

Introduction to Machine Learning

Gilles Gasso

INSA Rouen - ASI Departement
LITIS Lab

August 10, 2021

Machine Learning: introduction

- Machine Learning \equiv data-based programming
 - ability of computers to learn how to perform tasks (classification, detection, translation...) without being explicitly programmed
 - study of algorithms that **improve their performance** at **some task** based on **experience**

What is Machine Learning?



The rise

Data

- Big Data : continuous increase in data generated
 - Twitter : 50M tweets /day (=7 terabytes)
 - Facebook : 10 terabytes /day
 - Youtube : 50h of uploaded videos /minute
 - 2.9 millions of e-mails /second

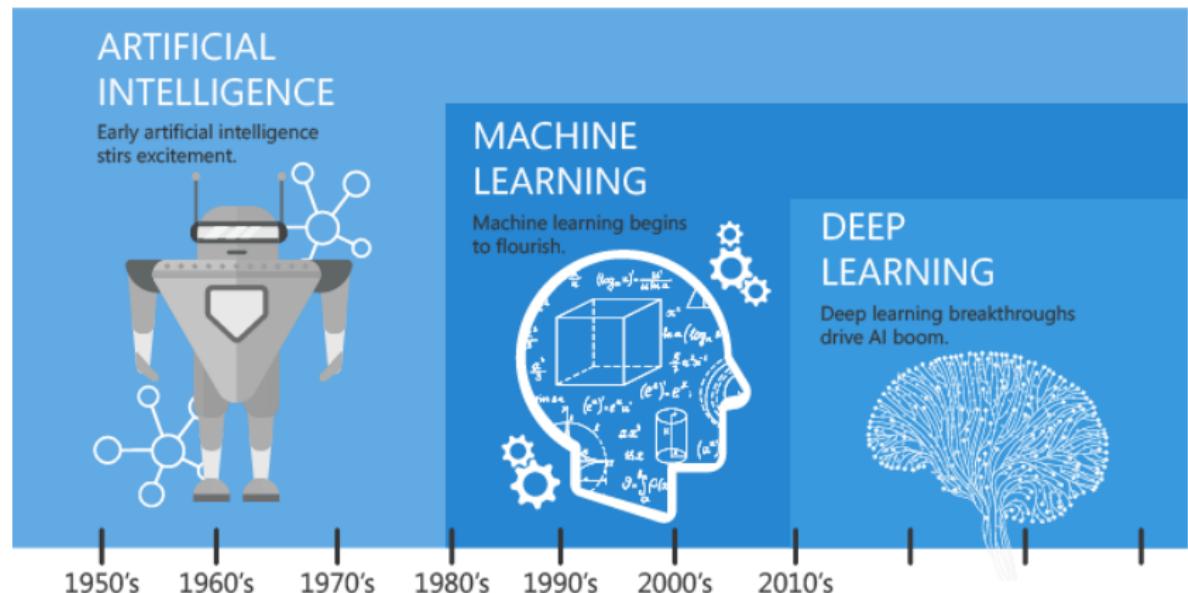
Computing Power

- Moore's law
- Massively distributed computing

The value-adding process

- Interest: from product to customers.
- Data Mining \equiv discovering patterns in large data sets

Historical perspective



Today's AI is Deep Learning (a technique of Machine Learning)

<https://blog.alore.io/machine-learning-and-artificial-intelligence/>

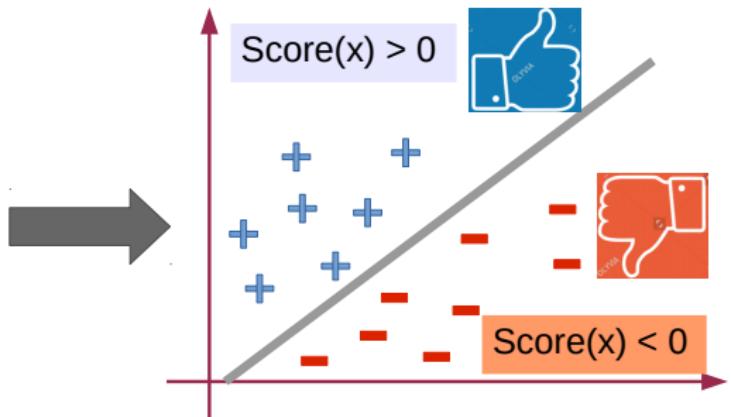
Applications: sentiment analysis

- Classify stored reviews according to users' sentiment

Here is my review about Tarantino's last movie I've watched last saturday. It was really awesome. I enjoy the characters, the script and the music. I fully recommend everyone to go to theater and enjoy the movie.

