

Machine Learning Vs Human Brain



Goals of today's event

Motivations:

- Scientific event with a subject transverse to the 4 departments
- Clear trend in our fields since already many years
- Already several use of ML in our researches

Goals:

- Inform about activities done @ IPHC
- Starting point to sharing expertise?
- Network ?
- Common interests – Structuration ?



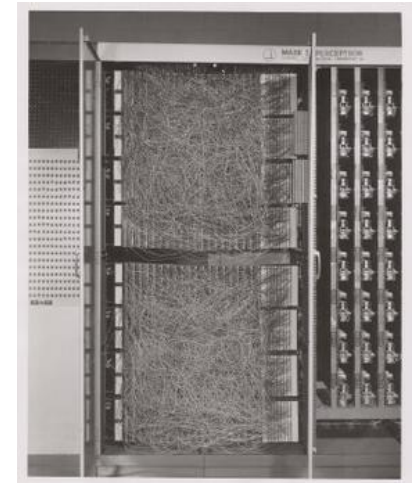
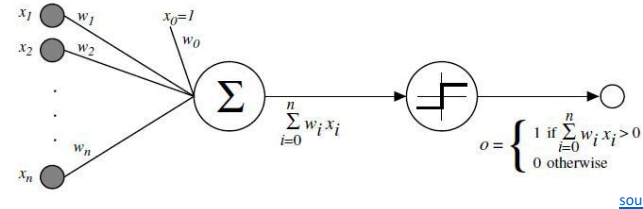
Morning's session

To be discussed during the round table

ML: an expanding field

- **Machine learning is not new**

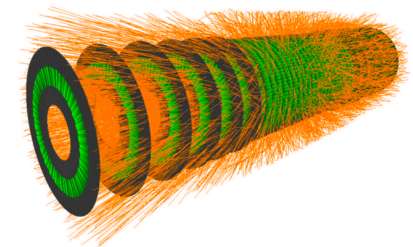
- First ideas already in 50's
- History of Neural Network
 - First paper (1943)
 - First algo still in use (Perceptron): 1958
 - First annual meeting "NN for computing" (1985)
 - AlphaGo (Google DeepMind) win a professional player (2015) ...



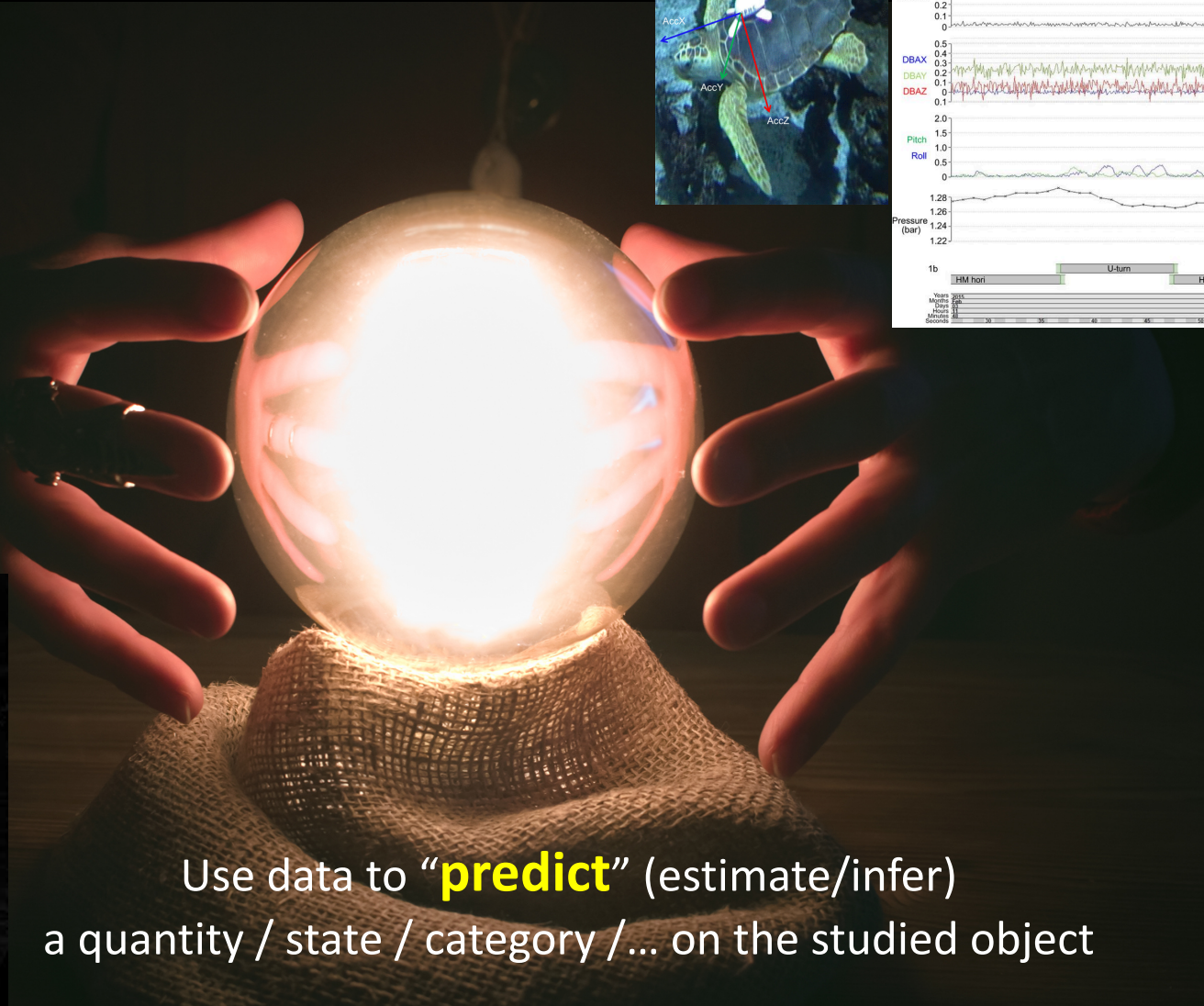
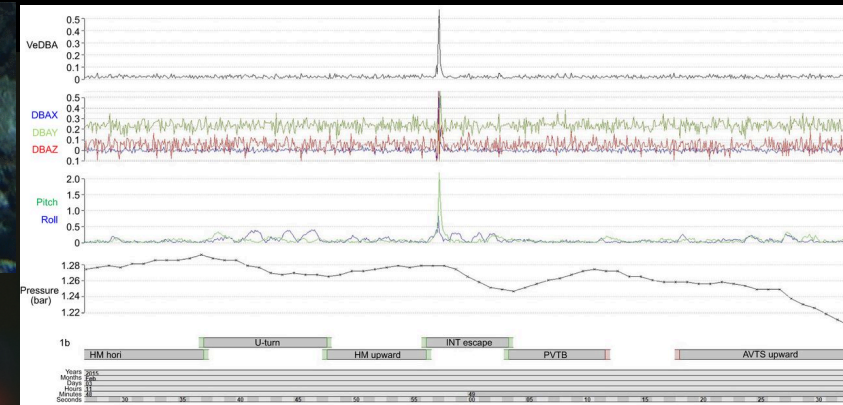
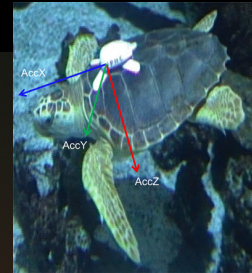
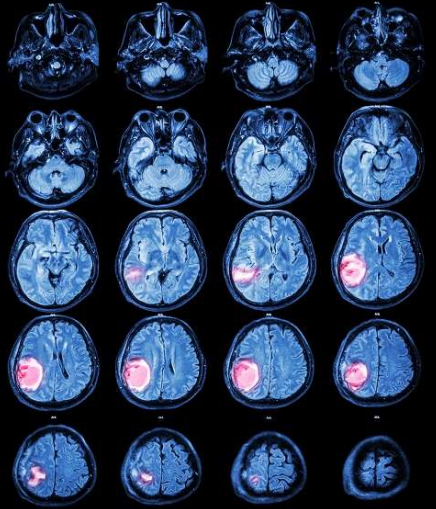
'Mark I Perceptron at the Cornell Aeronautical Laboratory', hardware implementation of the first Perceptron

- **Expansion of the ML field**

- Boost from **big data**
- Large increase of **computing** resources (massive parallelisation)
- ML implementation eased by user-friendly (open source) **packages**
- Variety of applications (interests in tech companies and industry)
- Dynamic field: many ML **contests**, ...



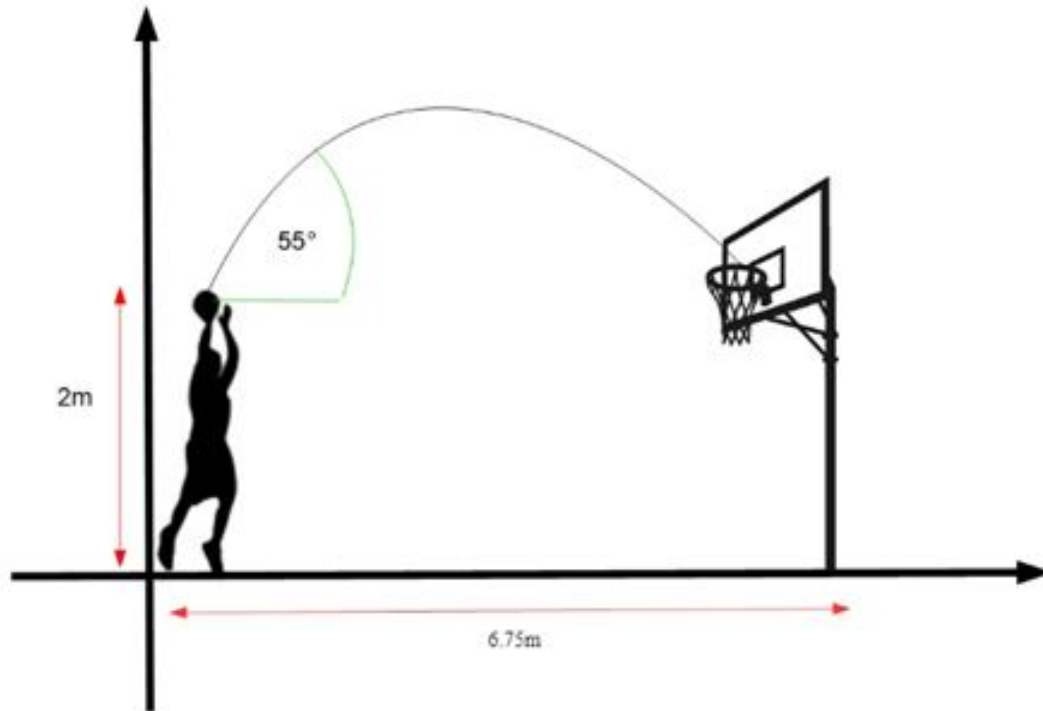
Machine Learning: *an e-crystal ball*



Use data to **“predict”** (estimate/infer)
a quantity / state / category / ... on the studied object

Machine Learning: *an e-crystal ball*

Do we need to know gravitational law to play basketball ?



