

Intelligence Artificielle et Analyse de données



Applications : Python-PyTorch

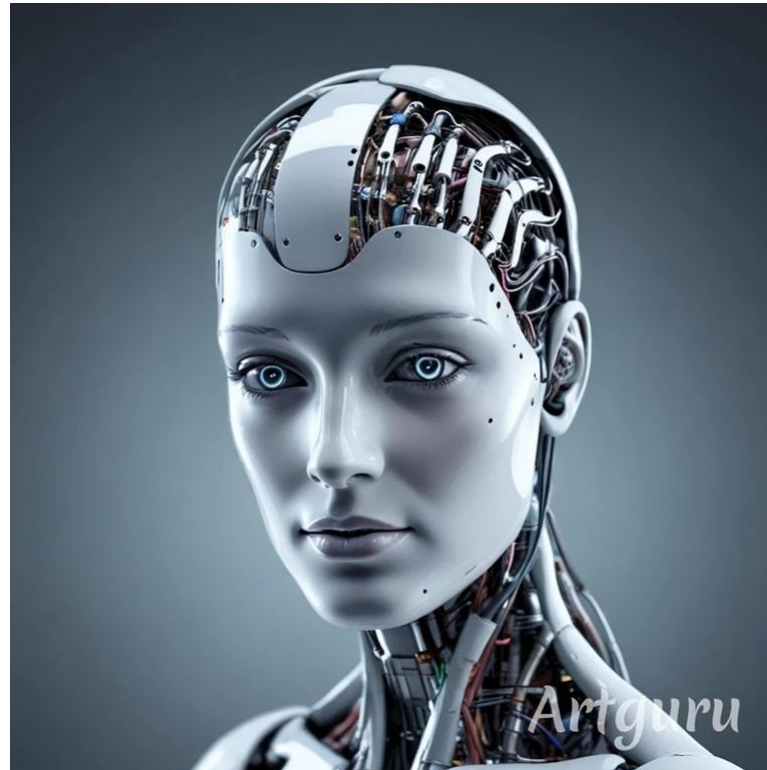


Johan Peralez

Plan du cours :

1. Définir l'Intelligence Artificielle (IA)
2. Prérequis -ressources (programmation-mathématiques)
3. Introduction aux réseaux de neurones (NNs, *Neural Networks*)
 - 3.1. Problèmes de régression
 - 3.2. Descente de gradient, dérivation automatique
 - 3.3. Un NN simple (MLP)
 - 3.4. Problèmes de classification
4. Les réseaux de neurones récurrents (RNNs)
 - 4.1. Données séquentielles
 - 4.2. Principaux RNNs
 - 4.3. Implémentation
5. Les réseaux de neurones convolutifs (CNNs)
 - 5.1. Problèmes de classification
 - 5.2. Principe des CNNs
 - 5.3. Implémentation des CNNs
6. Applications
 - 6.1. Apprentissage par renforcement
 - 6.2. Classification d'images (données CIFAR)

1. Définir l'Intelligence Artificielle (IA)



- L'IA est une notion floue et qui évolue rapidement \Rightarrow sa définition dépend:
 - du domaine (traitement d'image, contrôle, etc.)
e.g. distinguer des visages vs faire marcher un robot
 - de l'époque (depuis le milieu du XXe siècle)
e.g. Deep Blue vs Alpha Go
- Exemples:
 - « L'**automatisation** d'activités que nous **associons à la pensée humaine**, comme la prise de décision, la résolution de problème ou l'apprentissage. » (BELLMAN 1978)
subjectif (nous ?)
 - « L'étude de comment **programmer les ordinateurs** pour qu'ils réalisent des tâches pour lesquelles les êtres humains sont **actuellement meilleurs**. » (RICH & KNIGHT 1991)
paradoxal (jeu d'échec vs football ?)
 - « Ensemble de **théories et de techniques** mises en œuvre en vue de réaliser des machines capables de **simuler l'intelligence humaine**. » (Larousse 2024)
vague (définir intelligence ?)

- Difficulté à définir l'IA = difficulté à définir l'Intelligence.

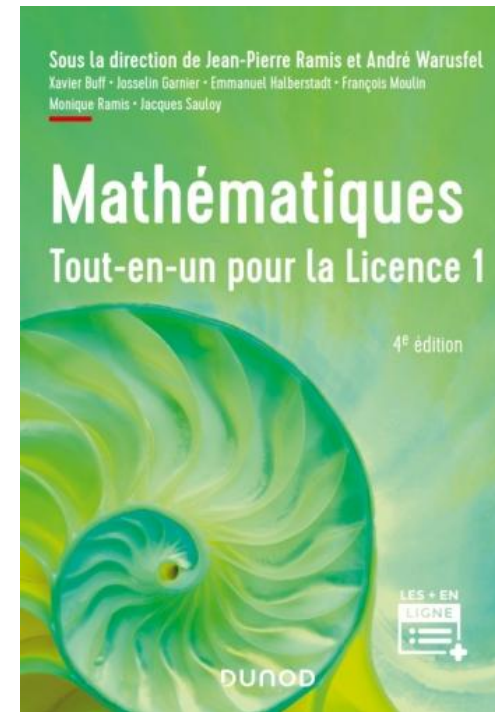
Selon [Larousse 2024], l'intelligence =

1. « Ensemble des fonctions mentales ayant pour objet la connaissance conceptuelle et rationnelle »
2. « Aptitude d'un être humain à s'adapter à une situation, à choisir des moyens d'action en fonction des circonstances »

- Approche “pragmatique”:

- IA = ensemble des méthodes qui sont habituellement classées dans l'IA par les gens de son domaine.
- En traitement de l'image : réseaux de neurones (NNs) convolutifs (CNNs), algorithmes de segmentation, etc.
- En contrôle : NNs récurrents (RNNs), algorithmes d'apprentissage par renforcement (RL), etc.

2. Prérequis -ressources (programmation-mathématiques)



- L'IA et l'analyse de données nécessitent des compétences :
 - en **programmation**
Python (langage le + utilisé en IA)
 - en **mathématiques**
Algèbre linéaire, calcul différentiel, probabilités, statistiques
- Ressources **Python** en ligne:
 - Cours et exos de base :
 - <https://www.learnpython.org/>
 - <https://www.france-ioi.org/algo/chapters.php>
 - <https://courspython.com/apprendre-numpy.html>
 - Installation (programmer en local)
 - <https://anaconda.org>
 - <https://pytorch.org/get-started/locally/>
 - Programmer en ligne :
 - <https://www.programiz.com/python-programming/online-compiler/>
 - <https://colab.research.google.com/>



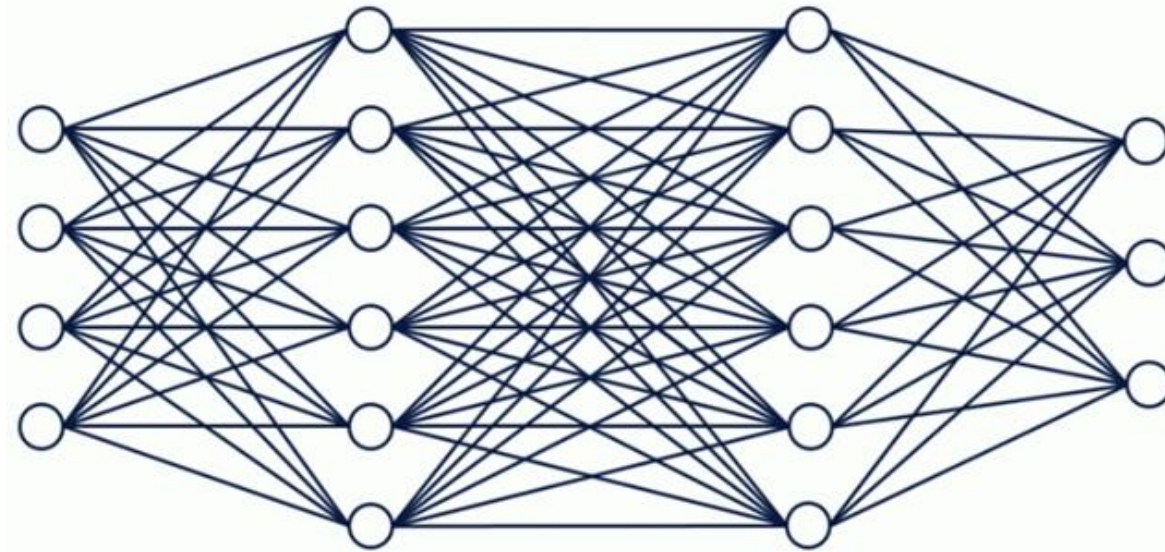
3. Introduction aux réseaux de neurones

3.1. Problèmes de régression

3.2. Descente de gradient, dérivation automatique

3.3. Un réseau de neurones simple (MLP)

3.4. Titanic



3.1 Problèmes de régression

Formulation : **apprendre à prédire** une valeur de sortie y à partir d'une donnée d'entrée x

Diagram illustrating the regression equation $\hat{y} = h(x, \theta)$. The components are annotated as follows:

- \hat{y} : valeur \hat{y} prédite
- h : fonction h choisie « à la main »
- θ : paramètre(s) θ à apprendre

afin de minimiser une fonction **coût** L (*loss*) :

Diagram illustrating the cost function minimization $\min_{\theta} L(\{y\}, \{\hat{y}\})$. The components are annotated as follows:

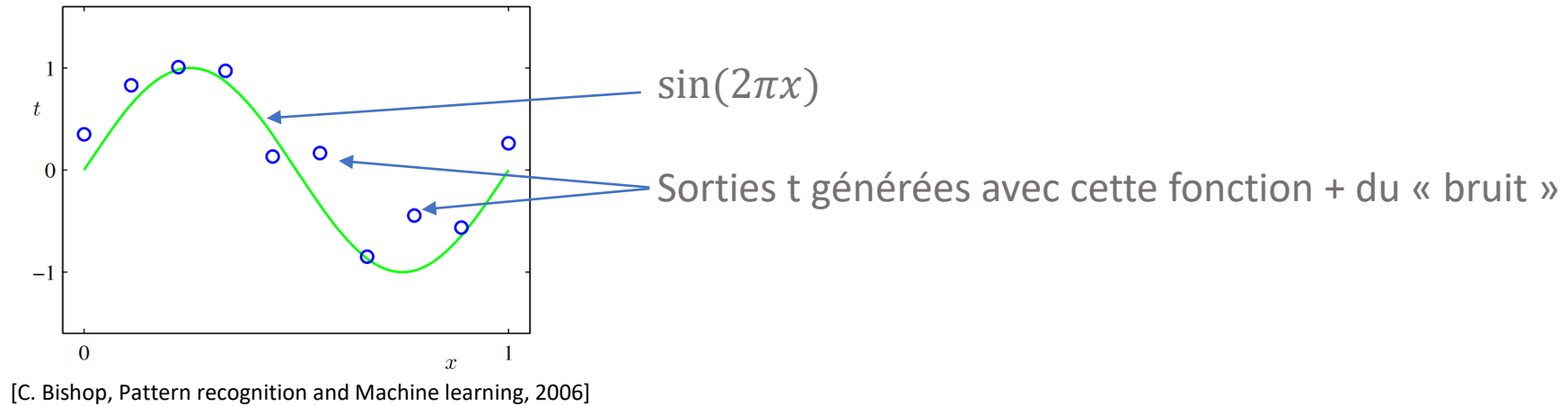
- $\{y\}$: données
- $\{\hat{y}\}$: prédictions

⇒ **problème d'optimisation**

Remarque : nous verrons plus tard que h peut-être un réseau de neurones ...

Exemple :

Données :



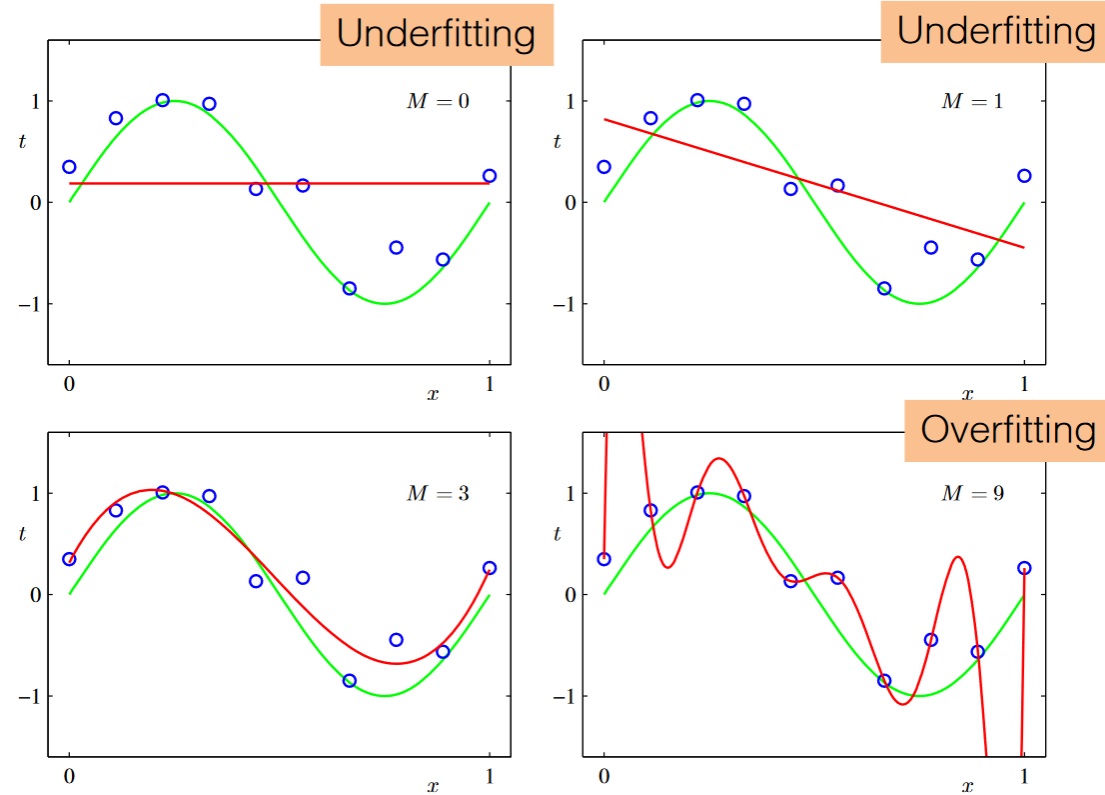
On choisit de prédire avec un polynôme d'ordre M :

$$\hat{t} = h(x, \theta) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_M x^M$$

Et une fonction de coût « moindres carrés » (moyenné):

$$L = \frac{1}{N} \sum_{i=1}^N (t_i - \hat{t}_i)^2, \text{ avec } N \text{ le nombre de données.}$$

Quel ordre M choisir pour notre polynôme ?



[C. Bishop, Pattern recognition and Machine learning, 2006]

- Pour M trop petit : problème de sous-apprentissage (***underfitting***).
- Pour M trop grand : problème de sur-apprentissage (***overfitting***).

Remarque (**polynômes de Lagrange**) : pour n données distinctes il existe un (unique) polynôme d'ordre $n-1$ qui passe exactement par chaque donnée.

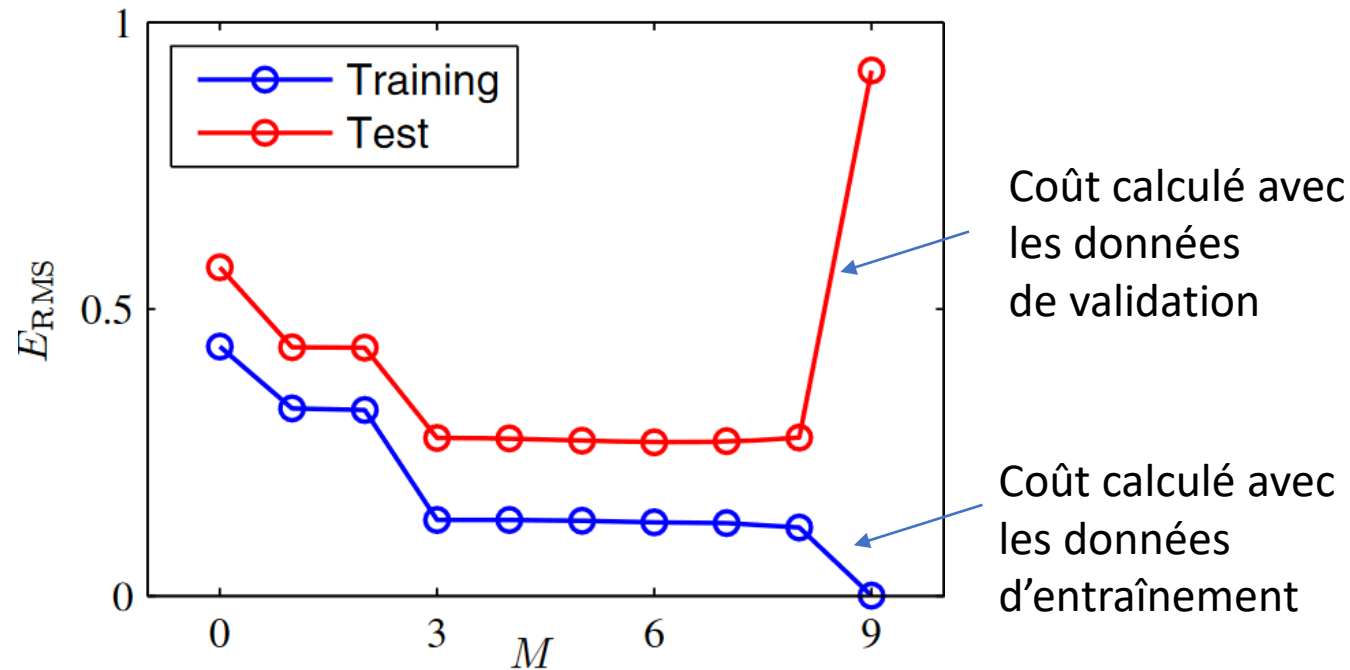
Séparation des données en, au moins, 2 ensembles :

- Données d'entraînement (*training set*).

Pour l'apprentissage de θ , i.e. la résolution du problème d'optimisation.

- Données de validation (*validation set*).

