

# Fundamentals of Machine learning for and with engineering applications

Enrico Riccardi<sup>1</sup>

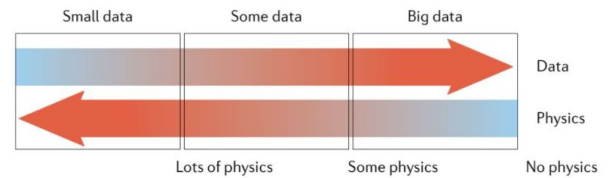
Department of Mathematics and Physics, University of Stavanger (UiS).<sup>1</sup>

Sep 2, 2025



© 2025, Enrico Riccardi. Released under CC Attribution 4.0 license

## Data vs Physics



## Uncertainty

- Def 1: Not knowing if an event is true or false. (Useful)
- Def 2: Things that cannot be measured. (Not useful)

Probability is how Uncertainty is quantified!

- Clarity test
- Assign a number between 0 and 1 to our degree of belief
- Error definition

Sentence also good for fortune cookies

Uncertainty is the only certainty

## Uncertainty and Probability

### Random quotes

- Probability: there is not science more worthy in our contemplations nor a more useful one for admission to our system of public education
- The theory of probabilities is at the bottom of nothing but common sense reduced to calculus.

### What is Statistics

**Clarity test.** Beer drinker?

Rain in Stavanger?

## Data properties

1 D

logs

2 D: maps

Quite limited but great for visualization

3 D

3d maps, seismic cubes. More informative, mostly ok in digital formats.

4 D

Trajectories

x D

Data realm

## Types of data

- Categorical / Nominal (classes)
- Categorical / Ordinal
- Continuous / Interval (e.g. Celsius)
- Continuous / Ratio
- Discrete: binned/grouped data
- Hard data: direct measurements
- Soft data: indirect measurements, very uncertain
- Primary data: variable(s) of interest
- Secondary data: descriptors
- Collective variables
- Latent variables

## Descriptive and Predictive statistics

### Estimation

- Process of obtaining the best value or range of a property in an unsampled location
- Local accuracy takes precedence over global spatial variability
- Not appropriate for forecasting

### Inference

- Predict unseen samples given assumptions about the population
- Test with a pre-trained model (ML definition)
- Generality versus Accuracy

## Variables and Features, Labels and Instances

### Population

Exhaustive, finite list of properties of interest over area of interest.

Generally the entire population is not accessible

### Samples/experiments/instances

The set of values and location that have been measured.

How many experiments are needed?

### Sampling distribution -(IID ?)

identical and independent distribution The set of values and location that have been measured.

### Features

The values to be measured for each sample/experiment/instance.

How many features are needed?

## Variables and Features, Labels and Instances

Predictors = input variables,  $X_1, \dots, X_M$

Response = output variables

### Error

Deviation from ... exact value (or expected value, mean value, trend...?)

Errors without definitions are just numbers.

### Error

Values, without error, are just number!!

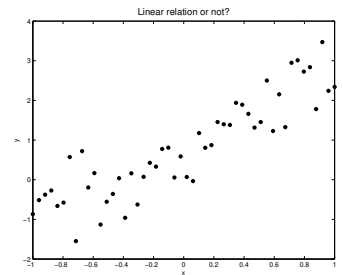
### Predictor and Response Features

Given a model  $Y = f(X_1, \dots, X_M) + e$

!Here and error! But is it even an error?

## Finding a suitable model

Soft modeling is in most cases based on **multivariate statistical methods**. Many of these methods may be viewed as sophisticated ways of performing curve fitting to data.



What would be the best model?

- Straight?:  $y(x) = ax + b$
- Parabolic?:  $y(x) = ax^2 + bx + c$
- Trigonometric?:  $y(x) = a\sin(x) + b\cos(x)$

## Uncertainty Modeling

**Given a model**, Generate multiple simulation to represent uncertainty

- Realizations: for the same input parameters, different random numbers.
- Scenarios: different input parameters.

**Sampling representative.**

### Random sampling

Each item of the population has an equal chance of being chosen.

- Very expensive
- Mostly not interesting
- Gives some global properties

### Bias sampling

Selection of data is (arbitrarily) distorted

- Sample probability bias has to be corrected for
- Might not capture the global picture
- It might distort the system under study -> false results

## Cognitive biases

- anchoring: The first bits are over-considered
- availability: over-estimating the importance of info
- bandwagon: P increases with the number of people holding a belief
- blind spot: not seen biases
- choice supporting: commitment/decision dependent
- clustering illusion: seeing patterns in random events
- confirmation bias
- conservatism bias
- Recency bias
- Supervision bias
- Many many more!

**Bias DO NOT cancel out! They sum up (or multiply?)**

## Simulations

Process of obtaining one or more values of a property

- Improved Global accuracy
- Better property distributions

### Why simulations then?

- We need to capture the full distribution of properties, extremes matter!
- We need more realistic models.

### Why not?

- High dimensionality level
- Computationally expensive
- Convergence limitations
- Constitutive equations need to be rather accurate.

## Representation

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 can 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).

## Data properties

- 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

## Data

DATA



SORTED



ARRANGED



PRESENTED VISUALLY



EXPLAINED WITH A STORY



## Comparability

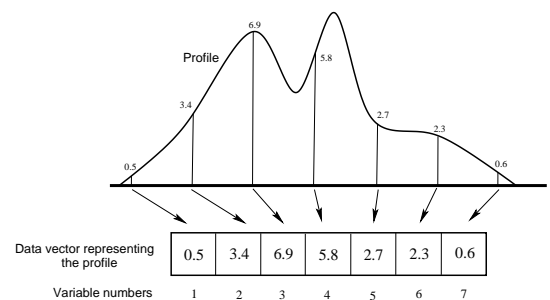
Same meaning **representations** for different objects (inputs).

### Discussion point!

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

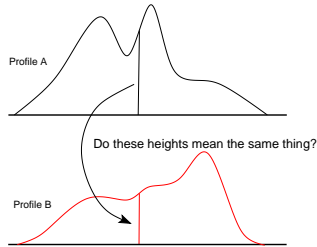
## Sampling point representation (SPR)

- 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.



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

## Data structures

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

## Workflow

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

## Spatial and Temporal Data

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.

## Actual 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

## Nomenclature Reminder

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.

## A quick example

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