# Machine Learning (MAL)

Mathilde Mougeot

ENSIIE

MAL 2019

# The "Data" phenomena

1. Data tsunami... Todays, data are everywhere.
   - Finance. Transactions data
   - Digital revolution in the Industry. Production data (Supply chain). physical data (Temperature, IR sensors)
   - Marketing/ consumption data. "Click" data
   - On your phone (GPS, mail, musique ...)
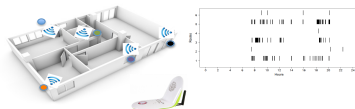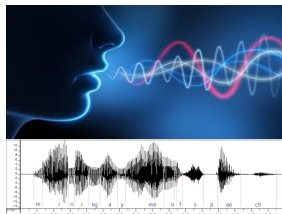
# The "Data" phenomena

1. Data tsunami... Todays, data are everywhere.
   - Finance. Transactions data
   - Digital revolution in the Industry. Production data (Supply chain). physical data (Temperature, IR sensors)
   - Marketing/ consumption data. "Click" data
   - On your phone (GPS, mail, musique ...)

2. Data zoology ... A large variety of data, well or no structured.
   - quantitative, qualitative, binary
   - synchronous, asynchronous, event data
   - ... image data, text data, speech data

# The "Data" phenomena

1. Data tsunami... Todays, data are everywhere.
   - Finance. Transactions data
   - Digital revolution in the Industry. Production data (Supply chain). physical data (Temperature, IR sensors)
   - Marketing/ consumption data. "Click" data
   - On your phone (GPS, mail, musique ...)

2. Data zoology ... A large variety of data, well or no structured.
   - quantitative, qualitative, binary
   - synchronous, asynchronous, event data
   - ... image data, text data, speech data

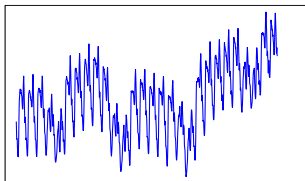3. Data base, data lakes available.
   From small data set to Big data set

# Several data sources

# Illustration. Several energy data



French electrical consumption



Energy Spot prices



Wind turbine power



Industrial equipment

# Potential questions on energy data



Exploratory analysis
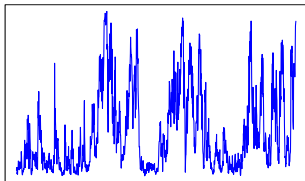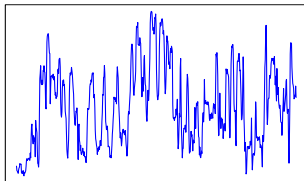


Virtual sensor



Monitoring



Forecasting electrical consumption

# Machine Learning. Statistical settings (1/3) :

Unsupervised learning. (Inputs $X$) :     $X \in \mathcal{X} \ (\mathbb{R}^p)$



Exploratory Analysis.



Clustering



Graph Analysis

# Machine Learning settings (2/3)

- Supervised regression learning $(Y, X)$ :

$$Y \in \mathbb{R}, \ X \in \mathcal{X} \ (\mathbb{R}^p) \qquad \boxed{Y = \mathcal{M}_{\text{data set}}(X) + \epsilon}$$

# Machine Learning settings (2/3)

- Supervised regression learning $(Y, X)$ :

  $Y \in \mathbb{R}$, $X \in \mathcal{X}$ $(\mathbb{R}^p)$ $\quad$ $\boxed{Y = \mathcal{M}_{\text{data set}}(X) + \epsilon}$

  - Example. Parametric models.

    $\mathcal{M}(X) : Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_p X_p + \epsilon$

    $\rightarrow$ Probabilistic or not, depending on the law assumption on $\epsilon$.

