

# Fundamentals of Machine learning for and with engineering applications

Enrico Riccardi<sup>1</sup>

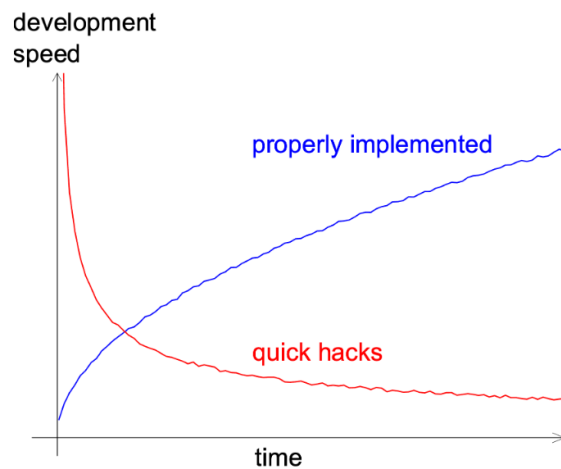
Department of Mathematics and Physics, University of Stavanger (UiS).<sup>1</sup>

Sep 10, 2025



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

## Flexible



## Developing approaches

Different code editors are available to interpret python language.

- jupyter notebooks are mostly dedicated to learning (Markdown)
- ipython is for interactive coding (similar to R, Matlab, etc)
- python packages (.py) developing suites (debug possibilities and git integration)

## Introducing code standards

When developing code, there are **guidelines** and best practices aimed at improving the **quality, readability, and maintainability** of a code.

There are different levels of coding quality, mostly depending on the code intended usage (and developer skills).

- Private codes can be whatever (Cpt. Obvious)
- Public packages shall use a 'Golden code standards' such to be used and eventually supported by communities.

## Community level standards

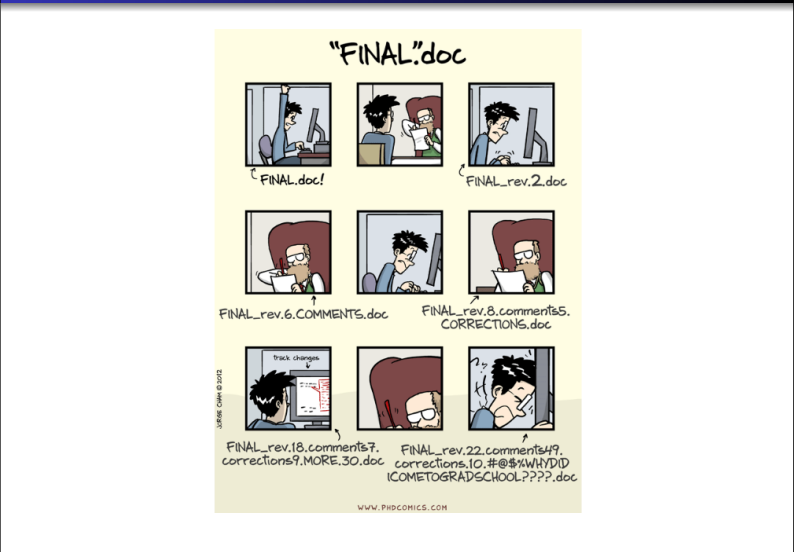
Principle of 'clean coding':

- 1 Readability and Clarity: A good code shall be possible to read as when reading a book
- 2 Structure and object oriented: A code shall be composed by objects, each of them connected in the less redundant way possible.
- 3 Consistency and Style: Variable naming, function naming and classes naming has to be consistent.
- 4 Documentation: Each file, each function and each class shall contain the relative description of its aim and its usage
- 5 Maintainability: Code dependencies have to be stated and consistently defined and updated, such that a suitable environment can be developed at any point in time.

## Community level standard

- 1 Testing: Unit testing shall cover the majority of the code
- 2 Error Handling: Each error shall be captured and properly identified.
- 3 Examples and benchmarks: Users shall be able to execute minimal examples of the code for computational checks.
- 4 Performance Optimization: Libraries shall be able to use the available computational power in the machine (e.g. GPU-CUDA)