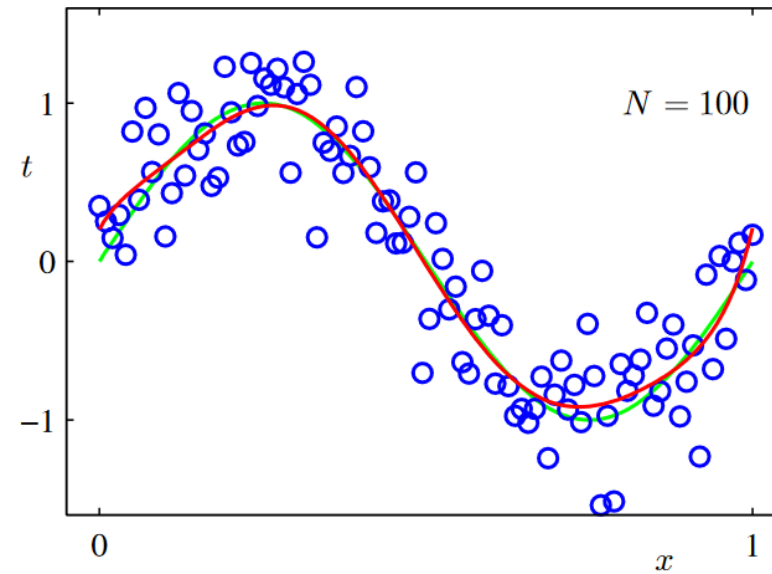
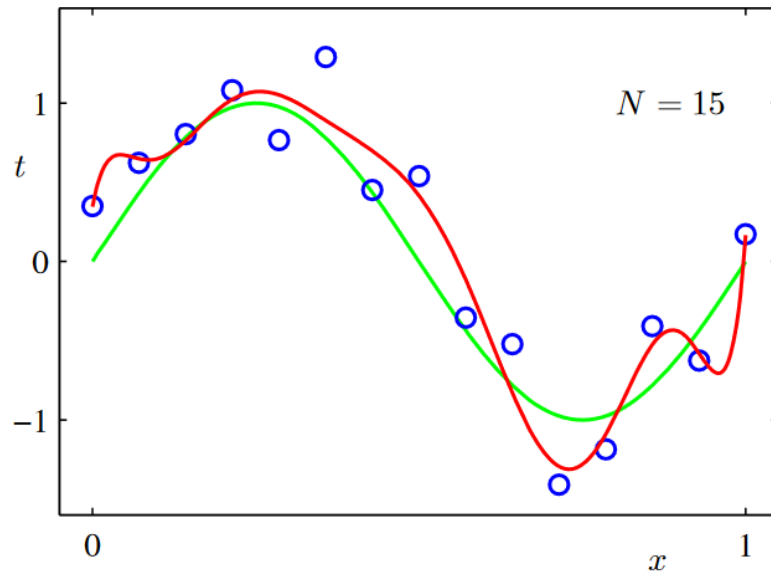
Permet de détecter le surapprentissage :



[C. Bishop, Pattern recognition and Machine learning, 2006]

Big Data ?

Augmenter (fortement) le nombre de données d'apprentissage permet de réduire le surapprentissage :



$M=9$

[C. Bishop, Pattern recognition and Machine learning, 2006]

Trois grandes difficultés des problèmes d'apprentissage :

1. Expressivité

Mon modèle, i.e. ma fonction h , peut-elle apprendre des phénomènes complexes ?

2. Difficulté à entraîner

Le problème d'optimisation, i.e. minimiser la différence entre prédictions et données, est-il difficile ?

e.g. si on résout par descente de gradient, la fonction coût est-elle dérivable ? lisse ?

3. Généralisation

Comment mon modèle se comporte sur des données qui ne sont pas dans le *training set* ?

Capacité à interpoler / extrapoler ?

e.g. le surapprentissage implique une mauvaise généralisation.

Exercice 3.a (Python, NumPy):

Générer et visualiser un jeu de données similaire à l'exemple ci-dessus.

1. Importer les librairies *numpy* et *matplotlib*

```
import numpy as np
from matplotlib import pyplot as plt
```

2. Paramètres

```
N = 30 # nombre de points
```

3. Générer (aléatoirement entre 0 et 1) les données d'entrées $\{x\}$, de taille N .

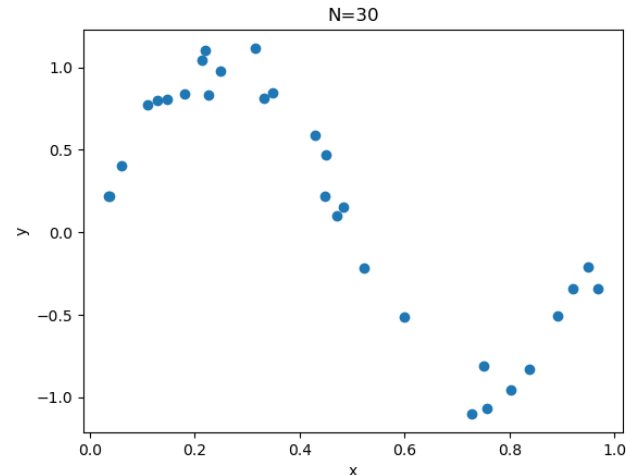
Utiliser la fonction `np.random.uniform(...)` pour utiliser une distribution uniforme.

4. Générer les données de sorties $\{y\}$, telles que $y = \sin(2\pi x) + w$.

Utiliser la fonction `np.random.normal(...)` pour générer w à partir d'une distribution normale centrée d'écart type (*standard deviation*) de 0.1

5. Visualiser le jeu de données

Utiliser `plt.plot`, `plt.legend`, `plt.xlabel`, etc.



3.2 Descente de gradient, dérivation automatique

Objectif : minimiser la fonction coût L .

Idée : mise à jour des paramètres, dans la direction qui fait diminuer L le plus fortement.

Mises à jour :

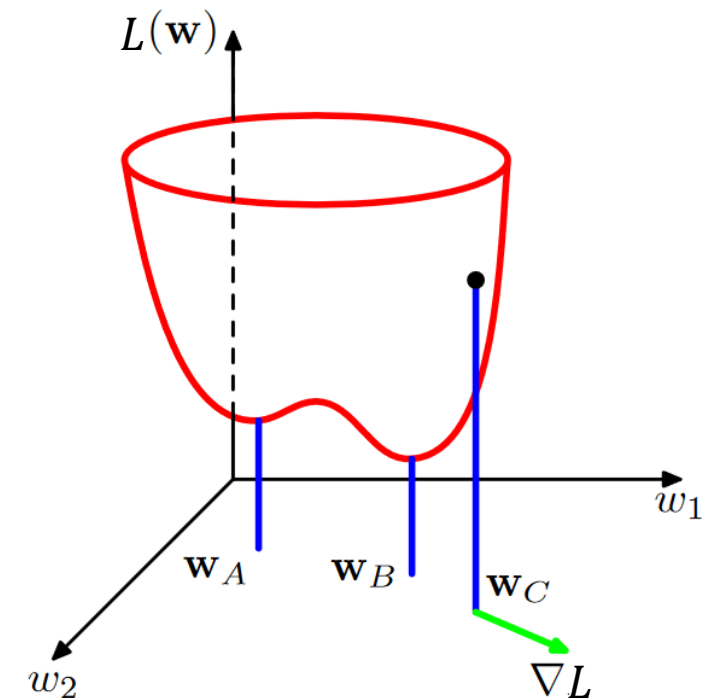
$$\theta^{k+1} = \theta^k - \eta \cdot \nabla_{\theta} L(h(x, \theta^k), y)$$

Paramètre à l'itération k Pas (*learning rate*) Gradient (se note aussi $\frac{\partial L}{\partial \theta}$)

Convergence : garantie pour un problème convexe, et

pour $0 < \eta < \frac{2}{k}$.

Illustration (non convexe): W_a est un minimum local.



Calculer le gradient $\nabla_{\theta} L = \left[\frac{\partial L}{\partial \theta_0}, \frac{\partial L}{\partial \theta_1}, \dots \right]^T$ revient à calculer les dérivées partielles $\frac{\partial L}{\partial \theta_i}$.

Ce calcul est à faire à chaque itération k (dépend de la valeur de θ) \Rightarrow calcul « à la main » prohibitif.

Calcul par différences finies :

$$\frac{\partial L}{\partial \theta_i} \approx \frac{L(h(x, \tilde{\theta}), y) - L(h(x, \theta), y)}{\epsilon}, \text{ avec } \tilde{\theta} = [\theta_0, \dots, \theta_{i-1}, \theta_i + \epsilon, \theta_{i+1}] \text{ et } \epsilon \text{ petit}$$

\Rightarrow **approximatif** et **coûteux** (nombreuses évaluations de L).

Calcul par dérivation automatique :

dérivation **exacte** (formelle) à l'aide d'un **logiciel**

\Rightarrow plus précis et plus rapide.

Exemple (bibliothèque Pytorch) :

```
import torch
x = torch.tensor([3.1], requires_grad = True) # initialise un scalaire (vecteur de taille 1)
y = x ** 2
grad = torch.autograd.grad(y, x) # calcule dy/dx = 2x
print(grad) # 6.2
```

```
x = torch.tensor([3.1, 2.0], requires_grad = True)
b = torch.tensor([5.5, 7.7])
y = b @ x # produit scalaire : y = b_1 * x_1 + b_2 * x_2
grad = torch.autograd.grad(y, x) # calcule dy/dx = [b_1, b_2]
print(grad) # [5.5, 7.7]
```

Exercice 3.a.suite (PyTorch):

Par descente de gradient, chercher le polynôme d'ordre 9 qui minimise la fonction coût L (moindres carrés). Compléter le programme:

```
# convertir données (numpy -> pytorch)
import torch
x, y = torch.tensor(x), torch.tensor(y)

# initialise theta
M = 3
theta = torch.zeros(M + 1, requires_grad = True)

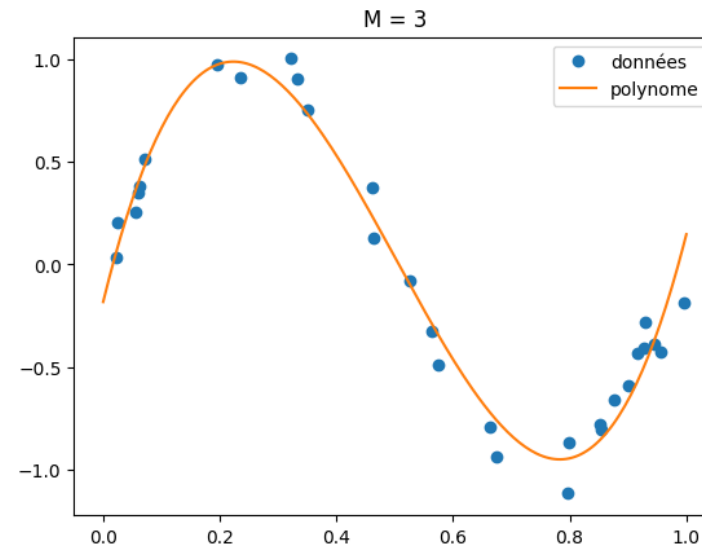
# fonction polynome
def get_predict(x, coefs):
    y_predict = torch.zeros_like(x)
    for i in range(len(theta)):
        y_predict = y_predict + coefs[i] * x**i
    return y_predict

# fonction coût
def get_loss(y_predict):
    return ((y_predict - y)**2).mean()

# descente de gradient
. . .
```

Tester pour:

- $M = 1, 2, 3, 8$.
- Différents *learning rate* (=0.01, 0.5, 2)



Pour aller plus loin : générer un 2^{ème} ensemble de données (*test set*) et tracer le coût en fonction de M .

