

Deep Learning

Perceptrons simples et multicouches

Ricco Rakotomalala

Université Lumière Lyon 2



Plan

1. Perceptron simple
2. Plus loin avec le perceptron simple
3. Perceptron multicouche
4. Plus loin avec le perceptron multicouche
5. Pratique des perceptrons (sous R et Python)
6. Références
7. Conclusion



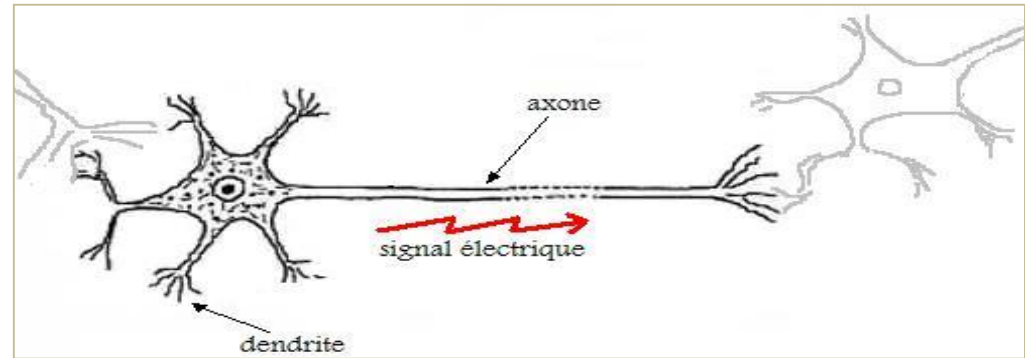
Métaphore biologique et transposition mathématique

PERCEPTRON SIMPLE

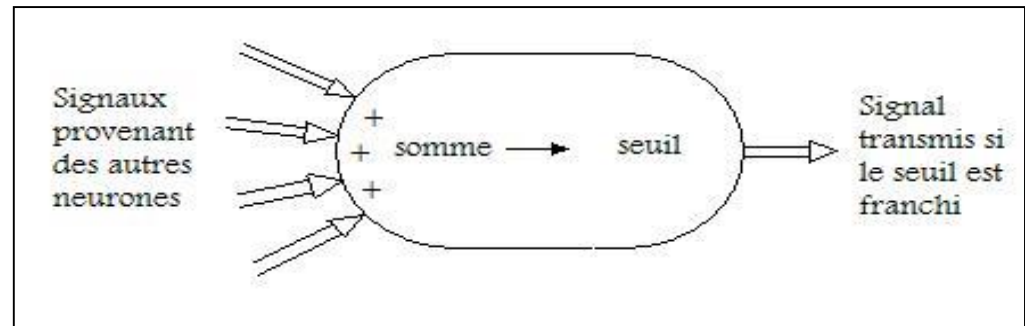
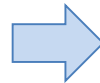


Métaphore biologique

Fonctionnement du cerveau
Transmission de l'information
et apprentissage



Idées maîtresses à retenir



Etapes clés :

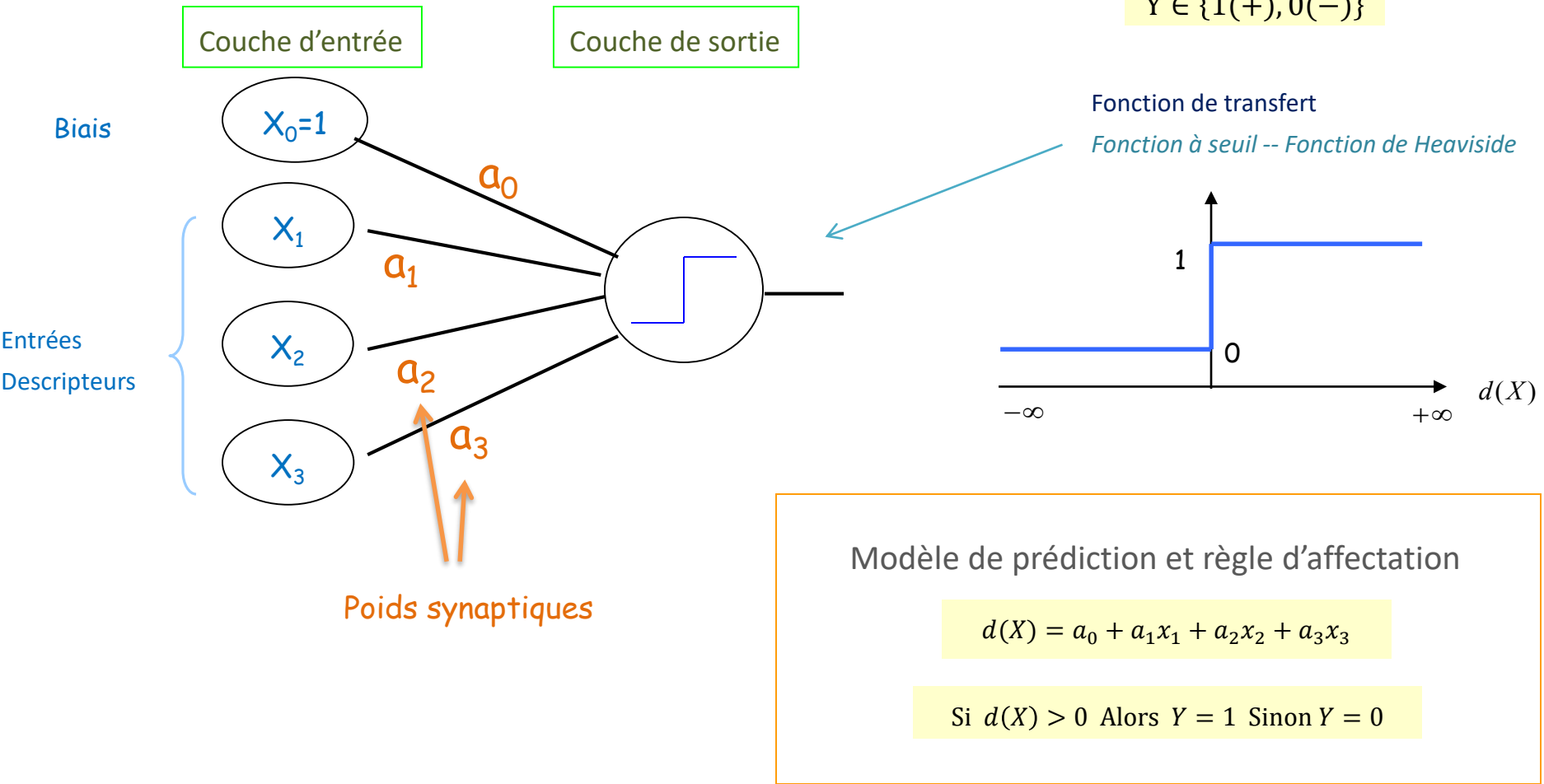
- Réception d'une information (signal)
- Activation + Traitement (simple) par un neurone
- Transmission aux autres neurones (si seuil franchi)
- A la longue : renforcement de certains liens → APPRENTISSAGE