## Version control



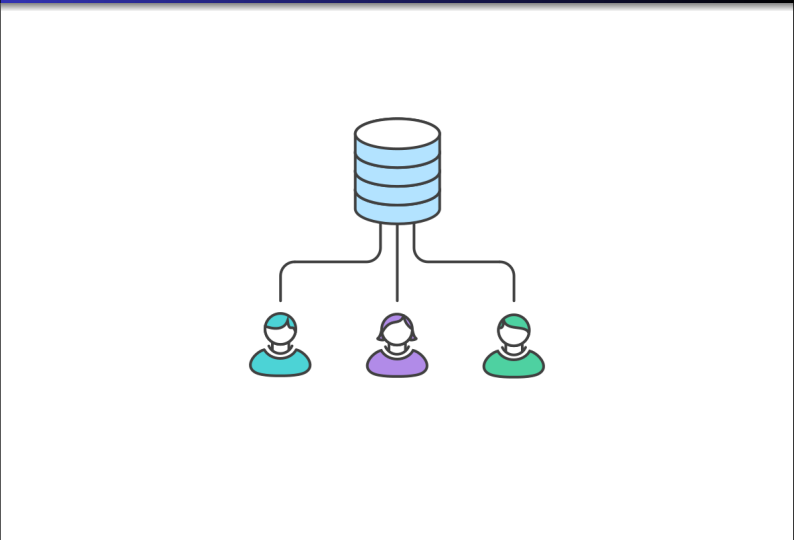
## Git

Git is a distributed version control system that tracks changes in any set of computer files, usually used for coordinating work among programmers who are collaboratively developing source code during software development. Its goals include speed, data integrity, and support for distributed, non-linear workflows (thousands of parallel branches running on different computers). [Wiki]

Let's try to be more accessible.

Git is a computer program/tool to save and download files on a hosting server (e.g. GitHub and GitLab).

Centralized workflow



## A distributed version control system

- Git facilitates users to track the various versions of files. It is not a necessary tool, but it can be very very helpful. Generally, the time spent to learn its syntax is well paid off

(do you remember to save some file like  
manuscript\_draft\_v4.02\_final\_definitive\_forreal\_lastcomments\_editedbyER\_submittedVe  
Exactly! Imagine to do that for a repository of files...)

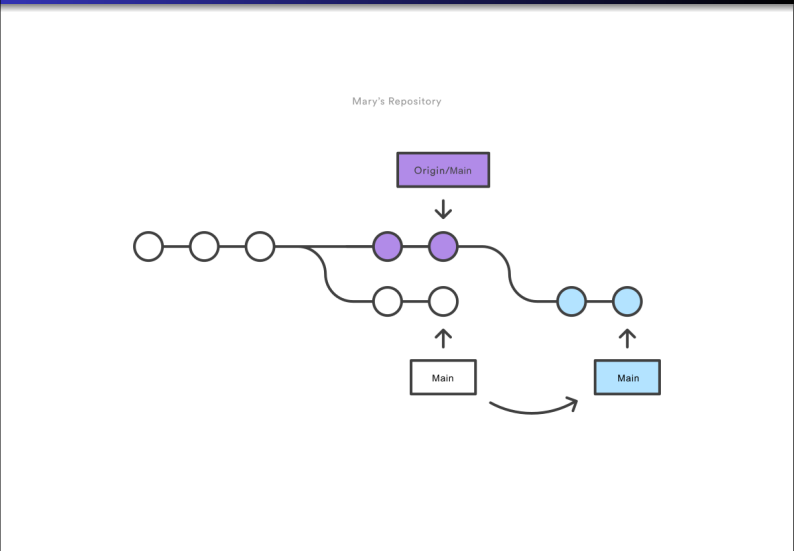
- It permits to save and share the intermediate stages of a work in progress (which software is complete and always up to date?) in an accessible, consistent and structured way, allowing an effective version tracking. It allows retrieval of previous working versions, limiting the risk to overwrite useful files.

## What is git actually for

The tool is particularly useful for programmers working in teams or in projects whose outcomes can be used by others.

- Git helps to co-develop a code, test its functions and the compatibility of the various code sections.
- A long list of further possibilities became possible by git.
- Different software integration on development platforms, based on git, will help you to develop and co-develop your code.
- The platform GitLab and GitHub have a large set of functionalities to further support code documentation and public releases.
- Files can be disclosed to the public, becoming a great integration of your CV, showing what you are able to do in an open and accessible way.

How does it work -in short-



## Why should I care?

As the open libraries are exploding in numbers, you might need some criteria to assert the reliability of a project.

Unit test driven development!

That is taking full advantage of python object oriented structure.

### Community

Good project are not only used by communities, but also **supported**

Git allows the development of projects without a clear lead. Community engagement is generally a desirable target to help develop to directly integrate feedbacks by users (and fix bugs).

## Statistics (recaps)

### Definition

Statistics is the science of acquiring and utilizing data

- It comprises tools for data collection, summarization, and interpretation.
- The aim is identifying the underlying structure, trends, and relationships inherent in the data.
- Is it all statistics then? Yes.
- **Numbers to data, data to information**

## Data properties

Before we talk about machine learning, we need to refresh some terminology.

### Population

The universe of all possible outcomes and events.