Yeah, yeah, what to say about this movie ? It is not that bad, but it is not the top Tarantino's movies. Has some good features (music, landscapes) but also bad features

As for the previous movies of Tarantino, I was bored all along the movie. The script was quite complicated to follow up with several references to past events you should know about before. Must not see movie, pass by your way.



Product recommendation

	Feature 1	Feature 2
User 1	?	?
User 2	?	?
User 3	?	?
User 4	?	?
User 5	?	?

User features



	Item 1	Item 2	Item 3	Item 4	Item 5
Feature 1	?	?	?	?	?
Feature 2	?	?	?	?	?

Product features

Matrix factorization
Purchase history of customers

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	0 ?	3	0 ?	3	0 ?
User 2	4	0 ?	0 ?	2	0 ?
User 3	0 ?	0 ?	3	0 ?	0 ?
User 4	3	0 ?	4	0 ?	3
User 5	4	3	0 ?	4	0 ?

https://katbailey.github.io/images/matrix_factorization.png

Purchased item



Digital Répéteur WiFi

Recommended products



AFTERSHOKZ trekz



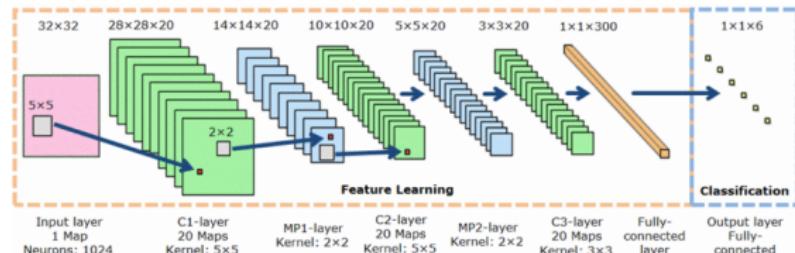
cowin E7 Casque Audio à



Image classification

Labeled training images

Deep classification architecture



https://www.researchgate.net/profile/Y_Nikitin/publication/270163511/figure/download/fig5/AS-295194831409153@1447391340221/MPCNN-architecture-using-alternating-convolutional-and-max-pooling-layers-13.png

Input images and predicted category



Medical diagnosis



Malignant melanoma detection
130 000 images including over 2 000 cases of cancer
Error rate 28 % (human 34 %)

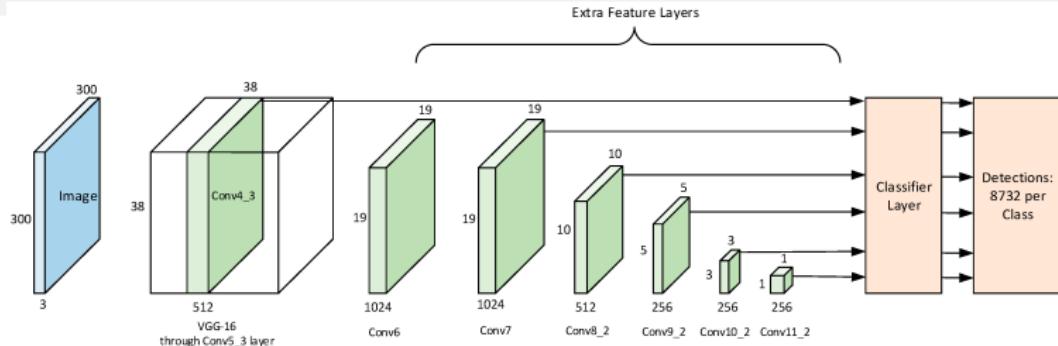


Digital Mammography DREAM Challenge
640 000 mammographies (1209 participants)
false-positive rate decreased by 5 %

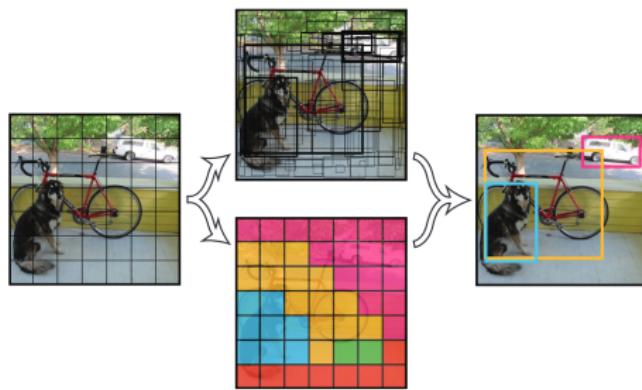


Heart rhythm analysis
500 000 ECG
accuracy 92.6 % (human 80.0 %) sensitivity of 97 %

