# Data Preprocessing

## Trung Nguyen

## Contents

# 1 Introduction

Machine learning "gives computers the ability to learn without being explicitly programmed" (Arthur Samuel, 1959).

There are five main tasks in the machine learning work flow.

- Data collection and preparation.

- Feature selection and feature engineering.

- Choosing the machine learning algorithm (choosing the model) and training the model.

- Evaluating the model

- Model tweaking, parameter tuning and prediction.

Data Prepossessing is an important step. It transforms the raw data into an understandable format for the machine.The phrase "Garbage In ,Garbage Out" is particularly applicable to machine learning and data mining.
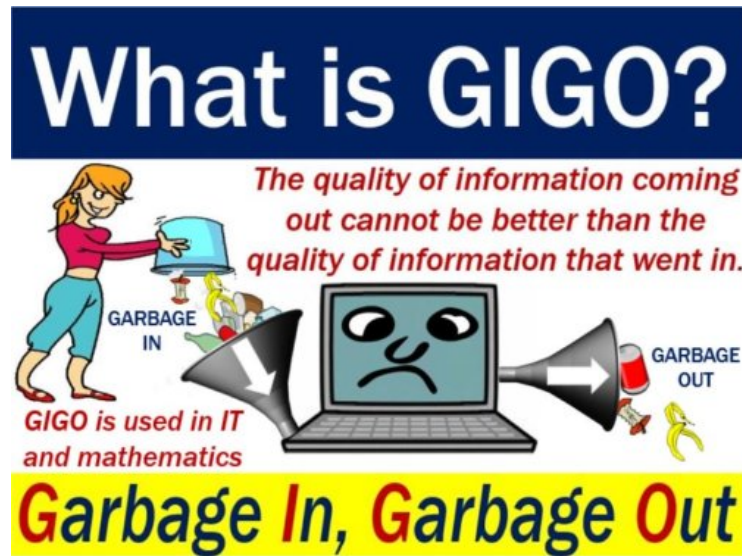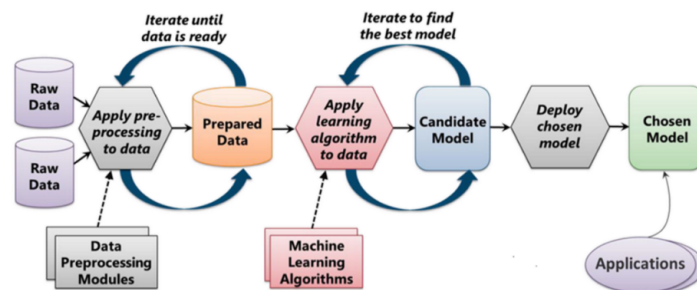
Figure 1: GIGO
source:[1]

If it is much irrelevant and redundant data or noisy and unreliable data then it is difficult to understand the pattern and our machine learning algorithm will not fit well to predict the correct estimations. (No Quality data, no correct predictions). In the next chapter, the discussion on the basic step of preparing the data will be further discussed.



Figure 2: Machine Learning Process
source:[2]

## 2 Data Reprocessing and Data Visualization

The visualization is not only important in the training steps but also important in the pre-processing steps.

The raw data are often incomplete, noisy, and inconsistent. Visualizing the data help users understand more in-depth of the data sets.

These are the most commonly used libraries in python.

- **Numpy:** the fundamental package for computing and mathematical operations with Python.

- textbfPandas: it is used for data manipulation and analysis.

- **Matplotlib :** python 2D plotting library.

- **Seaborn:** another python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

  Try an notebook example of data pre processing here.

## 3 Tips and Recommendation

A cheat sheet is a concise set of notes used for quick reference (wikipedia).

Looking into a cheat sheet will give a quick overview of the libraries function. An example of Numpy cheat sheet  is showing below:
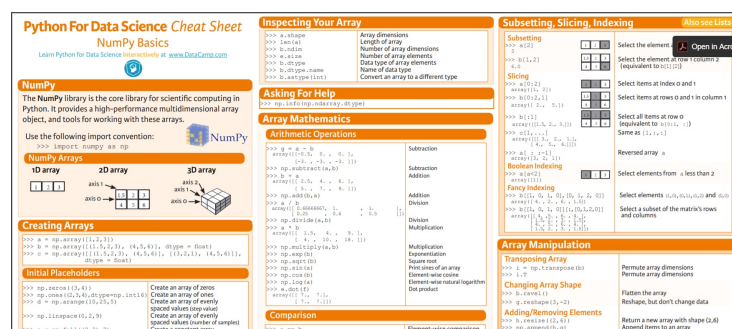


Figure 3: Numpy cheat sheet

Some useful cheat sheets for data manipulation:

- Numpy basic cheat sheet

- Panda basic cheat sheet

- Pandas Data Wrangling Cheat Sheet

- Matplotlib Cheat Sheet

- Seaborn Cheat Sheet

- Bokeh Cheat Sheet

There are a lot of tutorial from basic to advance with step to steps guideline which available from Medium and Towardsdatascience blog spot. Read and try out those examples not only help to improve knowledge in the field but also get up to with the current development.

# 4    Link to relevant information

Some more example of data preprocessing on Medium and Kaggle notebook.

Medium and towardsdatascience:

- Data Preprocessing- A significant step in Machine Learning

- Data Preprocessing : Concepts

- Data Pre Processing Techniques You Should Know

.

Kaggle Notebook:

- Full Preprocessing Tutorial Notebook

# 5    List of Documentation

1. Python Guideline for Beginner.

2. Setting up Online Python Notebook

3. Data Preprocessing

4. update...

# 6 Recommendation Courses

Machine learning by Andrew Ng.
Deep Learning using FastAI by Jeremy Howard.

# References

[1]  URL: https://marketbusinessnews.com/financial-glossary/gigo-garbage-in-garbage-out/.

[2]  David Chappelle. *Introducing Azure Machine Learning*.