Children learn by experiment how to solve this problem without determining explicitly any physics law !

Machine Learning goals



ML goals:

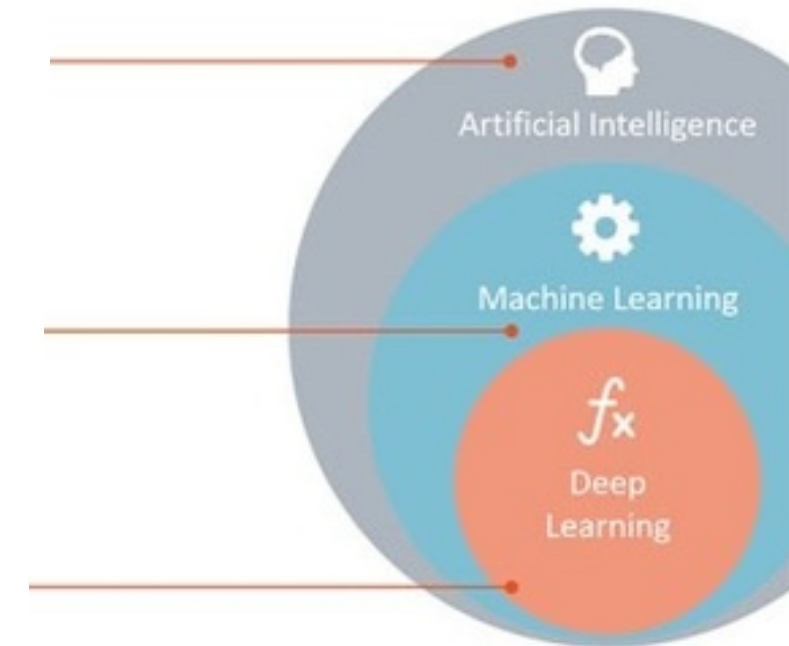
- Learning
- Finding patterns & trends
- Improving decision making
- Generating data
- Improving/optimizing “tasks”
- ...

A bit of clarification

Artificial Intelligence

Any technique which enables computers to mimic human behaviour. It uses ML but not only. Creates applications that react, adapt,

Ex: self-drive car, chatbot, ...



A bit of clarification

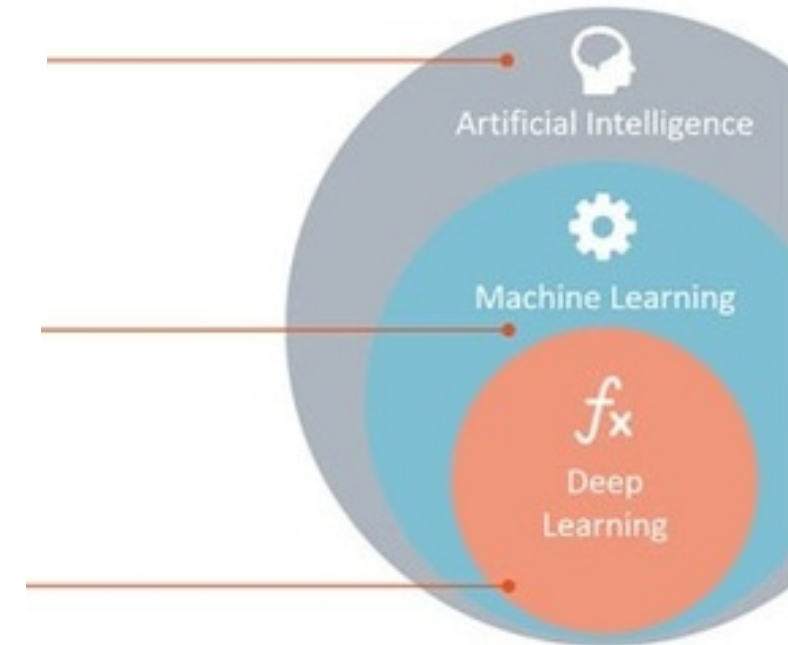
Artificial Intelligence

Any technique which enables computers to mimic human behaviour. It uses ML but not only. Creates applications that react, adapt,

Ex: self-drive car, chatbot, ...

Machine Learning

Algorithms able to find patterns/features **based on data** using statistics with the aim to perform **predictions**



A bit of clarification

Artificial Intelligence

Any technique which enables computers to mimic human behaviour. It uses ML but not only. Creates applications that react, adapt,

Ex: self-drive car, chatbot, ...

Machine Learning

Algorithms able to find patterns/features **based on data** using statistics with the aim to perform **predictions**

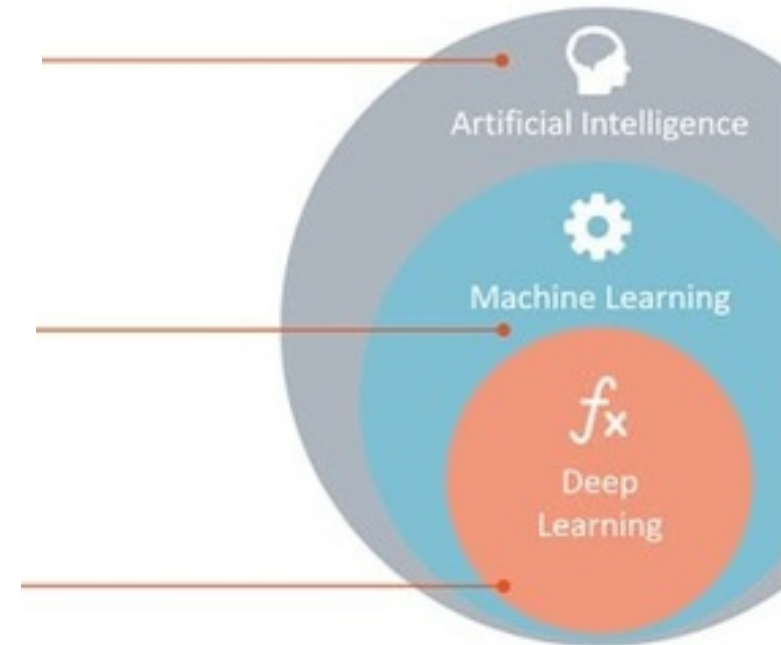
Deep Learning

Subset of ML which make use of multi-layer neural networks (Artificial NN, Convolutional NN, Recurrent NN, ...)

“**deep**” → number of layers & neurons

→ many parameters to be estimated

→ need large datasets and long training period



A bit of clarification

Maths/Statistics

Programming

Biological inspired
rules

Artificial Intelligence

Any technique which enables computers to mimic human behaviour. It uses ML but not only. Creates applications that react, adapt,

Ex: self-drive car, chatbot, ...

Machine Learning

Algorithms able to find patterns/features **based on data** using statistics with the aim to perform **predictions**

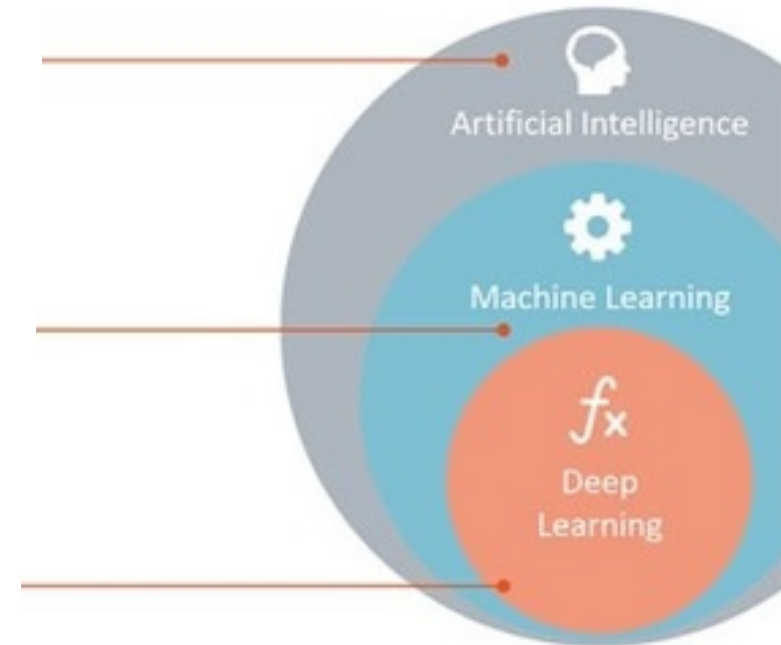
Deep Learning

Subset of ML which make use of multi-layer neural networks (Artificial NN, Convolutional NN, Recurrent NN, ...)

“**deep**” → number of layers & neurons

→ many parameters to be estimated

→ need large datasets and long training period



A bit of clarification

Maths/Statistics

Programming

Biological inspired rules

Big data

- Does not imply the use of ML
- ML complexity & performances depends on the size of the data used to learn

Artificial Intelligence

Any technique which enables computers to mimic human behaviour. It uses ML but not only. Creates applications that react, adapt,
Ex: self-drive car, chatbot, ...

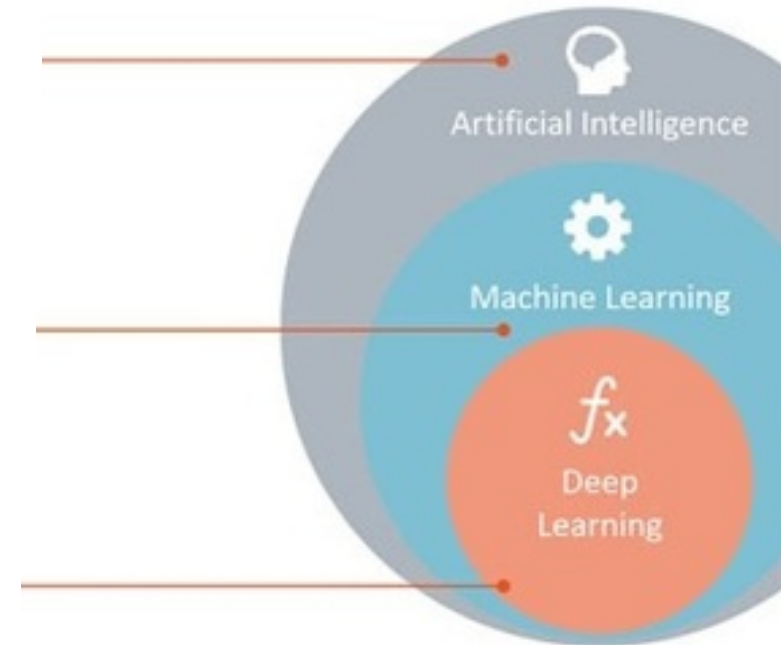
Machine Learning

Algorithms able to find patterns/features based on data using statistics with the aim to perform predictions