Object detection

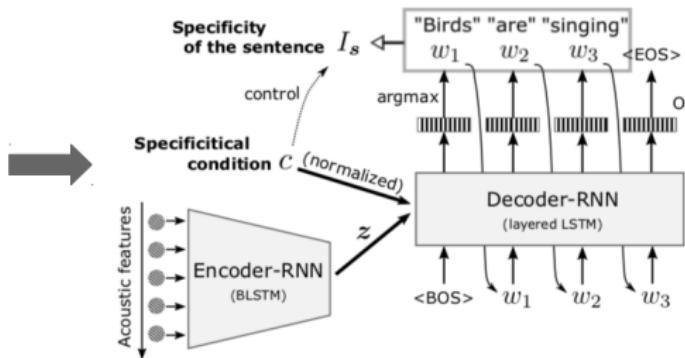
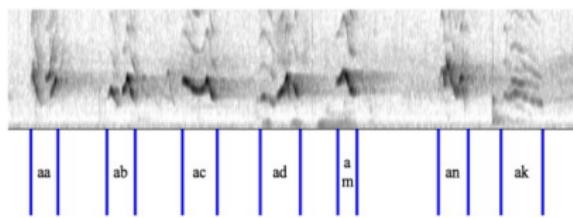


https://www.mdpi.com/applesci/applesci-09-01128/article_deploy/html/images/applesci-09-01128-g004.png



