

# Fundamentals of Machine learning for and with engineering applications

Enrico Riccardi<sup>1</sup>

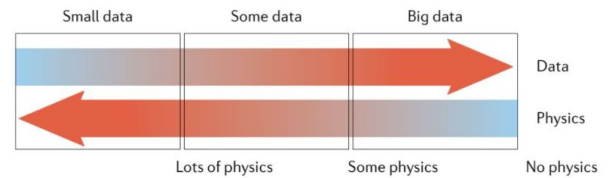
Department of Mathematics and Physics, University of Stavanger (UiS).<sup>1</sup>

Sep 1, 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!!

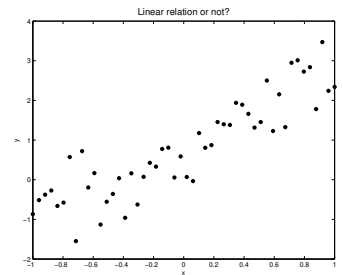
### Predictor and Response Features

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?

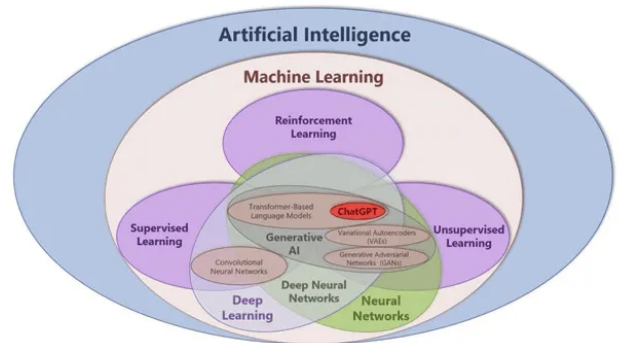
- We need to capture the full distribution of properties, extremes matter!
- We need more realistic models.

Why not?

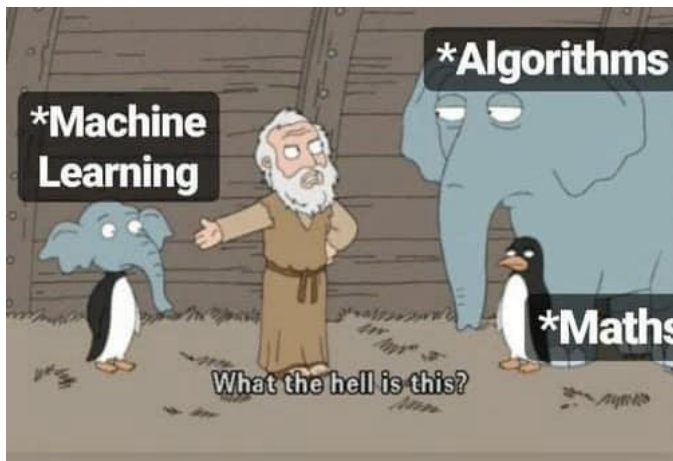
- High dimensionality level
- Computationally expensive
- Convergence limitations
- Constitutive equations need to be rather accurate.

## Statistics, Machine learning or Artificial intelligence?

What is the main difference between the three fields?



## How Machine Learning Started?



## Statistics

Let's start from the definition

- Statistics (origin "description of a state/country") is the discipline that concerns the collection, organization, analysis, interpretation, and presentation of data.
- It is conventional to begin with a statistical population or a statistical model to be studied. Populations can be diverse groups of people or objects such as "all people living in a country" or "every atom composing a crystal".
- Statistics deals with every aspect of data, including the planning of data collection in terms of the design of surveys and experiments.[Wikipedia]

## Machine learning

Definitions:

- Machine learning is a branch of artificial intelligence (AI) and computer science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy. [IBM]
- Machine learning (ML) is a field of study in artificial intelligence concerned with the development and study of statistical algorithms that can learn from data and generalize to unseen data, and thus perform tasks without explicit instructions. [WIKI]
- Machine learning is a subfield of artificial intelligence that uses algorithms trained on data sets to create models that enable machines to perform tasks that would otherwise only be possible for humans, such as categorizing images, analyzing data, or predicting price fluctuations. [Coursera]

## Machine learning

One technical definition

Machine learning is a set of computer based statistical approaches that aim to minimise the loss function to maximise inference accuracy. [Enrico, 5.2.2024]

**The loss function** is the actual engine in machine learning.

Loss function

It quantifies the difference between the predicted outputs of a machine learning algorithm and the actual target values.

