

# De la théorie de la compression d'échantillons aux algorithmes de méta-apprentissage

(Adaptation d'une présentation de Pascal Germain)

**Benjamin Leblanc**

<https://benthewhite.github.io/>

Université Laval, département d'informatique et de génie logiciel

22 mars 2024



## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

## Distribution génératrice des données

Chaque observation provient d'une **distribution**  $D$  sur  $\mathcal{Z}$ .

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

## Distribution génératrice des données

Chaque observation provient d'une **distribution**  $D$  sur  $\mathcal{Z}$ .

## Ensemble d'apprentissage

$$S := \{ z_1, z_2, \dots, z_n \} \sim D^n$$

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

## Distribution génératrice des données

Chaque observation provient d'une **distribution**  $D$  sur  $\mathcal{Z}$ .

## Ensemble d'apprentissage

$$S := \{ z_1, z_2, \dots, z_n \} \sim D^n$$

## Algorithme d'apprentissage

$$A(S) \longrightarrow h$$

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

## Distribution génératrice des données

Chaque observation provient d'une **distribution**  $D$  sur  $\mathcal{Z}$ .

## Ensemble d'apprentissage

$$S := \{ z_1, z_2, \dots, z_n \} \sim D^n$$

## Algorithme d'apprentissage

$$A(S) \longrightarrow h$$

## Prédicteur (où hypothèse)

$$h : \mathcal{X} \rightarrow \mathcal{Y}, \quad h \in \mathcal{H}$$

# Définitions

Une **observation**  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  est une paire **variables explicatives - étiquette**.

## Distribution génératrice des données

Chaque observation provient d'une **distribution**  $D$  sur  $\mathcal{Z}$ .

## Ensemble d'apprentissage

$$S := \{ z_1, z_2, \dots, z_n \} \sim D^n$$

## Algorithme d'apprentissage

$$A(S) \longrightarrow h$$

## Prédicteur (où hypothèse)

$$h : \mathcal{X} \rightarrow \mathcal{Y}, \quad h \in \mathcal{H}$$

## Fonction de perte

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$



**Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec**

**Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec**

S

| ID  | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | $y$<br>Prix |
|-----|----------------------|--------------------------|----------------------------|-----|----------------|-------------|
| 1   | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1000        |
| 2   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1230        |
| ⋮   | ⋮                    | ⋮                        | ⋮                          | ⋮   | ⋮              | ⋮           |
| ⋮   | ⋮                    | ⋮                        | ⋮                          | ⋮   | ⋮              | ⋮           |
| $n$ | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 829         |

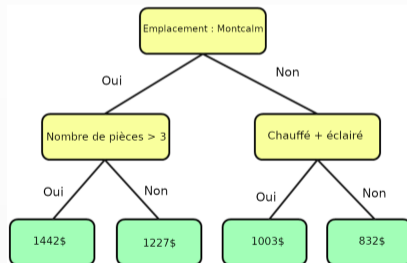
# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$S$

| ID  | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | $y$<br>Prix |
|-----|----------------------|--------------------------|----------------------------|-----|----------------|-------------|
| 1   | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1000        |
| 2   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1230        |
| ... | ...                  | ...                      | ...                        | ... | ...            | ...         |
| $n$ | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 829         |

$$A(S) = h$$



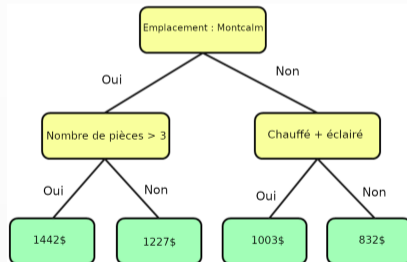
# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$S$

| ID  | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | $y$<br>Prix |
|-----|----------------------|--------------------------|----------------------------|-----|----------------|-------------|
| 1   | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1000        |
| 2   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1230        |
| ... | ...                  | ...                      | ...                        | ... | ...            | ...         |
| $n$ | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 829         |

$A(S) = h$



$\mathcal{H}$  : les différents arbres de décision qui existent

# Le problème de la généralisation

L'algorithme d'apprentissage se sert uniquement de l'ensemble d'entraînement  $S$  :

$$\hat{\mathcal{L}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i)$$

# Le problème de la généralisation

L'algorithme d'apprentissage se sert uniquement de l'ensemble d'entraînement  $S$  :

$$\hat{\mathcal{L}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i)$$

Le but : Minimiser la perte en moyenne sur  $D$

$$\mathcal{L}_D(h) := \mathbf{E}_{z \sim D} \ell(h(\mathbf{x}_i), y_i)$$

# Le problème de la généralisation

L'algorithme d'apprentissage se sert uniquement de l'ensemble d'entraînement  $S$  :

$$\hat{\mathcal{L}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i)$$

Le but : Minimiser la perte en moyenne sur  $D$

$$\mathcal{L}_D(h) := \mathbf{E}_{z \sim D} \ell(h(\mathbf{x}_i), y_i)$$

Bornes de généralisation PAC (Probablement Approximativement Correct)

# Le problème de la généralisation

L'algorithme d'apprentissage se sert uniquement de l'ensemble d'entraînement  $S$  :

$$\hat{\mathcal{L}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i)$$

Le but : Minimiser la perte en moyenne sur  $D$

$$\mathcal{L}_D(h) := \mathbf{E}_{z \sim D} \ell(h(\mathbf{x}_i), y_i)$$

Bornes de généralisation PAC (Probablement Approximativement Correct)

« Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ” »



# Le problème de la généralisation

L'algorithme d'apprentissage se sert uniquement de l'ensemble d'entraînement  $S$  :

$$\hat{\mathcal{L}}_S(h) := \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i)$$

Le but : Minimiser la perte en moyenne sur  $D$

$$\mathcal{L}_D(h) := \mathbf{E}_{z \sim D} \ell(h(\mathbf{x}_i), y_i)$$

Bornes de généralisation PAC (Probablement Approximativement Correct)

« Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ” »

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

- Classe finie d'hypothèses  $\mathcal{H}$  :

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

- Classe finie d'hypothèses  $\mathcal{H}$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

- Classe finie d'hypothèses  $\mathcal{H}$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$

- Classe dénombrable d'hypothèses  $h_i$ , avec avec probabilité *a priori*  $p(h_i)$  :

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

- Classe finie d'hypothèses  $\mathcal{H}$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$

- Classe dénombrable d'hypothèses  $h_i$ , avec avec probabilité *a priori*  $p(h_i)$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{p(h)\delta} \right)}$$

## PAC (Probablement Approximativement Correct)

Avec probabilité au moins “ $1-\delta$ ”, la perte de  $h$  sera au plus “ $\varepsilon(\cdot, \dots, \cdot)$ ”

$$\Pr_{S \sim D^n} \left( \mathcal{L}_D(h) \leq \varepsilon(\hat{\mathcal{L}}_S(h), n, \delta, \dots) \right) \geq 1-\delta$$

- Hypothèse unique  $h$  :

$$\mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{\delta} \right)}.$$

- Classe finie d'hypothèses  $\mathcal{H}$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$

- Classe dénombrable d'hypothèses  $h_i$ , avec avec probabilité *a priori*  $p(h_i)$  :

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \hat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{1}{p(h)\delta} \right)}$$

- Classe indénombrable d'hypothèses : Dimension VC, Complexité de Rademacher, PAC-Bayes...



## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples

- Borne de généralisation pour classificateurs binaires

- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage

- Résultats préliminaires - classification binaire

**Initialisation** : LITTLESTONE et WARMUTH 1986 : “Relating Data Compression and Learnability”

## The Set Covering Machine (SCM) et ses variants

MARCHAND et SHAWE-TAYLOR 2002 : “The Set Covering Machine”

MARCHAND et SOKOLOVA 2005 : “Learning with Decision Lists of Data-Dependent Features”

LAVIOLETTE et al. 2005 : “Margin-Sparsity Trade-Off for the Set Covering Machine”

DROUIN et al. 2014 : “Learning interpretable models of phenotypes from whole genome sequences with the Set Covering Machine”

GODON et al. 2022 : “RandomSCM: interpretable ensembles of sparse classifiers tailored for omics data”

## Autres

CAMPI et GARATTI 2023 : “Compression, Generalization and Learning”

PACCAGNAN et al. 2023 : “The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance”



- Classe indénombrable d'hypothèses : nouvel angle d'attaque

- Classe indénombrable d'hypothèses : nouvel angle d'attaque
- Classe dénombrable d'hypothèses : le terme  $|\mathcal{H}|$  peut être très pénalisant

- Classe indénombrable d'hypothèses : nouvel angle d'attaque
- Classe dénombrable d'hypothèses : le terme  $|\mathcal{H}|$  peut être très pénalisant

$$\forall h \in \mathcal{H}, \mathcal{L}_D(h) \leq \widehat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$

- Classe indénombrable d'hypothèses : nouvel angle d'attaque
- Classe dénombrable d'hypothèses : le terme  $|\mathcal{H}|$  peut être très pénalisant
- Permet de n'utiliser qu'un sous-ensemble des observations

$$\forall h \in \mathcal{H}, \quad \mathcal{L}_D(h) \leq \widehat{\mathcal{L}}_S(h) + \sqrt{\frac{1}{2n} \log \left( \frac{|\mathcal{H}|}{\delta} \right)}$$



Un prédicteur  $h$  peut être **compressé**  $h_i^\mu$  si obtenu par un nouvel algorithme  $R$  dépendant de deux sources d'informations complémentaires :

Un prédicteur  $h$  peut être **compressé**  $h_i^\mu$  si obtenu par un nouvel algorithme  $R$  dépendant de deux sources d'informations complémentaires :

- Un **ensemble de compression**  $S_i$ , un sous-ensemble de  $S$  :

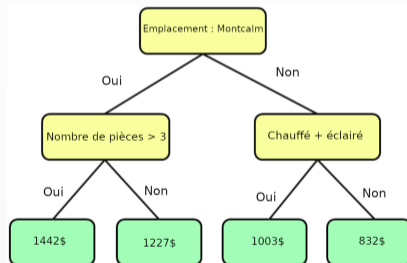
# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$S$

| ID  | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | $y$<br>Prix |
|-----|----------------------|--------------------------|----------------------------|-----|----------------|-------------|
| 1   | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1000        |
| 2   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1230        |
| ⋮   | ⋮                    | ⋮                        | ⋮                          | ⋮   | ⋮              | ⋮           |
| $n$ | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 829         |

$$A(S) = h$$

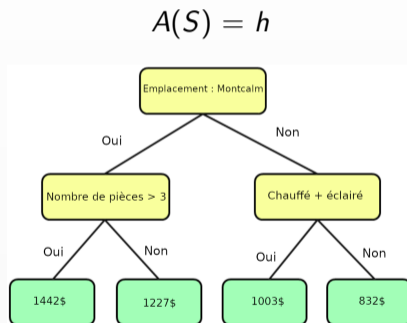


# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$$S_i = S_{\{1,29,263,1902\}}$$

| ID   | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | y<br>Prix |
|------|----------------------|--------------------------|----------------------------|-----|----------------|-----------|
| 1    | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1442      |
| 29   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1227      |
| 263  | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1003      |
| 1902 | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 832       |



