

Machine Learning Tutorial

Week1-Part 1

Introduction

Yasin Ceran

Anaconda Python

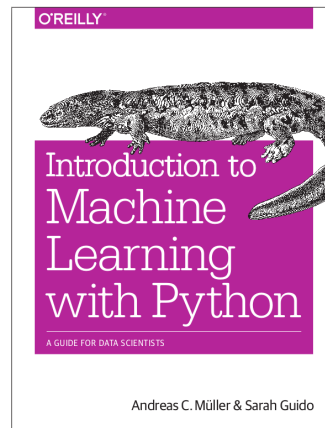
- First, download Anaconda. We recommend downloading Anaconda's latest Python 3 version.
<https://www.anaconda.com/distribution/>
- Second, install the version of Anaconda which you downloaded, following the instructions on the download page.
- Congratulations, you have installed Jupyter Notebook!

Python Knowledge

- Familiarity with Python programming and basic use of NumPy, pandas and matplotlib.
- A good reference is the Python Data Science Handbook by Jake VanderPlas.
- It's online for free and available as a notebook at the link below. I highly recommend going through it before starting the class.

<https://github.com/jakevdp/PythonDataScienceHandbook>

Book



Most of the slides and lecture content is based on the class notes and the textbook of Dr. Andreas Mueller.

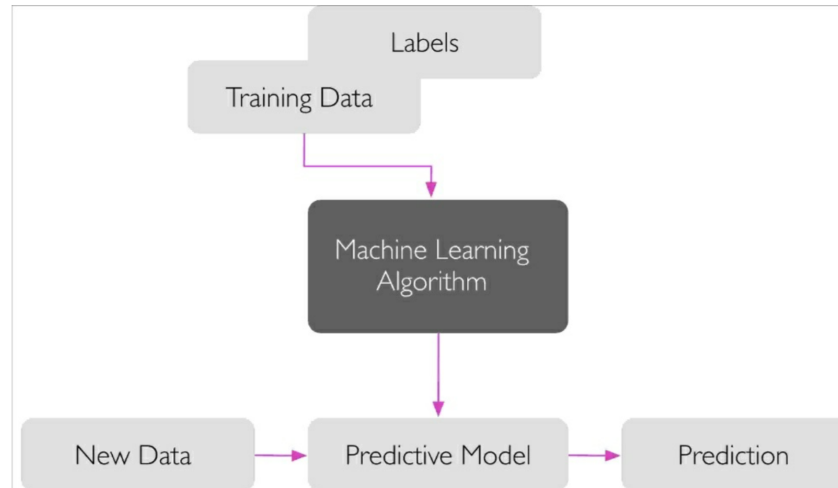
What is Machine Learning

- Large amount of structured and unstructured data
- Machine Learning helps capturing the knowledge from the data to improve the performance of predictive models and make data-driven decisions

Types of Machine Learning

Supervised Learning	<ul style="list-style-type: none">> Labeled data> Direct feedback> Predict outcome/future
Unsupervised Learning	<ul style="list-style-type: none">> No labels> No feedback> Find hidden structure in data
Reinforcement Learning	<ul style="list-style-type: none">> Decision process> Reward system> Learn series of actions

Supervised Learning



Supervised Learning

$$(x_i, y_i) \propto p(x, y) \text{ i.i.d.}$$

$$x_i \in \mathbb{R}^p$$

$$y_i \in \mathbb{R}$$

$$f(x_i) \approx y_i$$

Classification and Regression

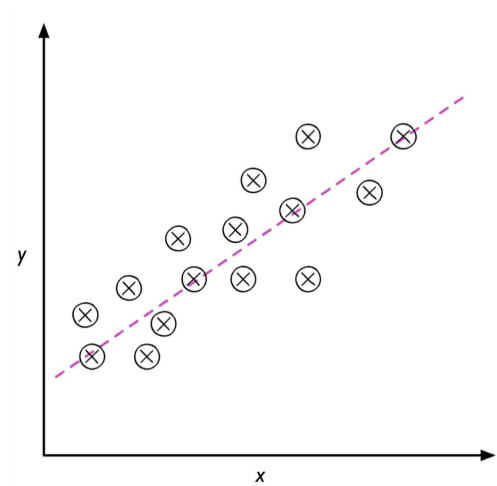
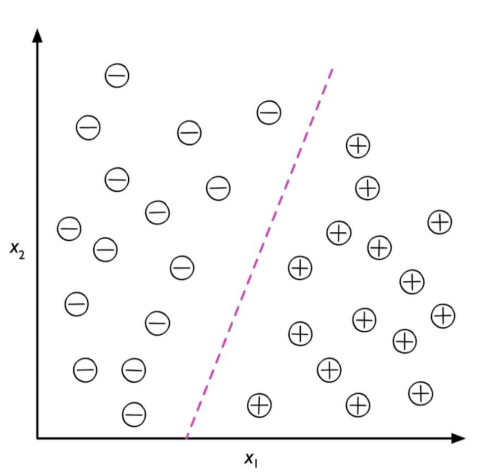
Classification

- target y discrete
- Will you pass?