## Artificial intelligences

And more definitions:

- Artificial intelligence is the intelligence of machines or software, as opposed to the intelligence of humans or other animals. It is a field of study in computer science that develops and studies intelligent machines. [WIKI]
- Artificial intelligence (AI) is the theory and development of computer systems capable of performing tasks that historically required human intelligence, such as recognizing speech, making decisions, and identifying patterns [Coursera]
- It is the science and engineering of making intelligent machines, especially intelligent computer programs. It is related to the similar task of using computers to understand human intelligence, but AI does not have to confine itself to methods that are biologically observable. [IBM]

## Data

DATA



SORTED



ARRANGED



PRESENTED  
VISUALLY



EXPLAINED  
WITH A STORY



## 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).

## Comparability

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

Discussion point!

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

## Data properties

- All starts from data: what are data-properties?
- Are there such things as good data and bad data?

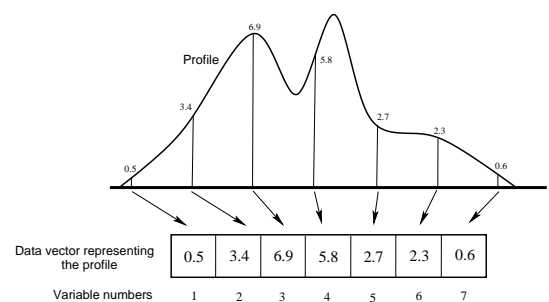
Life lesson (or exam question, same thing ;))

- Data **DO NOT** always have value.

- TRASH in TRASH out

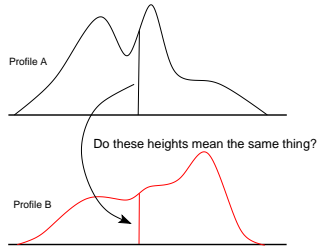
## Sampling point representation (SPR)

- 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

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

- the geo context of the data
- the spatial and time dependent relationship between data
- 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

## Version control



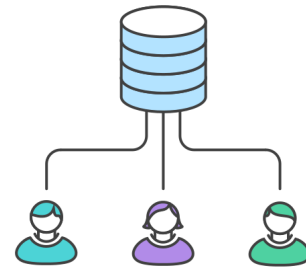
## Git

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

- 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_submittedV6` Exactly! Imagine to do that for a repository of files...)

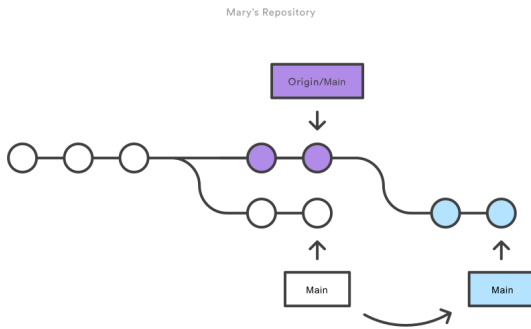
- It permits to save and share the intermediate stages of a work in progress (which software is complete and always up to date?) in an accessible, consistent and structured way, allowing an effective version tracking. It allows retrieval of previous working versions, limiting the risk to overwrite useful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

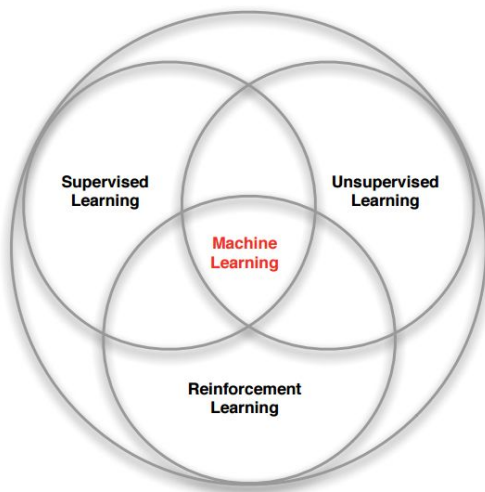
That is taking full advantage of python object oriented structure.