Modèle de Mc Colluch et Pitts – Le Perceptron Simple

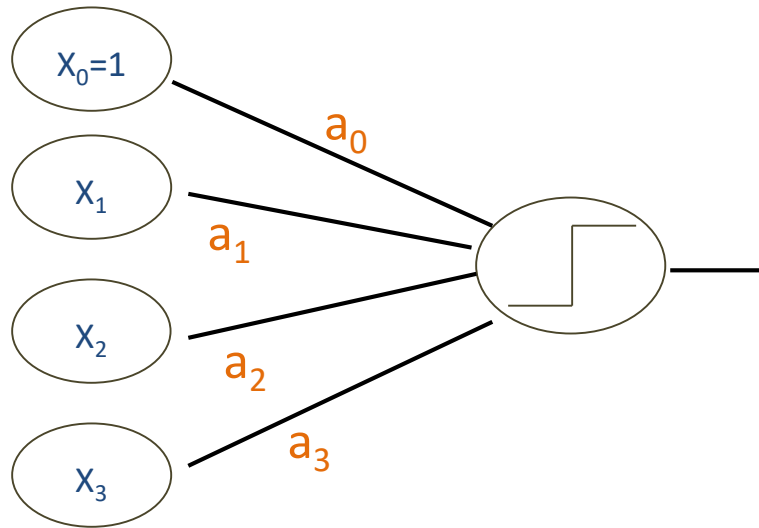
Problème à deux classes (positif et négatif)

$Y \in \{1(+), 0(-)\}$



Le perceptron simple est un modèle de prédiction (supervisé) linéaire





Comment calculer les poids synaptiques à partir d'un fichier de données ($Y ; X_1, X_2, X_3$)

Faire le parallèle avec la régression et les moindres carrés
Un réseau de neurones peut être utilisé pour la régression
(fonction de transfert avec une sortie linéaire)

- (1) Quel critère optimiser ?
- (2) Comment procéder à l'optimisation ?



- (1) Minimiser l'erreur de prédiction
- (2) Principe de l'incrémentalité (online)

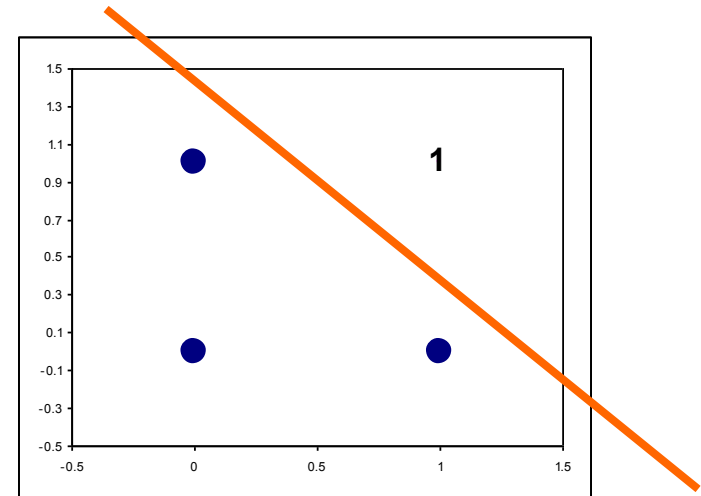


Exemple – Apprentissage de la fonction AND (ET logique)

Cet exemple est révélateur - Les premières applications proviennent de l'informatique

X1	X2	Y
0	0	0
0	1	0
1	0	0
1	1	1

Données



Représentation dans le plan

Principales étapes :

1. Mélanger aléatoirement les observations
2. Initialiser aléatoirement les poids synaptiques
3. Faire passer **les observations unes à unes**
 - Calculer l'erreur de prédiction pour l'observation
 - Mettre à jour les poids synaptiques
4. Jusqu'à **convergence** du processus

Une observation peut
passer plusieurs fois !