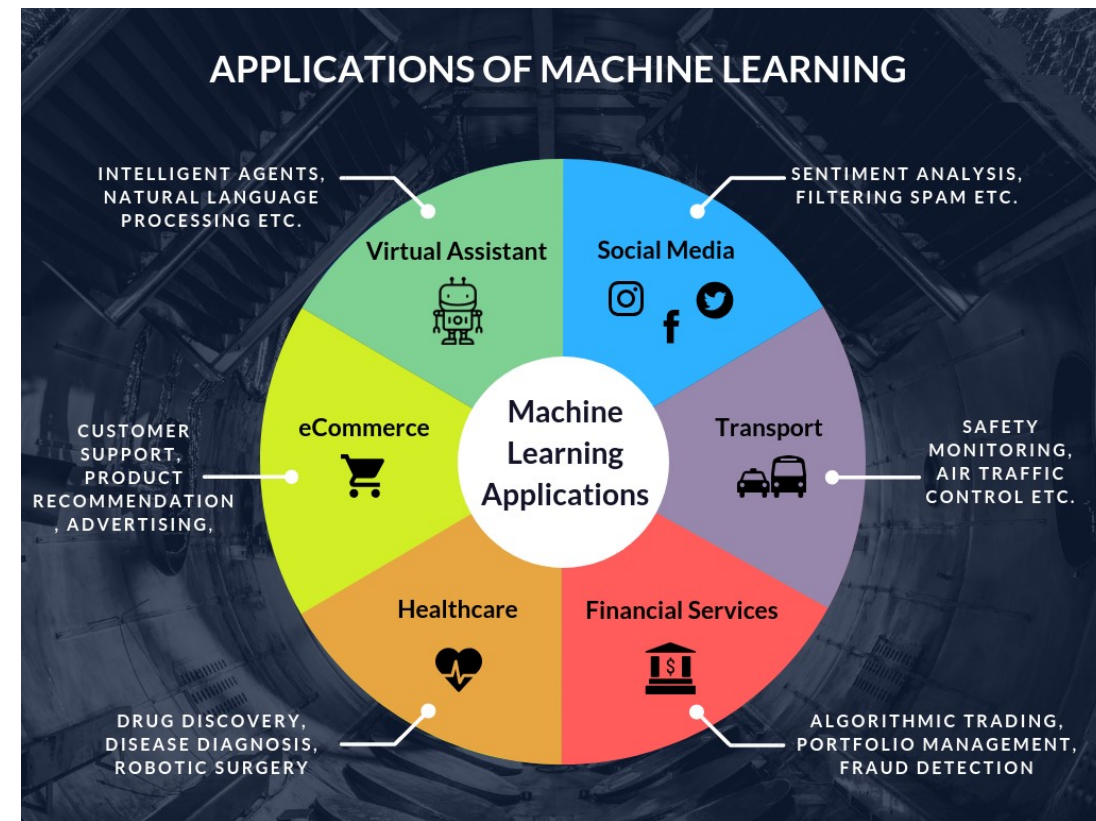
Deep Learning

Subset of ML which make use of multi-layer neural networks (Artificial NN, Convolutional NN, Recurrent NN, ...)

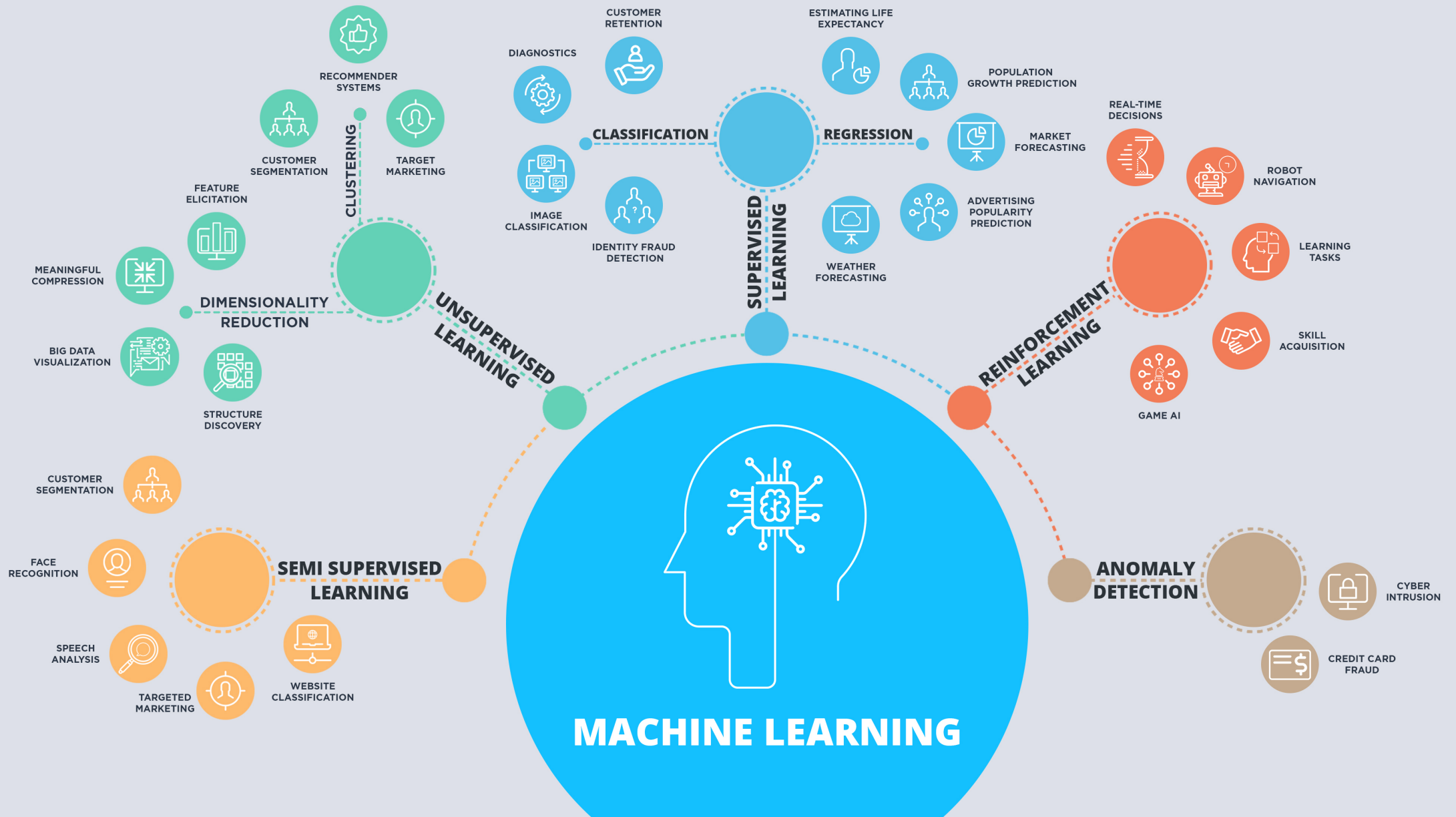
“deep” → number of layers & neurons
→ many parameters to be estimated
→ need large datasets and long training period



ML is everywhere in everyday's life



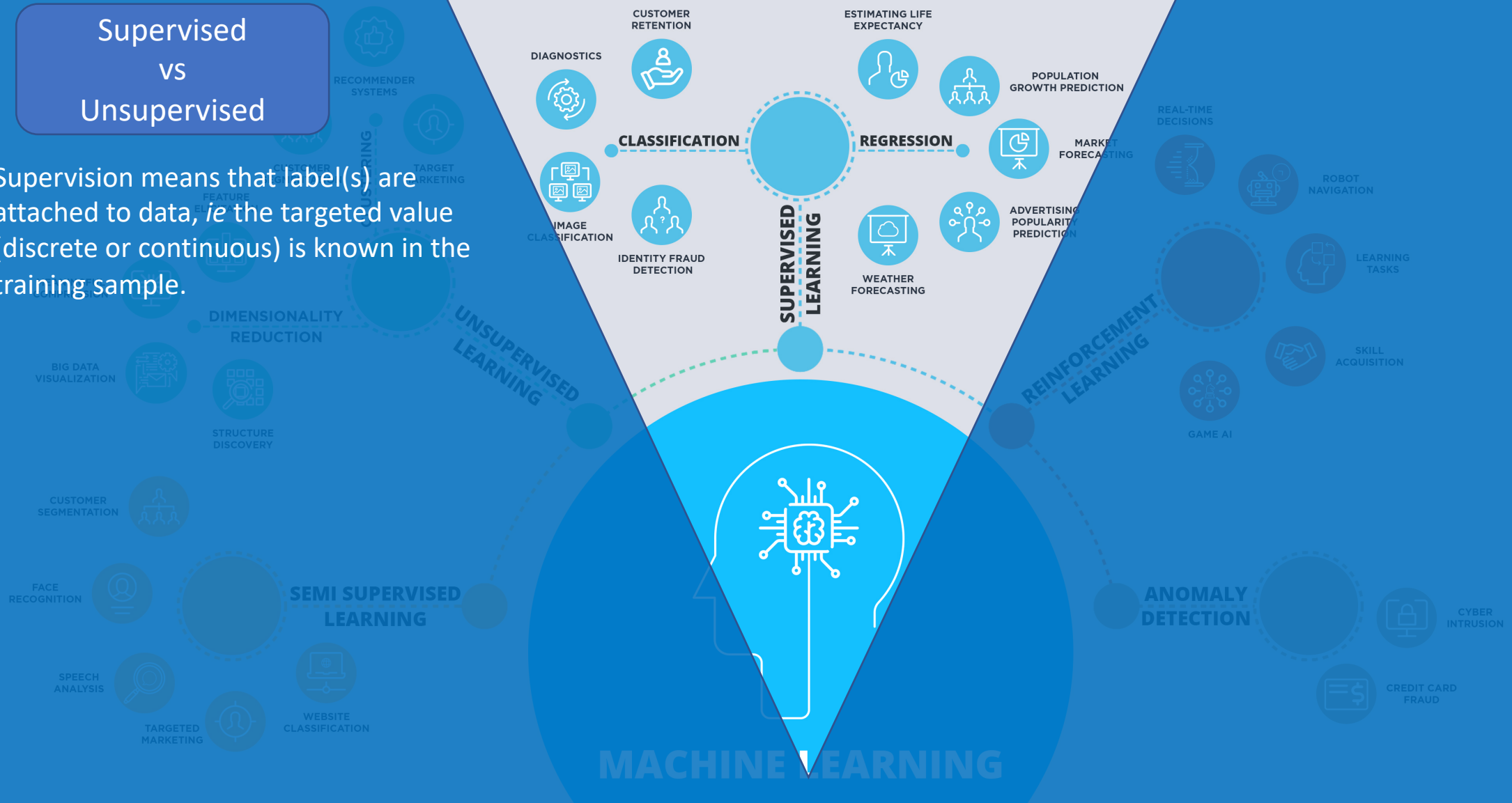
Class of problems solved by ML algorithms



Class of problems solved by ML algorithms

Supervised
vs
Unsupervised

Supervision means that label(s) are attached to data, i.e. the targeted value (discrete or continuous) is known in the training sample.



