


# Data Science for Actuaries (ACT6100)

Arthur Charpentier

Intro # 2 (Glossaire)

automne 2020

 <https://github.com/freakonometrics/ACT6100/>

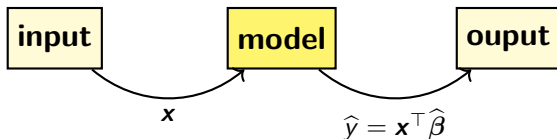
# Supervised vs. Non-Supervised

## Supervised Learning

In **supervised learning**, the algorithms learns a function that maps an input to an output based on example input-output pairs

**Example:** (Linear) regression (see [STT5100](#))

$$(y_i, \mathbf{x}_i), \begin{cases} y_i = \text{output} \in \mathbb{R} \\ \mathbf{x}_i = \text{input} \in \mathbb{R}^p \end{cases}$$



The **prediction**  $\hat{y}_i$  can be compared with the actual value  $y_i$  through some **loss function**  $\ell$ , e.g.  $\ell(y_i, \hat{y}_i) = [y_i - \hat{y}_i]^2$

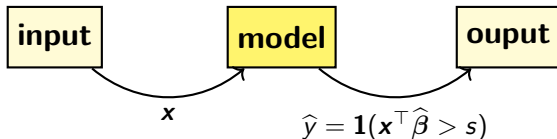
# Supervised vs. Non-Supervised

## Supervised Learning

In **supervised learning**, the algorithms learns a function that maps an input to an output based on example input-output pairs

**Example:** Logistic regression (see **STT5100**), or **classification**

$$(y_i, \mathbf{x}_i), \begin{cases} y_i = \text{output} \in \{0, 1\} \\ \mathbf{x}_i = \text{input} \in \mathbb{R}^p \end{cases}$$



The **prediction**  $\hat{y}_i$  can be compared with the actual value  $y_i$  through some **loss function**  $\ell$ , e.g.  $\ell(y_i, \hat{y}_i) = \mathbf{1}(\hat{y}_i \neq y_i)$

# Supervised vs. Non-Supervised

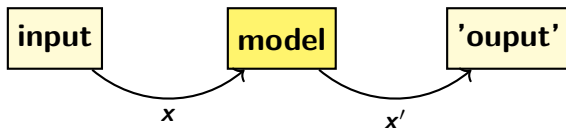
## Supervised Learning

In **unsupervised learning**, the algorithm looks for previously undetected patterns in a data set with no pre-existing labels

(Learn *useful properties of the structure of the data*)

**Example:** Dimension reduction

$$(\mathbf{x}_i), \mathbf{x}_i \in \mathbb{R}^p \rightarrow \mathbf{x}'_i \in \mathbb{R}^{p'}, p' < p$$



possibly with losing too much **information**  
(see **principal component analysis**)

# Supervised vs. Non-Supervised

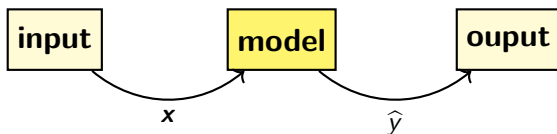
## Supervised Learning

In **unsupervised learning**, the algorithm looks for previously undetected patterns in a data set with no pre-existing labels

(Learn *useful properties of the structure of the data*)

**Example:** Clustering

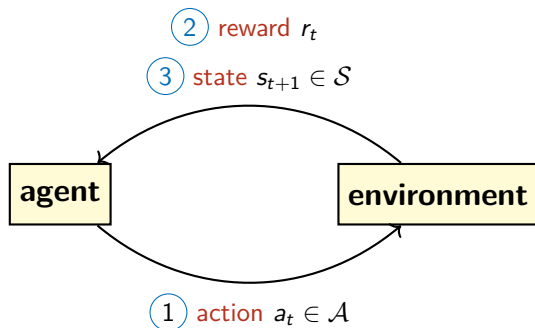
$$(\mathbf{x}_i), \mathbf{x}_i \in \mathbb{R}^p \rightarrow y'_i \in \{L_1, L_2, \dots, L_k\}$$



possibly with strong consistency of the groups (see **k-means** or **hierarchical clustering**)

# Sequential Learning

**Example:** Reinforcement Learning

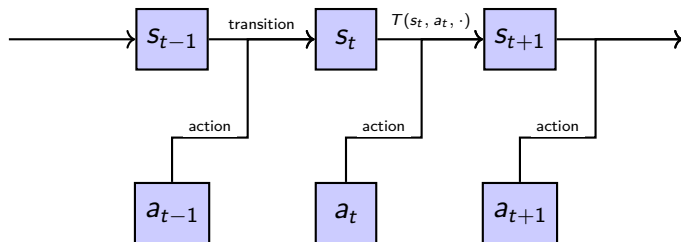


## Reinforcement Learning

In reinforcement learning, the algorithm ought to take actions in an environment in order to maximize the notion of cumulative reward