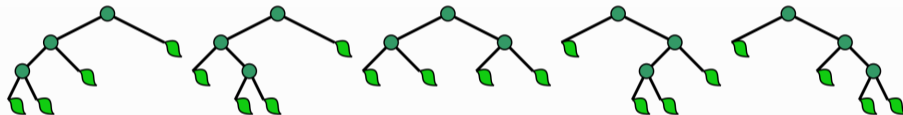
Un **prédicteur compressé**  $h_i^\mu$  est un prédicteur obtenu par un algorithme  $R$  dépendant de deux sources d'informations complémentaires :

- Un **ensemble de compression**  $S_i$ , un sous-ensemble de  $S$  :

Un **prédicteur compressé**  $h_i^\mu$  est un prédicteur obtenu par un algorithme  $R$  dépendant de deux sources d'informations complémentaires :

- Un **ensemble de compression**  $S_i$ , un sous-ensemble de  $S$  :
- Un **message**  $\mu \in \mathcal{M}_i$  qui contient de l'information additionnelle pour décrire  $h_i^\mu$ .

Profondeur?  
Nombre de feuilles?

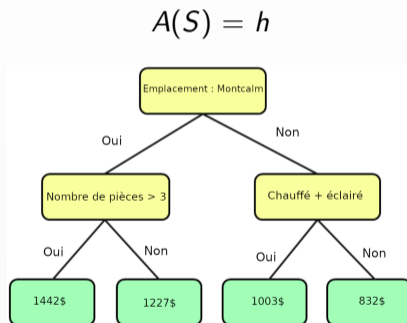


# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$$S_i = S_{\{1,29,263,1902\}}$$

| ID   | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | y<br>Prix |
|------|----------------------|--------------------------|----------------------------|-----|----------------|-----------|
| 1    | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1442      |
| 29   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1227      |
| 263  | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1003      |
| 1902 | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 832       |





Un **prédicteur compressé**  $h_i^\mu$  est un prédicteur obtenu par un algorithme  $R$  dépendant de deux sources d'informations complémentaires :

- Un **ensemble de compression**  $S_i$ , un sous-ensemble de  $S$  :
- Un **message**  $\mu \in \mathcal{M}_i$  qui contient de l'information additionnelle pour décrire  $h_i^\mu$ .

Un **prédicteur compressé**  $h_i^\mu$  est un prédicteur obtenu par un algorithme  $R$  dépendant de deux sources d'informations complémentaires :

- Un **ensemble de compression**  $S_i$ , un sous-ensemble de  $S$  :
- Un **message**  $\mu \in \mathcal{M}_i$  qui contient de l'information additionnelle pour décrire  $h_i^\mu$ .

Étant donné  $S_i \in \mathcal{Z}^{|\mathcal{I}|}$  et  $\mu \in \mathcal{M}_i$ , une **fonction de reconstruction**  $\mathcal{R}$  donne un prédicteur :

$$h_i^\mu = \mathcal{R}(S_i, \mu).$$

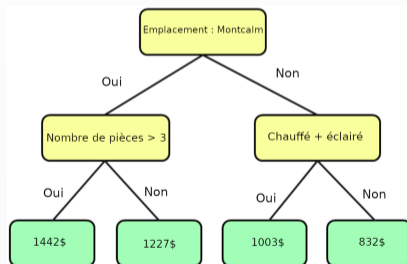
# Exemple

Tâche : prédire le loyer mensuel d'un appartement dans la ville de Québec

$$S_i = S_{\{1,29,263,1902\}}$$

| ID   | $x_1$<br>Emplacement | $x_2$<br>Nombre de pièce | $x_3$<br>Chauffé + Éclairé | ... | $x_d$<br>Étage | $y$<br>Prix |
|------|----------------------|--------------------------|----------------------------|-----|----------------|-------------|
| 1    | Montcalm             | 3 1/2                    | Oui                        | ... | 3e             | 1442        |
| 29   | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1227        |
| 263  | Montcalm             | 4 1/2                    | Oui                        | ... | 1er            | 1003        |
| 1902 | Sainte-Foy           | 2 1/2                    | Non                        | ... | 2e             | 832         |

$$A(S) = h$$
$$R(S_i, \mu) = h$$



## SVM : *Support Vector Machine* (marge rigide)

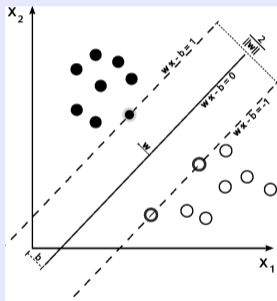


Image : Wikipedia

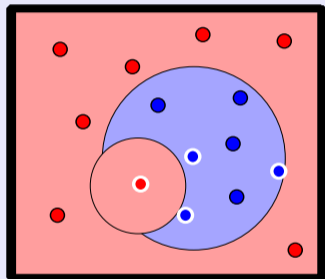
L'algorithme d'apprentissage du SVM agit comme sa propre fonction de reconstruction

