

Apprentissage statistique

Introduction: intelligence artificielle, machine learning, apprentissage statistique

Cours pour non-spécialistes

Aurélien Garivier

What is this course?

It is:

- an introduction to a field that is currently in fast growth
- possible to follow with no mathematical background
- necessary to listen to some technical elements
- interactive! please tell about your data problems

It is not:

- a comprehensive presentation of all problems and algorithms
- a course for perfectly mastering the algorithms and proofs
- where you will learn to master GPU clusters, etc.
- a general discussion about/over ML (it is ML, we present methods)

Prerequisite / Evaluation

Prerequisite: it will be easier with

- Mathematics: linear algebra, real analysis, elementary probability
- Statistics: some practice of data
- Computer science: programming

Evaluation: for choice

- Final exam
- Case study on self-provided data (10 pages report + oral presentation)

Outline

- Introduction: What is AI, data, ML, statistical learning?
- Unsupervised learning, clustering
 - Principal component analysis
 - Agglomerative Hierarchical Clustering
 - k-means, k-medoids and variants
 - overview of other methods: Affinity Propagation, dbscan, etc.
- Supervised learning: classification
 - k-nearest neighbors
 - Gaussian linear model, logistic regression, model selection
 - LASSO et variants
 - Support Vector Machines
 - Decision Trees
 - Ensemble methods: Bagging, Random Forests, Boosting
 - Neural networks, deep learning
- Introduction to reinforcement learning
- A few words on Big Data: what is new?

Machine Learning: when Artificial Intelligence meets Big Data

The Learning Models

Machine Learning Methodology

Artificial Intelligence (AI): Definition

Intelligence exhibited by machines

- emulate cognitive capabilities of humans
(big data: humans learn from abundant and diverse sources of data).
- a machine mimics "cognitive" functions that humans associate with other human minds, such as "learning" and "problem solving".

Ideal "intelligent" machine =

flexible rational agent that perceives its environment and takes actions that maximize its chance of success at some goal.

Founded on the claim that human intelligence

"can be so precisely described that a machine can be made to simulate it."

Operational goals

- Autonomous robots for not-too-specialized tasks
- In particular, vision + understand and produce language

Tension between operational and philosophical goals

- As machines become increasingly capable, facilities once thought to require intelligence are removed from the definition. For example, optical character recognition is no longer perceived as an exemplar of "artificial intelligence"; having become a routine technology.
- Capabilities still classified as AI include advanced Chess and Go systems and self-driving cars.

Arthur Samuel (1959)

Field of study that gives computers the ability to learn without being explicitly programmed

Tom M. Mitchell (1997)

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .

ML: Learn from and make predictions on data

- Algorithms operate by building a model from **example** inputs in order to make data-driven **predictions or decisions**...
- ...rather than following strictly static program instructions: useful when designing and programming explicit algorithms is unfeasible or poorly efficient.

Within Data Analytics

- Machine Learning used to devise complex models and algorithms that lend themselves to **prediction** - in commercial use, this is known as *predictive analytics*.
- www.sas.com: "Produce reliable, repeatable decisions and results" and uncover "hidden insights" through learning from historical **relationships and trends** in the data.
- evolved from the study of pattern recognition and computational learning theory in artificial intelligence.

Machine Learning: Typical Problems

- spam filtering, text classification
- optical character recognition (OCR)
- search engines
- recommendation platforms
- speech recognition software
- computer vision
- bio-informatics, DNA analysis, medicine
- etc.

For each of this task, it is possible but very inefficient to write an explicit program reaching the prescribed goal.

It proves much more succesful to have a machine infer what the good decision rules are.