Class of problems solved by ML algorithms

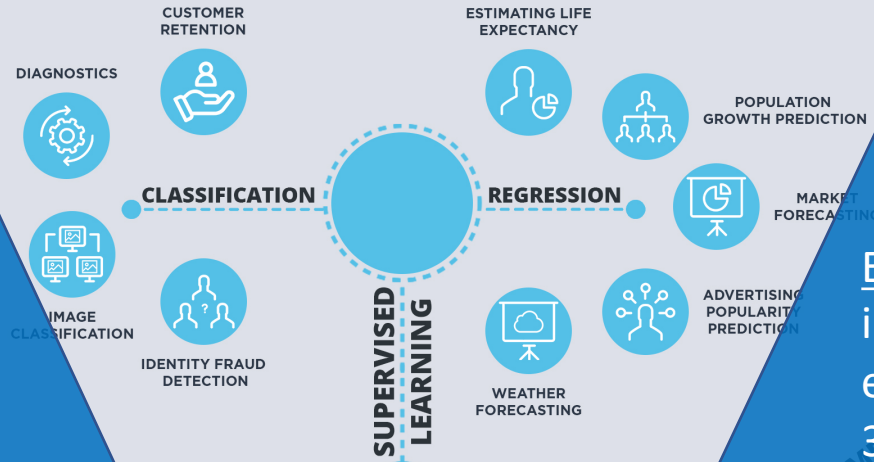
Supervised
VS
Unsupervised

Supervision means that label(s) are attached to data, *ie* the targeted value (discrete or continuous) is known in the training sample.

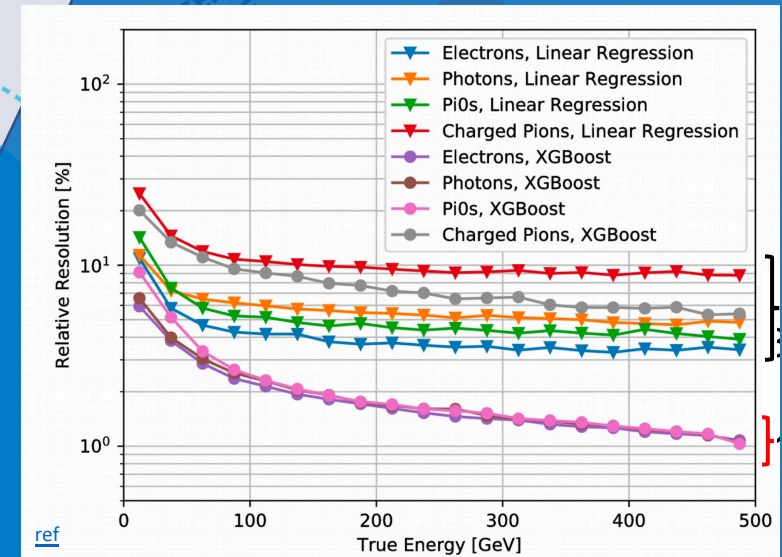
Regression
VS
Classification

Regression:
The target variable is continuous [ex: time, weight, mass, ...]

Classification:
The target variable is a discrete, *ie* a category.
The simplest case correspond to a binary classification.
[ex: signal vs bkg, species, ...]



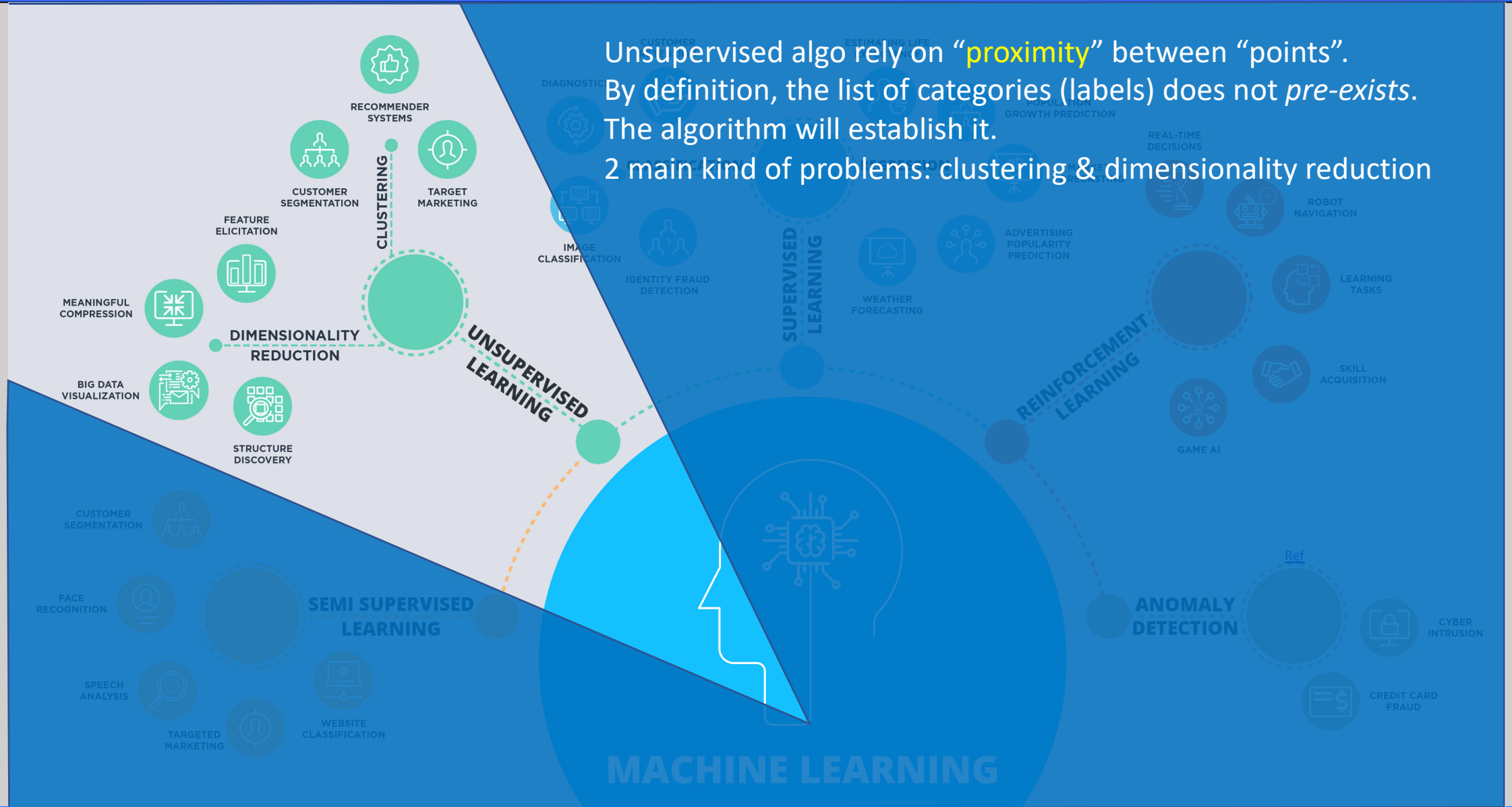
Example in HEP:
improvement of particle energy resolution
3-5% → 1%



Traditional methods
ML methods BDT

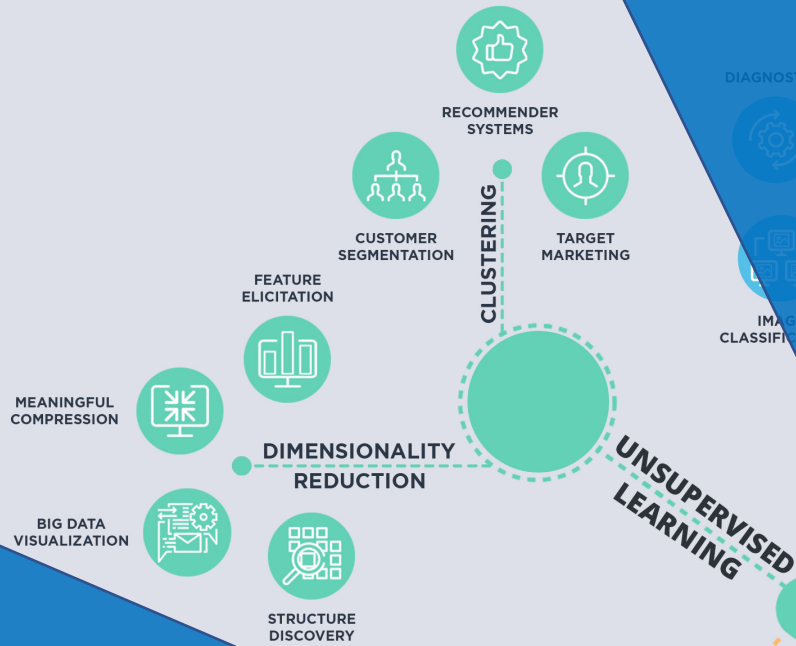
ML algo can improve “experimental performances”
Could be a cheaper/less demanding solution than building a new “apparatus”

Class of problems solved by ML algorithms



Unsupervised algo rely on “proximity” between “points”.
By definition, the list of categories (labels) does not *pre-exists*.
The algorithm will establish it.
2 main kind of problems: clustering & dimensionality reduction

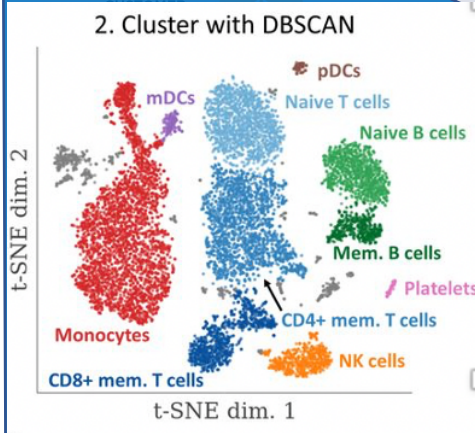
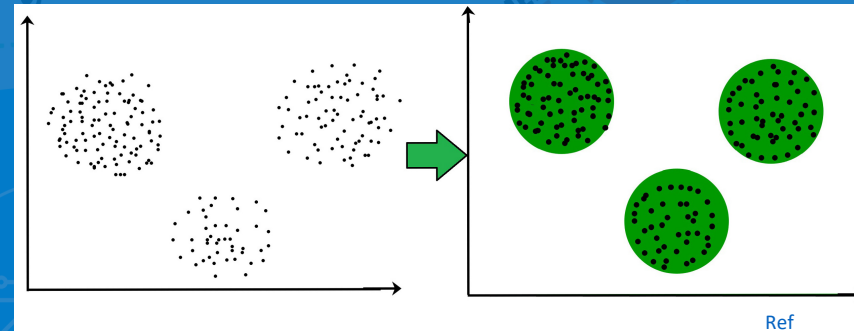
Class of problems solved by ML algorithms



Unsupervised algo rely on “**proximity**” between “points”.
By definition, the list of categories (labels) does not *pre-exists*.
The algorithm will establish it.

2 main kind of problems: clustering & dimensionality reduction

k-mean is an example of “geometric” algo.
Probabilistic approach also exists.



