# Fundaments of Machine learning for and with engineering applications

Enrico Riccardi[1]

Department of Energy Resources, University of Stavanger (UiS).[1]

Jan 17, 2025



University of Stavanger

DATA

SORTED

ARRANGED

PRESENTED
VISUALLY

EXPLAINED
WITH A STORY

A representation should **capture** the nature of the subject being studied.

Example: If you want to evaluate the 3D structure of a wind turbine, a set of descriptors an be:

1. Blade length
2. Turbine height
3. Geographical position
4. Output power
5. Wind direction

which are two decimal numbers, a 2d tuple, a 1D time series and a 2D time series (or 3D even).

Same meaning **represenations** for different objects (inputs).

### Discussion point!

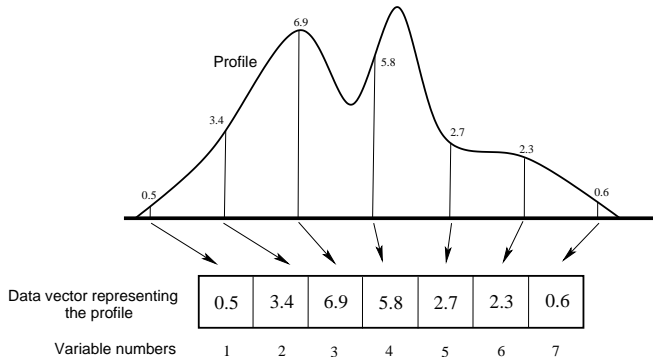How do we compare two wind turbines accounting for the 5 variables previously introduced?

- All starts from data: what are data-properties?

- Are there such things as good data and bad data?

**Life lesson (or exam question, same thing ;) )**
- Data DO NOT always have value.
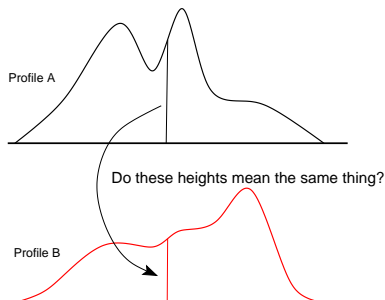
- TRASH in TRASH out

# Sampling point representation (SPR)

- An intuitive way to represent curves and spectra is the **sampling point representation**.
- We sample at regular intervals where each sample point is represented by a variable

# Sampling point representation (SPR)

- SPR is useful until point $i$ in a curve has the same meaning of the point $i$ in another curve.



Profile A

Do these heights mean the same thing?

Profile B

- Which parts of the profiles or shapes are comparable, i.e. have the same meaning?

Given a representation, it is then needed to decide on a suitable **data structure** for the problem.

### Definition
A data structure is a way of storing and organising data in a computer so that it can be used effectively.

Typical data structures used in data analysis are:

- Data points
- Arrays (vectors, matrices, N-mode (way) arrays)
- Graphs (trees)
- Databases

Data has to be prepared with these steps in mind

1. Plan experiments: Use experimental design to set up experiments in a *systematic* way
2. Pre-processing: Is there systematic variation in the data which should be removed Can cross-checking/validation procedures be designed?
3. Examine the data: Look at data (tables and plots). Strange behaviours? Smooth behaviour? WARNING!
4. Define desired model outcomes (speed, accuracy, false positive/negatives rate)
5. Estimate and validate model: What do the results tell us? Is the generated model general (valid for future sampling)?
6. Apply model to unknown samples

Statistics is collecting, organising, and interpreting data

Spatial and temporal statistics is a branch of applied statistics that emphasises:

1. the geo context of the data
2. the spatial and time dependent relationship between data
3. the different relative value and precision of the data.

The data matrix is an extremely common data structure.

$$X = \begin{bmatrix} 95 & 89 & 82 \\ 23 & 76 & 44 \\ 61 & 46 & 62 \\ 49 & 2 & 79 \end{bmatrix}$$

In python these can be saved as

- lists (vanilla python)
- numpy.arrays
- pandas dataframes

There are different conventions. Commonly we will construct data matrix such that:

- Rows are called instances, objects or samples.
- Columns are called features, variables.

One can think of each row to be an experiment, and the rows its properties. Each row (experiment, object, sample, ...) is thus a list of values, one for property.

**Note**

Mathematically speaking, this is just a notation. As long as one keeps track and is consistent, columns can be used as rows and vice versa.