### Sample

A finite subset extracted from the population.

### Exhaustivity

The samples covered the population spectra.

### Representativity

The population is properly described by the samples.

## Big data

We speak of big data when dataset are very large: i.e. many instances and features. Models have thus a large set of parameters (and often no one has a clue anymore of what is going on).

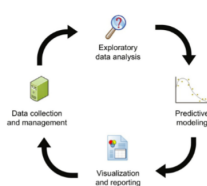
- Volume of data
- Variety different types of data sources with different length and scale.
- Frequency of data generation

## Sampling

Samples shall have no bias (to be randomly selected). If not, the bias has to be corrected for.

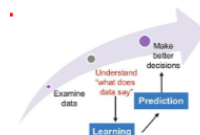
### Cycle of data

- 1 Data is collected
- 2 Checked upon
- 3 Some modelling
- 4 Analysis and visualization



## Data quality

- 1 Data has to be acquired and integrated
- 2 Data are passed to a quality analysis and control
- 3 Data cleaning, consistency check. Most of time goes here



## Preliminary Modeling

### Main tasks:

- 1 Hunt for redundancy
- 2 Reduce dimensionality
- 3 AnOmAlles removal

- Descriptive modeling (unsupervised learning)
- Predictive modeling (supervised learning)
- The model can be used to guide data acquisition (risky!)

## Visualization and reporting

- The data has to be condensed into a visualization to provide input for decisions.
- Depending on the goal, very very different visualizations are possible.
- Use a model to indicate what is undersampled or oversampled.

### Summarizing and visualizing data as a starting point for more analysis later on.

- Computing summary statistics (e.g. means and variance)
- Determining conditional probabilities of cause+effect relationships
- Calculating correlation and rank correlation coefficient between two variables
- Visualizing univariate, bivariate and multivariate data
- ...

## Exploratory data analysis

### Summarizing and visualizing data as a starting point for more analysis later on.

- ...
- Estimating probability coverage levels for different distributions
- Analyzing behaviour of normal distributions
- Calculating confidence interval and sampling distribution for the mean
- Testing for significance of difference in means
- Comparing two different distributions for statistical equivalence
- Developing a nonparametric regression model from given data
- Reducing data dimensionality
- Grouping data

## Random Variables

- A random variable is a real valued function that assigns a value to each outcome in the sample space
- A random variable (RV) can be either discrete or continuous
  - Discrete RV
  - Continuous RV
- The probability mass function (PMF),  $P$ , of a discrete RV,  $X$ , denotes the probability that the RV is equal to a specified value,  $a$ .

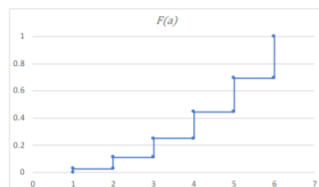
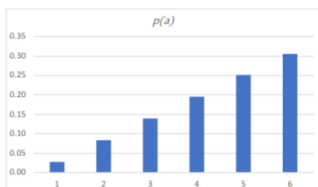
$$p(a) = p(X = a)$$

- The cumulative distribution function (CDF),  $F$ , denotes the sum

$$F(a) = P(X \leq a) = \sum_0^a f(x) dx$$

## Random Variables

$a$	1	2	3	4	5	6
$p(a)$	1/36	3/36	5/36	7/36	9/36	11/36
$F(a)$	1/36	4/36	9/36	16/36	25/36	1



## Sampling

- What are the effective sampling strategies? (Wind turbine example)
- Solar Panels to determine the efficiency of the source (Usage patterns, energy production forecast)
- Drilling (penetration rate)
- Corrosion extension
- Concrete Rigidity
- Experimental design!

Wind turbine example

Task (due 15.9.2025)

Turbine	Height	X	Y	Wind Speed	Air Density	Temperature	Power Output	Rotor Diameter	Hub Height	Air Pressure	Turbulence Intensity
WT-1	80	752.1	3945	7.5	1.225	15	1500	82	80	1013	0.1
WT-1	80	752.2	3945	8	1.223	15	1600	82	80	1012	0.12
WT-1	80	752.3	3945	7.8	1.224	16	1550	82	80	1013	0.11
WT-2	90	753.5	3946	6.5	1.226	14	1400	85	90	1012	0.15
WT-2	90	753.6	3946	7	1.225	14	1500	85	90	1011	0.13
WT-2	90	753.7	3946	7.2	1.227	14	1520	85	90	1012	0.14

- Task 1: make a histogram from a 2d random distribution
- Task 2: make a 2d heat map from a 2d random distribution
- Task 3: make a histogram for the source data you selected
- Task 4: convert the histogram into a discrete PMF
- Task 5: calculate the cumulative for each feature