# Sequential Learning

- ▶ the learner takes an **action**  $a_t \in \mathcal{A}$  (while at state  $s_t$ )
- ▶ the learner obtains a (short-term) **reward**  $r_t \in \mathcal{R}$
- ▶ then the **state** of the world becomes  $s_{t+1} \in \mathcal{S}$



see also **multi-armed bandits** in decision theory

# Loss & Risk for Supervised Learning

## Empirical Risk - for sample $\mathcal{D}_n$

The **empirical risk** associated with  $\hat{m}_n$  is  $\hat{\mathcal{R}}_n(\hat{m}_n)$

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_n(\mathbf{x}_i))$$

## Average Risk - under $\mathbb{P}$

The **average risk** associated with  $\hat{m}_n$  is  $\mathcal{R}_{\mathbb{P}}(\hat{m}_n)$

$$\mathbb{E}_{\mathbb{P}}[\ell(Y, \hat{m}_n(\mathbf{X}))]$$



# Cross-Validation for Supervised Learning

## Hold-Out Cross Validation

1. Split  $\{1, 2, \dots, n\}$  in  $T$  (training) and  $V$  (validation)
2. Estimate  $\hat{m}$  on training sample  $(y_i, \mathbf{x}_i)$ ,  $i \in T$ ,  $\hat{m}_T$

$$\operatorname{argmin}_{m \in \mathcal{M}} \frac{1}{|T|} \sum_{i \in V} \ell(y_i, m(\mathbf{x}_i))$$

3. Compute

$$\frac{1}{|V|} \sum_{i' \in V} \ell(y_{i'}, \hat{m}_T(\mathbf{x}_{i'}))$$

## Leave-one-Out CV

1. Estimate  $n$  models : estimate  $\hat{m}_{-j}$  on sample  $(y_i, \mathbf{x}_i)$ ,  $i \in \{1, \dots, n\} \setminus \{j\}$ ,  $\hat{m}_{-j}$

$$\operatorname{argmin}_{m \in \mathcal{M}} \frac{1}{n-1} \sum_{i \neq j} \ell(y_i, m(\mathbf{x}_i))$$

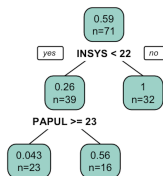
2. Compute

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}_{-i}(\mathbf{x}_i))$$

For a classification problem  $y \in \{0, 1\}$

## Trees (CART)

Each interior node corresponds to one of the input variables; there are edges to children for each of the possible values of that input variable. Each leaf represents a value of the target variable given the values of the input variables represented by the path from the root to the leaf.

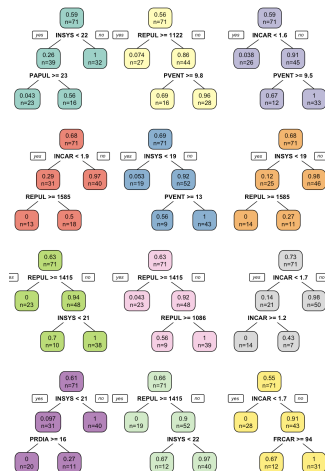


# Aggregation & Parallel Learning (bagging)

## Bagging (Bootstrap + Aggregation)

1. For  $k = 1, \dots$ 
  - (i) draw a bootstrap sample from  $(y_i, \mathbf{x}_i)$ 's
  - (ii) estimate a model  $\hat{m}_k$  on that sample
2. The final model is

$$m^*(\cdot) = \frac{1}{\kappa} \sum_{k=1}^{\kappa} \hat{m}_k(\cdot)$$



# Aggregation & Sequential Learning (boosting)

Learn slowly - **weak learner** - and learn from your mistakes

## Boosting & Sequential Learning

Starting from  $m_0(\cdot) = \bar{y}$ , update

$$m_k(\cdot) = m_{k-1}(\cdot) + \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \sum_{i=1}^n \ell(\underbrace{y_i - m_{k-1}(\mathbf{x}_i)}_{\varepsilon_i}, h(\mathbf{x}_i)) \right\}$$

# Overfit, Generalization, Bias & Variance

## Penalization

regularization is the process of adding information in order to solve an ill-posed problem or to prevent overfitting

$$\operatorname{argmin}_{m \in \mathcal{M}} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) + \lambda \operatorname{penalty}(m) \right\}$$

Usually, it yields **biased** estimates...

