

# Fondements de l'apprentissage machine

Automne 2014

Roland Memisevic

Leçon 1

## Qu'est-ce que l'apprentissage machine ("Machine Learning") ?

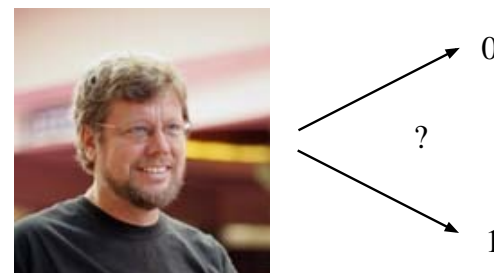
- ▶ Certaines tâches sont très difficiles ou fastidieuses pour programmer.
- ▶ L'apprentissage automatique remplace la programmation par l'apprentissage à partir d'exemples.

## IFT 3395/6390, Automne 2014

- ▶ Local :
  - ▶ cours théorique :  
Mercredi 15 : 30 - 17 : 30 S-144 Pav. Roger-Gaudry  
Jeudi 9 : 30 - 10 : 30 Z-210 Pav. Claire-McNicoll
  - ▶ démonstration :  
Jeudi 10 : 30 - 12 : 30 1340 Pav. André-Aisenstadt
- ▶ Professeur :  
Roland Memisevic,  
3349 Pav. Andre-Aisenstadt  
roland.memisevic@umontreal.ca
- ▶ Page Web du cours :

<http://www.iro.umontreal.ca/~memisevr/teaching/ift3395.2014/index.html>

## Que caractérise un visage ?



## Que caractérise le chiffre 2 ?

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

- Quelles sont les règles qui définissent le "2".

## Que caractérise la langue Française ?

*"Wikipedia.fr est un site de l'association  
Wikimedia France. Les resultats du  
moteur de recherche proviennent de  
Wikiwix"*

- Comment peut-on reconnaître la langue, étant donné un flux de caractères ?

## D'autres exemples

- Détection de pourriel
- Reconnaissance de l'identité faciale (eg., Facebook)
- Reconnaissance de sourire (par un appareil photo)
- Détection des anomalies des réseaux
- Reconnaissance de la musique (par téléphone mobile)
- Prédiction des cours boursiers
- Traduction automatique
- Système de recommandation
- Reconnaissance de parole
- Reconnaissance d'activités (console de jeu video)
- Surveillance automatique
- Construction de robots autonomes
- Des voiture qui peuvent conduire par elles même
- "Divination par télépathie"
- etc.

## L'analyse des données

- Dans beaucoup de domaines scientifiques on est confrontés à une rapide augmentation des données.
- L'apprentissage automatique fournit un moyen d'utiliser ces données.

```
01111110010000010011111101001111100000010000101  
01101010111101010001001000111000001010010100100  
001000000000001011111001111100110100100001001001  
110010100110010000101001001011000101000011101001  
000111010000110010101000111100000100100000001011  
10010100001011111001000000111011110010011011110  
10100111111011101100110101000110000101001100001
```

## Relation aux statistiques

- ▶ La statistique est la science classique pour l'analyse de données.
- ▶ Une tentative de distinguer les deux :
- ▶ L'accent dans la statistique est sur les outils qui permettent aux humains d'analyser les données. L'accent dans l'apprentissage machine est sur les systèmes qui comprennent des données par eux-mêmes.