Applications by example in biology for “species” identification

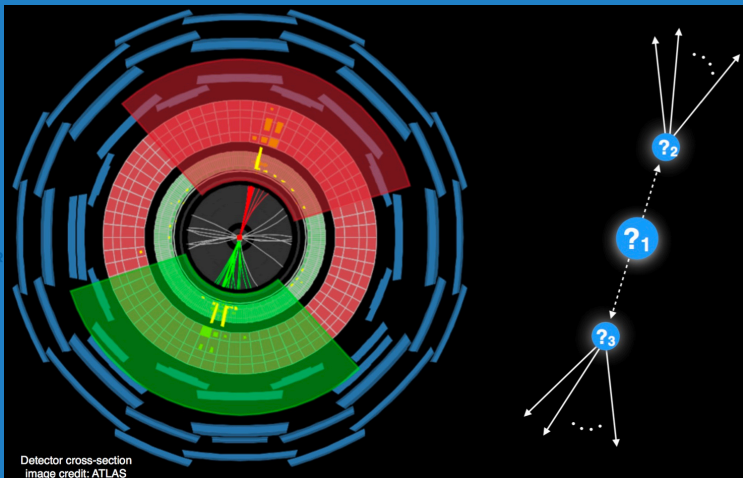
Class of problems solved by ML algorithms



LHC olympics 2020

BIG DATA
VISUALIZATION

<https://lhco2020.github.io/homepage/>



Detector cross-section
image credit: ATLAS

Anomaly = unexpected experimental signature (sign of new phenomena)

Example in HEP:

- Used to detector anomalies
- Attempt to use it to search for new phenomena (*example here*)

ESTIMATING LIFE EXPECTANCY

REGRESSION

WEATHER FORECASTING

SUPERVISED
LEARNING

POPULATION FORECASTING

MARKET FORECASTING

REINFORCEMENT
LEARNING

REAL-TIME DECISIONS

ROBOT NAVIGATION

LEARNING TASKS

SKILL ACQUISITION

GAME AI

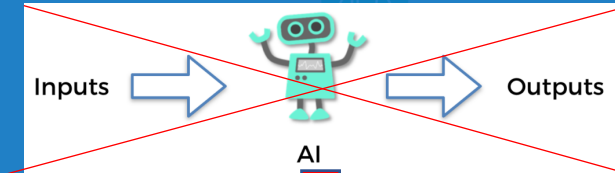
**ANOMALY
DETECTION**

CYBER
INTRUSION

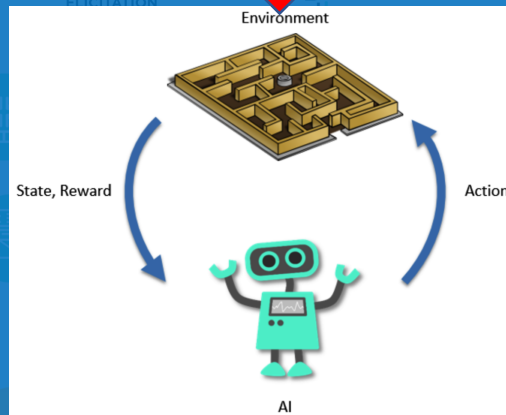
CREDIT CARD
FRAUD

MACHINE LEARNING

Class of problems solved by ML algorithms



- Feedback loop
- Introduce a reward
- Continuous improvement



REINFORCEMENT LEARNING

- REAL-TIME DECISIONS
- ROBOT NAVIGATION
- LEARNING TASKS
- SKILL ACQUISITION
- GAME AI

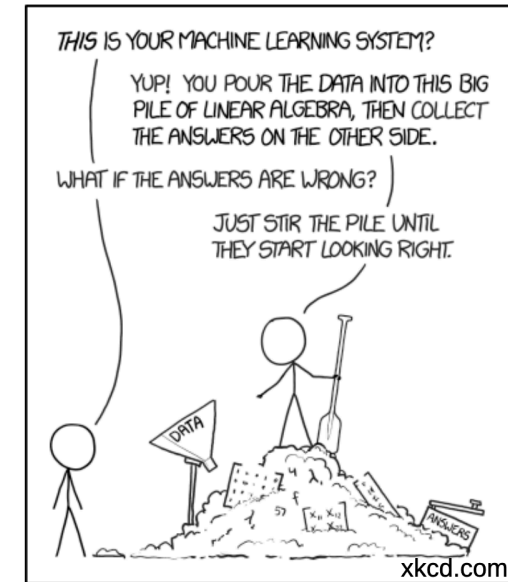
Example of applications in subatomic physics:
Finding optimal beam characteristics

Background collage of ML applications including: MEANINGFUL COMPRESSION, BIG DATA VISUALIZATION, CUSTOMER SEGMENTATION, FACE RECOGNITION, SPEECH ANALYSIS, TARGETED MARKETING, CLASSIFICATION, CUSTOMER RETENTION, ESTIMATING LIFE EXPECTANCY, POPULATION GROWTH PREDICTION, MARKET FORECASTING, ADVERTISING POPULARITY PREDICTION, WEATHER FORECASTING, IMAGE CLASSIFICATION, IDENTITY FRAUD DETECTION, SUPervised LEARNING, ANOMALY DETECTION, CYBER INTRUSION, CREDIT CARD FRAUD.

List of questions before starting using ML

1- Should I need to use ML ?

- Algorithm complexity, dataset size, computing resources used are **not guarantees** that the performances will be better than what you could achieve with more “**traditional methods**” !
- **ML is not the solution to all problems**: If you have an **analytical model** (*scientifically motivated*) describing your data, why should you spend resources to approximate it with ML ? (*discussion about GAN*)
- Useful to go beyond current knowledge in a data-driven way
- Gain vs effort
 - Can I expect a large improvement compare to current method ?
 - How much effort is needed to implement a ML solution?
 - Do I have enough resources to perform the training ?



List of questions before starting using ML

1- Should I need to use ML ?

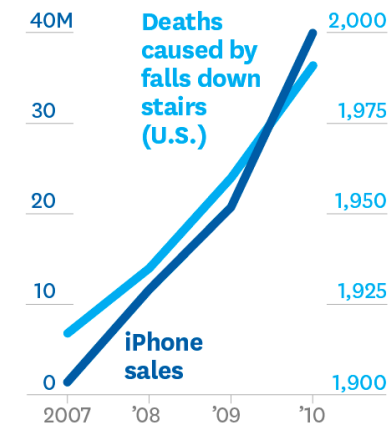
- Algorithm complexity, dataset size, computing resources used are **not guarantees** that the performances will be better than what you could achieve with more “traditional methods” !
- **ML is not the solution to all problems:**
- **Gain vs effort**

ML CANNOT be “blindly used”:

ML does not replace the human knowledge / judgment

- Needs human guidance in the process of applying ML techniques
 - Ex: “spurious correlation in large datasets” (p-score)
- For interpretation:
 - correlation does imply not causal relation
 - Remark: hard to interpret the “ML model”

**MORE IPHONES MEANS
MORE PEOPLE DIE FROM
FALLING DOWN STAIRS**



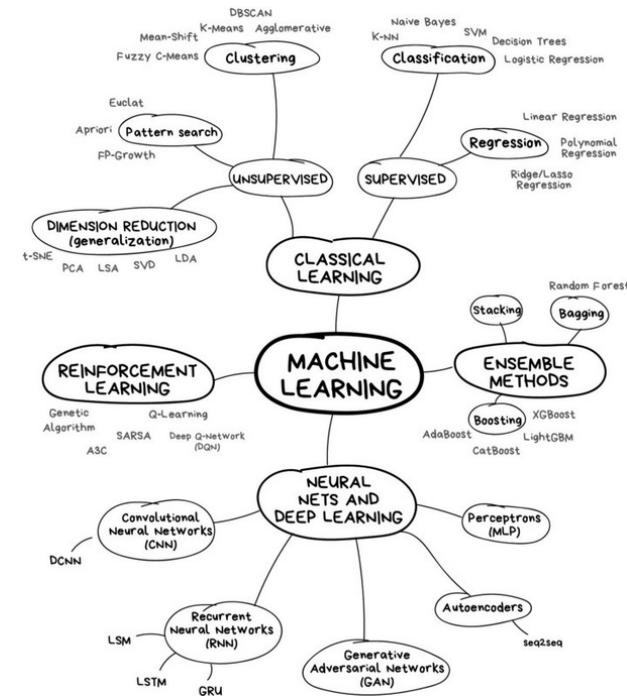
SOURCE TYLERVIGEN.COM
FROM “BEWARE SPURIOUS CORRELATIONS,” JUNE 2015



List of questions before starting using ML

2- Which category of ML algorithms should I use ?

- We have presented 5 class of algorithms
- Subcategories for some of the classes
- Many algorithms are devoted to a given class of problem



3- What is the “nature” of the data to analyze ?

4 types of data:

- Numerical data
- Categorical data
- Time series (video, sound, ...)
- Text

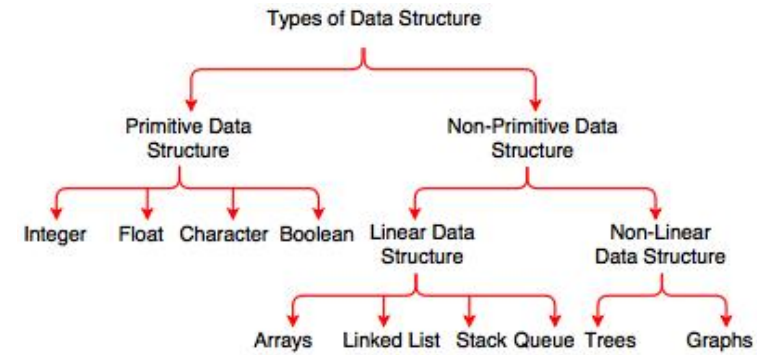


Fig. Types of Data Structure

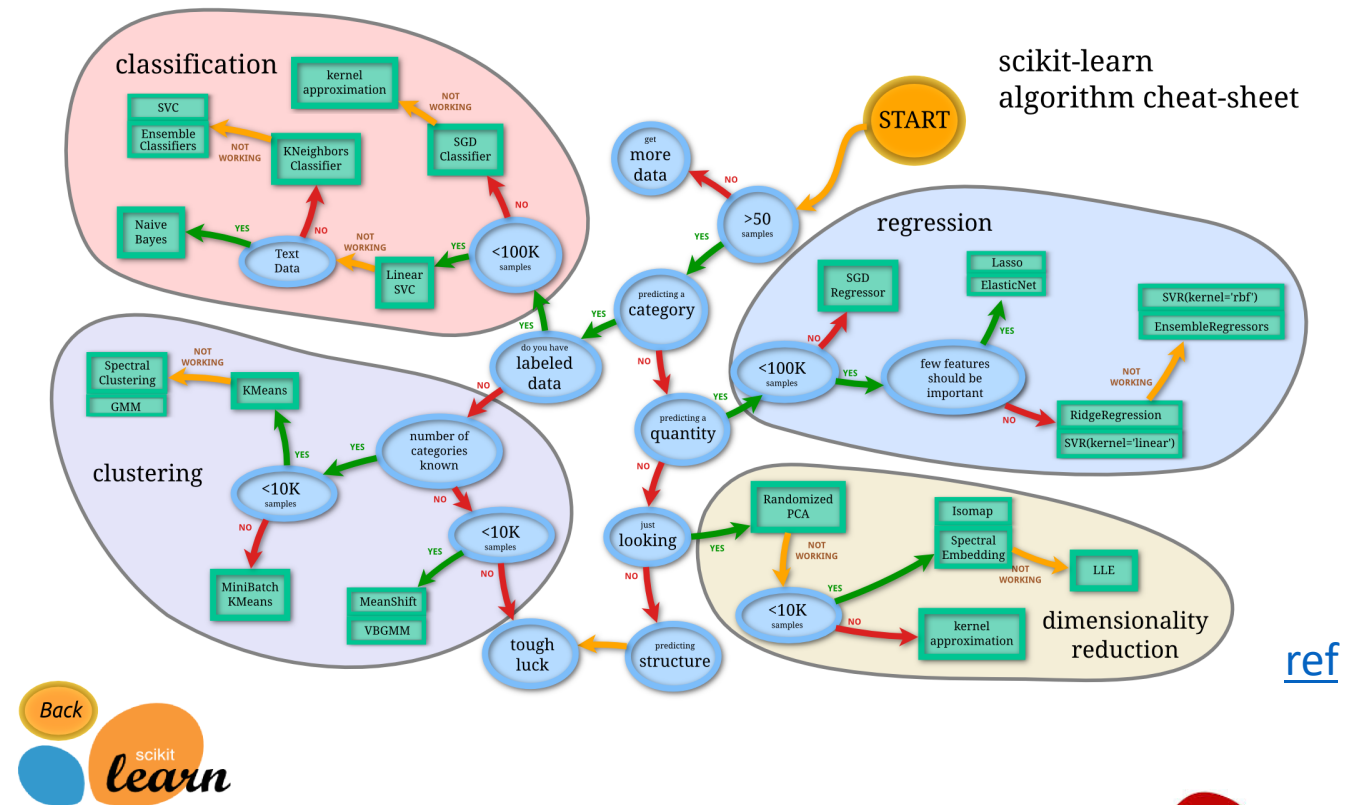


List of questions before starting using ML

4- How large is our sample ?

- Choose an algorithm adapted to our sample
- Deep learning requires very large datasets

You should now have a list of algorithms that may suit our needs.



What's next ?



ML implementation

Answer: finding a (free / open source / maintained) implementation !

- [\(Google\) Tensorflow](#) (2015)/[Keras](#)
 - Fast-growing easy-to-use python lib (but also in C++, ...)
 - Allows applications of deep-learning models
 - Interface to Tensorflow, Theano backends → GPU support
- [Scikit-learn](#) (2007)
 - Python lib that implements many (non-deep) techniques
 - A lot of data preprocessing & statistics tools
- [TMVA](#) (used in subatomic physics – ROOT based)

But also:

- [Torch](#) / [PyTorch](#)
- [Caffee](#) (C++/python)
- [Accord.net](#) (C++)
- [R](#)
- ...

Can be discussed
this afternoon



What's next ?



Hardware

Answer: where to run the ML training ?

- **Use our laptop:**
 - several ML-analyses can be done on (from few min to few hours)
 - Multi-threading
- **Use dedicated GPU(s)**
 - reasonable dataset – large computing time
 - Several algorithms profit a lot from GPU parallelization
 - From few 100's to few 1000's euros
- **Use server/computing center**
 - large data / complex models
 - Tier2
 - Mesocentre
 - CC-in2p3
 - CERN
 -



Can be discussed
this afternoon

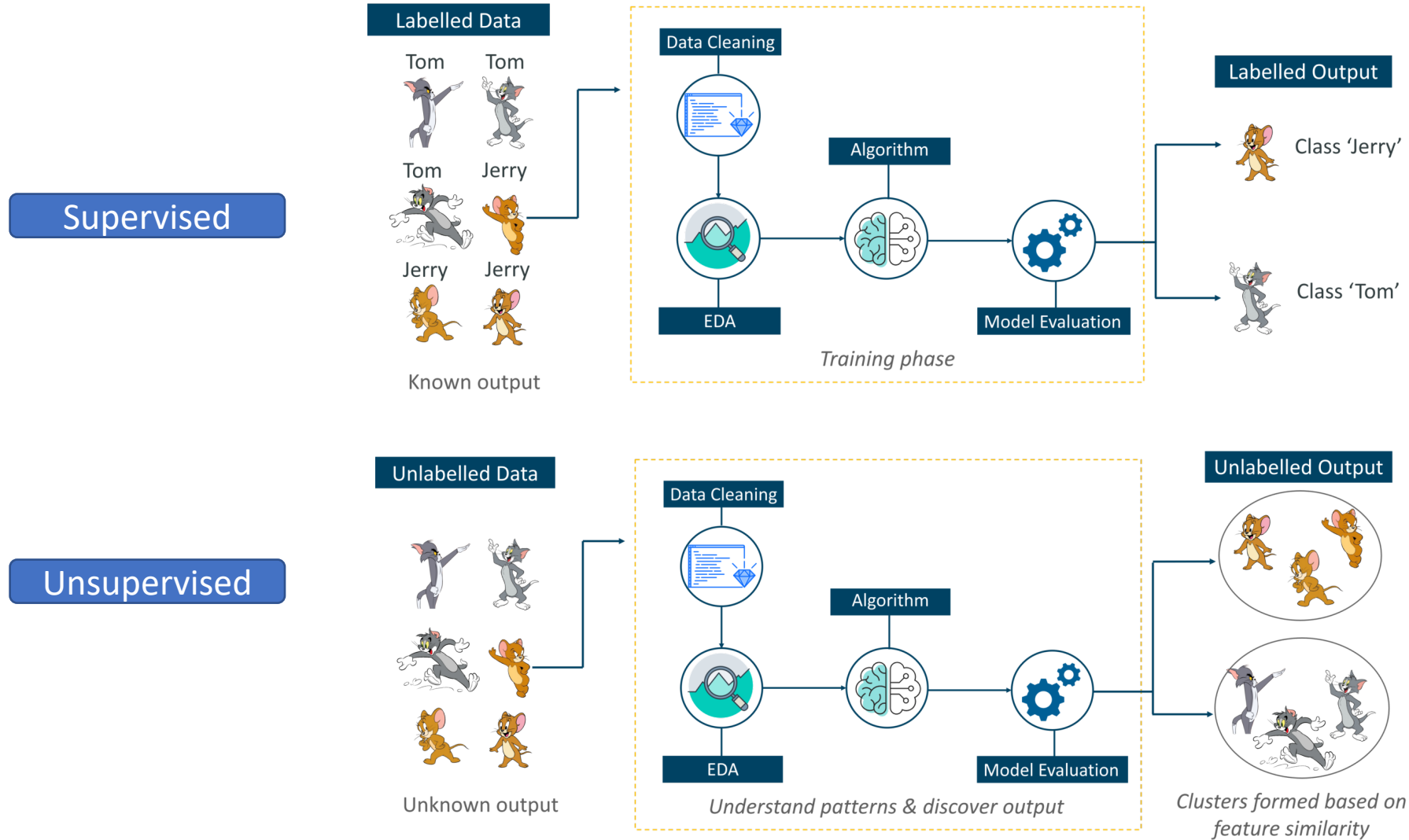


What's next ?

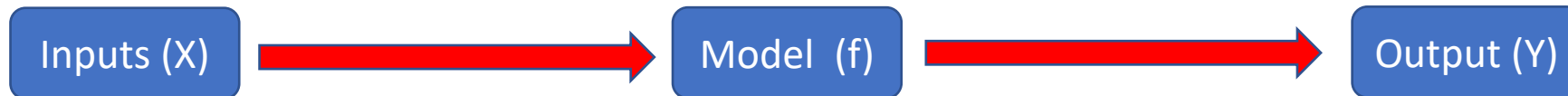
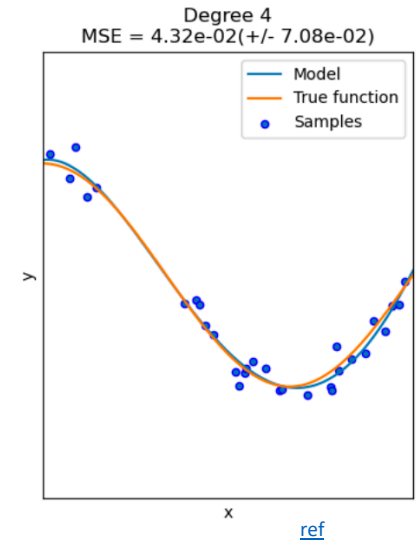
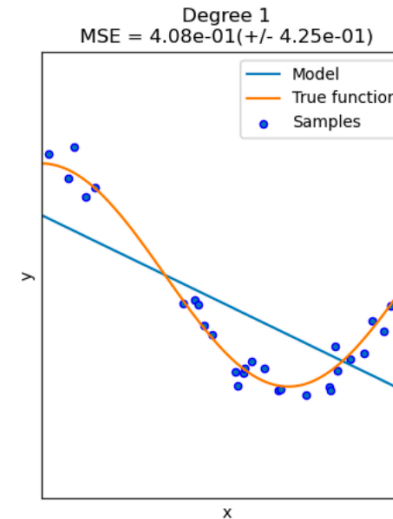
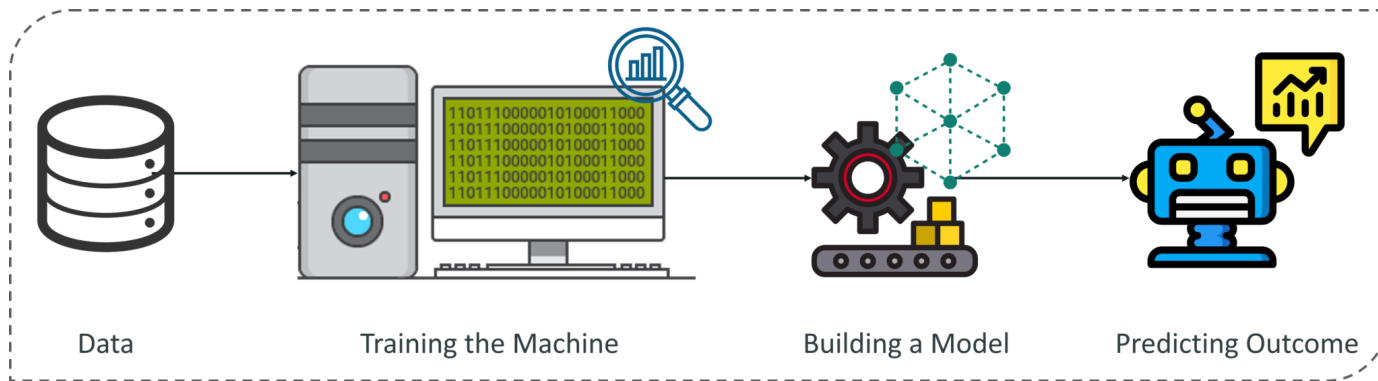


ML workflow

[Ref](#)