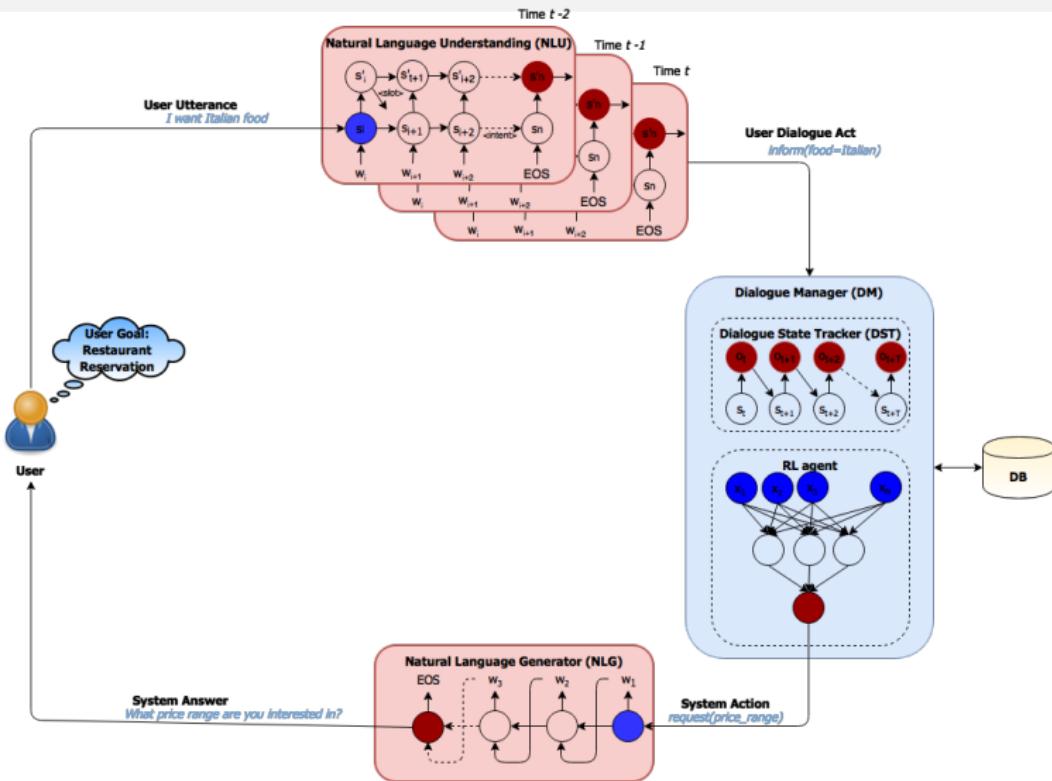
<https://arxiv.org/pdf/1506.02640.pdf>

Audio captioning



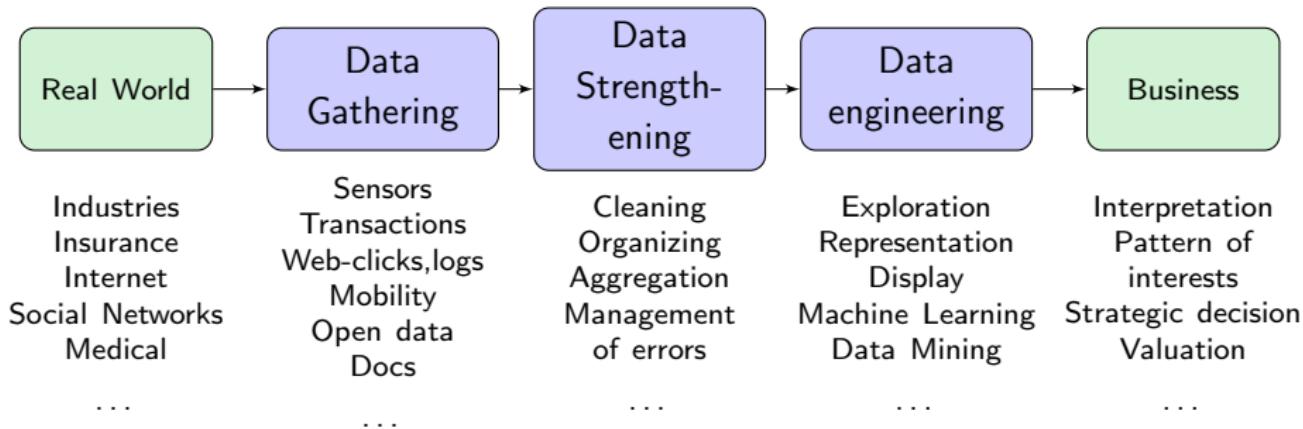
<https://ars.els-cdn.com/content/image/1-s2.0-S1574954115000151-gr1.jpg>

Chatbot

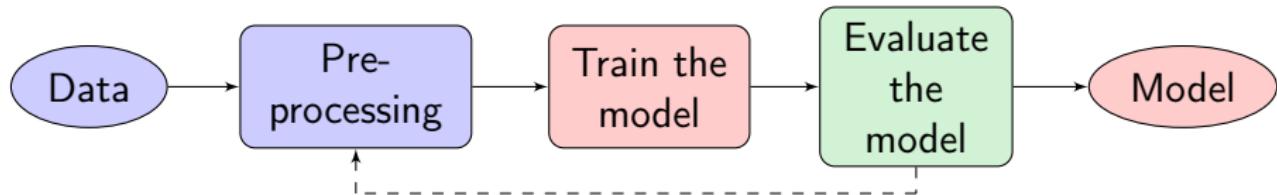


https://miro.medium.com/max/2058/1*LF5T9fsr4w2EqyFJkb-gng.png

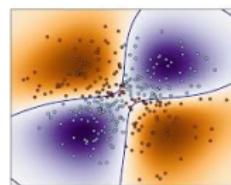
Implementing a Machine Learning project



Chain of the data engineering process



- ① Understand and specify project goals
- ② Pre-processing/visualize/analyze data
- ③ Which ML problem is it?
- ④ Design a solving approach
- ⑤ Evaluate its performance
- ⑥ Go to 2) if needed



Y



Course Goal: Study the steps from 2 to 5

The data

- Information (past experience) are examples with attributes
- Assume the data set consists of N samples

Attributes

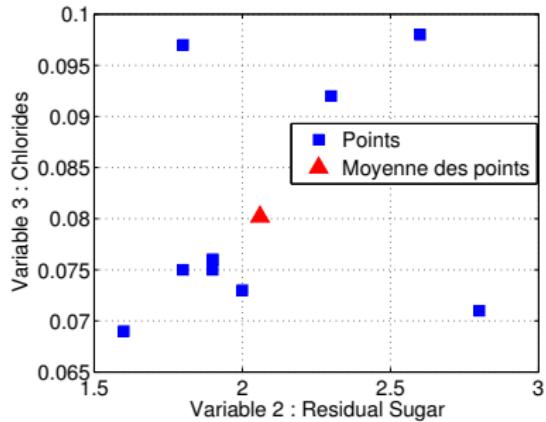
- An attribute is a **property** or **characteristic** of a phenomenon being observed. Also termed **feature** or **variable**

Sample

- It is an entity characterising an object; it is made up of attributes.
- Synonyms : **instance**, **point**, **vector** (usually in \mathbb{R}^d)