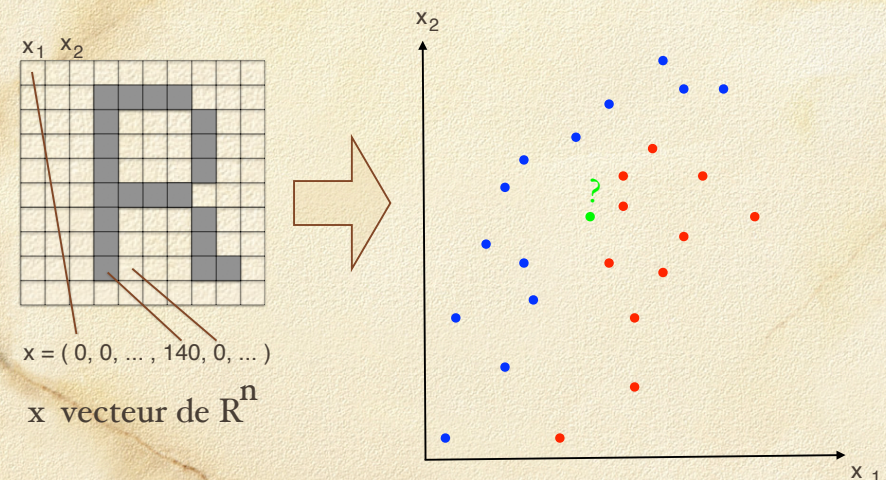
## Notation

- ▶  $\mathbf{x}$  : entrée,  $\mathbf{t}$  : sortie désirée,  $\mathbf{y}(\mathbf{x})$  : sortie de modèle
- ▶  $\mathbf{x}$  et  $\mathbf{t}$  peuvent être
  - ▶ des scalaire ou des vecteurs
  - ▶ pour des scalaire nous utiliserons aussi  $x$  and  $t$  (des caractères non-gras)
- ▶  $\mathbf{x}$  et  $\mathbf{t}$  peuvent être
  - ▶ continu : par exemple  $x = 1.73457$ , or  $\mathbf{x} \in \mathbb{R}^D$
  - ▶ discrète : par exemple  $x \in \{0, 1\}$ , or  $\mathbf{x} \in \{'a', 'b', 'c', \dots, 'z'\}^D$
- ▶ L'apprentissage revient souvent à ajuster des paramètres d'un modèle. Nous utilisons souvent le vecteur  $\mathbf{w}$  pour désigner tous les paramètres d'un modèle. Donc  $\mathbf{y}(\mathbf{x}) = \mathbf{y}(\mathbf{x}; \mathbf{w})$

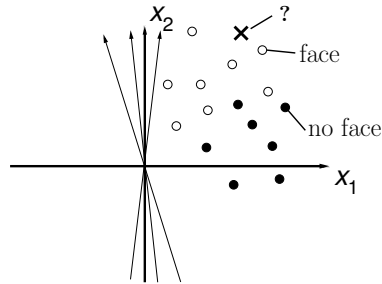
## Relation aux neuro-sciences

- ▶ Le meilleur système du monde pour l'analyse des données est le cerveau.
- ▶ Plusieurs techniques d'apprentissage automatique sont inspirées par ce que nous savons sur le traitement de l'information dans le cerveau.
- ▶ En retour, l'étude de la théorie de l'apprentissage peut fournir des conseils sur le traitement de l'information dans le cerveau.
- ▶ Exemple : "Le cerveau Bayésien" (The Bayesian Brain)

### Représentation des données



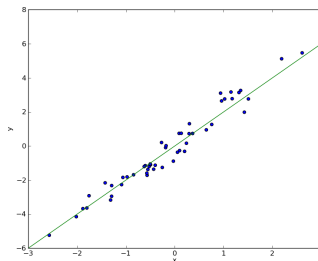
## Travailler avec des données de grande dimension



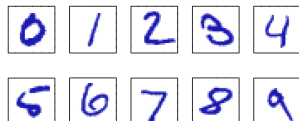
- ▶ Les exemples,  $\mathbf{x}$ , (surtout les entrées mais parfois aussi les sorties) sont souvent de grande dimension.
- ▶ Nous ne pouvons pas imaginer que trois dimensions.
- ▶ Nos intuitions tridimensionnelles fonctionnent souvent, mais parfois elles échouent !
- ▶ Des espaces de grande dimension ont des propriétés particulières ("malédiction de la dimension")

## Apprentissage supervisé

- ▶ **Regression** : Prévoir des sorties continues.



- ▶ **Classification** : Prévoir des sorties discrètes.

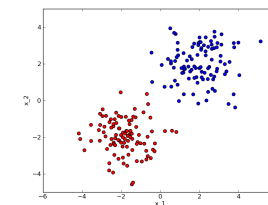


## Classification des tâches d'apprentissage

- ▶ **Apprentissage supervisé** : On a des exemples d'entraînement  $(\mathbf{x}, \mathbf{t})$ , et la tâche est d'apprendre une fonction  $\mathbf{y}(\mathbf{x})$  qui peut *prédire* les sorties,  $\mathbf{t}$ , pour des nouvelles entrées  $\mathbf{x}$ .
- ▶ **Apprentissage non supervisé** : On a des exemples d'entraînement non marqué  $\mathbf{x}$ , et la tâche est d'apprendre une meilleure représentation pour les données.
- ▶ **(Apprentissage par renforcement)** : La tâche est de prendre une série de mesures qui maximisent récompense (par exemple, dans les jeux. Ne sera pas vu en détail dans ce cours.)

## Apprentissage non supervisé

- ▶ **Regroupement ("clustering")** : Trouver des groupes. C'est comme mettre les données dans une nouvelle représentation (discrète).

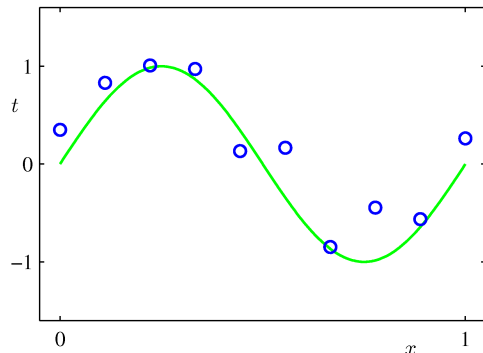


- ▶ **Réduction de la dimension** : Trouver une nouvelle représentation continue.
- ▶ L'apprentissage non supervisé est toujours une forme de *compression avec perte* : Les données sont représentées en utilisant moins de bits.