Exemple AND (1)

Règle de mise à jour des poids

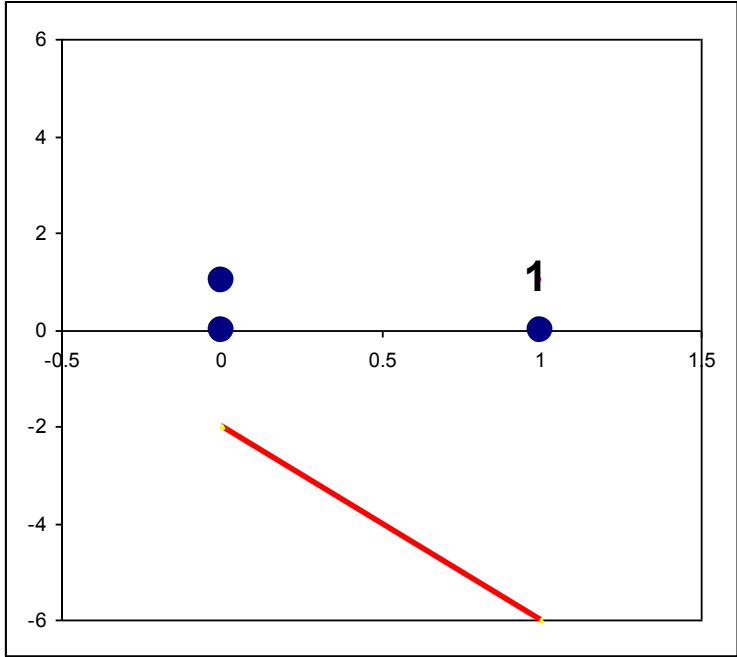
Pour chaque individu que l'on fait passer
(Principe de l'incrémentalité)

Initialisation aléatoire des poids : $a_0 = 0.1; a_1 = 0.2; a_2 = 0.05$

$$a_j \leftarrow a_j + \Delta a_j$$

Frontière :

$$0.1 + 0.2x_1 + 0.05x_2 = 0 \Leftrightarrow x_2 = -4.0x_1 - 2.0$$



S'assurer que les
variables sont sur la
même échelle
Force du signal

avec

$$\Delta a_j = \eta (y - \hat{y}) x_j$$

Erreur

Détermine s'il faut
réagir (corriger) ou non

Taux d'apprentissage

Détermine l'amplitude de l'apprentissage

Quelle est la bonne valeur ?

Trop petit \rightarrow lenteur de convergence

Trop grand \rightarrow oscillation

En général autour de 0.05 ~ 0.15 (0.1 dans notre exemple)

Ces 3 éléments sont au cœur du
mécanisme d'apprentissage



Exemple AND (2)

Observation à traiter

$$\begin{cases} x_0 = 1 \\ x_1 = 0 \\ x_2 = 0 \\ y = 0 \end{cases}$$



Appliquer le modèle

$$0.1 \times 1 + 0.2 \times 0 + 0.05 \times 0 = 0.1 \\ \Rightarrow \hat{y} = 1$$



Màj des poids

$$\begin{cases} \Delta a_0 = 0.1 \times (-1) \times 1 = -0.1 \\ \Delta a_1 = 0.1 \times (-1) \times 0 = 0 \\ \Delta a_2 = 0.1 \times (-1) \times 0 = 0 \end{cases}$$

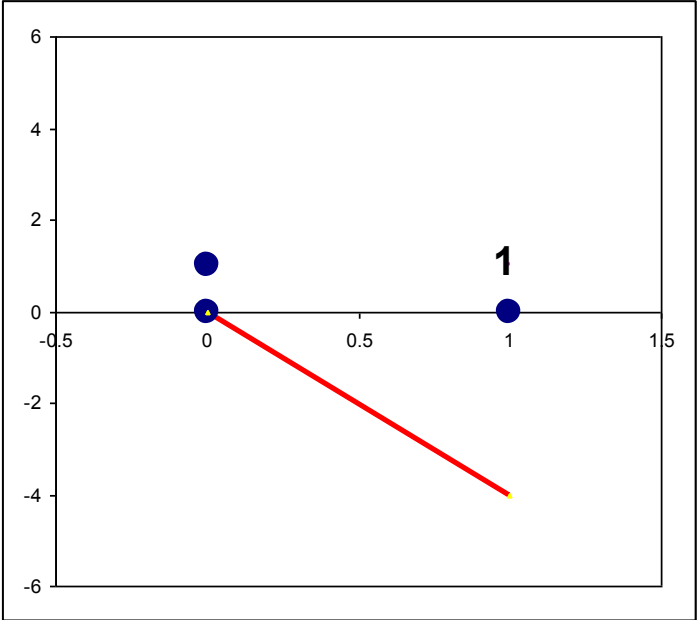


Signal nul ($x_1 = 0, x_2 = 0$), seule la constante a_0 est corrigée.

Valeur observée de Y et prédiction ne matchent pas, une correction des coefficients sera effectuée.

Nouvelle frontière :

$$0.0 + 0.2x_1 + 0.05x_2 = 0 \Leftrightarrow x_2 = -4.0x_1 + 0.0$$



Exemple AND (3)

Observation à traiter

$$\begin{cases} x_0 = 1 \\ x_1 = 1 \\ x_2 = 0 \\ y = 0 \end{cases}$$

Appliquer le modèle

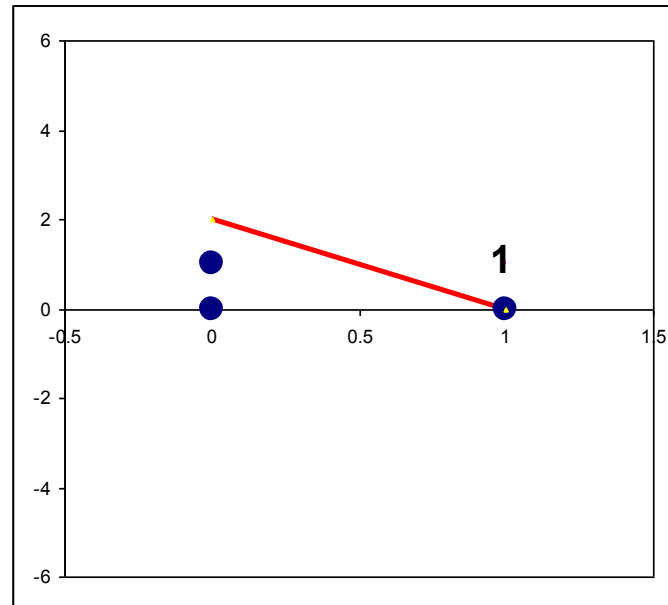
$$0.0 \times 1 + 0.2 \times 1 + 0.05 \times 0 = 0.2$$
$$\Rightarrow \hat{y} = 1$$

