# Python for Data Science

Sirakorn Lamyai
September 19, 2019

Student, Kasetsart U.

http://bit.ly/cpe-datascience

## Sirakorn Lamyai

- DAKDL Laboratory, Kasetsart University
- Research Assistant Intern, 2019, Vidyasirimedhi Institute of Science and Technology
- Research Assistant Intern, 2018, Vidyasirimedhi Institute of Science and Technology
- Love drinking tea
- Knows a little about Python

# I know a little about Python

When I say I know *a little* about Python…

When I say I know *a little* about Python…

- I think there's some better methods than I'm using

When I say I know *a little* about Python…

- I think there's some better methods than I'm using
- I think I do sometimes make mistakes

# I know a little about Python

When I say I know *a little* about Python...

- I think there's some better methods than I'm using
- I think I do sometimes make mistakes
- There are tons of people who know things much more than me

## I know a little about Python

When I say I know *a little* about Python…

- I think there's some better methods than I'm using
- I think I do sometimes make mistakes
- There are tons of people who know things much more than me
- I think there's much more for me to learn!

## Prerequisite

A basic Python knowledge will do!

# Your expectations from this talk

## Outline

# Data Science

- **Obtain** data from relevent sources

- **Obtain** data from relevent sources
- **Scrub**, sanitise, and clean the data into machine-understandable formats

## The Data Science Process: OSEMNI

- **Obtain** data from relevent sources
- **Scrub**, sanitise, and clean the data into machine-understandable formats
- **Explore** significant and meaningful patterns with statistical methods

- **Obtain** data from relevent sources
- **Scrub**, sanitise, and clean the data into machine-understandable formats
- **Explore** significant and meaningful patterns with statistical methods
- **Model** construction for prediction and forecast

- **Obtain** data from relevent sources
- **Scrub**, sanitise, and clean the data into machine-understandable formats
- **Explore** significant and meaningful patterns with statistical methods
- **Model** construction for prediction and forecast
- **iNterpret** and use the results obtained

# The Data Science Process: OSEMNI

- **Obtain** data from relevent sources
- **Scrub**, sanitise, and clean the data into machine-understandable formats
- **Explore** significant and meaningful patterns with statistical methods
- **Model** construction for prediction and forecast
- **iNterpret** and use the results obtained
- **Interate** and rethink about your outputs

# Why data?

# Data is the new oil

# Tools for data analysis

## With GUIs

- Spreadsheets
  - Excel
  - Google Spreadsheets
  - Lotus 1-2-3
- Modelling and Visualisation
  - RapidMiner Studio
  - Weka
  - Tableau

## As programming languages

- For data insights
  - R
  - Python
- For data retrieval
  - SQL

# Python

Courtesy: xkcd (https://xkcd.com/353/)

# I *loved* Python...

- Read it, understand it
- Multiparadigm
- Batteris included
- Lots of great, great libraries!

# Python Packages

`pip`

# pip

- PyPA (Python Packaging Authority)'s recommended package installer

pip

- PyPA (**Py**thon **P**ackaging **A**uthority)'s recommended package installer
- Obtains packages from PyPI (**Py**thon **P**ackaging **I**ndex)

# pip

- PyPA (**Py**thon **P**ackaging **A**uthority)'s recommended package installer
- Obtains packages from PyPI (**Py**thon **P**ackaging **I**ndex)
- Many useful packages for us to use!

# Anaconda Python Distribution

- Cross-platform Python Distribution

- Cross-platform Python Distribution
- Ships with its own package and environment manager

- Cross-platform Python Distribution
- Ships with its own package and environment manager
  - Its environment manager capability is not found in Python vanilla installation

- Cross-platform Python Distribution
- Ships with its own package and environment manager
  - Its environment manager capability is not found in Python vanilla installation
  - Fetches the packages from its own repository, not PyPI

- Cross-platform Python Distribution
- Ships with its own package and environment manager
  - Its environment manager capability is not found in Python vanilla installation
  - Fetches the packages from its own repository, not PyPI
- Aims for Data Science use

- Cross-platform Python Distribution
- Ships with its own package and environment manager
  - Its environment manager capability is not found in Python vanilla installation
  - Fetches the packages from its own repository, not PyPI
- Aims for Data Science use
- **Entirely separated Python**