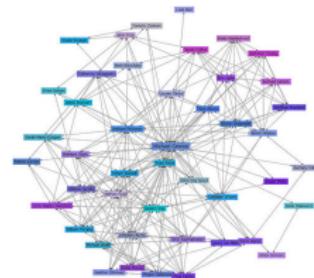
Data : visualization

Points x	Features	citric acid	residual sugar	chlorides	sulfur dioxide
	1	0	1.9	0.076	11
2	0	2.6	0.098	25	
3	0.04	2.3	0.092	15	
Point $x \in \mathbb{R}^4$	0.56	1.9	0.075	17	
5	0	1.9	0.076	11	
6	0	1.8	0.075	13	
7	0.06	1.6	0.069	15	
8	0.02	2	0.073	9	
9	0.36	2.8	0.071	17	
10	0.08	1.8	0.097	15	

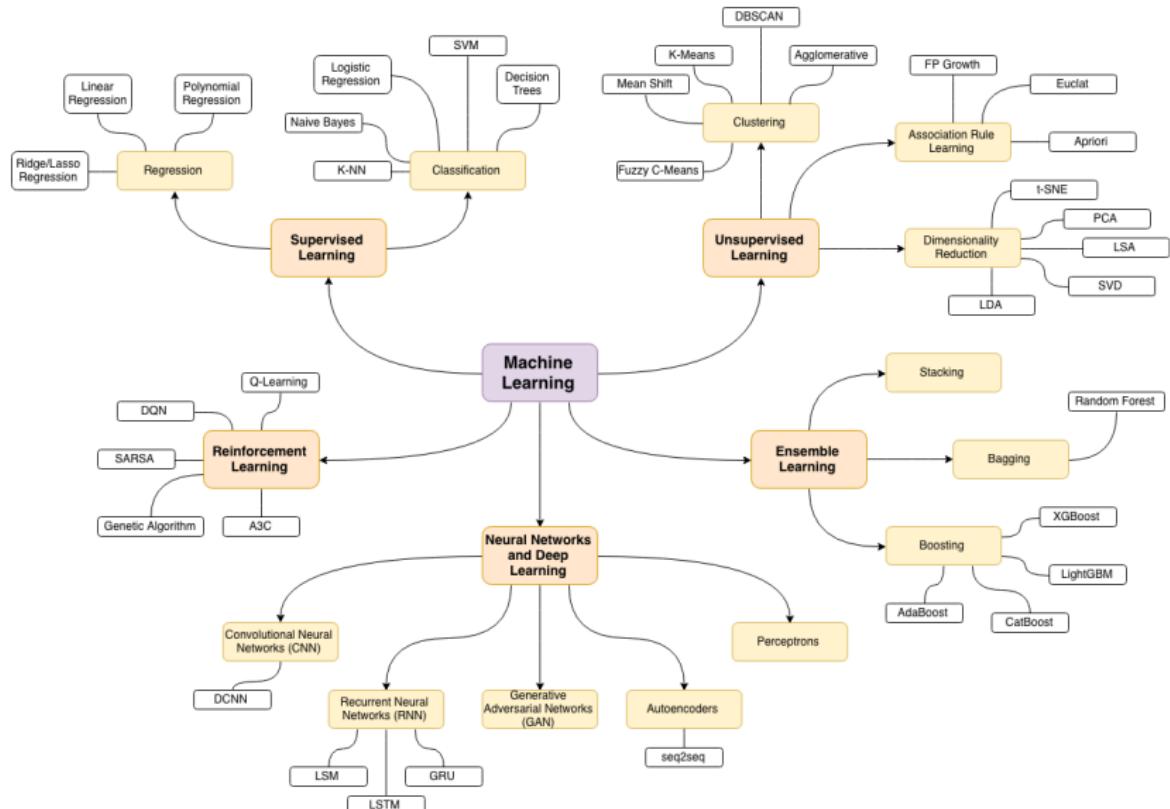


Data types

- Sensors → Quantitative and qualitative variables, ordinales, nominals
- Text → String
- Speech → Time Series
- Images → 2D Data
- Videos → 2D Data + time
- Networks → Graphs
- Stream → Logs, coupons...
- Labels → Expected output prediction



Approaches of Machine Learning



Supervised learning

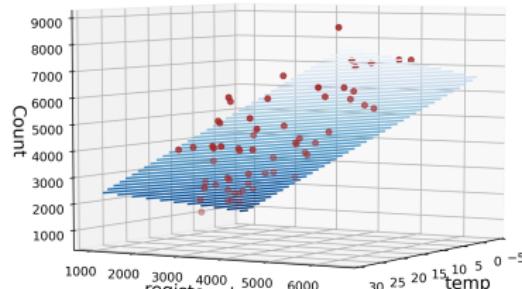
Principle

- Given a set of N training examples $\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}, i = \dots, N\}$, we want to estimate a prediction function $y = f(x)$.
- The *supervision* comes from the label knowledge

Examples

- Image classification, object detection, stock price prediction ...