Environmental measurements of rivers. The features (properties) can be:

- pH
- Temperature
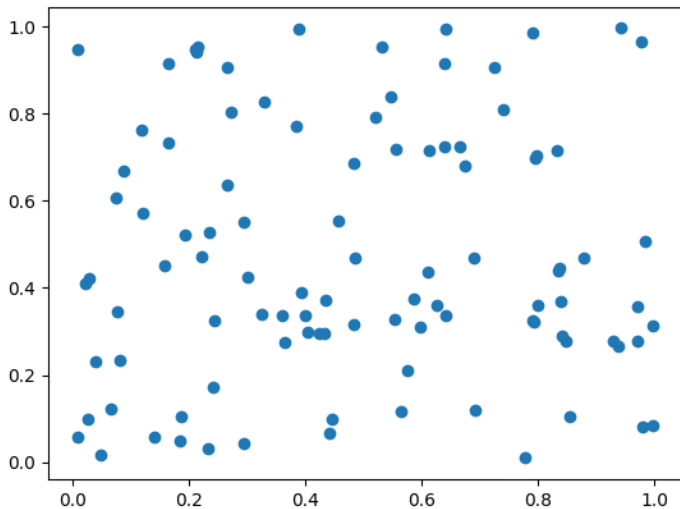- Concentration of pollutants
- Flow rate
- water speed

The experiments/observations/sample can be:
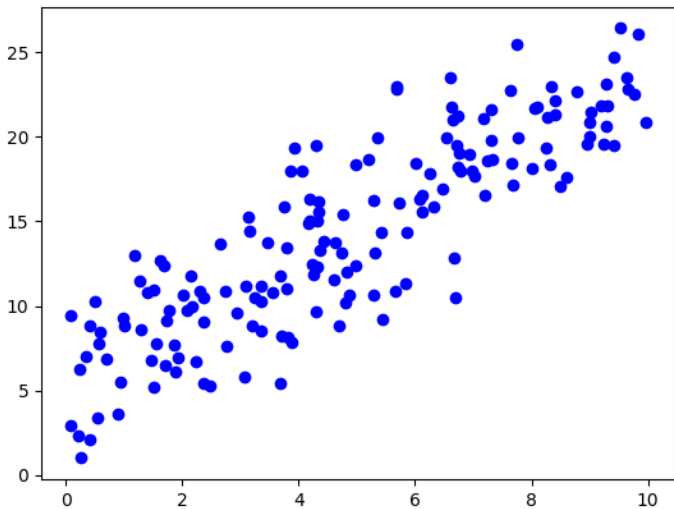
- Po
- Danube
- Rio delle Amazzoni
- Sjoa
- Atna

```python
import numpy as np
import matplotlib.pyplot as plt

# Generate some sample data
data = np.random.rand(100, 2)   # 100 data points with 2 features

plt.scatter(data[:, 0], data[:, 1])
plt.show()
```

```python
import numpy as np
import matplotlib.pyplot as plt

def generate_linear_data(n_random_points, noise=16):
    x = np.random.rand(n_random_points) * 10

    # Make 'perfect' data
    true_slope, true_intercept = 2, 5
    y = true_slope * x + true_intercept

    # Add noise
    y += np.random.randn(n_random_points)*noise

    return x, y, true_slope, true_intercept

# Use the function to generate data
x, y, true_slope, true_intercept = generate_linear_data(
        n_random_points=166,
        noise=3)

# Plot all
plt.scatter(x, y, color='blue', label='Data Points')
plt.show()
```
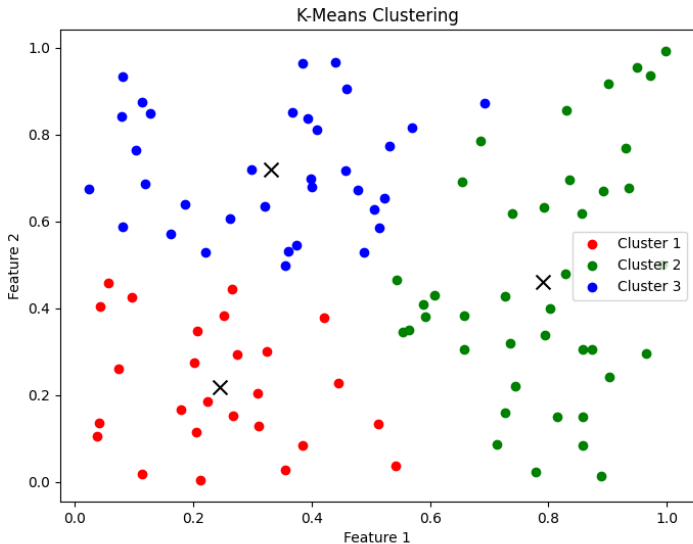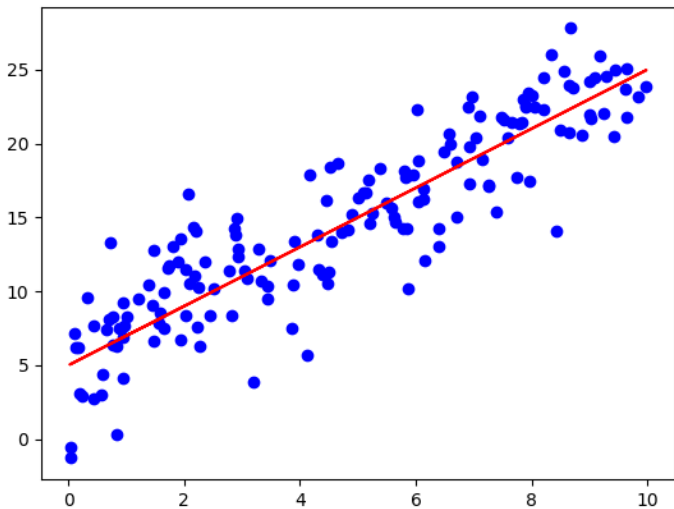
K-Means Clustering

This is why we focus so much on the data type.

