# Machine Learning in Python

Software Engineering for AI systems - DAT821
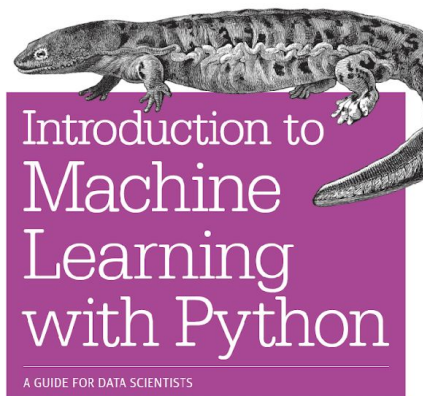
AI system lifecycle

# Objectives

**Learn** how subtasks of machine learning map unto Python

- **Understand** and **visualize** your training data before training
- **Apply** ML algorithm on your training data
- **Evaluate** your ML model on your test data

# **Great References** BOOKS & INTERNET



Andreas C. Müller & Sarah Guido

**_Python documentation_**
https://docs.python.org/3/

**_Python tutorials_**
- _Learn Python: step-by-step tutorial_ (Link)
- _Several others e.g., listed here_ (Link)

**_Machine learning with Python tutorials_**
- Scipy- lectures (Link)
- Others - _search machine learning with Python_
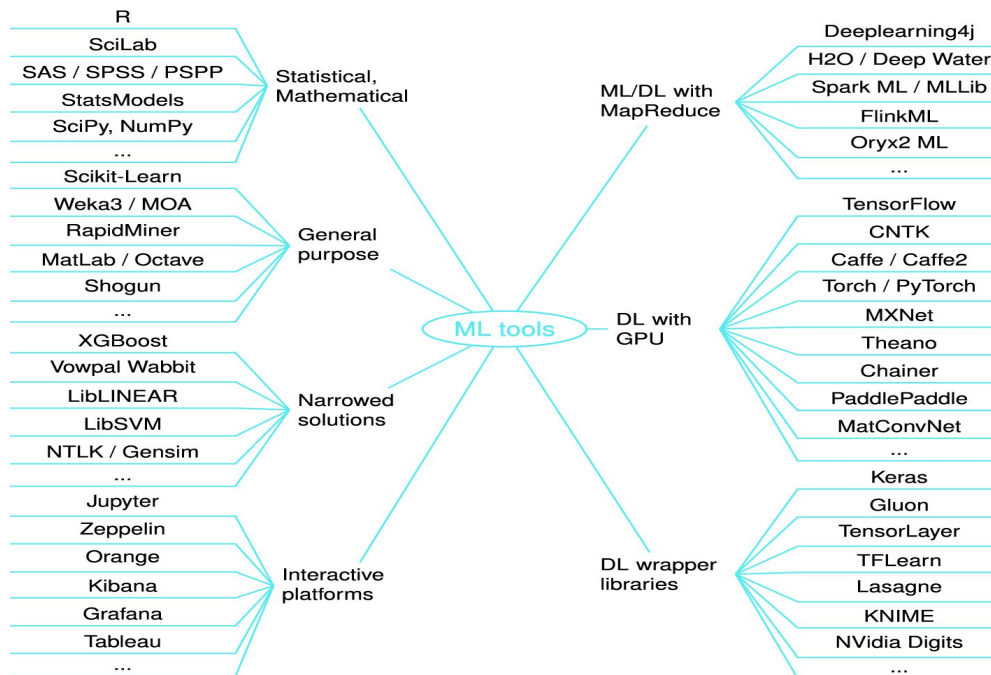
# What is, and why Python?