$$\text{SVM}(S) = h_i^\mu = \text{SVM}(S_i)$$

with  $S_i = \{\text{vecteurs de support}\}$   
and  $\mu = \emptyset$

## SCM : Set Covering Machine (MARCHAND et SHAWE-TAYLOR 2002)

→ conjonction variables explicatives booléennes créées



Le SCM apprend des *features*  $b_{i,j} : \mathcal{X} \rightarrow \{-1, +1\}$

Chaque variable explicative est une boule  $b_{i,j} \in \mathcal{B}$  définie par un centre  $(x_i, y_i)$  et une bordure  $(x_j, y_j)$  :

$$b_{i,j}(x) := \begin{cases} +y_i & \text{if } \|x - x_i\| \leq \|x_i - x_j\| + \epsilon \cdot y_i, \\ -y_i & \text{sinon.} \end{cases}$$

Le SCM apprend une conjonction comme classifieur :

$$h_i^\mu(x) := \bigwedge_{b_{i,j} \in \mathcal{B}} b_{i,j}(x) \quad (\text{où } +1 \equiv \text{Vrai} \quad -1 \equiv \text{Faux})$$

avec  $S_i = \{\text{les points « centre » et « bordure »}\}$   
and  $\mu = \{\text{les indices parmi } S_i \}$

## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

Étant donné  $h_i^\mu : \mathcal{X} \rightarrow \{-1, +1\}$ , et  $\mathcal{L}_D^{01}(h_i^\mu) = \mathbf{E}_{(x,y) \sim D} I(h_i^\mu(x) \neq y)$  la perte zéro-un.

Étant donné  $h_i^\mu : \mathcal{X} \rightarrow \{-1, +1\}$ , et  $\mathcal{L}_D^{01}(h_i^\mu) = \mathbf{E}_{(x,y) \sim D} I(h_i^\mu(x) \neq y)$  la perte zéro-un.

Théorème (MARCHAND et SOKOLOVA 2005 ; LAVIOLETTE et al. 2005)



Étant donné  $h_i^\mu : \mathcal{X} \rightarrow \{-1, +1\}$ , et  $\mathcal{L}_D^{01}(h_i^\mu) = \mathbf{E}_{(x,y) \sim D} I(h_i^\mu(x) \neq y)$  la perte zéro-un.

**Théorème (MARCHAND et SOKOLOVA 2005 ; LAVIOLETTE et al. 2005)**

Soit  $\mathcal{R}$  une fonction de reconstruction,  $P_{\mathcal{M}_i}$  une distribution sur les messages, et  $\delta \in (0, 1]$ . Avec forte probabilité ( $\geq 1 - \delta$ ) sur  $S \sim D^n$  :

# Bornes de généralisation pour des classifieurs binaires compressés

Étant donné  $h_i^\mu : \mathcal{X} \rightarrow \{-1, +1\}$ , et  $\mathcal{L}_D^{01}(h_i^\mu) = \mathbf{E}_{(x,y) \sim D} I(h_i^\mu(x) \neq y)$  la perte zéro-un.

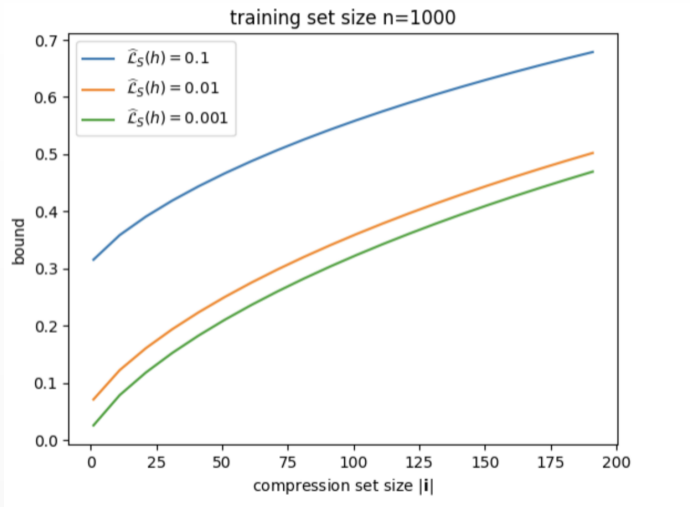
**Théorème (MARCHAND et SOKOLOVA 2005 ; LAVIOLETTE et al. 2005)**

Soit  $\mathcal{R}$  une fonction de reconstruction,  $P_{\mathcal{M}_i}$  une distribution sur les messages, et  $\delta \in (0, 1]$ . Avec forte probabilité ( $\geq 1 - \delta$ ) sur  $S \sim D^n : \forall i \in \mathcal{I}_n, \mu \in \mathcal{M}_i :$

$$\mathcal{L}_D^{01}(h_i^\mu) \leq 1 - \exp\left(\frac{-1}{n - |\mathbf{i}| - k_{S_{i^c}}} \left[ \ln \binom{n - |\mathbf{i}|}{k_{S_{i^c}}} + \ln \binom{n}{|\mathbf{i}|} + \ln \left( \frac{1}{P_{\mathcal{M}_i}(\mu) \cdot \xi(|\mathbf{i}|) \cdot \delta} \right) \right]\right)$$

où  $k_{S_{i^c}} := |\mathbf{i}^c| \widehat{\mathcal{L}}_{S_{i^c}}^{01}(h_i^\mu)$  est le nombre d'erreurs sur  $S_{i^c} := S \setminus S_i$  et  $\xi(a) := \frac{6}{\pi^2} (a + 1)^{-2}$ .