Related Fields

- **Computational Statistics:** focuses in prediction-making through the use of computers together with statistical models (ex: Bayesian methods).
- **Statistical Learning:** ML by statistical methods, with statistical point of view (probabilistic guarantees: consistency, oracle inequalities, minimax)
→ more focused on *correlation*, less on *causality*
- **Data Mining** (unsupervised learning) focuses more on exploratory data analysis: discovery of (previously) unknown properties in the data. This is the analysis step of Knowledge Discovery in Databases.
- Importance of **probability**- and **statistics**-based methods → **Data Science** (Michael Jordan)
- Strong ties to **Mathematical Optimization**, which delivers methods, theory and application domains to the field

Qu'est-ce qu'une (très grande) masse de données ?



VLDB
 XLDB
 Massive Data
 Big Data
 Very Big Data
 Data Deluge
 Data Masses

Data inflation

2

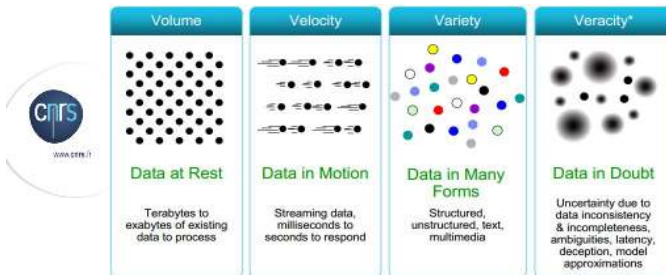
Unit	Size	What it means
Bit (b)	1 or 0	Short for "binary digit", after the binary code (1 or 0) computers use to store and process data
Byte (B)	8 bits	Enough information to create an English letter or number in computer code. It is the basic unit of computing
Kilobyte (KB)	1,000, or 2^{10} , bytes	From "thousand" in Greek. One page of typed text is 2KB
Megabyte (MB)	1,000KB; 2^{20} bytes	From "large" in Greek. The complete works of Shakespeare total 5MB. A typical pop song is about 4MB
Gigabyte (GB)	1,000MB; 2^{30} bytes	From "giant" in Greek. A two-hour film can be compressed into 1-2GB
Terabyte (TB)	1,000GB; 2^{40} bytes	From "monster" in Greek. All the catalogued books in America's Library of Congress total 15TB
Petabyte (PB)	1,000TB; 2^{50} bytes	All letters delivered by America's postal service this year will amount to around 5PB. Google processes around 1PB every hour
Exabyte (EB)	1,000PB; 2^{60} bytes	Equivalent to 10 billion copies of <i>The Economist</i>
Zettabyte (ZB)	1,000EB; 2^{70} bytes	The total amount of information in existence this year is forecast to be around 1.2ZB
Yottabyte (YB)	1,000ZB; 2^{80} bytes	Currently too big to imagine

The prefixes are set by an inter-governmental group, the International Bureau of Weights and Measures. Yotta and Zetta were added in 1991; terms for larger amounts have yet to be established.

Source: *The Economist*

Grandes Conf du domaine: VLDB, XLDB, ICDE, EDBT, ...

Complexité multidimensionnelle des Big Data



• Nouvelles archi. de stockage

• Nouvelles archi. d'interopérabilité

• Défi pour les réseaux de communication

• Nouveaux modèles de calcul sur des flux

• Nettoyage et transformation

• Fusion de données

Nouveaux modèles de qualité (données & processus de traitement)

<http://www.datasciencecentral.com/profiles/blogs/data-veracity>

Défis accompagnant les chgts

DÉFIS TRANSVERSES

- Passage à l'échelle
- Rapidité traitements
- Protection, sécurité
- Interaction

Acquisition

Extraction, nettoyage

Intégration

Analyse

Interpretation

Stockage

Accès/
Requêtage,
Raisonnement

Valeur

Véracité

Velocité

Variété

Volume

repenser les outils algorithmiques et mathématiques

inspired by "Big Data and Its Technical Challenges, Communications of the ACM, July 2014, vol 57, n°7", © H.V. Jagadish et al.