- Python is a dynamic, interpreted language
- Has no type declarations of variables, parameters, functions, or methods in source code
- It tracks the types of all values at runtime and flags code that does not make sense as it runs
- Has a growing and dominant ecosystem for machine learning

python Programming

# What is, and why Python?

## Overview of Machine Learning frameworks and libraries

*The number of ML algorithms and their software implementation is quite large*



Nguyen, et al., (2019). Machine Learning and Deep Learning frameworks and libraries for large-scale data mining: a survey. Artificial Intelligence Review, 52(1), 77-124
scipy-lectures

# What is, and why Python?

Overview of Machine Learning frameworks and libraries

*Findings from a survey:*

- Python is the most popular programming language for data mining, Machine Learning and Deep Learning applications

- The majority of frameworks and libraries are either Python based or support Python interfaces

| Tool | Licence | Written in | Computation graph | Interface | Popularity | Usage | Creator (notes) |
|---|---|---|---|---|---|---|---|
| TensorFlow (Numerical framework) | Open source, Apache 2.0 | C++, Python | Static with small support for dynamic graph | Python, C++[a], Java[a], Go[a] | Very High Growing very fast | Academic Industrial | – **Google** |
| Keras (Library) | Open source, MIT | Python | Static | Python Wrapper for TensorFlow, CNTK, DL4J, MXNet, Theano | High Growing very fast | Academic Industrial | F. Chollet |
| CNTK (Framework) | Open source, Microsoft permissive license | C++ | Static | Python, C++, BrainScript, ONNX | Medium Growing fast | Academic Industrial Limited mobile solution | – **Microsoft** |
| Caffe (Framework) | Open source, BSD 2-clause | C++ | Static | C++, Python, MatLab | High Growing fast | Academic Industrial | Y. Jia **BAIR** |
| Caffe2 (Framework) | Open source, Apache 2.0 | C++ | Static | C++, Python, ONNX | Medium-low Growing fast | Academic Industrial Mobile solution | Y. Jia **Facebook** |
| Torch (Framework) | Open source, BSD | C++, Lua | Static | C, C++, LuaJIT, Lua, OpenCL | Medium–low Growing low | Academic Industrial | R. Collobert, K. Kavukcuoglu, C. Farabet |
| PyTorch (Library) | Open source, BSD | Python, C | Dynamic | Python, ONNX | Medium Growing very fast | Academic Industrial | A. Paszke, S. Gross, S. Chintala, G. Chanan |
| MXNet (Framework) | Open source, Apache 2.0 | C++ | Dynamic dependency scheduler | C++, Python, Julia, MatLab, Go, R, Scala, Perl, ONNX | Medium Growing fast | Academic Industrial | – **Apache** |
| Chainer (Framework) | Open source, Owners permissive license | Python | Dynamic | Python | Low Growing low | Academic Industrial | – **Preferred Networks** |
| Theano (Numerical framework) | Open source, BSD | Python | Static | Python | Medium-low Growing low | Academic Industrial | Y. Bengio **University of Montreal** |

Nguyen, et al., (2019). Machine Learning and Deep Learning frameworks and libraries for large-scale data mining: a survey. Artificial Intelligence Review, 52(1), 77-124
scipy-lectures

# Python installation

- Check if Python is installed
- Installation guides Link
- Recommended installation
  - Anaconda Link
- Package installation
  - Using pip Link
  - Using conda  Link
- Virtual environment
  - using venv Link
  - using conda Link

```
$ pip install numpy scipy
matplotib scikit-learn pandas
```

```
$ conda install numpy scipy
matplotib scikit-learn pandas
```

Python Packaging index (https://pypi.org/)

# Python syntax

Example: Lab1 –
Linear regression
with one variable

```python
1   import os
2   import scipy
3   import sklearn
4   from sklearn.model_selection import train_test_split
5   from sklearn.linear_model import LinearRegression
6   from sklearn import metrics
7   import numpy as np
8   import pandas as pd
9   import matplotlib.pyplot as plt
10
11  # Step 1. load data
12  def load_data(data_path, file_name):
13      data_path = os.path.join(data_path, file_name)
14      return pd.read_csv(data_path, header = None)
15
16  my_path = "./lab1"
17  train_file = "ex1data1.txt"
18  train_data = load_data(my_path, train_file)
19  train_data.columns= ['Population', 'Profit']
20
21  # Step 2.  Explore and visualize data
22  print(train_data.head(5))
23  print(train_data.shape)
24  print(train_data.describe())
```

Pep8 – recommended style guideline for python code (https://www.python.org/dev/peps/pep-0008/#function-and-variable-names)

