \* scikit-learn
- \* Matplotlib and Seaborn
- \* Anaconda
- \* R

# IPython

- \* a shell for Python interactive programming
- \* interfaces with OS
  - \* file system
  - \* executing another program
- \* inline plot and other GUI
- \* a kernel for Jupyter

# spyder is an IDE for IPython

The image shows the Spyder IDE interface for Python 3.6. The main window is divided into three panes:

- Editor:** Contains a Python script named `practice1.py`. The script reads a CSV file, manipulates the data, and calculates scores. A red box labeled "scripts/code area" highlights the code editor.
- Variable explorer:** Displays the variables in the current namespace. A red box labeled "variables list" highlights this pane. It shows a table with columns: Name, Type, Size, and Value.
- IPython console:** Shows the execution of the code. A red box labeled "interactive shell" highlights this pane. It displays the output of the code, including the length of the `medals` DataFrame and the number of columns.

The Variable explorer table shows the following data:

Name	Type	Size	Value
col_total	Series	(16,)	Series object of pandas.core.series module
cols	list	16	['Country', 'SGames', 'SGold', 'SSilver', 'SBronze', 'STotal', 'WGames', 'WGold', 'WSilver', 'WBronze', 'WTotal', 'AllGames', 'AllGold', 'AllSilver', 'AllBronze', 'AllTotal']
medals	DataFrame	(147, 16)	Column names: Country, SGames, SGold, SSilver, SBronze, STotal, WGames ...

The IPython console shows the following output:

```
In [3]: len(medals)
Out[3]: 147

In [4]: len(medals.columns)
Out[4]: 16

In [5]: cols = ['Country', 'SGames', 'SGold', 'SSilver', 'SBronze', 'STotal', 'WGames',
File "<ipython-input-5-99a77ecafbce>", line 1
      cols = ['Country', 'SGames', 'SGold', 'SSilver', 'SBronze', 'STotal', 'WGames',
^
SyntaxError: unexpected EOF while parsing

In [6]:

In [6]: cols = ['Country', 'SGames', 'SGold', 'SSilver', 'SBronze', 'STotal', 'WGames',
...:           'WGold', 'WSilver', 'WBronze', 'WTotal', 'AllGames',
...:           'AllGold', 'AllSilver',
...:           'AllBronze', 'AllTotal']
...:

In [7]: medals.columns = cols

In [8]: col_total = medals.iloc[146]

In [9]:
```

The bottom status bar shows: Permissions: RW, End-of-lines: LF, Encoding: UTF-8, Line: 17, Column: 22, Memory: 65 %.

copyright Sethavidh Gertphol



# Jupyter Notebook

- \* web application that creates and shares documents that
  - \* has live code
  - \* has narrative text (in Markdown)
  - \* has visualization
- \* can view, run, and share codes through the web
- \* was part of IPython, but spinned out

web-based interface

markdown (text/bullets)

inserted picture

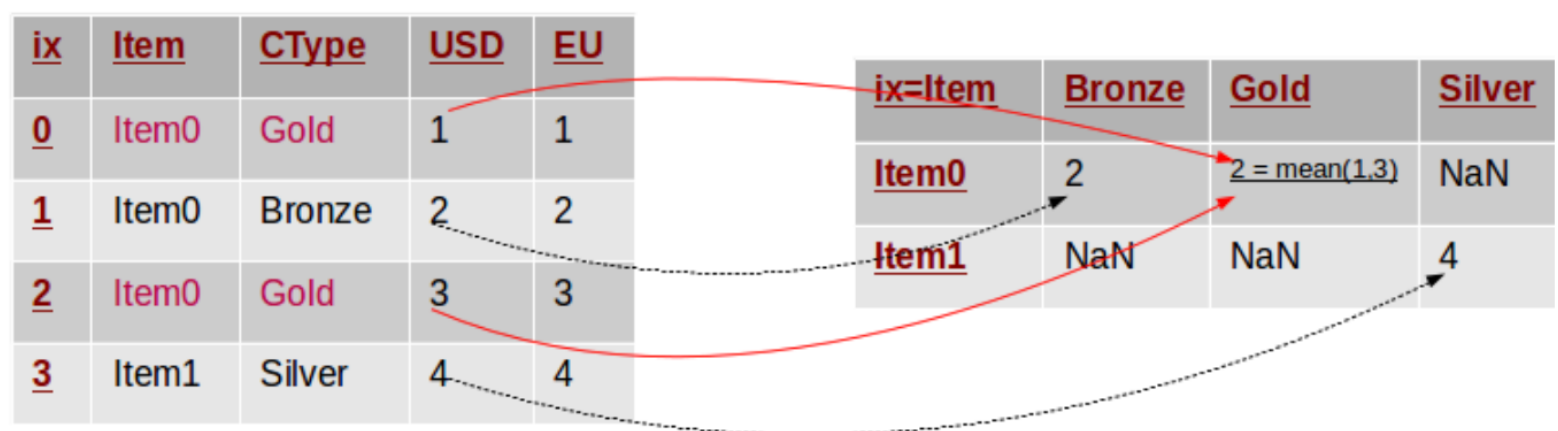
## Pivot Table

Pivot Table is used to create a new table where cells are summarized from another table.