4

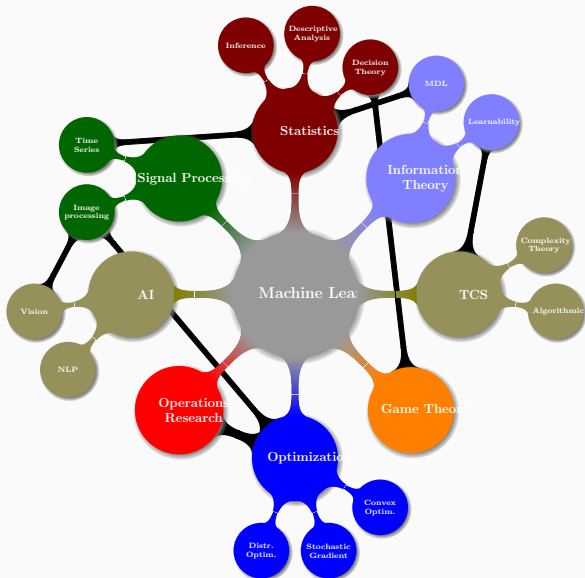
- Data analysis (inference, description) is the goal of statistics for long.
- Machine Learning has more **operational** goals (ex: consistency is important the statistics literature, but often makes little sense in ML).

Models (if any) are *instrumental*

Ex: linear model (nice mathematical theory) vs Random Forests.

- Machine Learning/big data: no separation between statistical modelling and optimization (in contrast to the statistics tradition).
- In ML, data is often here before (unfortunately)
- No clear separation (statistics evolves as well).

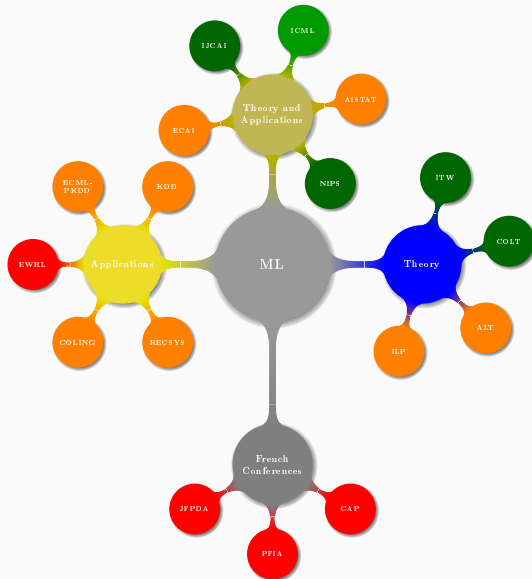
ML and its neighbors



ML journals



ML conferences



Machine Learning: when Artificial Intelligence meets Big Data

The Learning Models

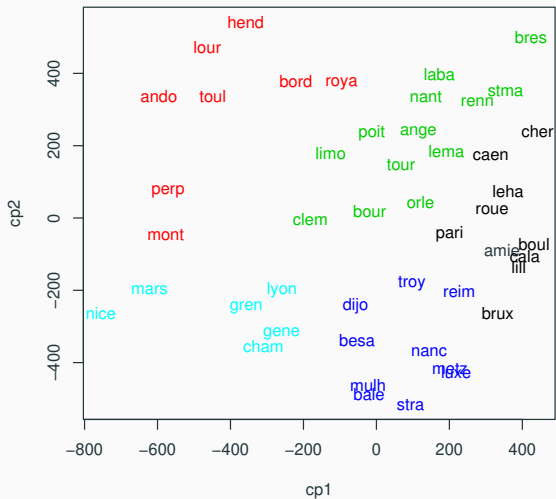
Machine Learning Methodology

What ML is composed of



- (many) observations on (many) individuals
- need to have a simplified, structured overview of the data
- *taxonomy*: untargeted search for *homogeneous clusters* emerging from the data
- Examples:
 - customer segmentation
 - image analysis (recognizing different zones)
 - exploration of data