# Python ecosystem for machine learning

- **NumPy** – For N-dimensional *array manipulation*
- **SciPy** – For advanced *mathematical routines*
- **Pandas** – For *data analysis* and *data structures*
- **Matplotlib** – For *2D plotting*
- **Scikit-learn** – For *machine learning* algorithms
- **PyTorch** – For *deep learning*

# Python ecosystem for machine learning

## Installation

### Pip

```
$ pip install numpy scipy matplotib scikit-learn pandas
```

### Conda

```
$ conda install numpy scipy matplotib scikit-learn pandas
```

# NumPy

- A fundamental package for scientific computing in Python

- It contains functionality for multidimensional arrays, high-level mathematical functions such as linear algebra operations

- NumPy is a fundamental data structure on scikit-learn

```
>>> import numpy as np
>>> a = np.array([0, 1, 2, 3])
>>> a
array([0, 1, 2, 3])
```

Scipy-lectures

# SciPy

- SciPy provides among other advanced linear algebra routines, mathematical function optimization, signal processing, special mathematical function and statistical distributions

- It is meant to operate efficiently on numpy arrays, so that numpy and scipy work hand in hand

- scikit-learn draws from SciPy's collection of functions for implementing its algorithms

scipy- lectures

# Pandas

- A library for data wrangling and analysing
- Built around a data structure called DataFrame, which is like a table (similar to Excel spreadsheet)
- Provides a wide range of methods to modify and operate on DataFrame
- Contrast to NumPy, pandas allows each column to have a separate data type
- Pandas gives the ability to ingest from a great variety of data files formats

```python
11  # Step 1. load data
12  def load_data(data_path, file_name):
13      data_path = os.path.join(data_path, file_name)
14      return pd.read_csv(data_path, header = None)
15
16  my_path = "./lab1"
17  train_file = "ex1data1.txt"
18  train_data = load_data(my_path, train_file)
19  train_data.columns= ['Population', 'Profit']
20
21  # Step 2.  Explore and visualize data
22  print(train_data.head(5))
23  print(train_data.shape)
24  print(train_data.describe())
25
26  X = train_data.Population.values.reshape(-1, 1)
27  y = train_data.Profit.values.reshape(-1,1)
28
29  print(X)
30  print(y)
```

https://pypi.org/project/pandas/

# Matplotlib

- Matplotlib is the primary scientific plotting library in Python
- It provides functions for making good quality visualizations, such as histograms, scatter plots etc.,
- Visualizations are important for giving insights into your data analysis

```python
31
32    # Plotting the data
33    plt.scatter(X, y, c='blue')
34    plt.title('Scatterplot of training data')
35    plt.xlabel('Population')
36    plt.ylabel('Profit in $10,000')
37    plt.show()
38
```

# Scikit-learn

- Scikit-learn contains a number of machine learning algorithms for classification, regression, clustering ,as well as pre-processing
- Scikit-learn depends on two other Python packages NumPy and sciPy
- ML algorithms implemented in scikit-learn expect data to be stored in a **two-dimensional array or matrix** e.g., numpy arrays

```
38
39    # Step 3. Train linear regression mode
40    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
      =0.2, random_state=0)
41    lr = linear_model.LinearRegression()
42    model = lr.fit(X,y)
43
44    # Theta as attributes of model using intercept_ for theta_zero and
      coef_ for theta_one
45    print('Intercept: ', model.intercept_)
46    print('Slope: ', model.coef_)
47
```

# Scikit-learn

- Provides metrics functions for evaluating ML model performance

```
48    # Step 4. Evaluate trained model
49    y_pred = model.predict(X_test)
50
51    df = pd.DataFrame({'Actual': y_test.flatten(), 'Predicted':
      y_pred.flatten()})
52    print(df)
53
54    #Plotting linear model
55    plt.scatter(X_test, y_test, color='blue')
56    plt.plot(X_test, y_pred, color='red', linewidth=2)
57    plt.show()
58
59    # Evaluate with mean squared error to evaluate the model
60    print('Mean squared error: ', metrics.mean_squared_error(y_test,
      y_pred))
```