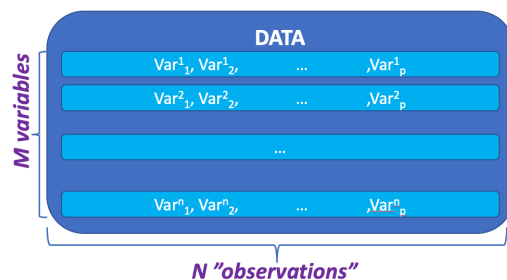
Supervised ML workflow



- N “samples”.
- M “features” per sample [dimension]
- (*predicator variables*)

- Predictive function: $y = f(x)$
- Depends on the choice of the algo
- Configuration (ex: NN: #layers, #nodes, ...)
- Have **many** parameters to be determined

- One (or few) response (*target*) variable
- Can be discrete or continuous



Find an function $f(x)$
that approximates the function
F such that
 $Y = F(x)$



Supervised ML workflow

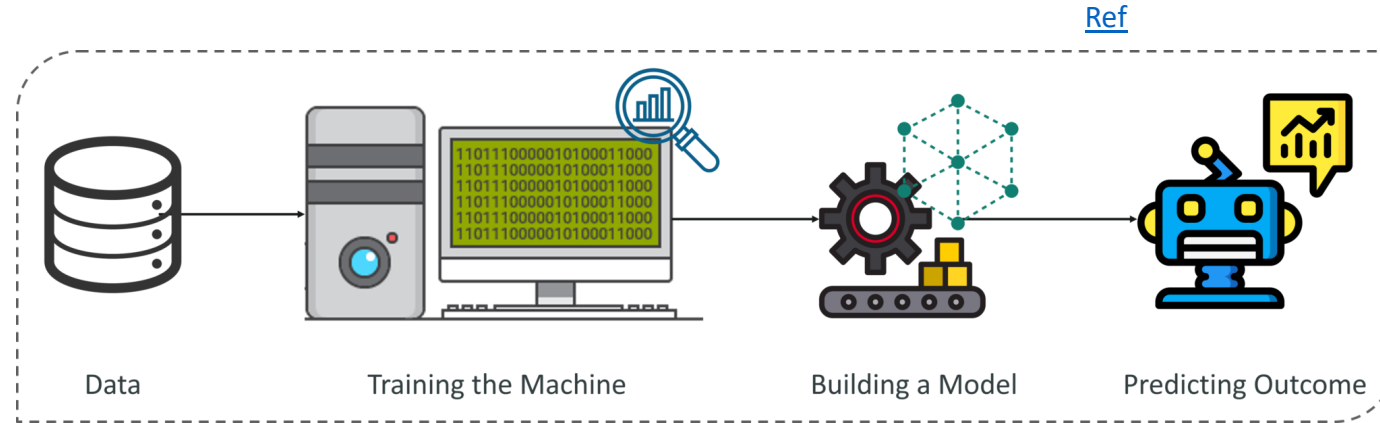
Requirements:



Step 1: Define the objective

- choice of algo/implementation
- choice of hardware chosen

Step 2: Data gathering



Step 3: Data preparation

data cleaning, filtering, binning, transformation ...

Step 4: Exploratory Data Analysis (EDA)



Understand patterns, trends, correlations, ..

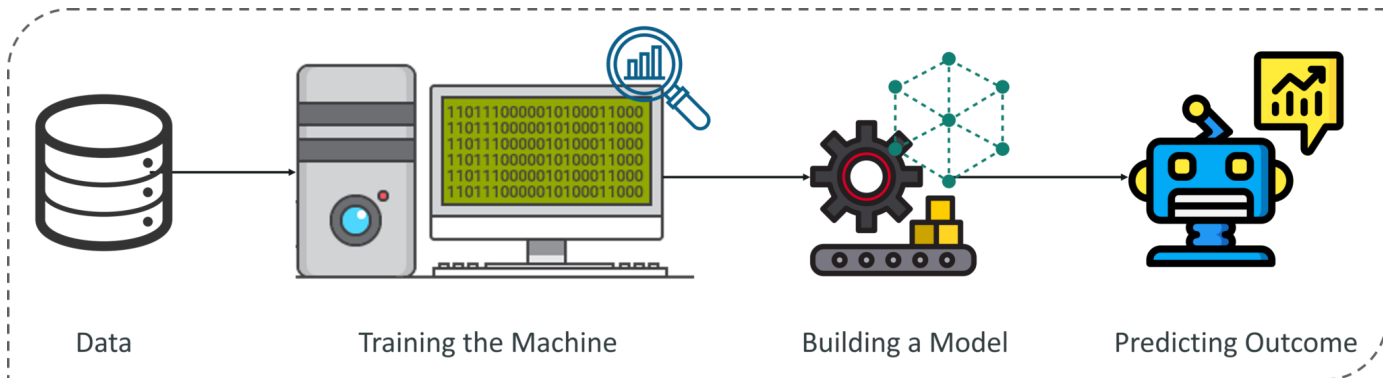
Goal: Selection input variables (features) & “samples”

Advices:

- Avoid variables too sensitive to noise, uncertainties, ...
- Avoid variable badly described by our simulation if the dataset is simulation based
- Use a representative “samples”



Supervised ML workflow

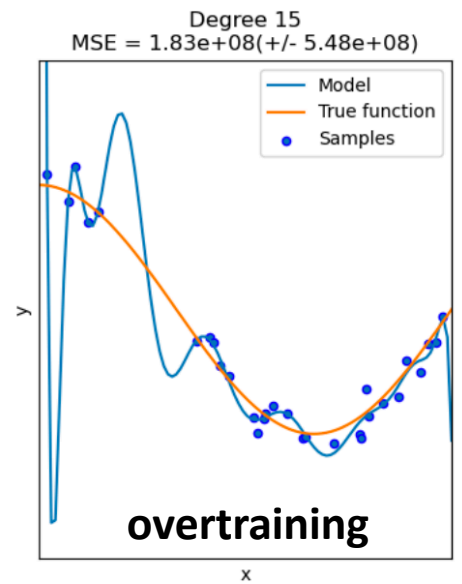


- CPU**
 - Small models
 - Small datasets
 - Useful for design space exploration
- GPU**
 - Medium-to-large models, datasets
 - Image, video processing
 - Application on CUDA or OpenCL
- TPU**
 - Matrix computations
 - Dense vector processing
 - No custom TensorFlow operations
- FPGA**
 - Large datasets, models
 - Compute intensive applications
 - High performance, high perf./cost ratio

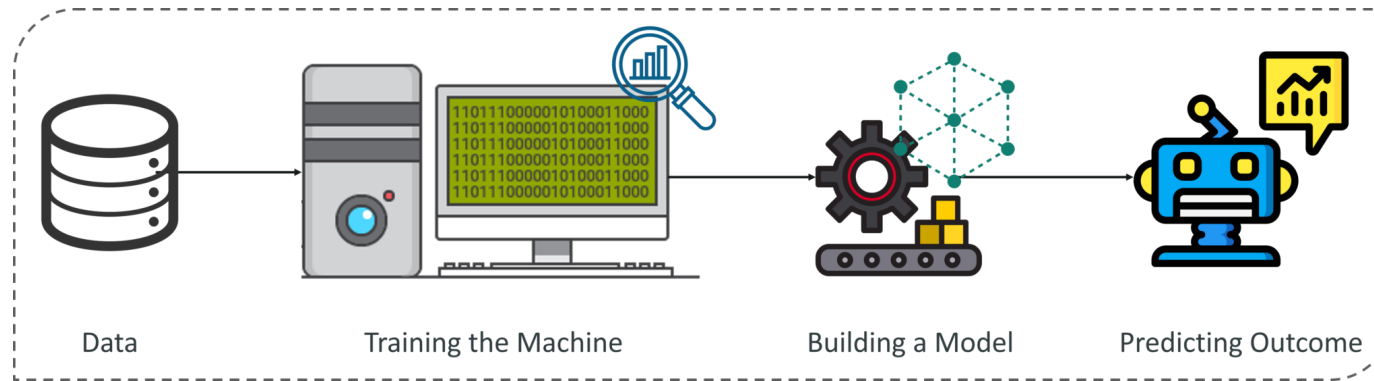
Tensor PU
Dev. By google
2016

Step 5: Building a Machine Learning Model == Model training

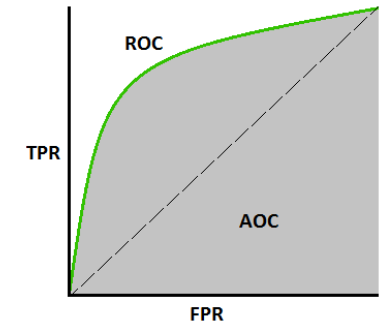
- Define a (large) subset of data == **training sample**
- Most CPU/GPU demanding task
- **Goal:** determining the parameters of the model
- **Advices:**
 - sample should be large enough (depends on model complexity)
 - Computation time can be speed-up depending on both the hardware architecture chosen & the model
 - Remark: processing does not necessarily linearly scale with #threads !
 - Check that there is no **overtraining**



Supervised ML workflow



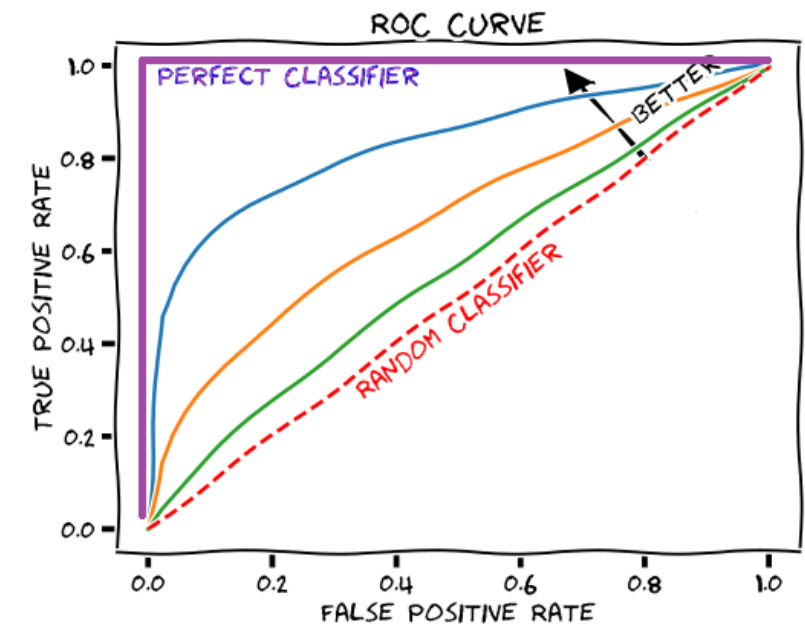
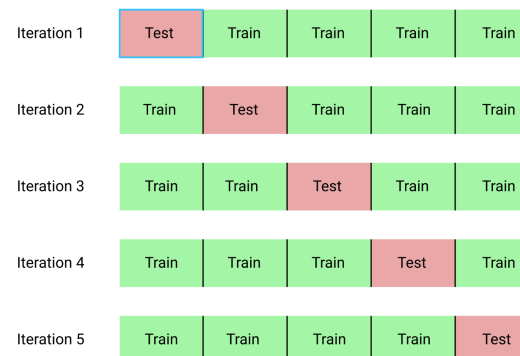
[Ref](#)