## Sujets (prévu)

- ▶ Introduction
- ▶ Régression linéaire
- ▶ Classification linéaire
- ▶ Les réseaux de neurones
- ▶ Méthodes à noyaux
- ▶ Raisonnement Bayésien
- ▶ Clustering, modèles de mélange, l'algorithme EM
- ▶ Séquences et modèles de Markov cachés (HMM)
- ▶ Les modèles graphiques, inférence approximative, l'échantillonnage
- ▶ Prédictions structurées
- ▶ L'apprentissage de traits, l'apprentissage profond

## Exemple : Estimation d'une courbe (1-D)



- ▶ La vraie relation sous-jacente entre  $x$  et  $t$  est en vert.
- ▶ Les observations sont bruyantes (en bleue).

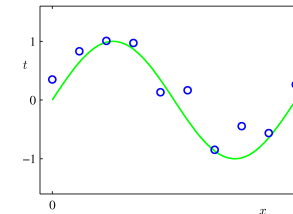
## Généralisation et sur-apprentissage ("overfitting")

- ▶ Formellement, l'apprentissage est d'adapter des paramètres du modèle sur la base de données d'entraînement.
- ▶ Une question centrale est de savoir si le modèle va fonctionner sur des données futures : "Généralisation".
- ▶ Ceci est une fonction (a) de comment bien le modèle fait sur les données d'entraînement, et (b) de la taille de l'espace d'hypothèses.

### Exemple :

Il est très facile de *memoriser* chaque exemple d'entraînement. Mais cela sera inutile sur de nouvelles données !

## Exemple : Estimation d'une courbe (1-D)



- ▶ Modèle polynomial :

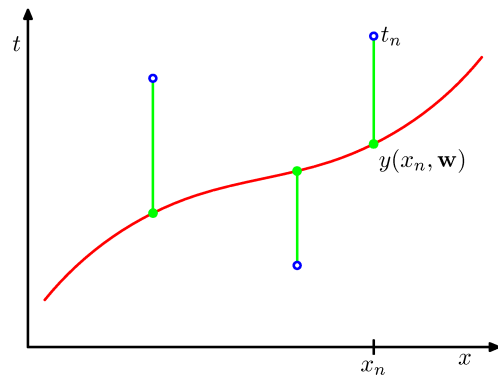
$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M$$

avec des paramètres

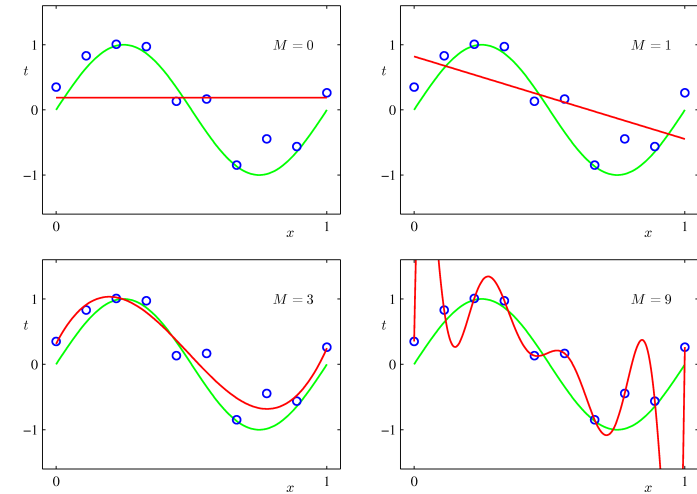
$$\mathbf{w} = (w_0, w_1, \dots, w_M)^T$$

- ▶ Pour estimer les paramètres, on peut minimiser l'erreur quadratique :  $E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$

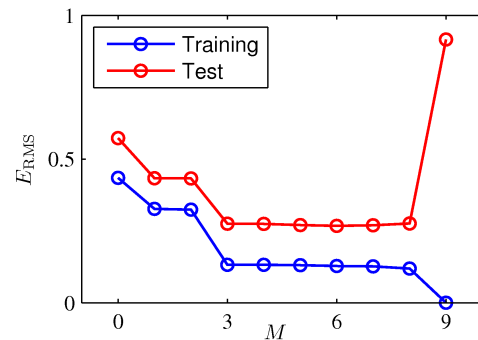
## Exemple : Estimation d'une courbe (1-D)