The row (index) and column of the pivot table can be chosen from columns of a another table.

aggfunc parameter is the name of function to summarized data

<https://nikgrozev.com/2015/07/01/reshaping-in-pandas-pivot-pivot-table-stack-and-unstack-explained-with-pictures/>



d.pivot\_table(index='Item', columns='CType', values='USD', aggfunc=np.mean)

In [8]: pt = players.pivot\_table(index='team', columns='position', values='passes', aggfunc='mean').head(5)  
pt

Out[8]:

position	defender	forward	goalkeeper	midfielder
team				
Algeria	102.000000	10.000000	23.5	68.222222
Argentina	138.166667	112.666667	47.0	151.714286
Australia	78.750000	30.000000	51.0	63.090909
Brazil	190.000000	73.000000	69.0	111.750000
Cameroon	77.375000	41.000000	54.0	77.666667

code

result

# NumPy

numeric

- \* NumPy is a python library for scientific computing

- \* add supports for **large multi-dimensional array object**

nd array

วัตถุหลายมิติ (เวกเตอร์)

- \* includes operations like sorting array, linear algebra and Fourier transform

- \* basis for other scientific packages

# SciPy

Science

→ numpy numpy ၂၀၁၁  
(extension numpy)

- \* provides many more numerical operations on NumPy array such as
  - \* linear programming, optimization
  - \* linear algebra
  - \* statistics
  - \* signal processing
  - \* etc.
- \* used for scientific computing and much more



# Pandas

- \* provides high-performance **data structure and data analysis**
- \* Series and DataFrame objects
- \* data manipulation on those objects
- \* parallel operations for performance
- \* some compare Pandas to a command-line Excel
- \* basis for data science works

# scikit-learn

- \* Python package for **Machine Learning**
- \* basic ML algorithms
  - \* SVM, random forest, KNN classification
  - \* ridge/lasso regression
  - \* K-means, spectral clustering
- \* also
  - \* dimensionality reduction
  - \* model selection

# Matplotlib

- \* Python package for **plotting data and creating charts**
- \* charts can be output to IPython or Jupyter Notebook
- \* can be used as quick, exploratory plot
- \* or can be fine-tune to create publication quality charts

# Seaborn

new visualization → plot {dimension}!

- \* Python library for quickly plotting charts
- \* Use matplotlib underneath
- \* has default settings and standard interface to matplotlib
- \* users can create chart quickly with **single call** to seaborn



# Anaconda

- \* a Distribution of Python and packages
  - \* has everything we talked about in one distro
- \* also a package manager that allows searching and installing of new packages
- \* also an environment manager that allows multiple configurations of Python to exist in the same system



Home



Environments



Learning



Community

Documentation

Developer Blog



Applications on

base (root)

Channels

Refresh



JupyterLab

0.35.3

An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.

Launch



Notebook

5.7.2

Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.

Launch



Qt Console

4.3.1

PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.

Launch

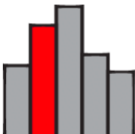


Spyder

3.3.2

Scientific PYTHON Development EnviRonment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection features

Launch



Glueviz

0.13.3

Multidimensional data visualization across files. Explore relationships within and among related datasets.

Install



Orange 3

3.17.0

Component based data mining framework. Data visualization and data analysis for novice and expert. Interactive workflows with a large toolbox.

Install



RStudio

1.1.456

A set of integrated tools designed to help you be more productive with R. Includes R essentials and notebooks.

Install



VS Code

1.30.1

Streamlined code editor with support for development operations like debugging, task running and version control.

Install

# R

- \* a language for statistical computing
- \* has GUI frontend (RStudio)
- \* package management (CRAN)
- \* has huge package selections

django

# Summary

Data  
Gathering

Data  
Cleaning

Data  
Transform

Data  
Visualization

Data  
Modeling

Python

Pandas

Matplotlib

scipy/  
scikit-  
learn

R / ggplot2