Step 6: Model Evaluation

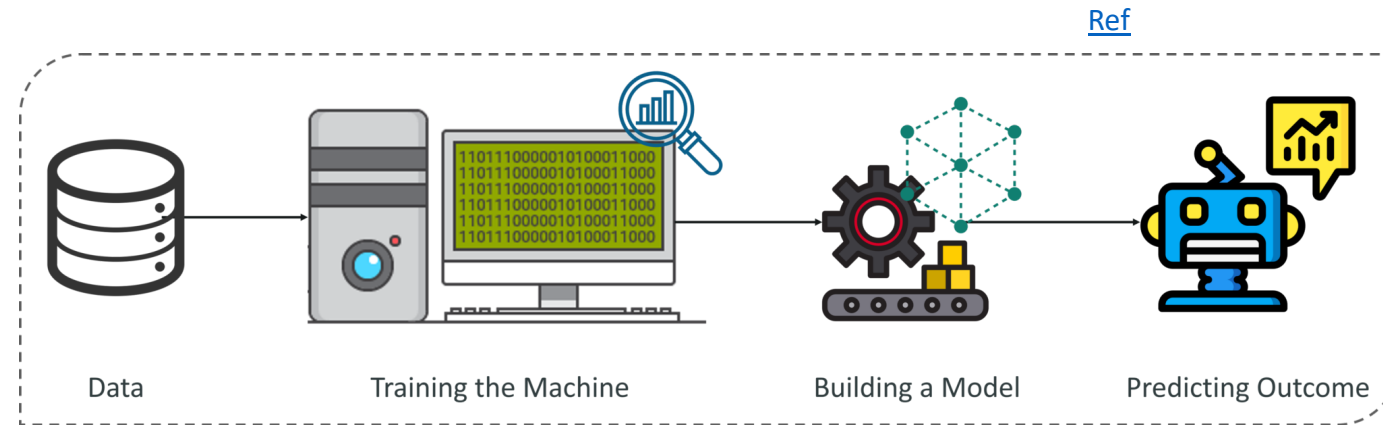
- Define a (small / independent) subset of data == **test sample**
- Evaluate the performances
 - True vs False Positive Rate
 - ROC: Receiver Operating Characteristic
 - AUC: Area under the curve
- Cross-validation
 - Avoid overtraining

Cross-validation



[ref](#)

Supervised ML workflow



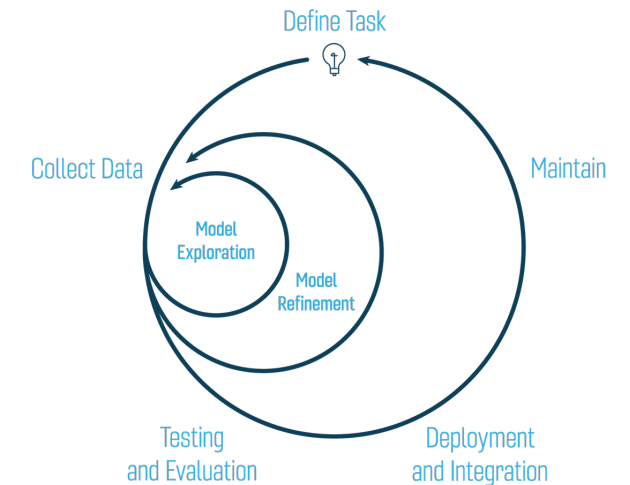
Step 7: Model optimization (refinement)

- (Hyper-)parameter tuning (ex with NN: #layers, #nodes, ..)
- Test/Optimize feature selection (ex: change/add input var)
- **Imply retraining** (CPU/GPU demanding)



→ This **loop** can be done several times ...

Machine Learning Development Lifecycle



Step 8: Predictions ... !

→ algo deployment/integration/maintenance

Remark: Retraining can be needed (even if inputs/parameters are fixed) if the model is applied on data collected in new conditions (new sensors/calibration/alignment etc ...)

Networks & training

Networks

- ML collaborations @ IN2P3
- IML (LHC ML working group)
- ...

Training

- Master “Big data and ML” @ unistra
- [IN2P3 school of statistics](#)
- Ecole doctorale (ED 182)
- MOOC
- ...

Can be discussed
this afternoon

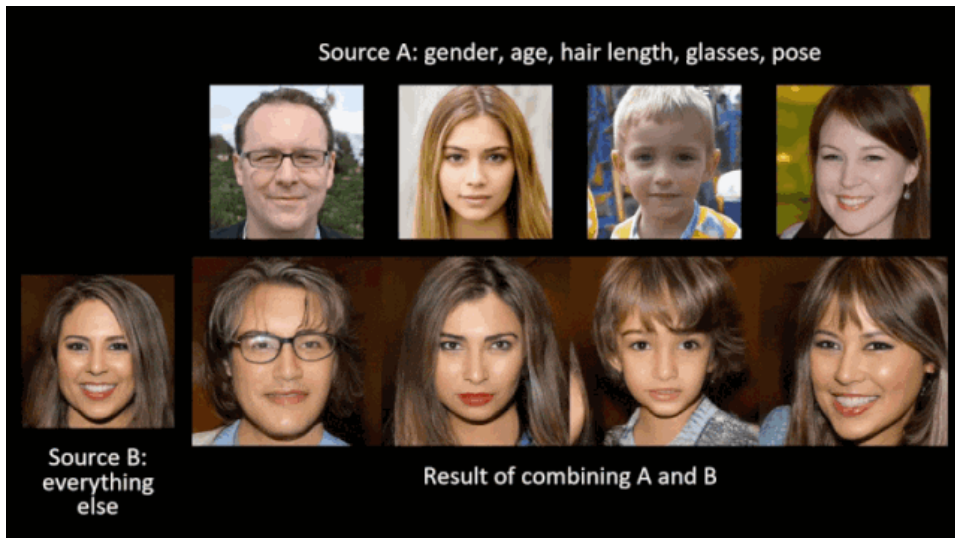


GAN: Generative Adversarial Network

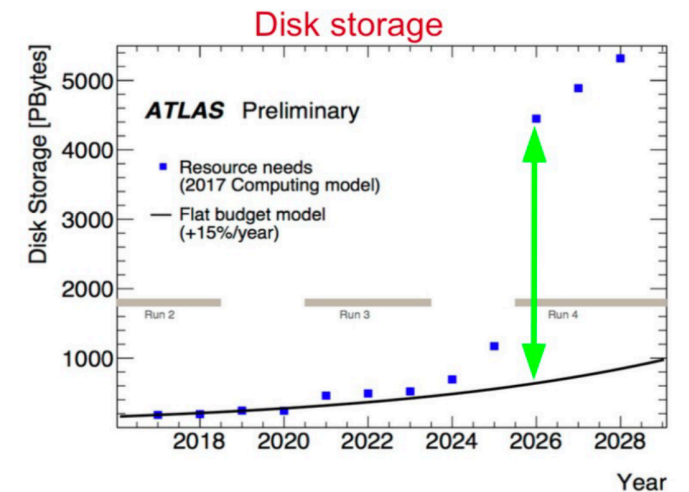
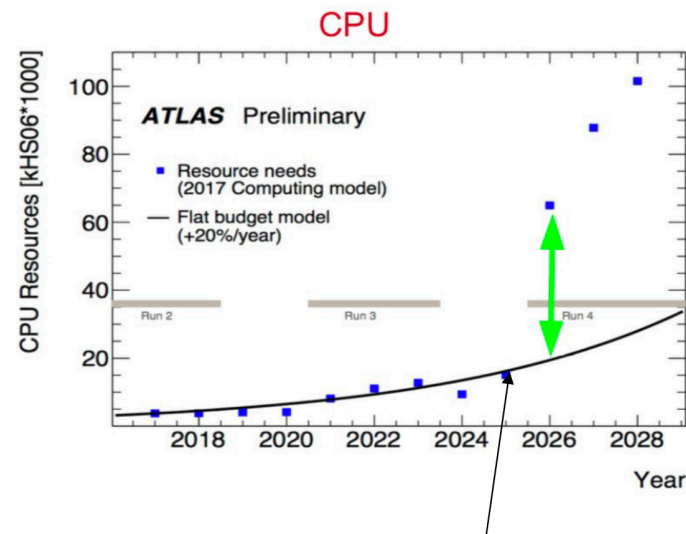
Technique used to generate “realistic human face”

Can have application in HEP to **speed-up simulation**

- Simulation of proton-proton collision @ LHC
- Simulation of detector response



[Ref](#)



Dominated by : calorimeter simulation and tracking

09:00	→ 09:20	Accueil et introduction Orateurs: Eric CHABERT (IPHC/UDS), Pierre Van Hove (CNRS)	🕒 20m
09:20	→ 10:20	Exploitation des données de détecteurs	
09:20		Bio-loggers, utilisation de "Convolutional Neural Network" Orateur: Lorène Jeantet	🕒 15m
09:35		Exploitation des données des bio-loggers Orateur: Marie-Amélie Forin-Wiart	🕒 25m
10:00		Capteurs CMOS et réseau de neurones Orateur: auguste besson (Institut Pluridisciplinaire Hubert Curien)	🕒 20m
10:20	→ 11:00	Particules et collisions	
10:20		Reconstruction des événements à l'aide de ML en physique des particules Orateur: Emery Nibigira (LPC Clermont)	🕒 25m
10:45		Traitement global d'événements sur collisionneur: utilisation de DNN Orateur: Giulio Dujany (CNRS - IPHC)	🕒 15m
11:00	→ 11:20	Pause pommes	🕒 20m
11:20	→ 11:55	Traitement d'image	
11:20		Reconnaissance d'espèces: image et vidéo Orateur: Claire Saraux	🕒 20m
11:40		Projet dessin Orateurs: Cédric Sueur (Institut Pluridisciplinaire Hubert Curien), Marie Pelé	🕒 15m
11:55	→ 12:15	Spectrométrie de masse	
11:55		Utilisation de ML et études mathématiques pour le traitement des données de spectrométrie de masse Orateurs: Alexandre Burel (CNRS), Marie Chion	🕒 20m
12:15	→ 12:35	Support	
12:15		Informations, conseils et développements du service informatique Orateurs: Jerome Pansanel (IPHC - CNRS), sebastien geiger (IPHC IN2P3)	🕒 15m
12:35	→ 14:00	Déjeuner	🕒 1h 25m
14:00	→ 16:00	Table ronde Fils de discussions envisagés: -Prochaines étapes -Retour d'expériences -Conseils/suggestions du service info Présidents de session: Eric CHABERT (IPHC/UDS), Pierre Van Hove (CNRS)	

Beyond activities presented this morning, additional are done

Non exhaustive list:

- Clustering in CMOS pixels with NN (*DRHIM, Finck & al*)
- Use of AI to improve detector performances in the context of PET (tomography) (*DRHIM, Brasse & al*)
- Automatic ML-based localisation of radioactive contamination zone with an autonomous drone (learning based on MC) (*DRS, Arbor & al*)