3.3 Un réseau de neurones simple : le **MLP** (multi-layers perceptron)

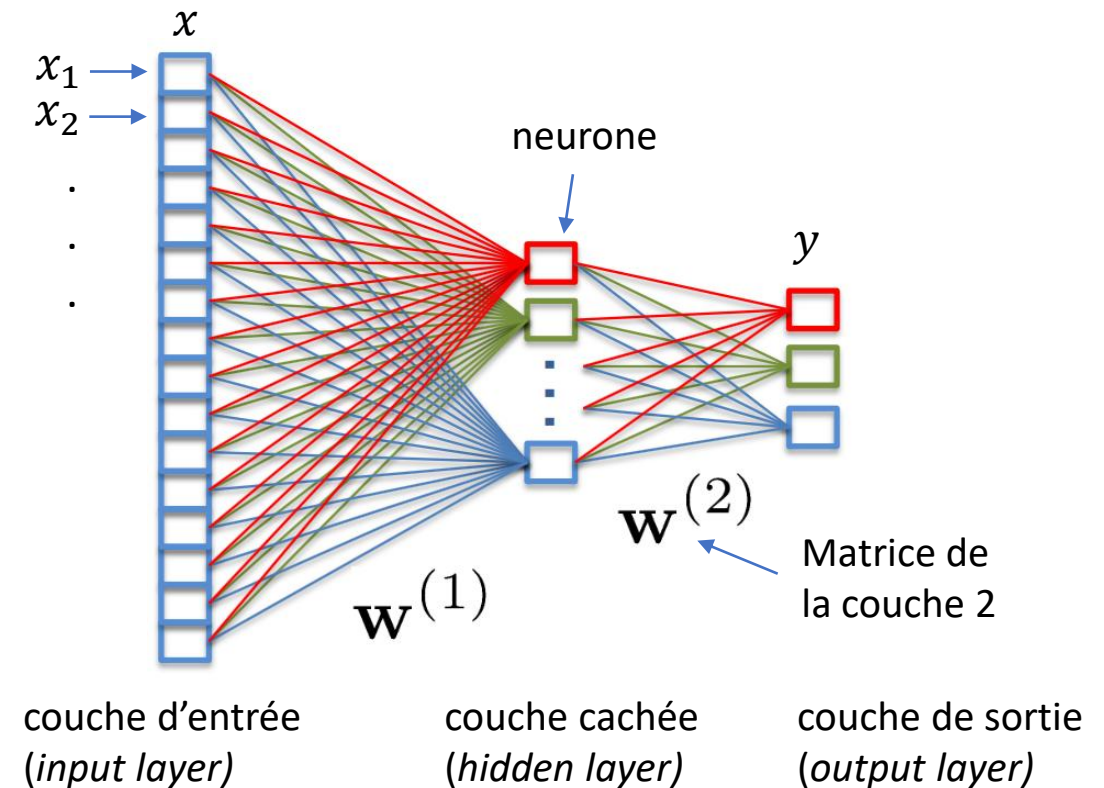
Principe du MLP (réseau dense) :

- Plusieurs couches (*layers*)
- Une couche transforme un vecteur d'entrée e en vecteur de sortie s .
- s est une combinaison linéaire de e , suivie (pour les couche cachées) par une non-linéarité σ .

⇒ le MLP est une fonction :

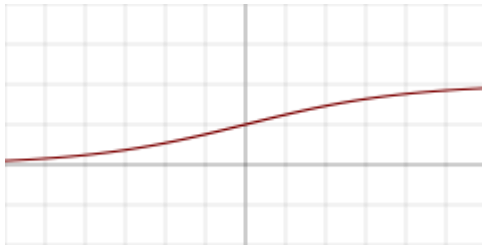
- Sortie d'une couche i : $s^{(i)} = \sigma(W^{(i)} \times e^{(i)} + b^{(i)})$
- Sortie du réseau : $y = W^{(2)} \times \sigma(W^{(1)} \times x + b^{(1)}) + b^{(2)}$
- Sortie d'un neurone : $s_j^{(i)} = w_j^{(i)} * e^{(i)} + b_j^{(i)}$

Entraîner un MLP = apprendre (identifier) ses paramètres,
i.e. les matrices $W^{(i)}$ (les poids) et les vecteurs $b^{(i)}$ (les biais).



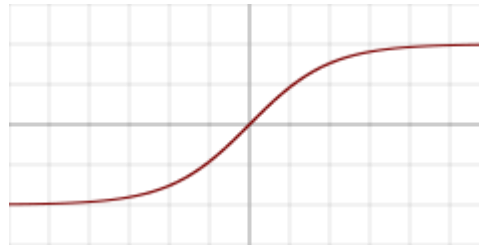
Non-linéarités σ (fonction d'activation):

- Fonction scalaire, i.e. $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
- Sur un vecteur, s'applique élément par élément, i.e. $\sigma([x_1, x_2, \dots]) = [\sigma(x_1), \sigma(x_2), \dots]$
- Les activations les + courantes :



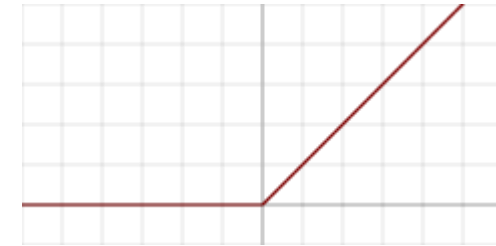
Sigmoïde

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Tangente hyperbolique

$$\sigma(x) = \frac{2}{1 + e^{-2x}} - 1$$



ReLU

$$\sigma(x) = \max(0, x)$$

Théorème d'approximation universelle :

Un MLP peut approcher d'aussi près que l'on veut n'importe quelle fonction continue.

Pour cela une seule couche cachée suffit (à condition de prendre une taille suffisante, i.e. une matrice $W^{(1)}$ suffisamment grande) et une fonction d'activation non polynomiale (les fonctions ci-dessus respectent cette condition). Et ... il faut trouver les bons paramètres !

Implémentation du MLP (PyTorch):

Hérite de la classe *torch.nn.Module*.

Exemple pour une entrée de taille 3, une matrice $W^{(1)}$ (couche cachée) de taille 3×32 , et une sortie de taille 2.

```
class MLP(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.fc1=torch.nn.Linear(3, 32) #  $W^1$ ,  $b^1$ 
        self.fc2=torch.nn.Linear(32, 2) #  $W^2$ ,  $b^2$ 

    def forward(self, x): # méthode appelée par MLP.forward(x) ou MLP(x)
        y = self.fc2(torch.functional.F.relu(self.fc1(x)))
        return y

model = MLP()
print(model)
print(list(model.parameters()))
```

Question: de quelle taille sont les vecteurs de biais ?

Entraînement du MLP :

PyTorch propose des outils pour faciliter la chaîne d'actions (*pipeline*) de l'entraînement :

- Classe `DataSet` pour la génération des données (*train set*, et *valid set*).
- L'implémentation de NNs.
- Des *optimizer* pour gérer la descente de gradient et la mise du *learning rate*.

Dans le code liés à ce cours, vous trouverez un exemple de pipeline. La descente de gradient utilisée est stochastique (*stochastic gradient descent*, **SGD**).