Example



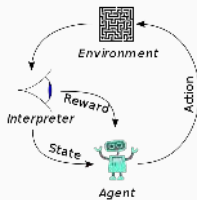
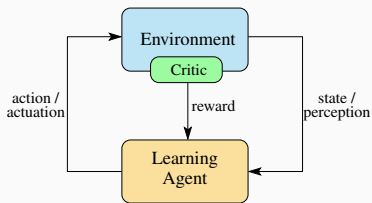
Supervised Learning

- observations = pairs (X_i, Y_i)
- goal = learn to *predict* Y_i given X_i
- regression (when Y is continuous)
- classification (when Y is discrete)
- statistical technique: linear models

Example: Character Recognition

Input space \mathcal{X}	64×64 images
Output space \mathcal{Y}	$\{0, 1, \dots, 9\}$
Joint distribution $P(x, y)$?
Prediction function $h \in \mathcal{H}$	
Risk $R(h) = P(h(X) \neq Y)$	
Sample $\{(x_i, y_i)\}_{i=1}^n$	MNIST dataset
Empirical risk $\hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{h(x_i) \neq y_i\}$	
Learning algorithm $\phi_n : (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$	NN, boosting...
Expected risk $R_n(\phi) = \mathbb{E}_n[R(\phi_n)]$	
Empirical risk minimizer $\hat{h}_n = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h)$	
Regularized empirical risk minimizer $\hat{h}_n = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h) + \lambda C(h)$	

Reinforcement Learning



[Src: https://en.wikipedia.org/wiki/Reinforcement_learning]

- area of machine learning inspired by behaviourist psychology
- how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.
- Model: random system (typically : Markov Decision Process)
 - agent
 - state
 - actions
 - rewards
- sometimes called approximate dynamic programming, or neuro-dynamic programming

Markov decision process

A **Markov Decision Process** is defined as a tuple $M = (X, A, p, r)$:

- X is the **state** space,
- A is the **action** space,
- $p(y|x, a)$ is the **transition probability** with

$$p(y|x, a) = \mathbb{P}(x_{t+1} = y | x_t = x, a_t = a),$$

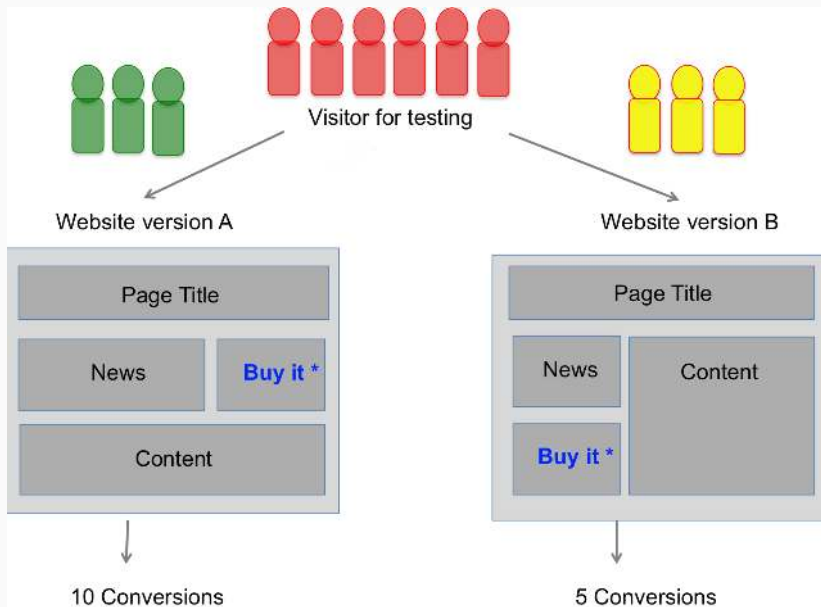
- $r(x, a, y)$ is the **reward** of transition (x, a, y) .