Régression linéaire avec 2 variables



Unsupervised learning

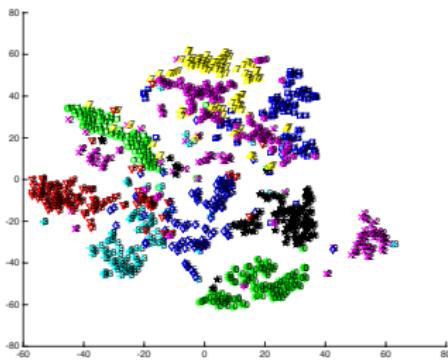
Principle

- Only the $\{x_i \in \mathcal{X}, i = \dots, N\}$ are available. We aim to **describe how data is organized and extract homogeneous subsets from it.**

Examples

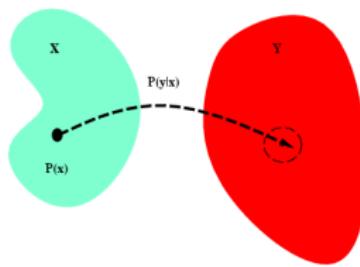
- Applications: Customers segmentation, image segmentation, data visualization, categorization of similar documents ...

0 0
1 1
2 2
3 3
4 4
5 5
6 6
7 7
8 8
9 9



Supervised learning : concept

- Let \mathcal{X} and \mathcal{Y} be two sets. Assume $p(X, Y)$ the joint probability distribution of $(X, Y) \in \mathcal{X} \times \mathcal{Y}$.



- Goal** : find a prediction function $f : \mathcal{X} \rightarrow \mathcal{Y}$ which correctly estimates the output y corresponding to x .
- f belongs to a space \mathcal{H} called **hypothesis class**. Example of \mathcal{H} : set of polynomial functions

Supervised learning: principle

- **Loss function $L(Y, f(X))$**

- evaluates how "close" is the prediction $f(x)$ to the true label y
- it penalizes errors: $L(y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ \geq 0 & \text{if } y \neq f(x) \end{cases}$

- **True risk (expected prediction error)**

$$R(f) = E_{(X,Y)}[L(Y, f(X))] = \int_{\mathcal{X} \times \mathcal{Y}} L(y, f(x)) p(x, y) dx dy$$

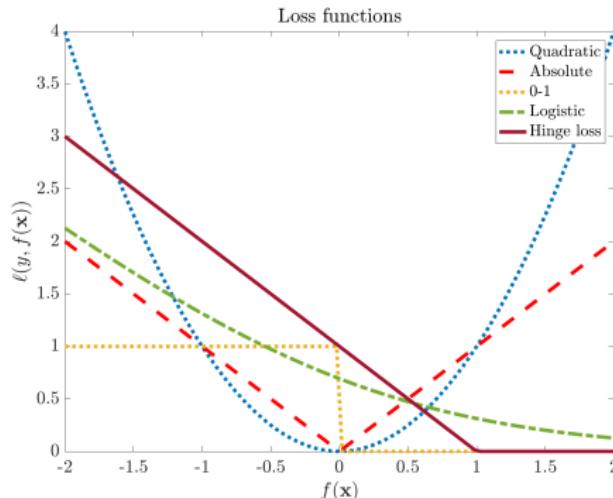
Objective

Identify the prediction function which minimizes the true risk i.e.

$$f^* = \arg \min_{f \in \mathcal{H}} R(f)$$

Loss functions

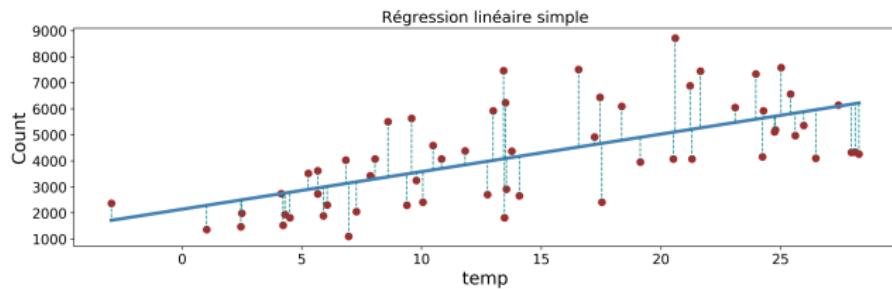
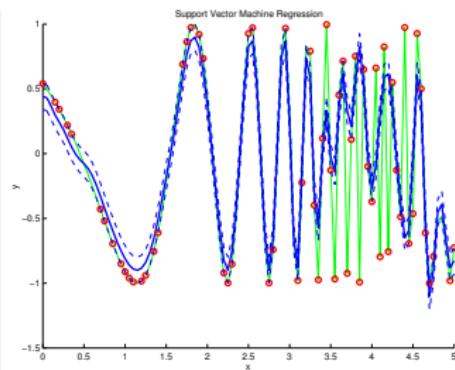
- Quadratic loss : $L(Y, f(X)) = (Y - f(X))^2$
- ℓ_1 loss (absolute deviation): $L(Y, f(X)) = |Y - f(X)|$
- 0 – 1 loss: $L(y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ 1 & \text{otherwise} \end{cases}$
- Hinge loss: $L(y, f(x)) = \max(0, 1 - yf(x))$