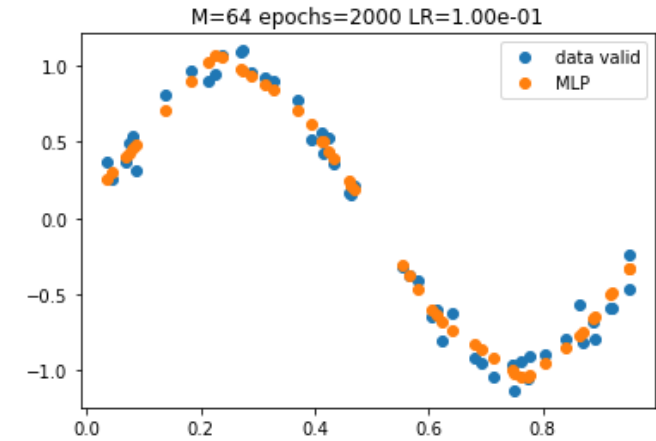
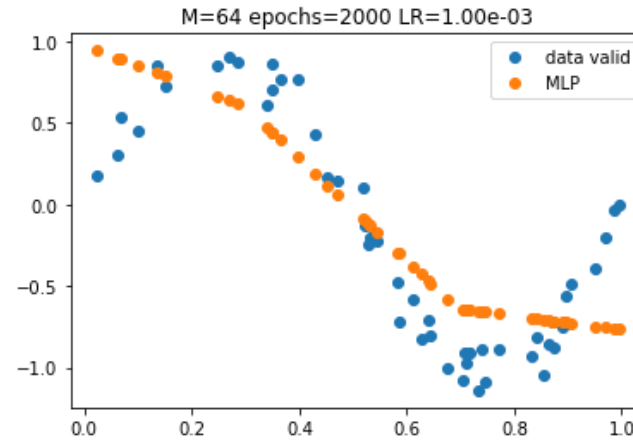
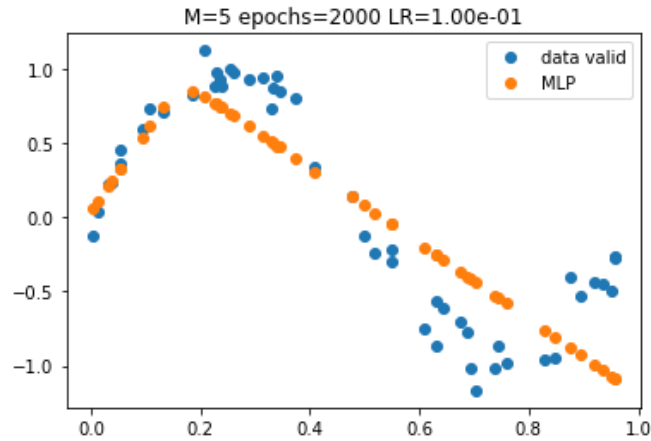
Avec SGD, à chaque itération, les données *train* sont découpées en sous-ensemble (***mini-batch***) et un calcul de gradient et une mise à jour des paramètres sont faits sur chaque *mini-batch*. Le découpage est fait aléatoirement à chaque itération (i.e. stochastique). Les avantages sont:

- Apprendre plus vite sur les grandes quantités de données (mises à jour plus fréquentes).
- Moins de risque d'être bloqué dans un minimum local.

Résultats (problème de régression simple):

Soit N le nombre de données d'entraînement. Nous fixerons $N = 200$ (ce sont des données du problème).

Soit M le nombre de colonnes de $W^{(1)}$, $epochs$ le nombre d'itérations, LR le *learning rate*. $M, epochs, LR$ sont des réglages (**hyperparamètres**) de l'entraînement.

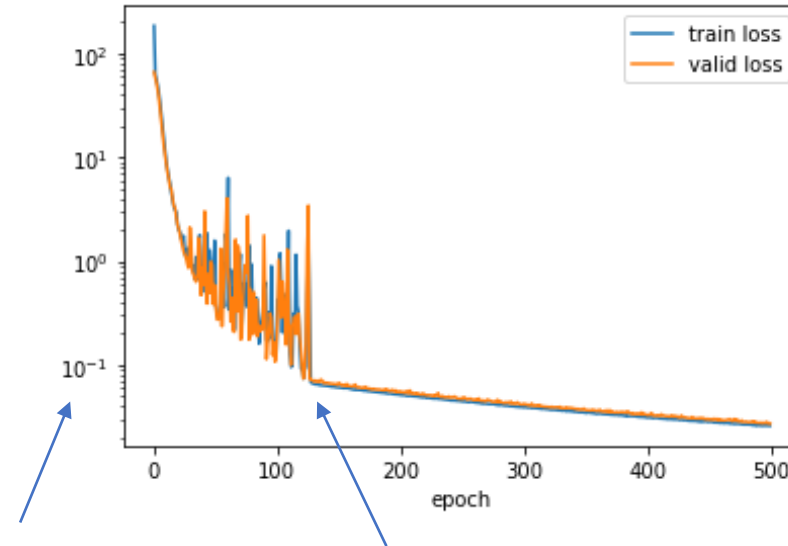
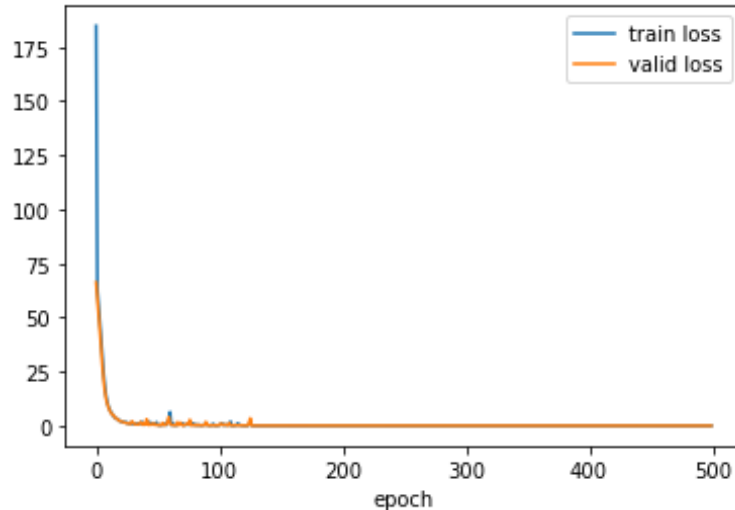


Exercice 3.b : Résoudre le problème de régression avec les données (vectorielle) générées ainsi:

- $x \in \mathbb{R}^2, x_1 \in [-1; 1], x_2 \in [1; 4]$ (aléatoire uniforme)
- $y = [\sqrt{x_1 + x_2}; x_2^3] + N(0, .1)$ (bruit gaussien)

Résultats pour

- Nombre de données : 2000 (entraînement) et 500 (validation)
- Nombre d'itérations (*epochs*): 500
- Learning rate : $1e^{-3}$ puis $2e^{-4}$



Echelle semilog

Réduction du *learning rate*

3.4 Problèmes de classification

Les données de sortie (qu'on cherche à prédire) sont des valeurs discrètes, i.e. un échantillon y appartient à un ensemble finis de classes $\{1, \dots, C\}$.

Ex: à partir d'une image d'un animal, prédire s'il s'agit d'un chat, d'un chien, etc.

Sorties du modèle :

La sortie du modèle est un vecteur $\hat{y} \in \mathbb{R}^C$.

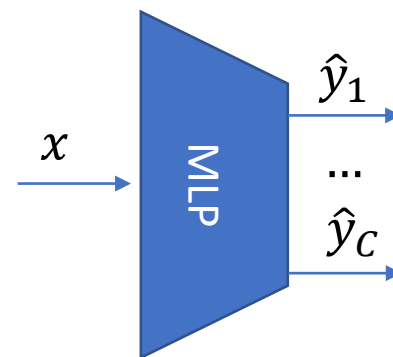
A chaque \hat{y}_i on associe la probabilité $\hat{P}(i|x) = \frac{\exp(\hat{y}_i)}{\sum_{c=1}^C \exp(\hat{y}_c)}$

La prédiction du modèle est la classe avec la plus grande probabilité.

Remarque : on vérifie qu'il s'agit d'une loi de proba : $\hat{P}(i|x) > 0$ et $\sum_{i=1}^C \hat{P}(i|x) = 1$.

Fonction coût (**cross entropy**) : Soit $y = c$ pour l'entrée x (i.e. la *vraie* classe est c), alors

$$L = -\log \hat{P}(c|x)$$



Analyse des données :

- Type d'entrées : numérique, catégorie (binaire, multiple).
- Statistique simple : moyennes, écart-types, min, max.
- Corrélations :

Coefficient de Pearson (r ou R) : mesure de la force et de la direction de la relation entre deux variables.

$$R(x, y) = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

Propriétés :

1. $-1 \leq R \leq 1$ (avec, e.g., relation linéaire positive parfaite pour $r = 1$, aucune relation pour $R = 0$).
2. $R(x, y) = R(y, x)$
3. $R(x, x) = 1$

Matrice de corrélation : pour des variables $\{x, y, z, \dots\}$

matrice carrée représentant les coeffs de Pearson.

Remarques :

propriété 2 \Rightarrow la matrice est symétrique.

propriété 3 \Rightarrow diagonale remplies de 1.

| | x | y | z |
|---|----------|----------|----------|
| x | $R(x,x)$ | $R(x,y)$ | $R(x,z)$ |
| y | $R(y,x)$ | ... | |
| z | ... | | |

Pré-traitement des données :

Certaines données brutes sont mal adaptées aux NN.

Pour faciliter leur entraînement, i.e. rendre la descente de gradient plus simple, il faut :

- Eviter des valeurs numériques extrêmes (problèmes de conditionnement, explosion ou annulation du gradient).

Ex : représentation d'une pression en *hPa* plutôt qu'en Pascal (pression atmosphérique $> 1e^5$ Pa)

- Favoriser une représentation informative des entrées.

Ex : la représentation d'une image est sous forme de 3 matrices RGB (*red, green, blue*) est plus informative que sous la forme d'un vecteur (notions de voisinage, couleur).

Normalisation (valeurs numériques) : Afin de s'assurer que les données respectent des valeurs centrées autour de zéro, et un écart-type de 1, on peut normaliser les données ainsi:

$$x^{norm} := \frac{x - \bar{x}}{\sigma_x} \text{ où } \bar{x} \text{ représente la moyenne et } \sigma_x \text{ l'écart-type.}$$

One-hot encoding (classes) : plutôt que de représenter la classe par un seul nombre entier, on représente par un vecteur de longueur le nombre de classe; le vecteur est alors formé de zéros et d'un seul 1.