Màj des poids

$$\begin{cases} \Delta a_0 = 0.1 \times (-1) \times 1 = -0.1 \\ \Delta a_1 = 0.1 \times (-1) \times 1 = -0.1 \\ \Delta a_2 = 0.1 \times (-1) \times 0 = 0 \end{cases}$$

Nouvelle frontière :

$$-0.1 + 0.1x_1 + 0.05x_2 = 0 \Leftrightarrow x_2 = -2.0x_1 + 2.0$$



Exemple AND (4) – Définir la convergence

Observation à traiter

$$\begin{cases} x_0 = 1 \\ x_1 = 0 \\ x_2 = 1 \\ y = 0 \end{cases}$$

Appliquer le modèle

$$-0.1 \times 1 + 0.1 \times 0 + 0.05 \times 1 = -0.05$$
$$\Rightarrow \hat{y} = 0$$

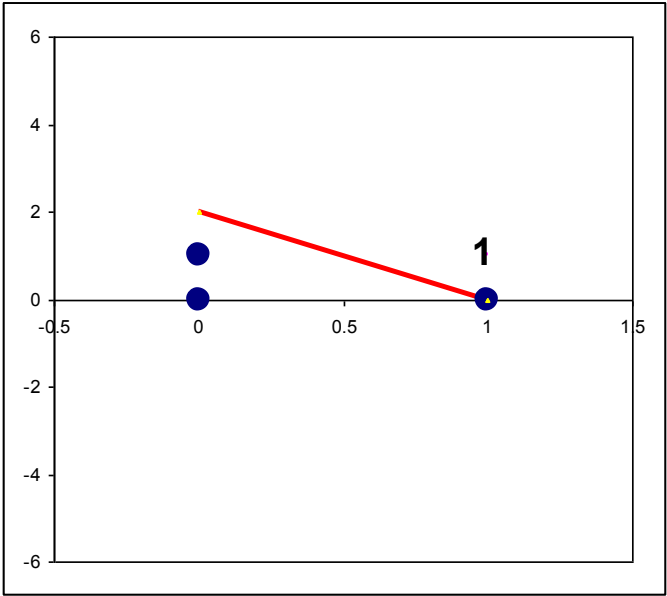
Màj des poids

$$\begin{cases} \Delta a_0 = 0.1 \times (0) \times 1 = 0 \\ \Delta a_1 = 0.1 \times (0) \times 0 = 0 \\ \Delta a_2 = 0.1 \times (0) \times 1 = 0 \end{cases}$$

Pas de correction ici ? Pourquoi ?
Voir aussi la position du point par rapport à la frontière dans le plan !

Frontière inchangée :

$$-0.1 + 0.1x_1 + 0.05x_2 = 0 \Leftrightarrow x_2 = -2.0x_1 + 2.0$$



Remarque : Que se passe-t-il si on repasse l'individu ($x_1=1$; $x_2=0$) ?

Convergence ?

- (1) Plus aucune correction effectuée en passant tout le monde
- (2) L'erreur globale ne diminue plus « significativement »
- (3) Les poids sont stables
- (4) On fixe un nombre maximum d'itérations
- (5) On fixe une erreur minimale à atteindre

(2), (4) et (5) deviennent des « paramètres » de l'algorithme à considérer
(attention aux valeurs par défaut) dans les logiciels !!! Il y en aura d'autres...