# Bornes non-triviales !



## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

## La question

Peut-on apprendre une fonction de reconstruction  $\mathcal{R}(\mathbf{i}, \mu)$  pour construire un réseau de neurones ?

## La question

Peut-on apprendre une fonction de reconstruction  $\mathcal{R}(\mathbf{i}, \mu)$  pour construire un réseau de neurones ?

Il faudrait alors adapter nos résultats de telle sorte que :

## La question

Peut-on apprendre une fonction de reconstruction  $\mathcal{R}(\mathbf{i}, \mu)$  pour construire un réseau de neurones ?

Il faudrait alors adapter nos résultats de telle sorte que :

- 1 Les pertes à valeurs continues sont considérées

## La question

Peut-on apprendre une fonction de reconstruction  $\mathcal{R}(\mathbf{i}, \mu)$  pour construire un réseau de neurones ?

Il faudrait alors adapter nos résultats de telle sorte que :

- 1 Les pertes à valeurs continues sont considérées
- 2 Une quantité indénombrable de message



# Nouvelle borne pour une quantité indénombrable de messages

## PAC-Bayes à la rescousse !

On considère maintenant une distribution *a posteriori*  $Q_{\mathcal{M}}$  sur l'ensemble (potentiellement continu) des messages  $\mathcal{M}$ .

# Nouvelle borne pour une quantité indénombrable de messages

## PAC-Bayes à la rescousse !

On considère maintenant une distribution *a posteriori*  $Q_{\mathcal{M}}$  sur l'ensemble (potentiellement continu) des messages  $\mathcal{M}$ .

## Théorème

Soit  $\mathcal{R}$ , une fonction de reconstruction,  $P_{\mathcal{M}}$  une distribution *a priori* sur les messages, et  $\delta \in (0, 1]$ . Avec forte probabilité ( $\geq 1 - \delta$ ) sur  $S \sim D^n$  :

# Nouvelle borne pour une quantité indénombrable de messages

## PAC-Bayes à la rescousse !

On considère maintenant une distribution *a posteriori*  $Q_{\mathcal{M}}$  sur l'ensemble (potentiellement continu) des messages  $\mathcal{M}$ .

## Théorème

Soit  $\mathcal{R}$ , une fonction de reconstruction,  $P_{\mathcal{M}}$  une distribution *a priori* sur les messages, et  $\delta \in (0, 1]$ . Avec forte probabilité ( $\geq 1 - \delta$ ) sur  $S \sim D^n$  :

$\forall \mathbf{i} \in \mathcal{I}_n$ ,  $Q_{\mathcal{M}}$  sur  $\mathcal{M}$  :

$$\mathbf{E}_{\mu \sim Q_{\mathcal{M}}} \mathcal{L}_D(h_{\mathbf{i}}^{\mu}) \leq \mathbf{E}_{\mu \sim Q_{\mathcal{M}}} \widehat{\mathcal{L}}_{S_{\mathbf{i}c}}(h_{\mathbf{i}}^{\mu}) + \frac{1}{\sqrt{n - |\mathbf{i}|}} \left[ \text{KL}(Q_{\mathcal{M}} \| P_{\mathcal{M}}) + \frac{\sigma^2}{2} + \ln \binom{n}{|\mathbf{i}|} + \ln \left( \frac{1}{\xi(|\mathbf{i}|) \cdot \delta} \right) \right].$$

## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

## 1 Préambules

## 2 Théorie de la compression des données

- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

# Le méta-apprentissage ?

**Idée**

# Le méta-apprentissage ?

## Idée

- Considérer plusieurs tâches simultanément

# Le méta-apprentissage ?

## Idée

- Considérer plusieurs tâches simultanément

## Motivations



# Le méta-apprentissage ?

## **Idée**

- Considérer plusieurs tâches simultanément

## **Motivations**

- Partage de connaissance entre les tâches

# Le méta-apprentissage ?

## Idée

- Considérer plusieurs tâches simultanément

## Motivations

- Partage de connaissance entre les tâches
- Davantage d'exemple par entraînement

## Méta-entraînement

$$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\} \text{ tel que } S^{(i)} \sim (D^{(i)})^{n_i}.$$

## Méta-entraînement

$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\}$  tel que  $S^{(i)} \sim (D^{(i)})^{n_i}$ .

## Algorithme

$\mathcal{A}(\mathcal{S}) \rightarrow (\mathcal{C}, \mathcal{R})$

## Méta-entraînement

$$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\} \text{ tel que } S^{(i)} \sim (D^{(i)})^{n_i}.$$

## Algorithme

$$\mathcal{A}(\mathcal{S}) \longrightarrow (\mathcal{C}, \mathcal{R})$$

## Fonction de compression

$$\mathcal{C}(\mathcal{S}) \longrightarrow (\mathbf{i}, \mu)$$

## Méta-entraînement

$$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\} \text{ tel que } S^{(i)} \sim (D^{(i)})^{n_i}.$$

## Algorithme

$$\mathcal{A}(\mathcal{S}) \longrightarrow (\mathcal{C}, \mathcal{R})$$

## Fonction de compression

$$\mathcal{C}(S) \longrightarrow (\mathbf{i}, \mu)$$

## Fonction de reconstruction

$$\mathcal{R}(\mathbf{i}, \mu) \longrightarrow \theta$$

## Méta-entraînement

$$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\} \text{ tel que } S^{(i)} \sim (D^{(i)})^{n_i}.$$