## Environments 101: $PATH

```
$ echo $PATH
/home/srakrn/.pyenv/plugins/pyenv-virtualenv/shims:/home/
    srakrn/.pyenv/shims:/home/srakrn/.pyenv/bin:/home/srakrn
    /.local/bin:/usr/local/bin:/usr/local/sbin:/home/srakrn/.
    local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/
    local/games:/snap/bin
```

## Different machines, different Pythons

On my laptop…

```
srakrn@epsilon-ubuntu:~$ which python
/home/srakrn/.pyenv/shims/python
```

On my https://charles.srakrn.me/ server…

```
srakrn@charles:~$ which python
/usr/bin/python
srakrn@charles:~$ which python3
/usr/bin/python3
```

# Installed pip

```
$ pip -V
pip 8.1.1 from /usr/lib/python2.7/dist-packages (python 2.7)
$ pip3 -V
pip 8.1.1 from /usr/lib/python2.7/dist-packages (python 3.6)
```

MY PYTHON ENVIRONMENT HAS BECOME SO DEGRADED
THAT MY LAPTOP HAS BEEN DECLARED A SUPERFUND SITE.

Courtesy: xkcd (https://xkcd.com/1987/)

Interactive computing environment

- Thinks of a more *dynamic* coding environment.

- Thinks of a more *dynamic* coding environment.
- Inserts snippets of codes alternately with texts, maths, and images.

- Thinks of a more *dynamic* coding environment.
- Inserts snippets of codes alternately with texts, maths, and images.
- A wonderful tool for coding *documented code.*

# Google Colaboratory (Colab)

- Think of an online Jupyter Notebook provided by Google

- Think of an online Jupyter Notebook provided by Google
- The runtime relies on Google's server fram

- Think of an online Jupyter Notebook provided by Google
- The runtime relies on Google's server fram
  - In other words, your code are remotely executed

# Google Colaboratory (Colab)



- Think of an online Jupyter Notebook provided by Google
- The runtime relies on Google's server fram
  - In other words, your code are remotely executed
- Could be more powerful for some tasks (like Deep Learning) than your computer

- Think of an online Jupyter Notebook provided by Google
- The runtime relies on Google's server fram
  - In other words, your code are remotely executed
- Could be more powerful for some tasks (like Deep Learning) than your computer
- Free!

# https: //colab.research.google.com/

# Caveats 1: Execution order

```
[2]  a = 10

[1]  a = 5

[3]  print(a)

     10
```

# Caveats 1: Execution order

```
[2]  a = 10

[1]  a = 5

[3]  print(a)

     10
```

You'll do a lot of out-of-order code execution!

# Caveats 2: Cell edits

```
[1]  a = 10

[2]  a *= 2

[4]  print(a)
```

60

# Caveats 2: Cell edits

```
[1]  a = 10

[2]  a *= 2

[4]  print(a)
```
⮕  60

You might sometimes remove a cell, and that shows no visible trace without explicit query.

# Caveats 2: Cell edits

```
[1]  a = 5

[2]  a = 20

[3]  print(a)
```

10

# Caveats 2: Cell edits



```
[1]  a = 5

[2]  a = 20

[3]  print(a)

     10
```

Jupyter Notebook offers no cell edited marks, while Colab offers them

# Caveats 2: Cell edits

```
[1]  a = 5

[2]  a = 20

[3]  print(a)

     10
```

Jupyter Notebook offers no cell edited marks, while Colab offers them (note: observe the greyed out cell number)

## Caveats 3: Be neat and tidy

Jupyter Notebook and Colab, unlike IDE and code editors, offers a relatively poor
**clean code** tools

- Syntax error highlighting
- Autocomplete
- Linting
- Code formatter

Sirakorn's Workflow Demo

# Python Data Structures

# Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
    ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

· **Lists** are a compilation of objects.

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
    ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

## Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
    ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

- **Lists** are a compilation of objects.
- Can store multiple data types.

# Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
    ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

- **Lists** are a compilation of objects.
- Can store multiple data types.
  - This includes storing lists in a list

# Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
     ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

- **Lists** are a compilation of objects.
- Can store multiple data types.
  - This includes storing lists in a list
  - So-called a **nested list**

# Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
     ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

- **Lists** are a compilation of objects.
- Can store multiple data types.
  - This includes storing lists in a list
  - So-called a **nested list**
- Can be resized.