Mort et résurrection du perceptron

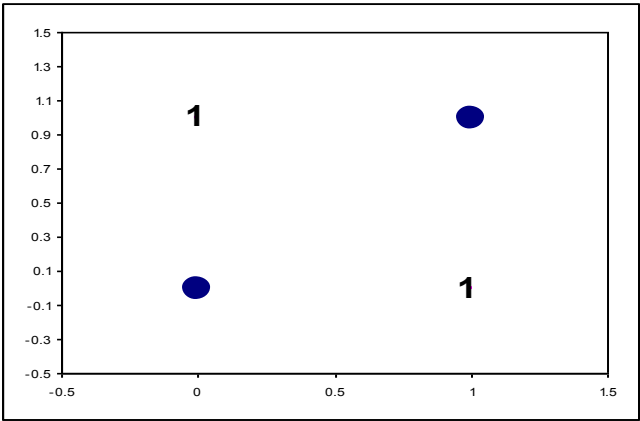
PERCEPTRON MULTICOUCHE



Problème du XOR – L'impossible séparation linéaire

X1	X2	Y
0	0	0
0	1	1
1	0	1
1	1	0

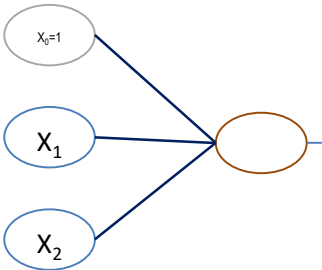
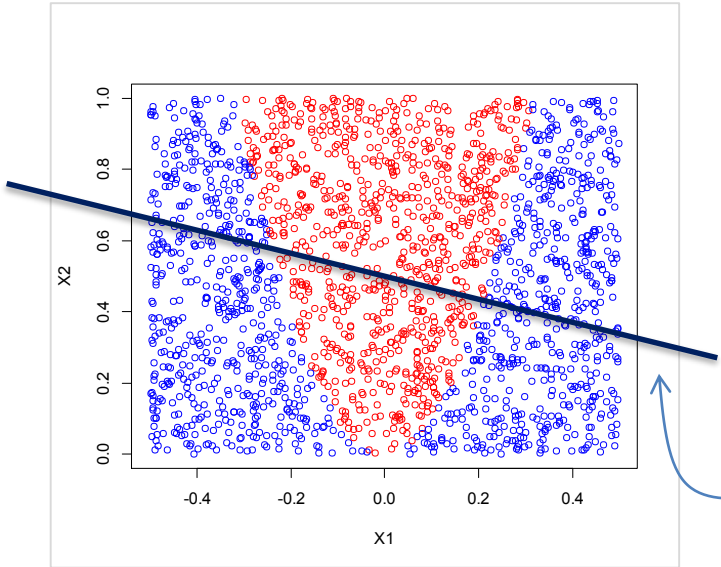
Données



Non séparable linéairement
(Minsky & Papert, 1969)



Un perceptron simple ne sait
traiter que les problèmes
linéairement séparables.

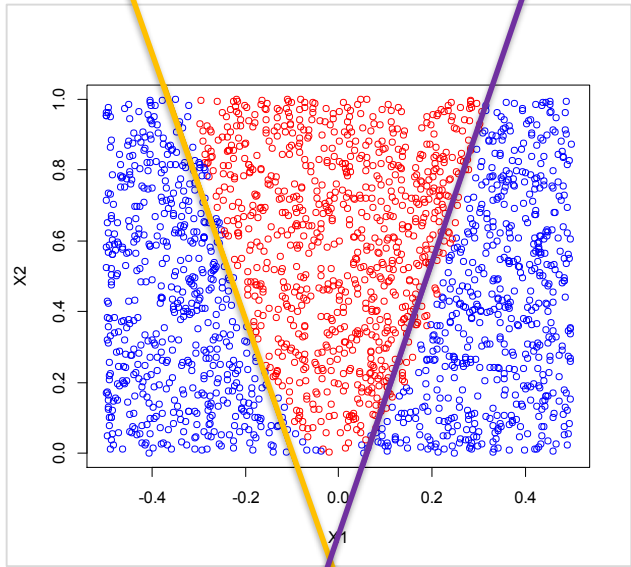


Trouver une droite de séparation
efficace n'est pas possible.

Perceptron multicouche - Principe

$26.0 + 266.76 \times X1 + 63.08 \times X2 = 0$

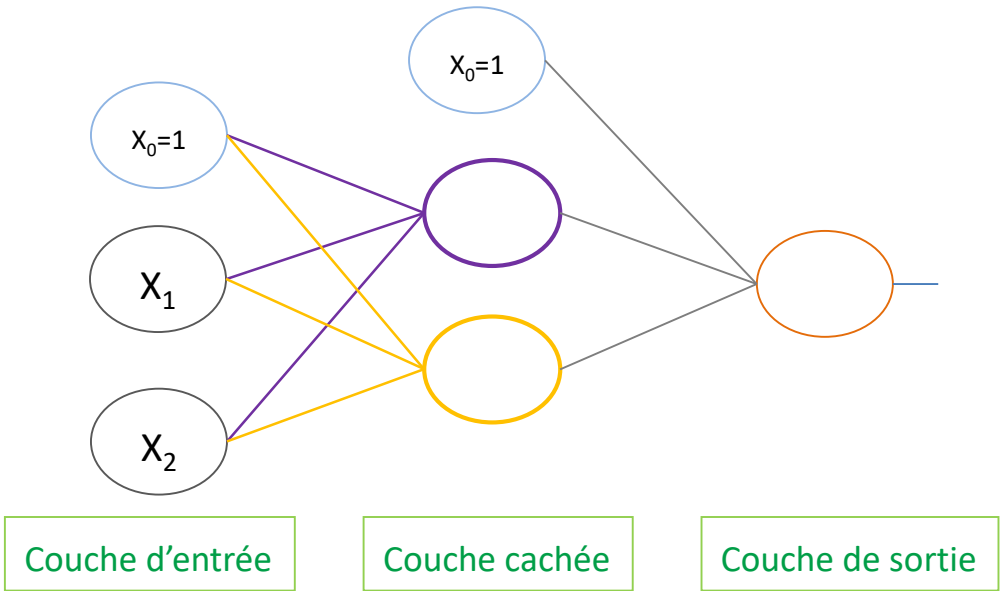
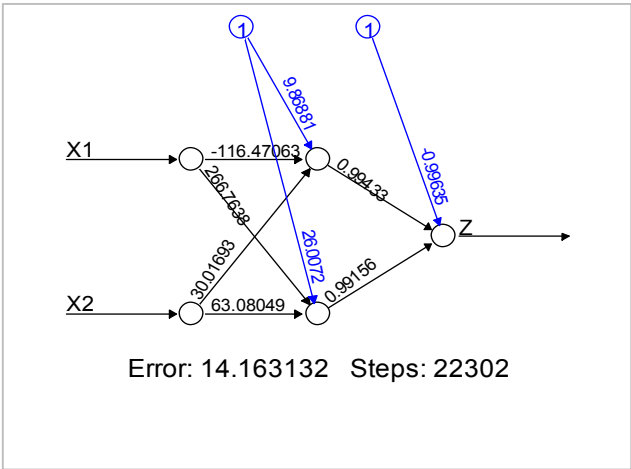
$9.86 - 116.47 \times X1 + 30.01 \times X2 = 0$



Perceptron Multicouche (PMC)

Une combinaison de séparateurs linéaires permet de produire un séparateur global non-linéaire (Rumelhart, 1986).