Regression

- target y continuous
- How many points will you get in the exam?

Classification and Regression



Generalization

Not only

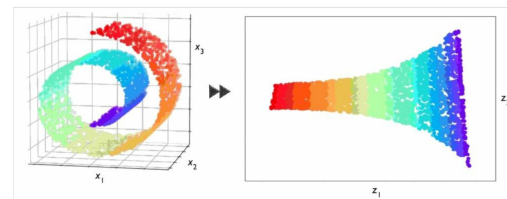
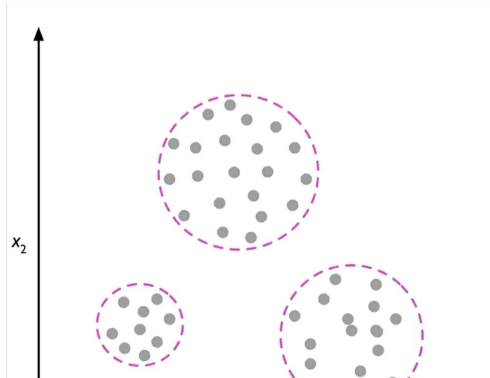
$$f(x_i) \approx y_i,$$

also for new data:

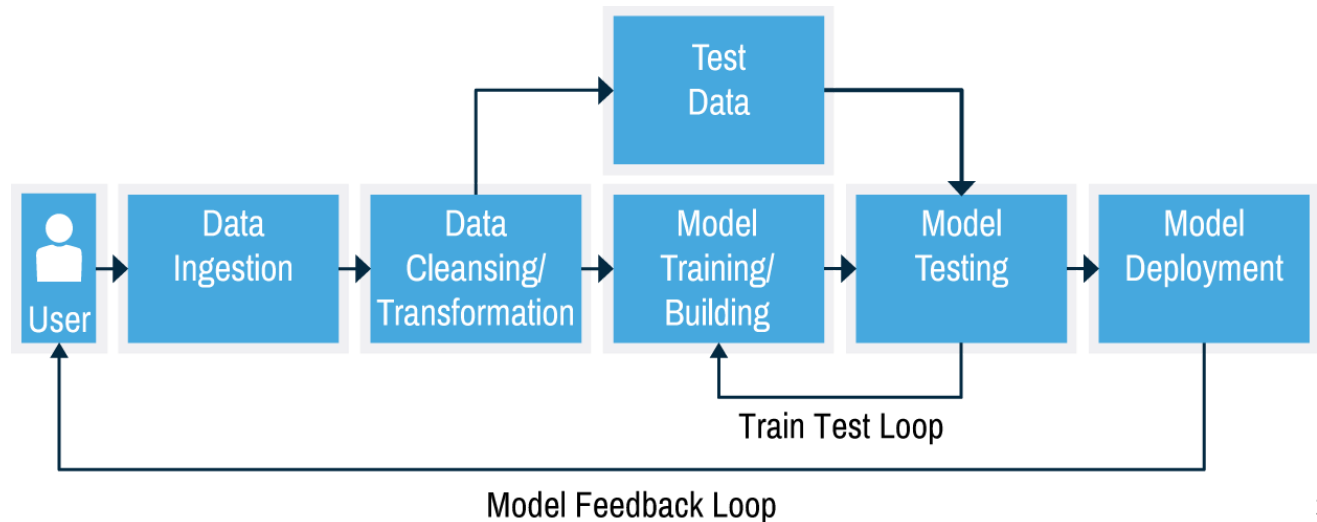
$$f(x) \approx y$$

Unsupervised Learning

- Clustering is an explanatory data analysis technique
- Dimensionality reduction is used to remove noise and compress data



The Machine Learning Work-Flow



Representing Data

Training and Test Data

training set

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix}$$

test set

$$y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

Jupyter Notebook

Part 1- Data Loading

Questions ?