    Data Set : $n$ observations : $(y_i, x_i)$

    $\rightarrow$ to estimate (to compute) the parameters $\hat{\beta}$ of the model

# Machine Learning settings (2/3)

- Supervised regression learning $(Y, X)$ :

  $Y \in \mathbb{R}$, $X \in \mathcal{X}$ $(\mathbb{R}^p)$    $\boxed{Y = \mathcal{M}_{\text{data set}}(X) + \epsilon}$

  - Example. Parametric models.
    $\mathcal{M}(X) : Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_p X_p + \epsilon$
    $\rightarrow$ Probabilistic or not, depending on the law assumption on $\epsilon$.

    Data Set : $n$ observations : $(y_i, x_i)$
    $\rightarrow$ to estimate (to compute) the parameters $\hat{\beta}$ of the model

  - Non Parametric models No analytical (complex) expression, infinity of parameters



$\rightarrow$ Models : Kernels, decision trees, neural networks,...

<span style="color:red">Crucial question : How to find the correct model ?</span>

# Machine Learning settings (3/3)

- Supervised Classification learning $(Y, X)$ :

  $Y \in G_1, ..., G_K,\ X \in \mathcal{X}\ (\mathbb{R}^p)$ $\quad\boxed{Y = \mathcal{M}_{\text{data set}}(X)}$

# Machine Learning settings (3/3)

- Supervised Classification learning $(Y, X)$ :

  $Y \in G_1, ..., G_K, \; X \in \mathcal{X} \; (\mathbb{R}^p)$  $\boxed{Y = \mathcal{M}_{\text{data set}}(X)}$

  - Example. Parametric model. Logistic regression
    $\mathcal{M}(X) : P(Y = G_k / X = x) = \frac{1 + e^{x\beta}}{1 - e^{x\beta}}$

    Data Set : $n$ observations : $(y_i, x_i)$
    $\rightarrow$ to estimate (to compute) the parameters $\hat{\beta}$ of the model

# Machine Learning settings (3/3)

- Supervised Classification learning $(Y, X)$ :

  $Y \in G_1, ..., G_K,\ X \in \mathcal{X}\ (\mathbb{R}^p)$ $\boxed{Y = \mathcal{M}_{\text{data set}}(X)}$
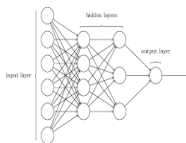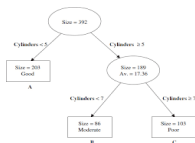
  - Example. Parametric model. Logistic regression
    $\mathcal{M}(X) : P(Y = G_k / X = x) = \frac{1 + e^{x\beta}}{1 - e^{x\beta}}$

    Data Set : $n$ observations : $(y_i, x_i)$
    $\rightarrow$ to estimate (to compute) the parameters $\hat{\beta}$ of the model

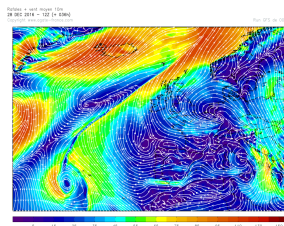  - Non Parametric models No analytical (complex) expression, infinity of parameters



$\rightarrow$ Models : Kernels, decision trees, neural networks,...

  Crucial question : How to find the correct classification model ?

# Machine Learning or Statistical Modeling vs Physical, simulation models

- **Physical Modelling** mosly based on physical equations :

  - Explicit equation
    $PV = nRT...$

  - Simulation models. Partial differential equations,...
    $\rightarrow$ Need of a numerical model to study the dynamic and the evolution of a model. Ex : Navier-Stokes equation

# Outline

Machine learning. The 2019 choice...

**1** Supervised setting. Classification.

- Parametric models. Bayes model. LDA. QDA
- Performance criteria.
- Non Parametric models. Classification trees.
- Ensemble methods.
  Bagging. Random Forest. Boosting. Stacking.

**2** Supervised setting. Regression.

- Non Parametric models. Regression trees.
- Ensemble methods.
  Bagging. Random Forest. Boosting. Stacking.

**3** Unsupervised setting. Clustering.

- K-means.
- Spectral clustering.