## Lists

```
1 a = [1, 2, 3, 4, 5]
2 b = ["Cats", "Dogs", "Penguins
    ", "Tonkatsu Pieces"]
3 c = [1, "1", True]
```

· **Lists** are a compilation of objects.
· Can store multiple data types.
  · This includes storing lists in a list
  · So-called a **nested list**
· Can be resized.
  · No need to declare its size on the first declaration.

# Lists

```python
a = [1, 2, 3, 4, 5]
a[0]    # Accessing elements
a[1:3]  # Slicing
```

Accessing list

# Lists

```
1 a = [1, 2, 3, 4, 5]
2 a[0]     # Accessing elements
3 a[1:3]  # Slicing
```

Accessing list

```
1 a = [1, 2, 3, 4, 5]
2 a[0]    # Accessing elements
3 a[1:3]  # Slicing
```

Accessing list

- **Elementwise**: accessing one elements at a time)

# Lists

```python
a = [1, 2, 3, 4, 5]
a[0]    # Accessing elements
a[1:3]  # Slicing
```

Accessing list

- **Elementwise**: accessing one elements at a time)
- **Slicing**: accessing a sublist

# List Functions

```python
vowels = ["a", "e", "o", "u"]

# Get a's length
len(a)
# Append the new element to the end of a
a.append("y")
# Deletes the first occurence of the element from a
a.remove("y")
# Inserts the item into a list with a specified index
a.insert(2, "i")
```

# List Functions

```
1 vowels = [1, 3, 2, 5, 4]
2 # Get the first index of a specified element
3 a.index(4)
4 # Sort a list and store into a new list
5 sorted_a = sorted(a)
6 # Sort a list, making changes directly to the old one
7 a.sort()
```

## Dictionary

```python
names = {
    "Cherprang": "Cher",
    "Manipa": "Khamin",
    "Jiradapa": "Pupe"
}

kami_nickname = names["Manipa"
    ]
```

## Dictionary

```
1 names = {
2     "Cherprang": "Cher",
3     "Manipa": "Khamin",
4     "Jiradapa": "Pupe"
5 }
6
7 kami_nickname = names["Manipa"
      ]
```

- **Dictionaries** store values in a **key-pair** format.

## Dictionary

```
1 names = {
2     "Cherprang": "Cher",
3     "Manipa": "Khamin",
4     "Jiradapa": "Pupe"
5 }
6
7 kami_nickname = names["Manipa"
      ]
```

- Dictionaries store values in a key-pair format.
- From the positional index, dictionary takes the key as an index instead.

# NumPy

- NumPy is a powerful library for mathematical computation in Python

- NumPy is a powerful library for mathematical computation in Python
- It offers wide range of tools from data structures to advanced functions and operations

- NumPy is a powerful library for mathematical computation in Python
- It offers wide range of tools from data structures to advanced functions and operations
- A very strong ease in mathematical computation, **don't reinvent the wheels!**

*Courtesy: Rebellious Professor,*

# Let's go to Notebook!

*I'm too lazy to cover the contents twice...*

# Pandas

# Pandas

- Tabular-like structure

- Tabular-like structure
- High performance
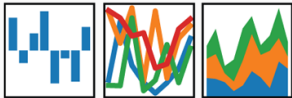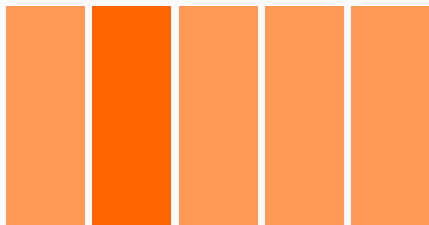
- Tabular-like structure
- High performance
- Easy to use

$$y_i t = \beta' x_{it} + \mu_i + \epsilon_{it}$$

- Tabular-like structure
- High performance
- Easy to use
- Helps a lot in data preparation

- Tabular-like structure
- High performance
- Easy to use
- Helps a lot in data preparation
- Considered as a wrapper for Numpy, although there's a lot more