Some supervised learning problems

Regression

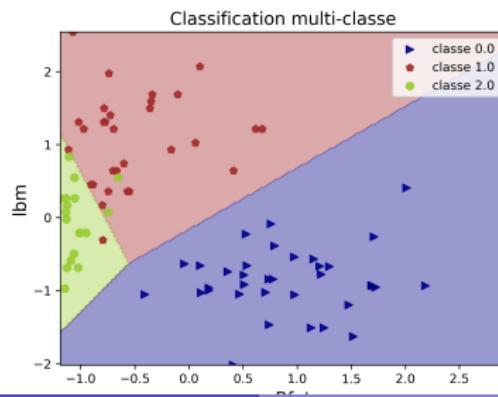
- We talk about **regression** when \mathcal{Y} is a subset of \mathbb{R}^d .
- Usual related loss function:
quadratic loss $(y - f(x))^2$



Some supervised learning problems

Classification

- Output space \mathcal{Y} is an un-ordered discrete set
- Binary classification: $\text{card}(\mathcal{Y}) = 2$
 - Example: $\mathcal{Y} = \{-1, 1\}$
 - Loss functions: 0-1 loss, hinge loss
- Multiclass classification: $\text{card}(\mathcal{Y}) > 2$
 - Example: $\mathcal{Y} = \{1, 2, \dots, K\}$



From true risk to empirical risk

Minimizing the true risk is not doable (in allmost all practical applications)

- The joint distribution $p(X, Y)$ is unknown!
- Only a finite **training set** $\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^N$ is available
- **Empirical risk**

$$R_{emp}(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$$

- Empirical risk minimization

$$\hat{f} = \arg \min_{f \in \mathcal{H}} R_{emp}(f)$$

Overfitting

Overfitting

Empirical risk is not appropriate for model selection: if \mathcal{H} is large enough, $R_{emp}(f) \rightarrow 0$ but the generalized error (true risk) is high.