Ex: soit une donnée appartenant à $\{1,2,3,4,5\}$. L'appartenance à la classe 1 est codée $[1,0,0,0,0]$, à la 3 $[0,0,1,0,0]$.

Interprétation des résultats / utilisation du modèle:

Une fois entraîné, le modèle nous intéresse pour :

- La **prédiction d'une classe** = index de la sortie avec la plus grande valeur, i.e.

$$classe\ prédite = arg \max_{i \in \{1, \dots, C\}} y_i$$

- La **confiance** du modèle en sa prédiction = $P(y = i|x) = \frac{\exp(y_i)}{\sum_{c=1}^C \exp(y_c)}$

Les fonctions de coût utilisé en classification (e.g. *cross entropy*) sont des **métriques indirectes** de la qualité du modèle.

La métrique (quantité qui évalue notre modèle) qui nous intéresse vraiment est le **pourcentage d'erreur dans la prédiction des classes**. C'est une métrique dont l'évolution est à regarder pendant l'entraînement et pour interpréter la qualité du modèle. Mais cette métrique ne peut pas être utilisé pour l'entraînement par descente de gradient car pas dérivable.

Exemple: Titanic

Données : pour 891 passagers, prix de leur billet, genre (H/F), tarif de leur billet, classe de leur cabine, un identifiant et si ils ont survécu.

Problème : entraîner un modèle qui prédise si un passager a survécu.



Type de données :

Affichons les premières lignes du fichier de données (*titanic.csv*).

Chaque ligne correspond à un passager.

- Sortie :
 - « *Survived* » (survécu vs pas survécu) = classe binaire
- Entrées :
 - « *p_id* » (identifiant) = classe multiples,
 - « *cabin_class* » = classe multiples,
 - « *female* » (femme vs homme) = classe binaire,
 - « *fare* » (tarif du billet) = numérique.

```
# p_id, survived, cabin_class, female, fare
0,0,3,0,7.250
1,1,1,1,71.283
2,1,3,1,7.925
3,1,1,1,53.100
4,0,3,0,8.050
5,0,3,0,8.458
6,0,1,0,51.862
7,0,3,0,21.075
8,1,3,1,11.133
```

Traçons les valeurs des champs « p_id » et « cabin_class ».

p_id attribue une valeur unique à chaque passager.

Cabin_class est une classe multiple; ses éléments sont {1,2,3}:

```
print(np.unique(cabin_class))
```

Combien de passagers ont survécu ? `print(sum(survived))`

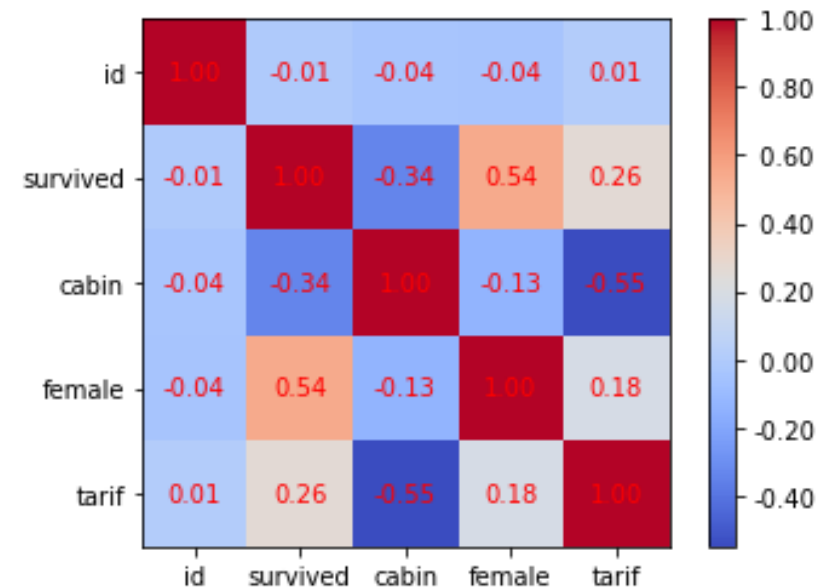
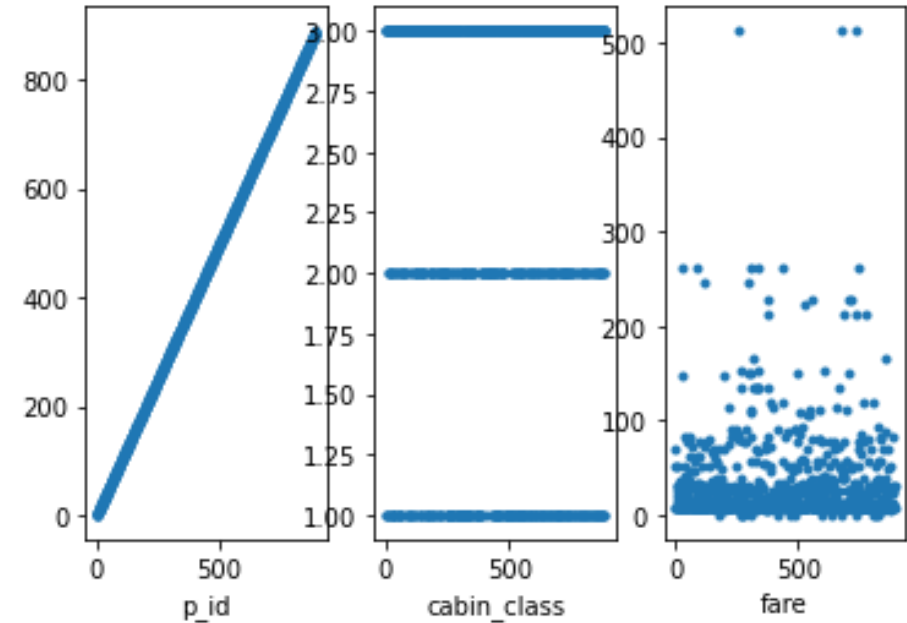
Statistiques simples :

```
print("fare: min %.3f max %.3f écart-type %.3f" %  
      (fare.min(), fare.max(), fare.std()))
```

```
>> fare: min 0.000 max 512.329 écart-type 49.666
```

Corrélation :

Forte corrélation négative entre « cabin_class » et « survived »,
i.e. les passagers de 1^{ère} classe ont plus de chance de survie que
ceux en 3^{ème} classe.



Pré-traitement des données:

Normalisation de « fare » (valeur numérique).

Encodage *one-hot* de « p_id », « cabin_class » et « female » (classes).

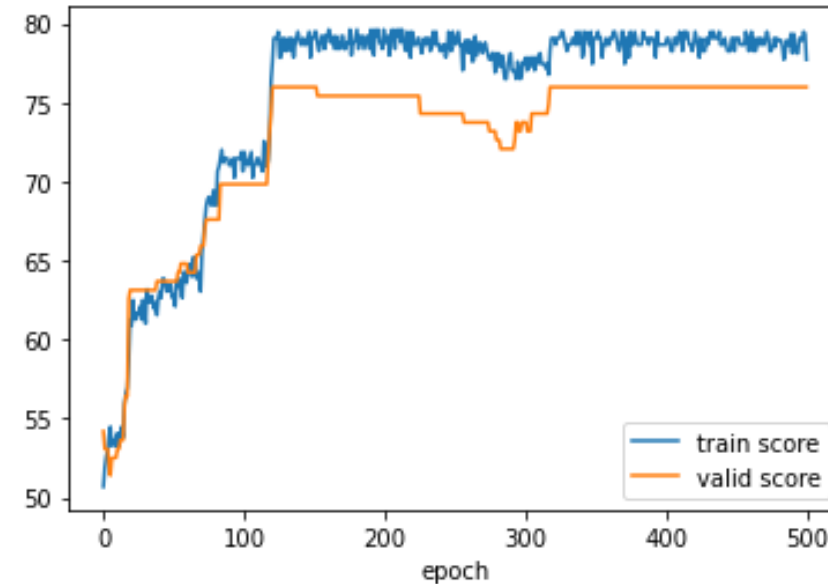
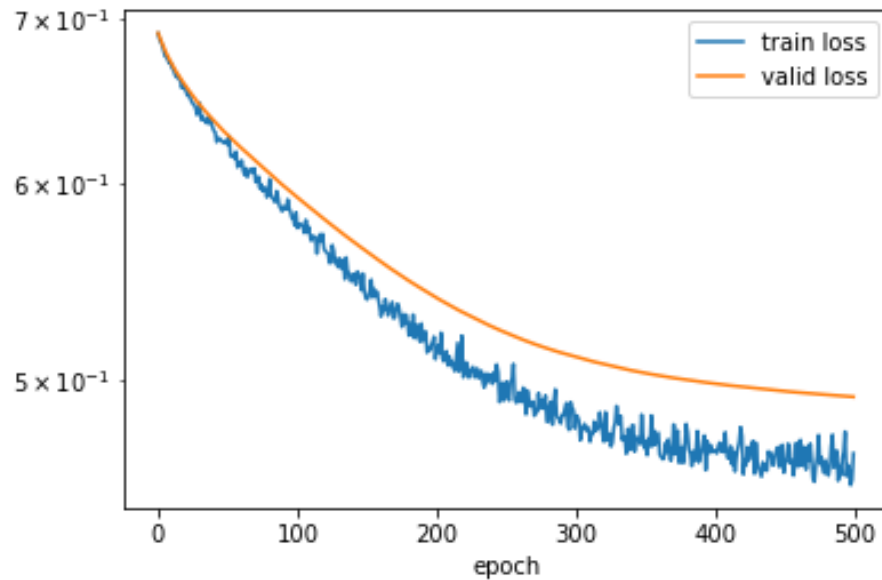
Il est préférable d'**éliminer l'entrée non pertinente** « p_id » : n'a pas de lien avec ce qu'on cherche à prédire (ce que confirme l'analyse de corrélation). Au contraire cette entrée peut :

- Complexifier le modèle (plus d'entrées impliquent plus de paramètres).
- Favoriser le sur-apprentissage (en apprenant une relation causale qui n'existe pas).