Example: the Retail Store Management Problem

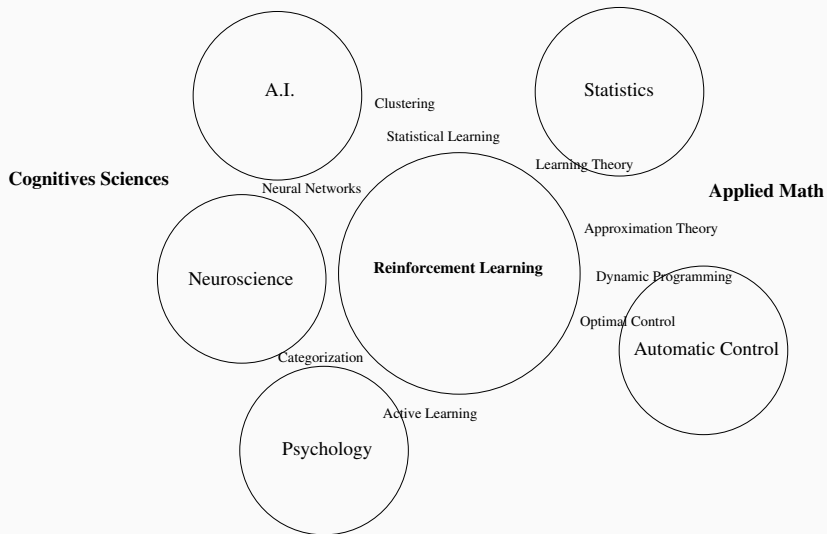
At each month t , a store contains x_t items of a specific goods and the demand for that goods is D_t . At the end of each month the manager of the store can order a_t more items from his supplier. Furthermore we know that:

- The **cost** of maintaining an inventory of x is $h(x)$.
- The **cost** to order a items is $C(a)$.
- The **income** for selling q items is $f(q)$.
- If the demand D is bigger than the available inventory x , customers that cannot be served leave.
- The **value of the remaining inventory** at the end of the year is $g(x)$.
- **Constraint**: the store has a maximum capacity M .

Example: A/B testing



Reinforcement Learning and the others



Machine Learning: when Artificial Intelligence meets Big Data

The Learning Models

Machine Learning Methodology

n -by- p matrix X

- n examples = points of observations
- p features = characteristics measured for each example

Questions to consider:

- Are the features centered?
- Are the features normalized? bounded?

In `scikitlearn`, all methods expect a 2D array of shape (n, p) often called

`X (n_samples, n_features)`

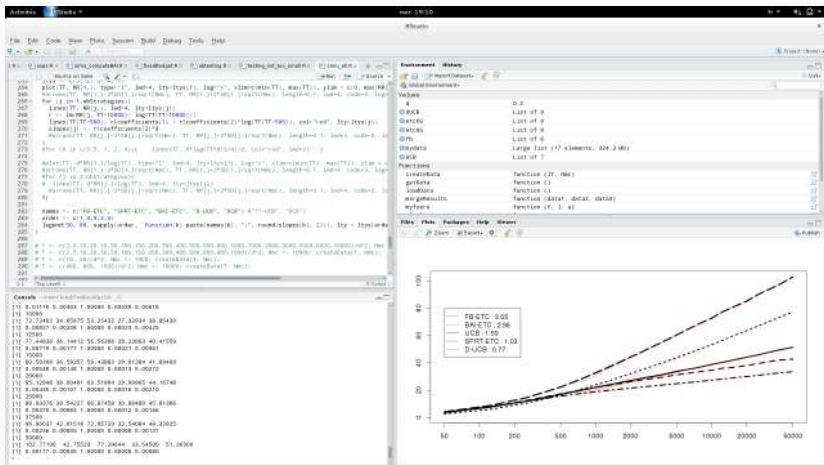
- Inside R: package datasets
- Inside scikitlearn: package sklearn.datasets
- UCI Machine Learning Repository
- Challenges: Kaggle, etc.