The data properties dictate what statistical model can be adopted.

An statistical model has leverages our understanding of the data structure to improve its **predictions** (inference).

The numerical recipe that we used to generate the data is defined the **truth**

Psychology or data science?

Most Machine learning tools are aimed to find the truth. In most cases, we are happy to not find lies.

Unsupervised learning, a term that resonates with the autonomy of machine intelligence, operates on the principle of identifying patterns and structures in datasets without labelled responses.

This branch of machine learning is distinguished by its lack of explicit guidance, where algorithms are tasked with uncovering hidden structures from unlabeled data.

The most common clustering strategies are :

- filtering
- clustering
- dimensionality reduction
- association learning

It is a bit of a holy grail: a computer that finds patterns without guidance. (Yes, it doesn't work, most of the time)
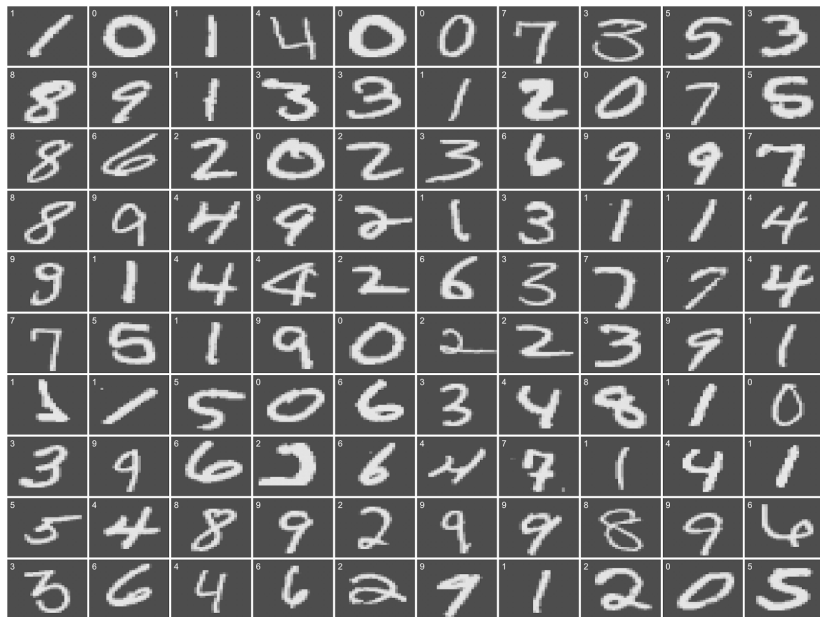
Still, it has been shown efficient for:

- Computer vision
- Anomaly detection
- Exploratory data analysis

## Main challenge
The right result is quite undefined, Uncertain goal.

We will demonstrate it with a famous problem.

A less popular type of machine learning problem is when labels are assigned to groups of instances.

The group of instances is called **bag**.

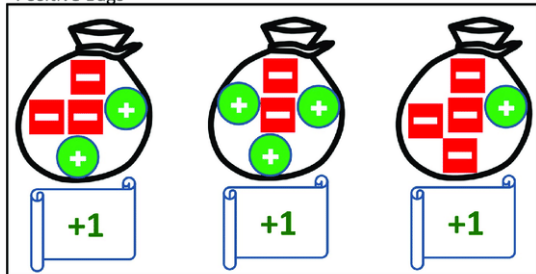The question is, what is the level of a previously unforeseen bag?

This data structure and question type request a hybrid treatment between supervised and supervised learning.
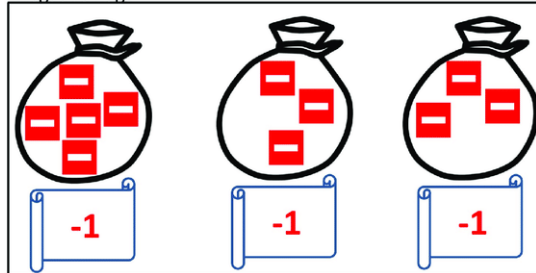
**Multiple instance learing**

Multiple instances are needed to learn (quite clear name)

Finally, there is a further approach.

## Reinforcement learning (RL)

It aims to train an intelligent agent to take actions in a dynamic environment in order to maximise the cumulative reward.

It learns from outcomes and decides which action to take next. After each action, the algorithm receives feedback that helps it determine whether the choice it made was correct, neutral or incorrect.

It is a self-teaching system that essentially learns by trial and error.

It is a dependable tool for automated decision making.