## Algorithme

$$\mathcal{A}(\mathcal{S}) \longrightarrow (\mathcal{C}, \mathcal{R})$$

## Fonction de compression

$$\mathcal{C}(\mathcal{S}) \longrightarrow (\mathbf{i}, \mu)$$

## Fonction de reconstruction

$$\mathcal{R}(\mathbf{i}, \mu) \longrightarrow \theta$$

## Prédicteur

$$h_{\theta}(\mathbf{x}) \longrightarrow y$$

## Méta-entraînement

$$\mathcal{S} := \{S^{(1)}, S^{(2)}, \dots, S^{(m)}\} \text{ tel que } S^{(i)} \sim (D^{(i)})^{n_i}.$$

## Algorithme

$$\mathcal{A}(S) \longrightarrow (\mathcal{C}, \mathcal{R})$$

## Fonction de compression

$$\mathcal{C}(S) \longrightarrow (\mathbf{i}, \mu)$$

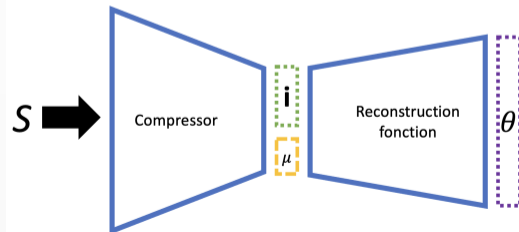
## Fonction de reconstruction

$$\mathcal{R}(\mathbf{i}, \mu) \longrightarrow \theta$$

## Prédicteur

$$h_{\theta}(\mathbf{x}) \longrightarrow y$$

**Combiner meta-learning et sample compression :**





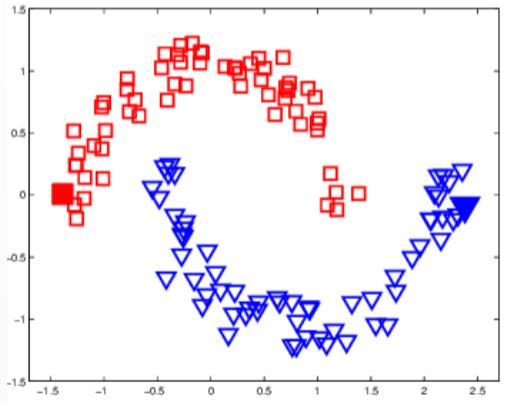
## 1 Préambules

## 2 Théorie de la compression des données

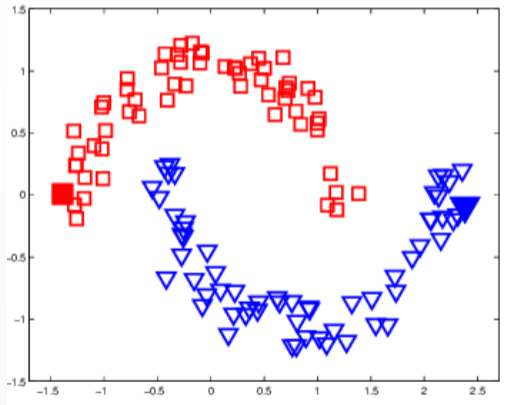
- Définition et exemples
- Borne de généralisation pour classificateurs binaires
- Nouvelle borne de généralisation

## 3 Deep Reconstruction Machine

- Le paradigme du méta-apprentissage
- Résultats préliminaires - classification binaire