On peut avoir plusieurs couches cachées, cf. plus loin





Classifieur très précis (si bien paramétré)

Incrémentalité

Scalabilité (capacité à être mis en œuvre sur de grandes bases)



Modèle boîte noire (causalité descripteur – variable à prédire)

Difficulté de paramétrage (nombre de neurones dans la couche cachée)

Problème de convergence (optimum local)

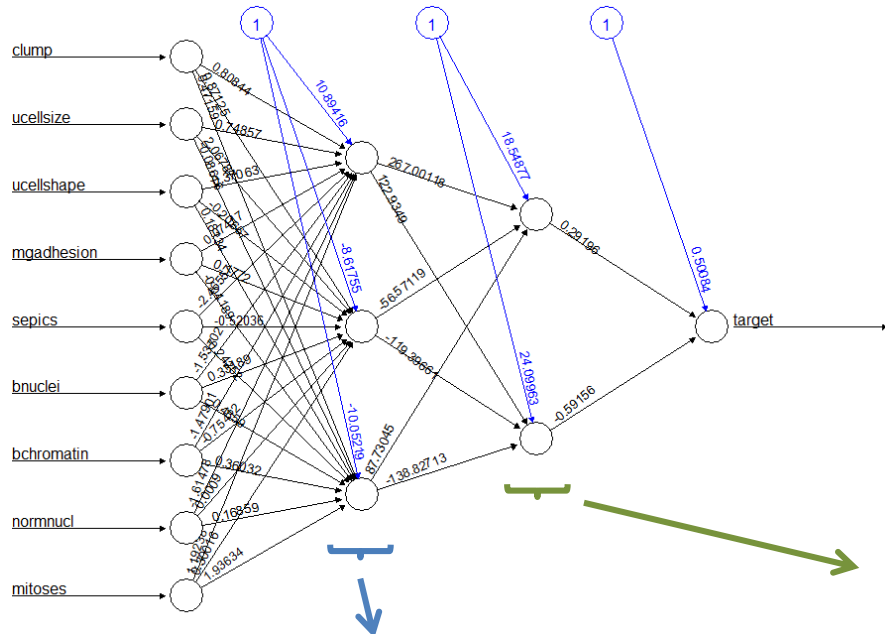
Danger de sur-apprentissage (trop de neurones dans la couche cachée)



PMC avec plusieurs couches cachées

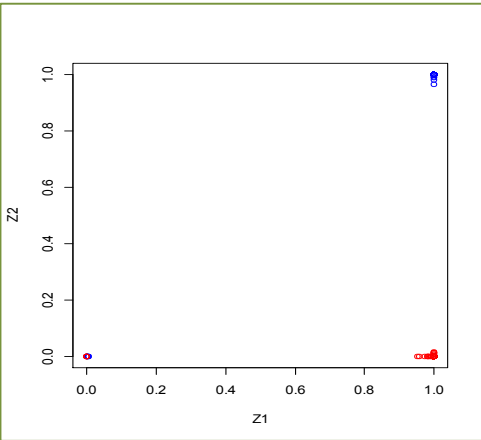
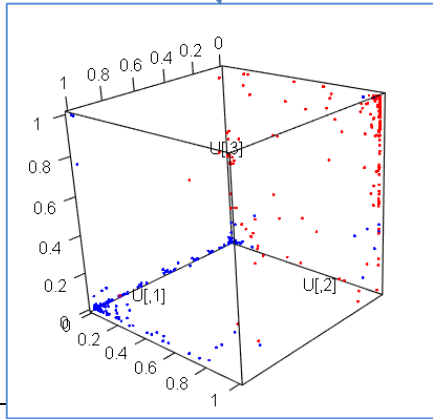
Pourquoi plusieurs couches cachées ? Démultiplie le pouvoir explicatif du modèle, mais : (1) le paramétrage est quasi inextricable (sauf à l'utiliser les couches intermédiaires comme filtres, cf. les réseaux de convolution) ; (2) la lecture des couches (des systèmes de représentation) intermédiaires n'est pas évidente.

« Breast cancer »
dataset



A noter : (A) Dès l'espace U, on avait une bonne discrimination ;
(B) Z2 suffit à la discrimination ;
(C) saturation des valeurs à 0 ou 1.

Espace de représentation à 3 dimensions.



Espace de représentation à 2 dimensions. Séparation linéaire possible.

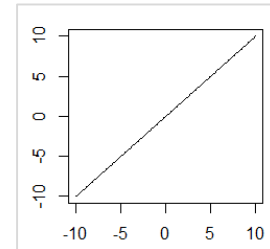
Fonctions d'activation

Différents types de fonctions d'activation sont utilisables. En fonction du problème à traiter, de la définition des espaces intermédiaires, du filtrage que l'on veut effectuer sur les données.... Il faut être attentif à ce que l'on veut obtenir.

Fonction linéaire

Pour la régression, pas de transformation

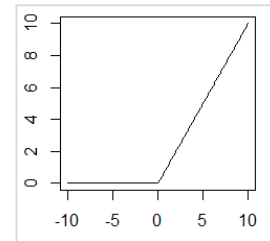
$$u = x$$



Fonction ReLu (Rectified Linear units)

Filtre les valeurs négatives

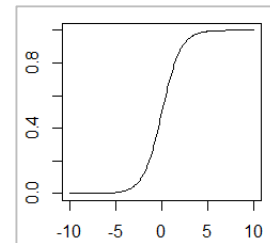
$$u = \max(0, x)$$



Fonction Sigmoide

Ramène les valeurs entre [0, 1]. Utilisé pour le classement. Y codé {0,1}.