Remarque : puisque les passagers sont identifiés par un nombre unique, « p_id » il est possible d'apprendre à prédire de façon parfaite, sur les données d'entraînement, si un passager a survécu seulement avec cette entrée. Mais sur de nouvelles données, la prédiction sera totalement aléatoire (cas extrême de surapprentissage).

Résultats :

Pour un NN avec 1 couche caché de taille 32, learning rate à $1e^{-3}$, des mini-batch de taille 32 :



La courbe de droite correspond au % de réussite dans la prédiction de la classe de sortie.

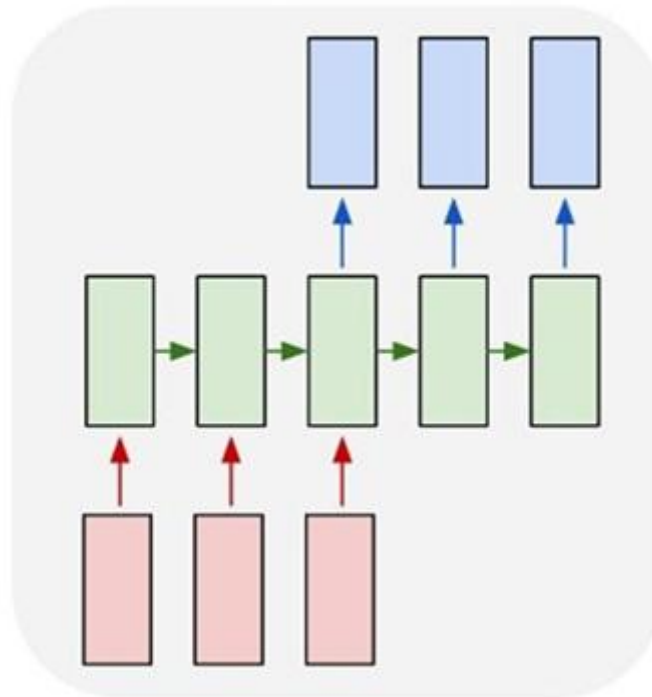
Remarque : dans cet exemple (Titanic) le nombre de classe en sortie est de 2 (survécu / pas survécu). Il est donc logique qu'un modèle non entraîné ait un score proche de 50% (réponse « au hasard » en début d'entraînement).

4. Les réseaux de neurones récurrents (RNNs)

4.1. Données séquentielles

4.2. Principaux RNNs

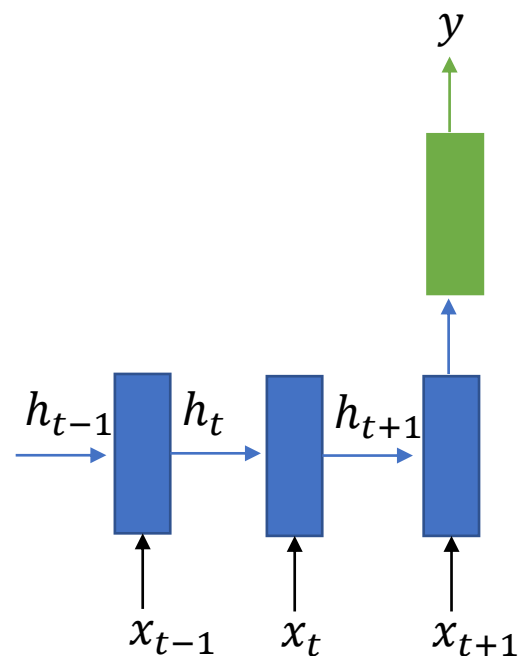
4.3. Implémentation



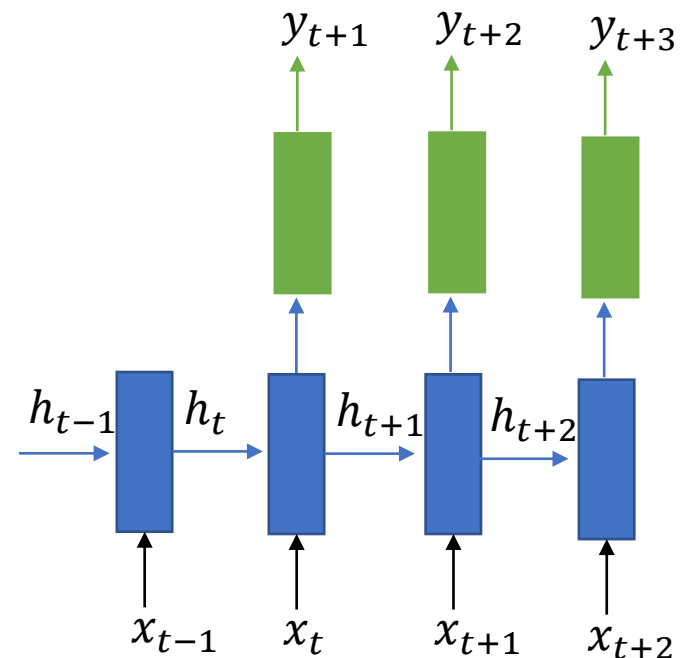
4.1 Données séquentielles



Données simples
 $y = f(x)$



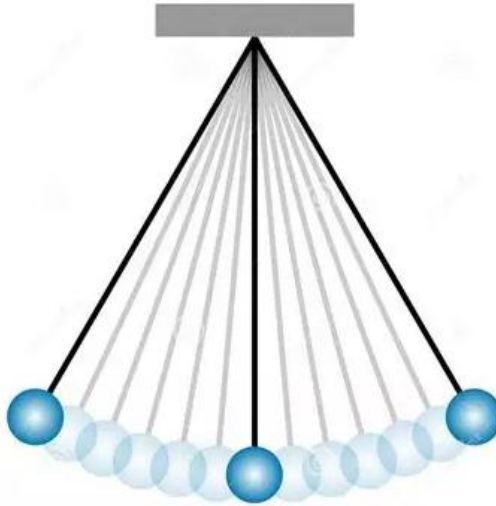
Entrées séquentielles
 $y = f_2(h_{t+1}, x_{t+1})$
 $h_{t+1} = f_1(h_t, x_t)$



Entrées et sorties séquentielles
 $y_{t+1} = f_2(h_t, x_t)$
 $h_{t+1} = f_1(h_t, x_t)$

Récurrence

Exemple (pendule) :



Problème 1 :

On enregistre la position angulaire θ au cours du temps (données d'entrées) pour différentes longueurs l et masse m du pendule (données de sortie) \Rightarrow apprendre un modèle pour estimer les valeurs de l et m .

Problème 2 :

On enregistre la position angulaire θ au cours du temps (données d'entrées x) et on cherche à prédire les futures valeurs de θ .

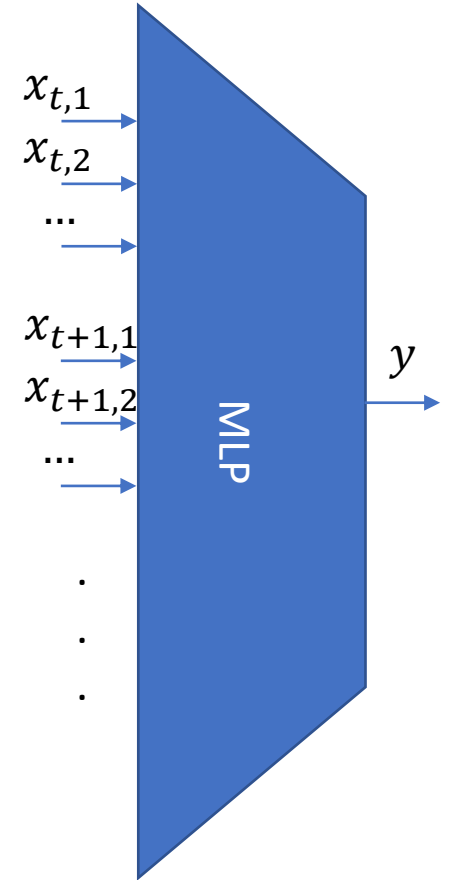
Approche simple :

- Représentation de la séquence en simple vecteur (*flat*).
- Utilisation d'un MLP.

Pas adapté :

- La structure du MLP est mal adaptée à la structure des données (e.g. la similitude entre un vecteur x_t et x_{t+1} n'est pas exploitée).
- Mal adapté à des séquences de longueurs variables.