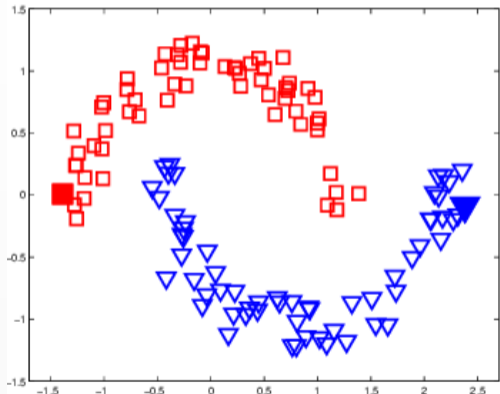


**Tâche #1**



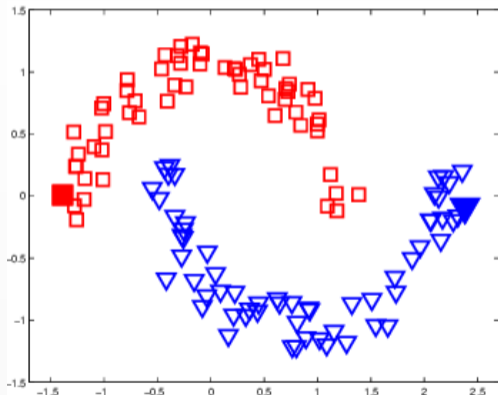
## Tâche #1

- Non-linéairement séparable



## Tâche #1

- Non-linéairement séparable
- Translations / rotations / homothéties

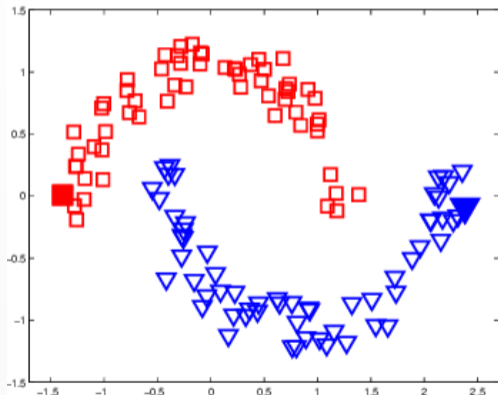


## Tâche #1

- Non-linéairement séparable
- Translations / rotations / homothéties

## Le prédicteur

- $h$  : réseau de neurones simple



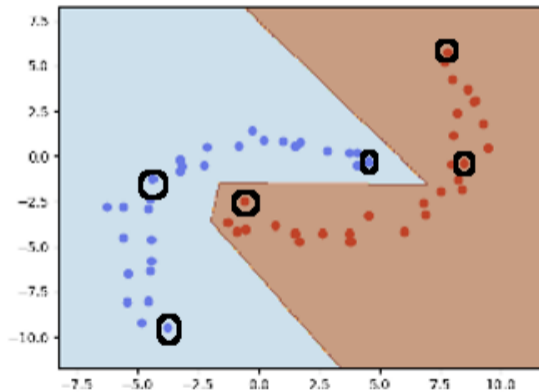
## Tâche #1

- Non-linéairement séparable
- Translations / rotations / homothéties

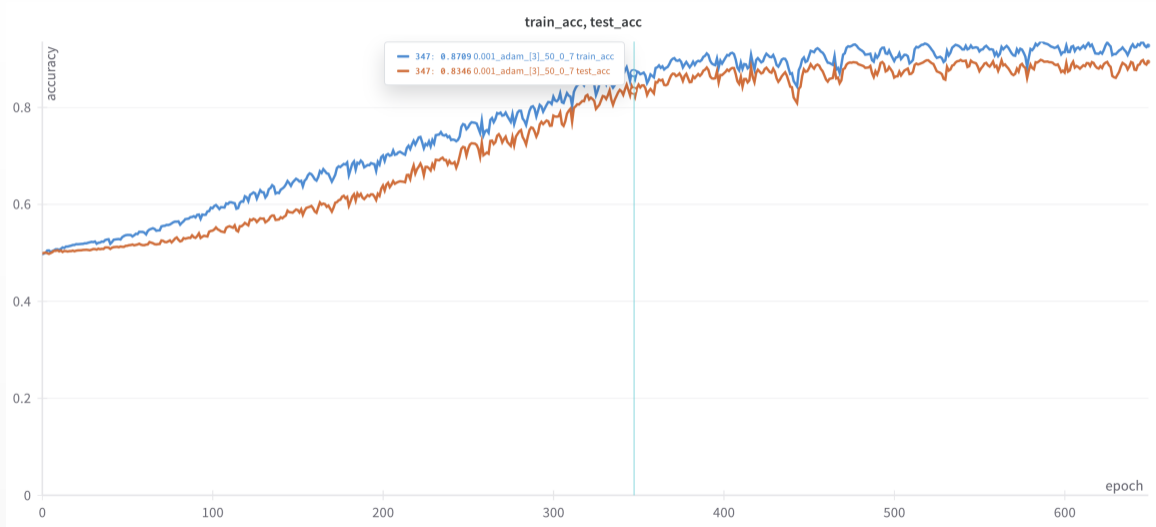
## Le prédicteur

- $h$  : réseau de neurones simple
- $\theta = (\mathbf{w}, b)$

# Résultats préliminaires

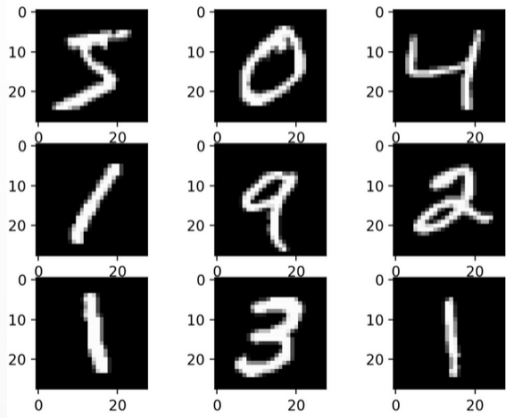


# Résultats préliminaires

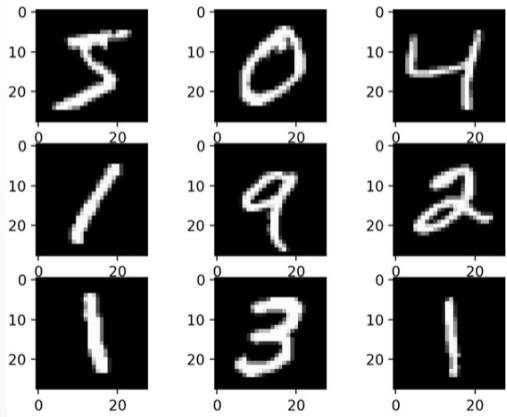






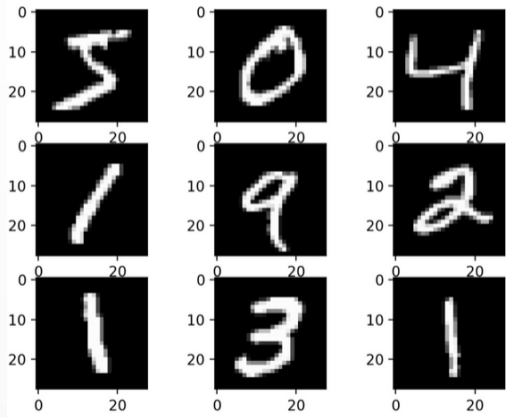


Tâche #2



## Tâche #2

- Entrée de dimension 28x28

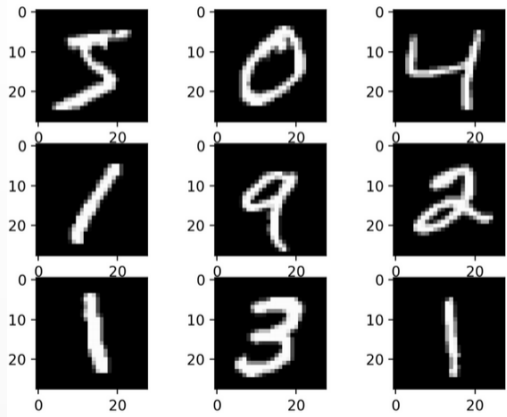


## Tâche #2

- Entrée de dimension 28x28

## Le prédicteur

- $h$  : séparateur linéaire
- $h(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$



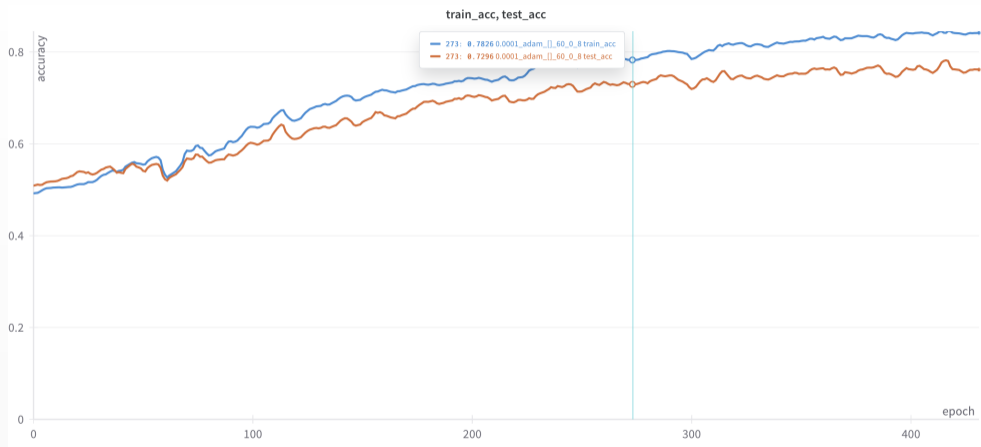
## Tâche #2

- Entrée de dimension 28x28

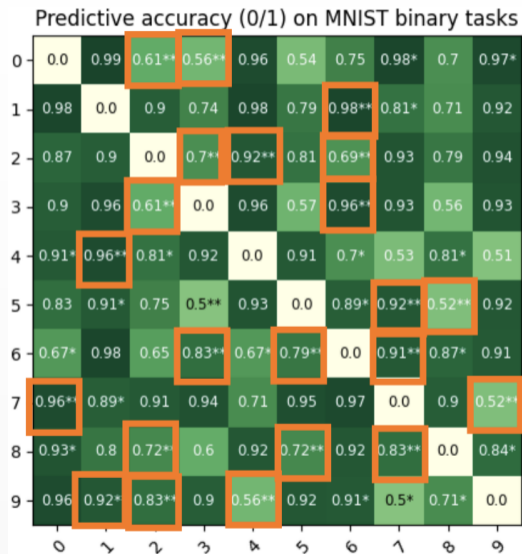
## Le prédicteur

- $h$  : séparateur linéaire
- $h(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$
- $\theta = (\mathbf{W}_i, \mathbf{b}_i)_{i=1}^2$

# Résultats préliminaires



# Résultats préliminaires



Merci de votre écoute :) !