## Exemple : Estimation d'une courbe (1-D)

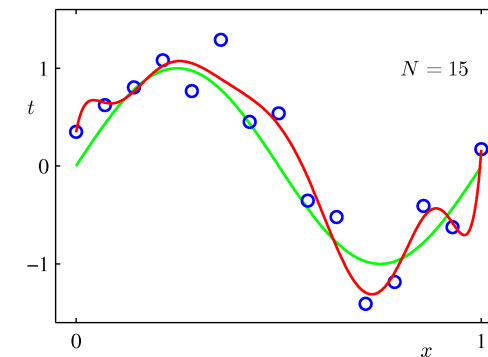


## Sur-apprentissage ("overfitting")



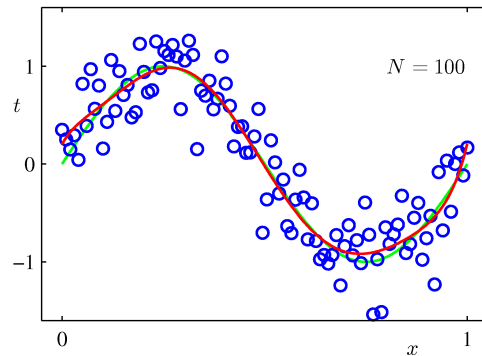
- La diminution de l'erreur sur les données d'entraînement ne signifie pas que celle-ci diminue sur des données test.
- La réduction de la capacité (du  $M$ ) permet de généraliser.

## Sur-apprentissage



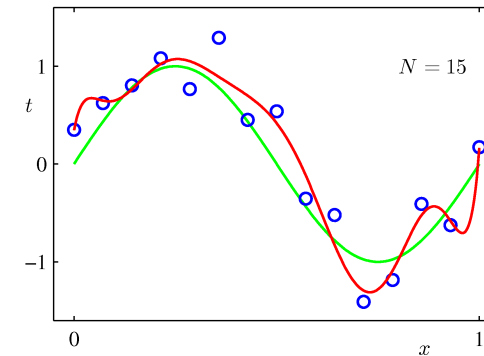
- L'augmentation des données permet aussi de généraliser.

## Sur-apprentissage



- L'augmentation des données permet aussi de généraliser.

## Exemple : Estimation d'une courbe (1-D)



- Que faire si la quantité de données d'entraînement est limitée ?

## Prévenir sur-apprentissage

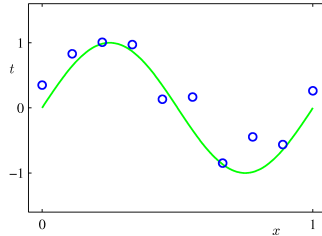
- **Selection de modèles** : Choisissez la complexité correcte.
- **Régularization** : Pénaliser les grands coefficients :
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$
- **Modelisation Bayésienne** : Ne pas essayer d'adapter un modèle à tous ! Déduire une distribution de probabilité conditionnelle  $p(\text{modele}|\text{data})$ . Cette approche est beaucoup plus difficile à faire dans la pratique, mais il est un moyen naturel de prévenir la sur-apprentissage, et cela fonctionne très bien.
- Toutes ces approches sont liées.

## Validation croisée

- Pour effectuer la sélection du modèle, il est courant de diviser les données d'entraînement dans un sous-ensemble de données **d'entraînement** approprié et un sous-ensemble de données de **validation**.
- Ensuite, nous estimons plusieurs modèles différents, par exemple en choisissant différentes valeurs pour  $M$  ou  $\lambda$ , et choisissons le modèle qui fait le mieux sur les données de validation.
- On peut échanger les rôles de formation et de validation des sous-ensembles pour obtenir une estimation plus stable. C'est s'appelle validation croisée (**cross-validation**).
- Cas extrême : Prendre tous les sous-ensembles avec  $N - 1$  de cas d'entraînement comme ensemble d'apprentissage. "**Leave-one-out**" **cross validation**.



## No free lunch



- ▶ Si tout ce que nous avons sont des données d'entraînement, il est *rien* de façon de généraliser.
- ▶ Pour avoir en mesure d'apprendre, nous *doit* faire des supposition. Biais inductif **Inductive bias**.
- ▶ Une supposition très commun : La douceur de la fonction sous-jacente.

## Beaucoup de questions de recherche

- ▶ Amélioration de la vitesse, précision, généralité des méthodes.
- ▶ Trouvez le biais inductif correct pour la tâche à accomplir.
- ▶ L'apprentissage de bout en bout de modèles complexes.
- ▶ Des applications.

### Conférences et revues principales

- ▶ NIPS : Neural Information Processing Systems, ICML, UAI, AISTATS
- ▶ PAMI : Pattern Analysis and Machine Intelligence, Journal of Machine Learning Research, Journal Machine Learning, Neural Computation