### Community

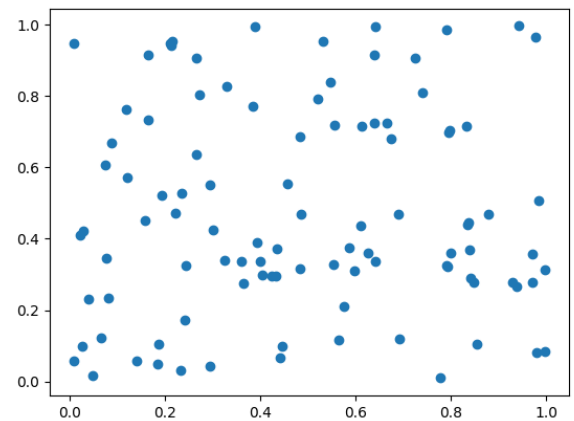
Good project are not only used by communities, but also **supported**

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).

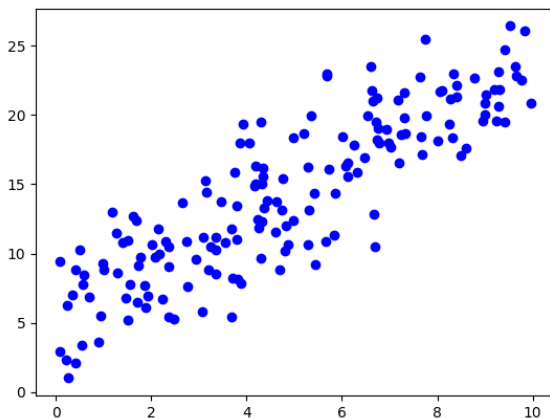
## Families of Machine learning



## What can we do with that?



## What about in this case?



## Python Source code 1

```
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 Source code 2

```
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

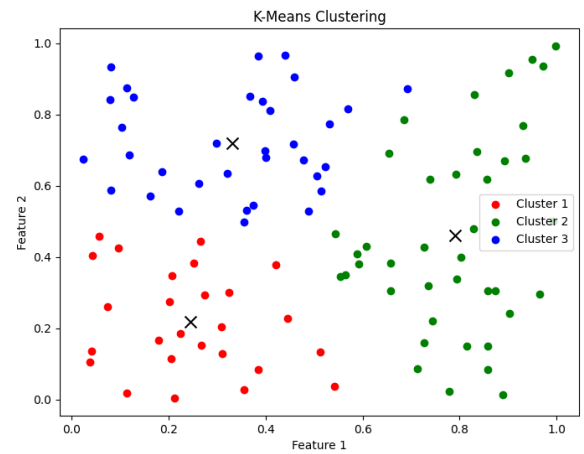
    # Add noise
    y += np.random.randn(n_random_points)*noise

    return x, y, true_slope, true_intercept

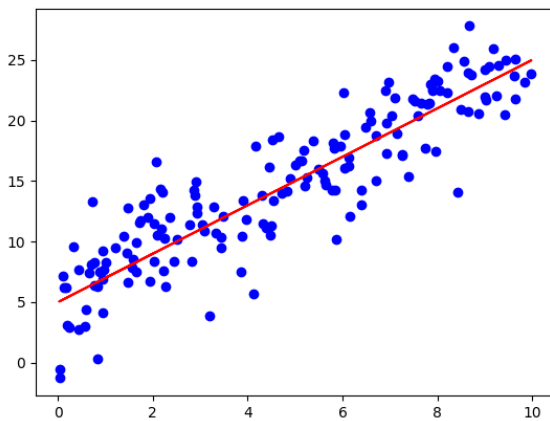
# Use the function to generate data
x, y, true_slope, true_intercept = generate_linear_data(
    n_random_points=166,
    noise=3)