# References I

- CAMPI, Marco C. et Simone GARATTI (2023). "Compression, Generalization and Learning". In : *JMLR abs/2301.12767*.
- DROUIN, Alexandre, Sébastien GIGUERE, Vladana SAGATOVICH, Maxime DÉRASPE, François LAVIOLETTE, Mario MARCHAND et Jacques CORBEIL (2014). "Learning interpretable models of phenotypes from whole genome sequences with the Set Covering Machine". In : *arXiv preprint arXiv:1412.1074*.
- GODON, Thibaud, Pier-Luc PLANTE, Baptiste BAUVIN, Elina FRANCOVIC-FONTAINE, Alexandre DROUIN et Jacques CORBEIL (2022). "RandomSCM: interpretable ensembles of sparse classifiers tailored for omics data". In : *CoRR abs/2208.06436*. DOI : 10.48550/ARXIV.2208.06436. arXiv : 2208.06436. URL : <https://doi.org/10.48550/arXiv.2208.06436>.
- LAVIOLETTE, François, Mario MARCHAND et Mohak SHAH (2005). "Margin-Sparsity Trade-Off for the Set Covering Machine". In : *ECML*. T. 3720. Lecture Notes in Computer Science. Springer, p. 206-217.
- LITTLESTONE, Nick et Manfred K. WARMUTH (1986). "Relating Data Compression and Learnability". In : *Technical Report*.
- MARCHAND, Mario et John SHAWE-TAYLOR (2002). "The Set Covering Machine". In : *JMLR 3*.
- MARCHAND, Mario et Marina SOKOLOVA (2005). "Learning with Decision Lists of Data-Dependent Features". In : *JMLR 6*.
- PACCAGNAN, Dario, Marco C. CAMPI et Simone GARATTI (2023). "The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance". In : *NeurIPS*.