# Let's go to Notebook!
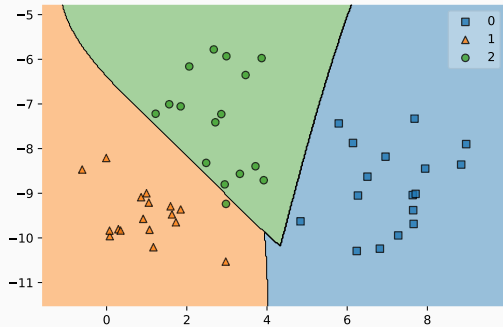
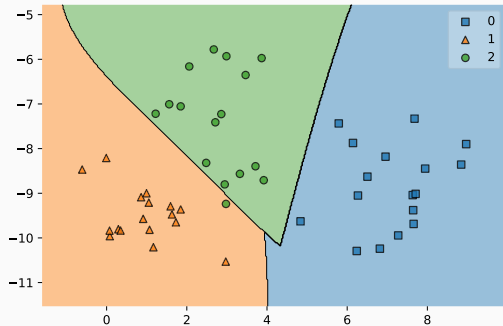Why Visualise?

Why Visualise?
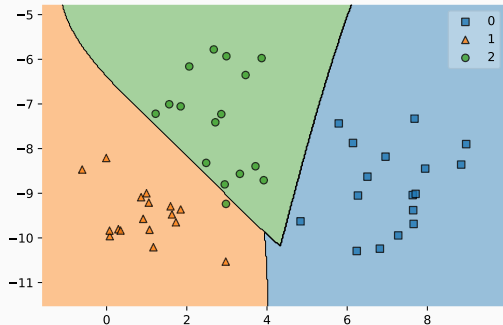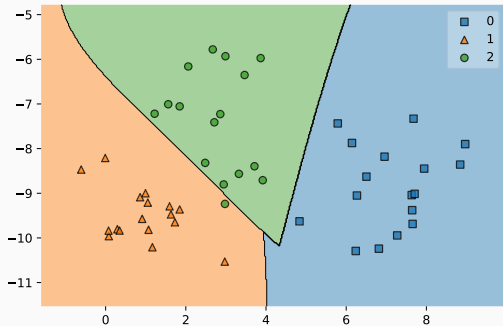
· Our visual system is so great!
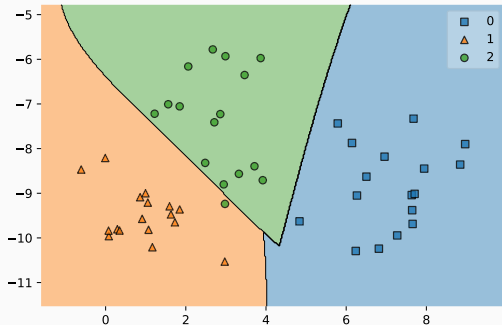
Why Visualise?

- Our visual system is so great!
- A much better way to represent data than statistical values

Why Visualise?

- Our visual system is so great!
- A much better way to represent data than statistical values
- Meaningful plots show meaningful insights without needing to do much

# The Datasaurus Dozen

`https://www.autodeskresearch.com/publications/samestats`
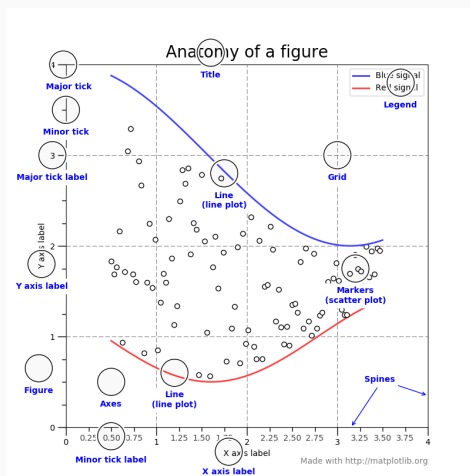
# Matplotlib

- feature-rich plotting library in python

- feature-rich plotting library in python
- plots multiple types of charts

- feature-rich plotting library in python
- plots multiple types of charts
- extensive support for jupyter notebook

*Source:* `https://matplotlib.org/tutorials/introductory/usage.html`

# Let's go to Notebook!