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

## Unsupervised learning



## Supervised learning



## The data decides

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

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

## Application of unsupervised 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

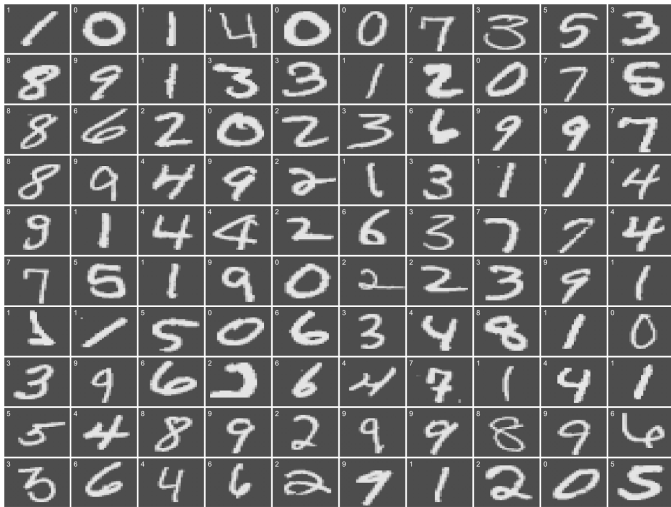
Main challenge

The right result is quite undefined, Uncertain goal.

We will demonstrate it with a famous problem.



## Uncertain goal



## Weak supervised learning

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

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

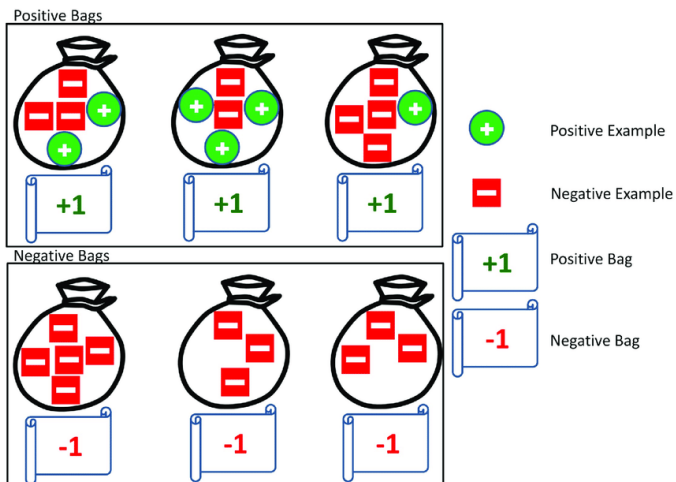
The question is, what is the level of a previously unforeseen bag?

This data structure and question type request a hybrid treatment between supervised and supervised learning.

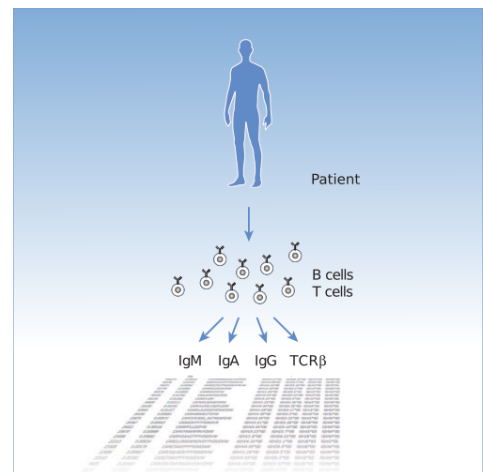
### Multiple instance learning

Multiple instances are needed to learn (quite clear name)

## Weak Supervised learning



## Weak Supervised learning



## Reinforcement learning

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.