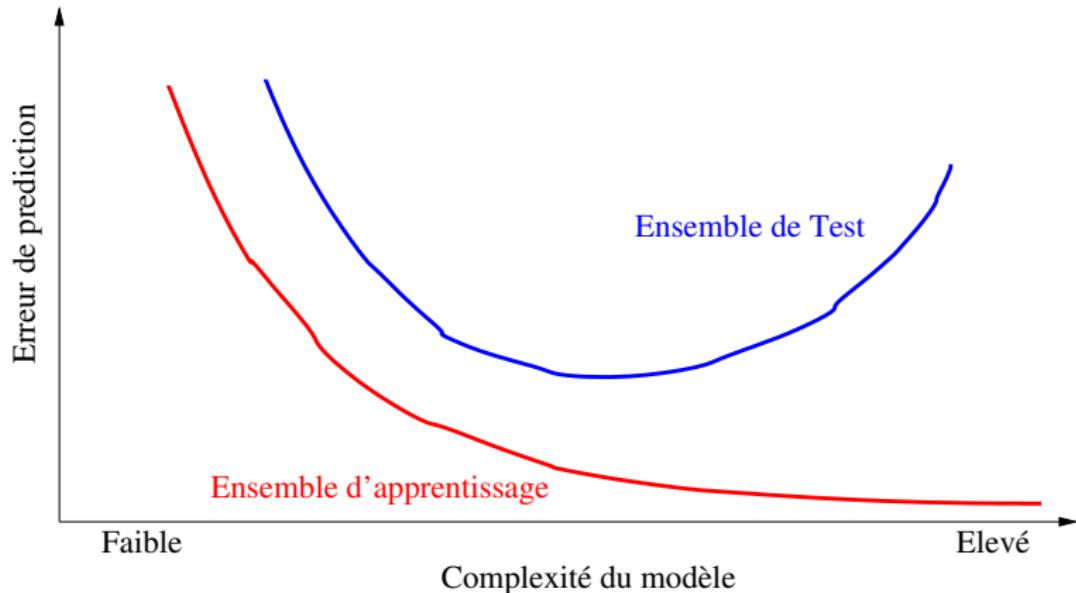
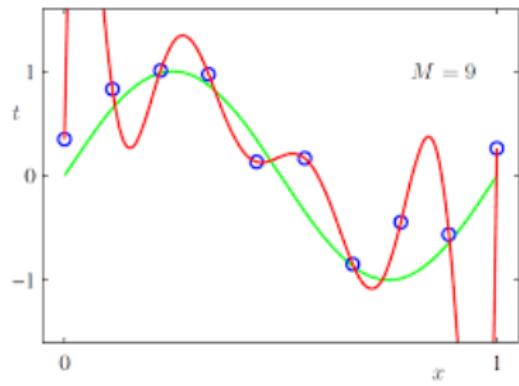
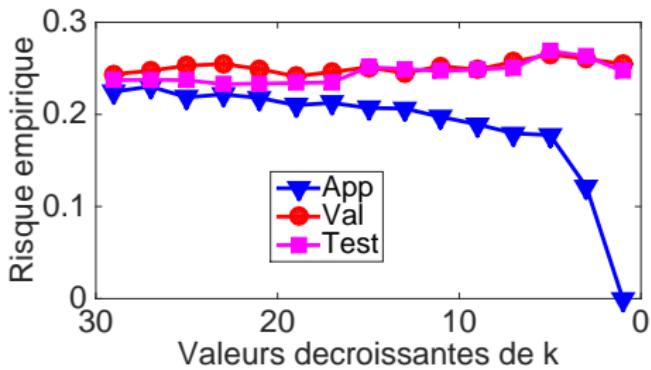


Illustration of overfitting



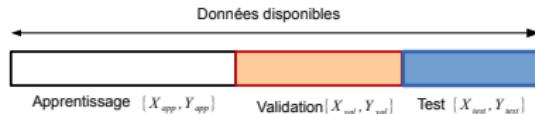
https://www.cs.princeton.edu/courses/archive/spring16/cos495/slides/ML_basics_lecture6_overfitting.pdf

Model selection

- Find in \mathcal{H} the best function f that learned based on the training set will well generalize (low true risk)
- Example : We are looking for a polynomial function of degree α minimizing the risk : $R_{emp}(f_\alpha) = \sum_{i=1}^N (y_i - f_\alpha(x_i))^2$.
- Goal :
 - ① propose a model estimation method in order to choose (approximately) the best model belonging to \mathcal{H} .
 - ② once the model is selected, estimate its generalization error.

Model selection : basic approach

Case 1 : N is really big (large scale \mathcal{D}_N)



- ① Randomly split $\mathcal{D}_N = \mathcal{D}_{train} \cup \mathcal{D}_{val} \cup \mathcal{D}_{test}$
- ② For each α , train f_α based on \mathcal{D}_{train}
- ③ Evaluating its performance on \mathcal{D}_{val} $R_{val} = \frac{1}{N_{val}} \sum_{i \in \mathcal{D}_{val}} L(y_i, f(x_i))$
- ④ Select the model with the best performance on \mathcal{D}_{val}
- ⑤ Test selected model on \mathcal{D}_{test}

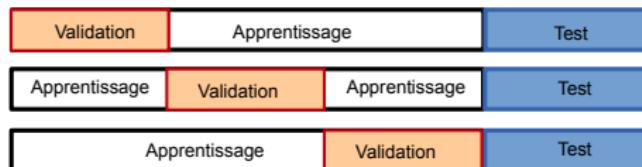
Note

- \mathcal{D}_{test} is used once!

Model selection : Cross-validation

Case 2 : Small or medium scale \mathcal{D}_N

- Estimate the generalization error by re-sampling.
- Principle
 - ① Split \mathcal{D}_N into K sets of equal size.
 - ② For each $k = 1, \dots, K$, train a model by using the $K - 1$ remaining sets and evaluate the model on the k -th part.
 - ③ Average the K error estimates obtained to have the cross-validation error.



Conclusions

To successfully carry out an automatic data processing project

- Clearly identify and spell out the needs.
- Create or obtain data representative of the problem
- Identify the context of learning
- Analyze and reduce data size
- Choose an algorithm and/or a space of hypotheses
- Choose a model by applying the algorithm to pre-processed data
- Validate the performance of the method

Au final ...

Find your way...