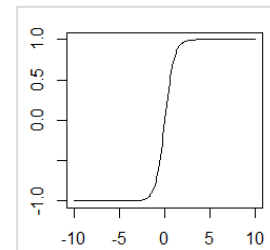
$$u = \frac{1}{1 + e^{-x}}$$



Fonction Tangente hyperbolique

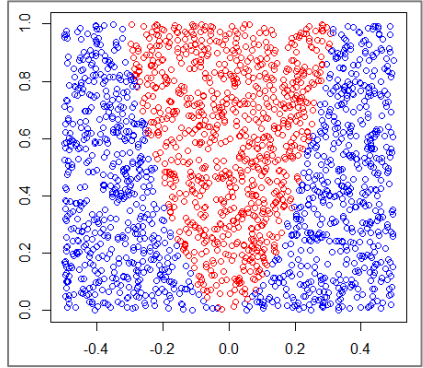
Ramène les valeurs entre [-1, 1].
Alternative à sigmoïde pour le classement. Y codé {-1,+1}.

$$u = \frac{e^{2x} - 1}{e^{2x} + 1}$$



Fonctions d'activation

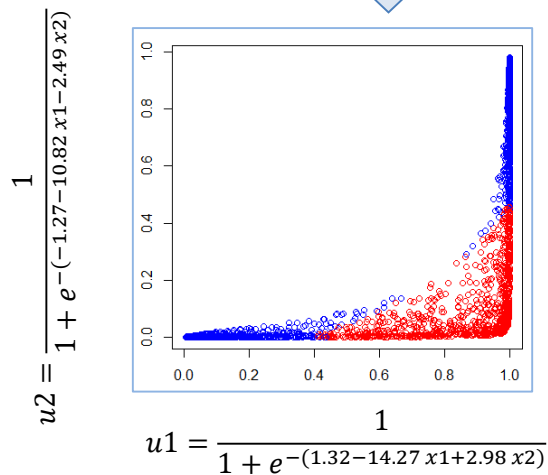
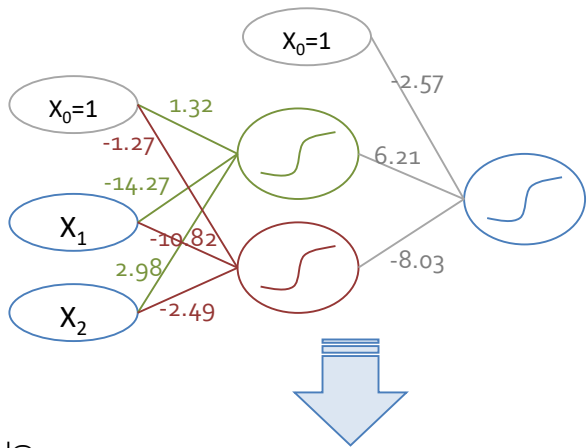
Mixer les fonctions dans un PMC



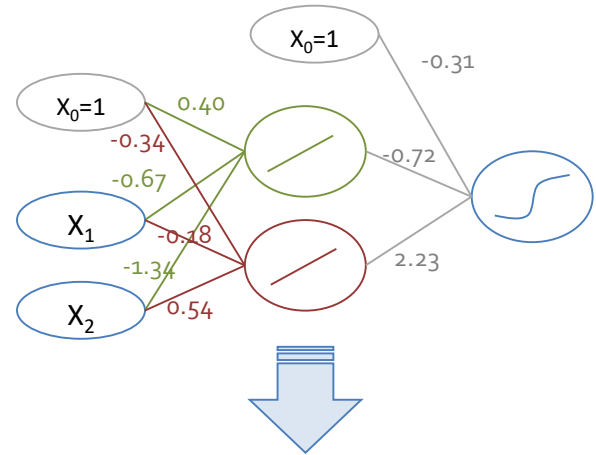
Problème à deux classes. A l'évidence une couche cachée avec 2 neurones suffit... mais cela est-il valable pour tous types de fonctions d'activations ?



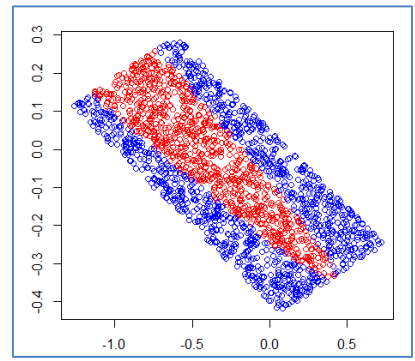
Fonction sigmoïde dans toutes les couches.



Fonction linéaire dans la couche cachée.
Fonction sigmoïde dans la sortie.



Avec la fonction d'activation linéaire, l'espace intermédiaire n'est pas discriminant (linéairement s'entend).

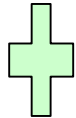


$z1 = 0.40 - 0.67 x_1 - 1.34 x_2$

$z2 = -0.34 - 0.18 x_1 + 0.54 x_2$



De nombreuses améliorations ces dernières années



Notamment avec des algorithmes d'optimisation (au-delà du gradient stochastique) performants

Librairies de calcul puissantes disponibles sous R et Python

Attention au paramétrage



Bien lire la documentation des packages pour ne pas s'y perdre

Faire des tests en jouant sur les paramètres « simples » (ex. architecture du réseau) dans un premier temps, affiner par la suite.