The big steps of data analysis

1. Extracting the data to expected format
2. Exploring the data
 - detection of outliers, of inconsistencies
 - descriptive exploration of the distributions, of correlations
 - data transformations
3. Random partitioning of the data: (see also: cross-validation)
 - learning sample
 - validation sample
 - test sample
4. For each algorithm: parameter estimation using training and validation samples
5. Choice of final algorithm using testing sample, risk estimation

Machine Learning tools: R



Machine Learning tools: python

The screenshot shows a Python IDE with a code editor on the left, a console on the right, and a plot window in the center. The code defines a linear regression model with parameters w_0 and w_1 , and a function `predict` to calculate the predicted value \hat{y} for a given input x .

```
def predict(x):  
    w0 = 0.0  
    w1 = 0.0  
    return w0 + w1 * x
```

The plot, titled "Figure 2", shows the predicted values \hat{y} versus the input values x . The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from 0.0 to 1.4. The data points are blue dots, and the predicted values are shown as a solid blue line. The plot is labeled "RM_mathsuite".

x	y
0.0	0.0
0.1	0.1
0.2	0.2
0.3	0.3
0.4	0.4
0.5	0.5
0.6	0.6
0.7	0.7
0.8	0.8
0.9	0.9
1.0	1.0

The screenshot shows the scikit-learn website homepage. At the top, there is a navigation bar with the scikit-learn logo and links for Home, Tutorials, User guides, and Examples. A search bar is also present. The main header features a grid of colorful images and the text "scikit-learn Machine Learning in Python". Below this, there are three bullet points: "Simple and efficient tools for estimating and applying models", "Accessible to practitioners and researchers in various domains", and "Full community, fully open-source".

The main content is organized into several sections:

- Classification**: Focusing on which category an object belongs to. Applications: Spam detection, image recognition. Also there: Support vector machines, random forests, ...
- Regression**: Predicting a continuous value (stock prices, house prices, ...). Applications: Drug response, stock prices. Algorithms: SVM, Ridge regression, ...
- Clustering**: Automatic grouping of similar objects into sets. Applications: Customer segmentation, anomaly detection, ... Also there: K-Means, hierarchical clustering, ...
- Dimensionality reduction**: Finding the number of important features in a set of data. Applications: SVM, k-NN, ...
- Model selection**: Choosing the best model from a set of candidates. Goal: Minimize your model's performance. Models: Logistic regression, ...
- Preprocessing**: Preparing of data and feature scaling. Applications: Imputation, ...
- News**: Keeping developers: What's new.
- Community**: About us: We're looking for contributors. Note: Machine Learning is a skill.
- Who uses scikit-learn?**: AWeber logo.

Knime, Weka and co: integrated environments

The screenshot displays the Weka Explorer interface. The 'Classify' tab is active, showing the 'Classifier' dropdown set to 'J48 -C 0.25 -M 2'. Under 'Test options', 'Use training set' is selected. The 'Classifier output' window shows the following summary:

```
==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances      144      96 %
Incorrectly Classified Instances     6        4 %
Kappa statistic                    0.94
Mean absolute error                 0.035
Root mean square error
```

The 'Tree View' window shows a decision tree structure:

```
graph TD
    Root((petalwidth)) -- "<= 0.6" --> Node1[Iris-setosa (50.0)]
    Root -- "> 0.6" --> Node2((petalwidth))
    Node2 -- "<= 1.7" --> Node3((petalength))
    Node2 -- "> 1.7" --> Node4[Iris-virginica (46.0/1.0)]
    Node3 -- "<= 4.9" --> Node5[Iris-versicolor (48.0/1.0)]
    Node3 -- "> 4.9" --> Node6((petalwidth))
    Node6 -- "<= 1.5" --> Node7[Iris-virginica (3.0)]
    Node6 -- "> 1.5" --> Node8[Iris-versicolor (3.0/1.0)]
```

The 'Result list' on the left includes options like 'View in main window', 'View in separate window', 'Load model', and 'Visualize classifier tree'. The 'Visualize' window shows a scatter plot with 'Y: petalwidth (Num)' and 'Select Instance' controls. A legend at the bottom identifies 'Iris-versicolor' in red and 'Iris-virginica' in green.