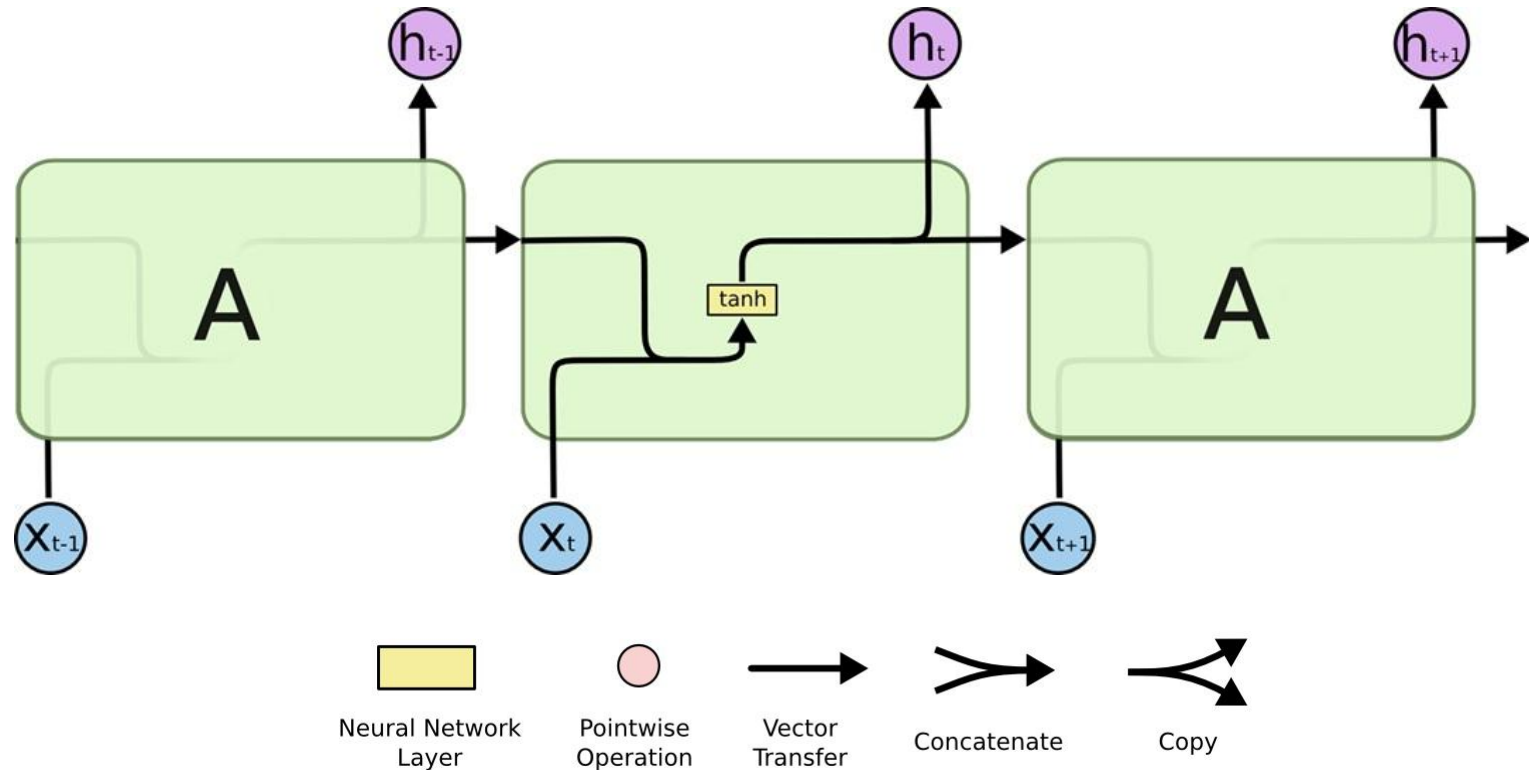
⇒ expressivité limitée ⇒ prédictions peu précises et/ou nécessite de nombreux paramètres (beaucoup de neurones) ⇒ risque de surapprentissage.



4.2 Principaux RNNs

Réseau Elman (= vanilla):

$$h_t = \tanh(W_x x_t + W_h h_{t-1} + b_h)$$

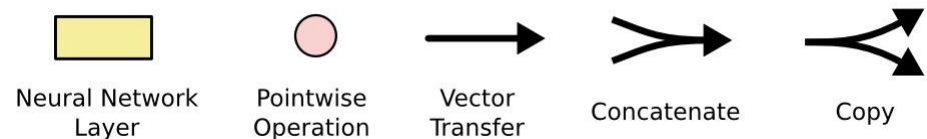
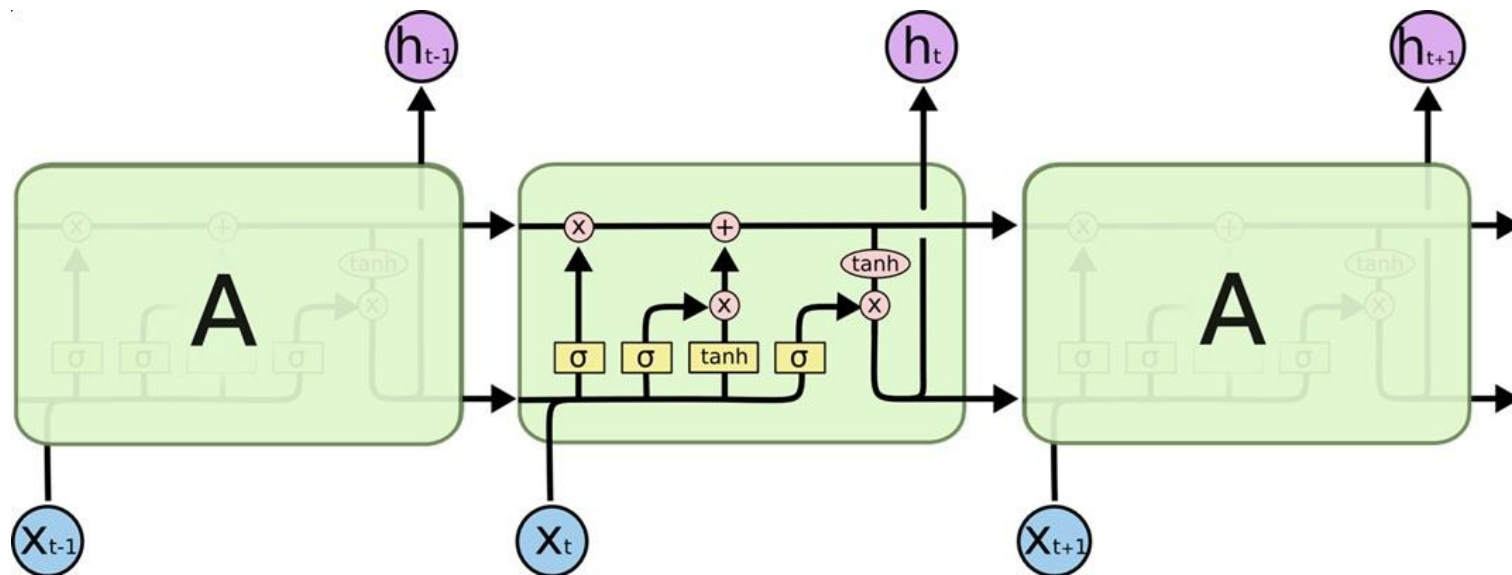


y_t est une fonction de h_t , par exemple $y_t = W h_t + b_y$ (relation linéaire).

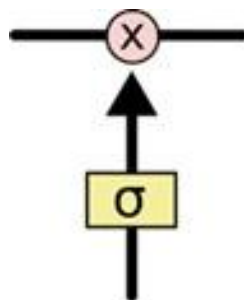
Amélioration par rapport au MLP : ajout d'une mémoire \Rightarrow partage des paramètres entre les pas de temps $t \Rightarrow$ plus facile à entraîner (nécessite moins de paramètres).

Limitations : temps-invariant, difficile à entraîner sur de longues séquences (perte de sensibilité = *vanishing gradient*).

LSTM :



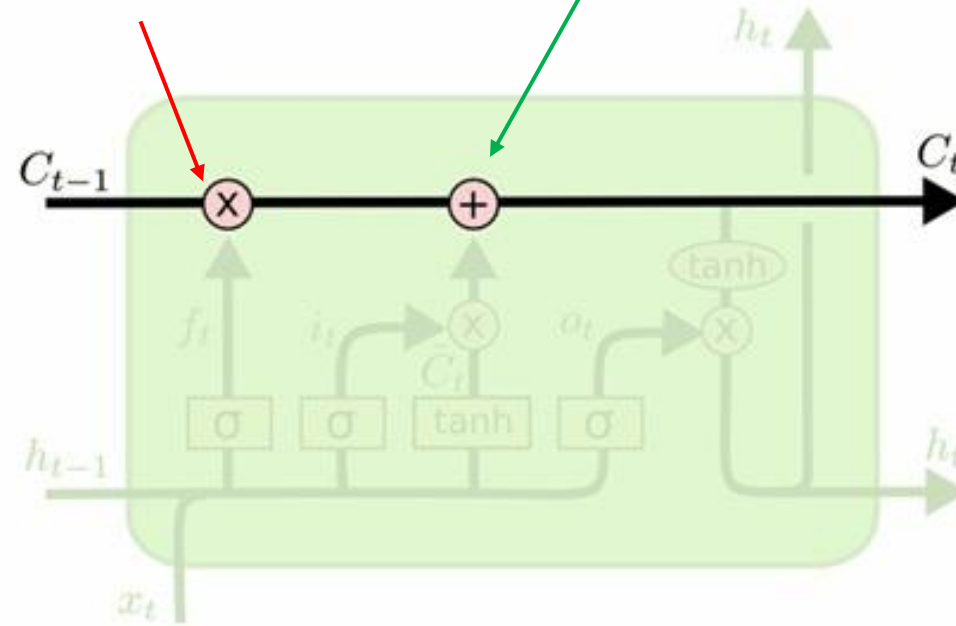
Mécanisme de porte (**gate**) pour décider si on laisse passer une information :
 σ représente la fonction sigmoïde, avec en sortie un signal entre 0 (ne laisse rien passer) et 1 (laisse tout passer).



LSTM :

Oublier (ou pas) une partie
de l'ancienne information C_{t-1}

Ajouter de nouvelles informations
à C_t



Expérimentalement, les LSTM ont montré leur efficacité pour modéliser des systèmes complexes.
En théorie, ils sont difficiles à analyser et restent encore mal compris.

4.3. Implémentation