Librairies sous Python et R

PRATIQUE DES PERCEPTRONS



[Scikit Learn](https://scikit-learn.org/) est une librairie de machine learning puissante pour Python.

```
#importation des données
import pandas
D = pandas.read_table("artificial2d_data2.txt",sep="\t",header=0)

#graphique
code_couleur = D['Y'].eq('pos').astype('int')
D.plot.scatter(x="X1",y="X2",c=pandas.Series(['blue','red'])[code_couleur])

#séparer cible et descripteurs
X = D.values[:,0:2]
Y = D.values[:,2]

#subdivision en apprentissage et test
from sklearn import model_selection
XTrain,XTest,YTrain,YTest = model_selection.train_test_split(X,Y,train_size=1000)

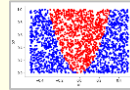
#initialisation du classifieur
from sklearn.neural_network import MLPClassifier
rna = MLPClassifier(hidden_layer_sizes=(2,),activation="logistic",solver="lbfgs")

#apprentissage
rna.fit(XTrain,YTrain)

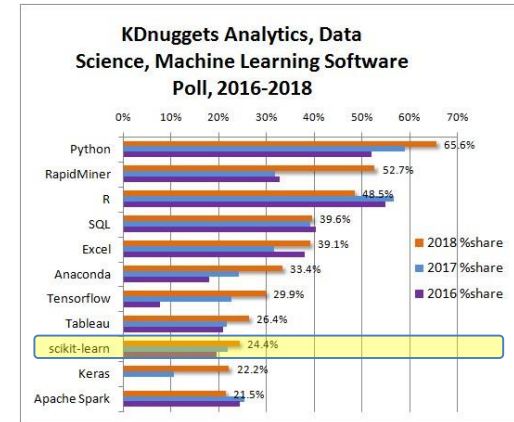
#affichage des coefficients
print(rna.coefs_)
print(rna.intercepts_)

#prédiction sur l'échantillon test
pred = rna.predict(XTest)
print(pred)

#mesure des performances
from sklearn import metrics
print(metrics.confusion_matrix(YTest,pred))
print("Taux erreur = " + str(1-metrics.accuracy_score(YTest,pred)))
```



1000 TRAIN, 1000 TEST



solver : {'lbfgs', 'sgd', 'adam'}, default 'adam'

The solver for weight optimization.

- 'lbfgs' is an optimizer in the family of quasi-Newton methods.
- 'sgd' refers to stochastic gradient descent.
- 'adam' refers to a stochastic gradient-based optimizer proposed by Kingma, Diederik, and Jimmy Ba

Note: The default solver 'adam' works pretty well on relatively large datasets (with thousands of training samples or more) in terms of both training time and validation score. For small datasets, however, 'lbfgs' can converge faster and perform better.

De l'importance du paramétrage. cf. [LBFGS](#).

	neg	pos
neg	565	5
pos	2	428

$\epsilon = 0.007$



```
#importer le package
library(keras)

#construire l'architecture du perceptron
rna <- keras_model_sequential()
rna %>%
  layer_dense(units = 2, input_shape = c(2), activation = "sigmoid") %>%
  layer_dense(units = 1, activation = "sigmoid")
```

```
#paramétrage de l'algorithme
rna %>% compile(
  loss="mean_squared_error",
  optimizer=optimizer_sgd(lr=0.15),
  metrics="mae"
)
```

```
#codage de la cible - éviter la saturation
yTrain <- ifelse(DTrain$Y=="pos",0.8,0.2)
```

```
#apprentissage avec son paramétrage
rna %>% fit(
  x = as.matrix(DTrain[,1:2]),
  y = yTrain,
  epochs = 500,
  batch_size = 10
)
```

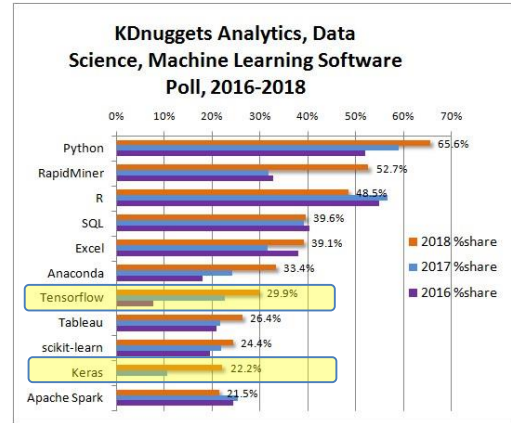
```
#affichage des poids calculés
print(keras::get_weights(rna))
```

```
#prédiction sur l'échantillon test
pred <- rna %>% predict_classes(as.matrix(DTest[,1:2]))
```

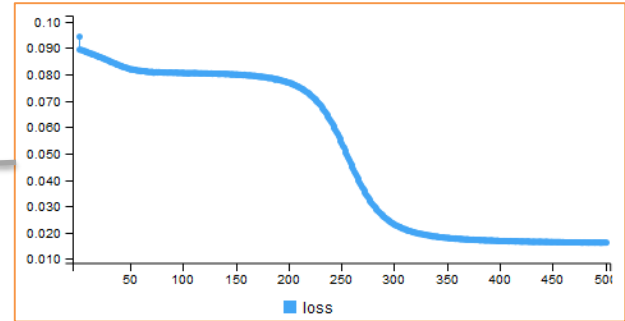
```
#matrice de confusion
print(mc <- table(DTest$Y,pred))
print(paste("Taux erreur =", 1-sum(diag(mc))/sum(mc)))
```

Keras sous R

[Keras](#) repose sur la technologie [Tensorflow](#). Ces deux bibliothèques de « deep learning » sont très populaires.



Pour éviter que les valeurs des combinaisons linéaires soient élevées en valeurs absolues, et que la transformation sigmoïde sature à 0 ou 1 (dans la zone où la dérivée est nulle), on peut ruser en codant Y de manière à se situer plus vers la partie linéaire de la fonction d'activation.



L'évolution de la perte est assez singulière quand-même. Si nous avions fixé (epochs = nombre de passages sur la base ≤ 200), nous n'aurions pas obtenu le bon résultat !

	neg	pos
neg	565	10
pos	1	424

$\epsilon = 0.011$