

## DÉMYSTIFICATION DE L'IA

Prise en main du computer vision



- Avoir une vision technique de ce qu'est l'I.A
- Savoir explorer / comprendre / interpréter un jeu de données
- Savoir entraîner / test / évaluer un modèle de Deep Learning
- Connaître les principales architectures de Deep Learning
- Savoir utiliser la bonne architecture
- Identifier la faisabilité d'un projet de DL
- Connaître les technologies du marché
- Identifier les compétences recherchées



**01**

**INTRODUCTION AU DEEP  
LEARNING**

## CONTENU



- Qu'est-ce que le DL
- Pourquoi le DL
- Définition d'un perceptron
- Backpropagation

## COMPÉTENCES ATTENDUES

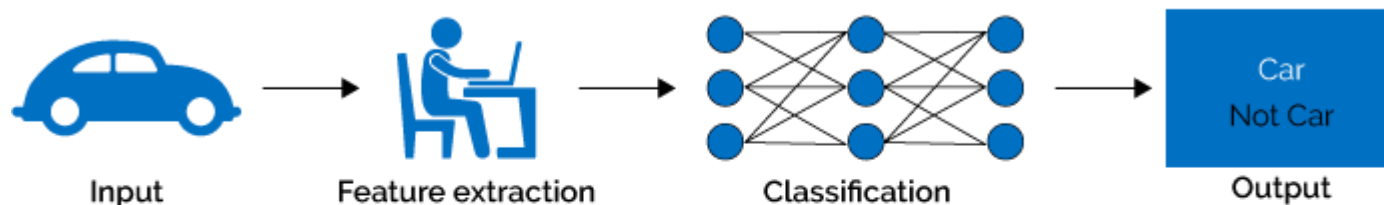


- Définir le Deep Learning (DL)
- Différencier l'approche DL de l'approche ML
- Comprendre la logique de l'apprentissage
- Evaluer un réseau de neurones

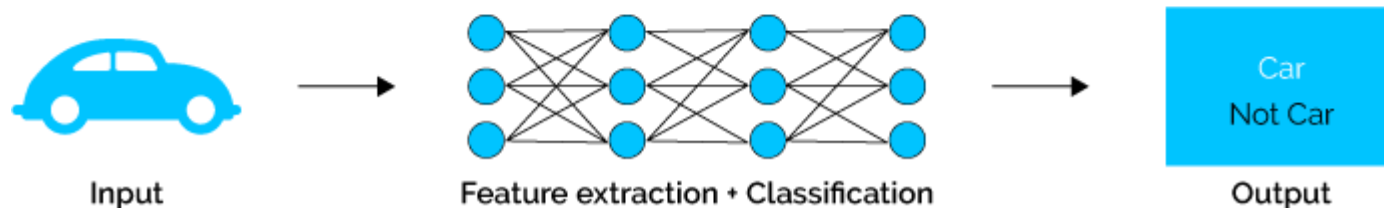
“L'apprentissage profond est un ensemble de méthodes d'apprentissage automatique tentant de modéliser avec un haut niveau d'abstraction des données grâce à des architectures articulées de différentes transformations non linéaires.”

Wikipedia

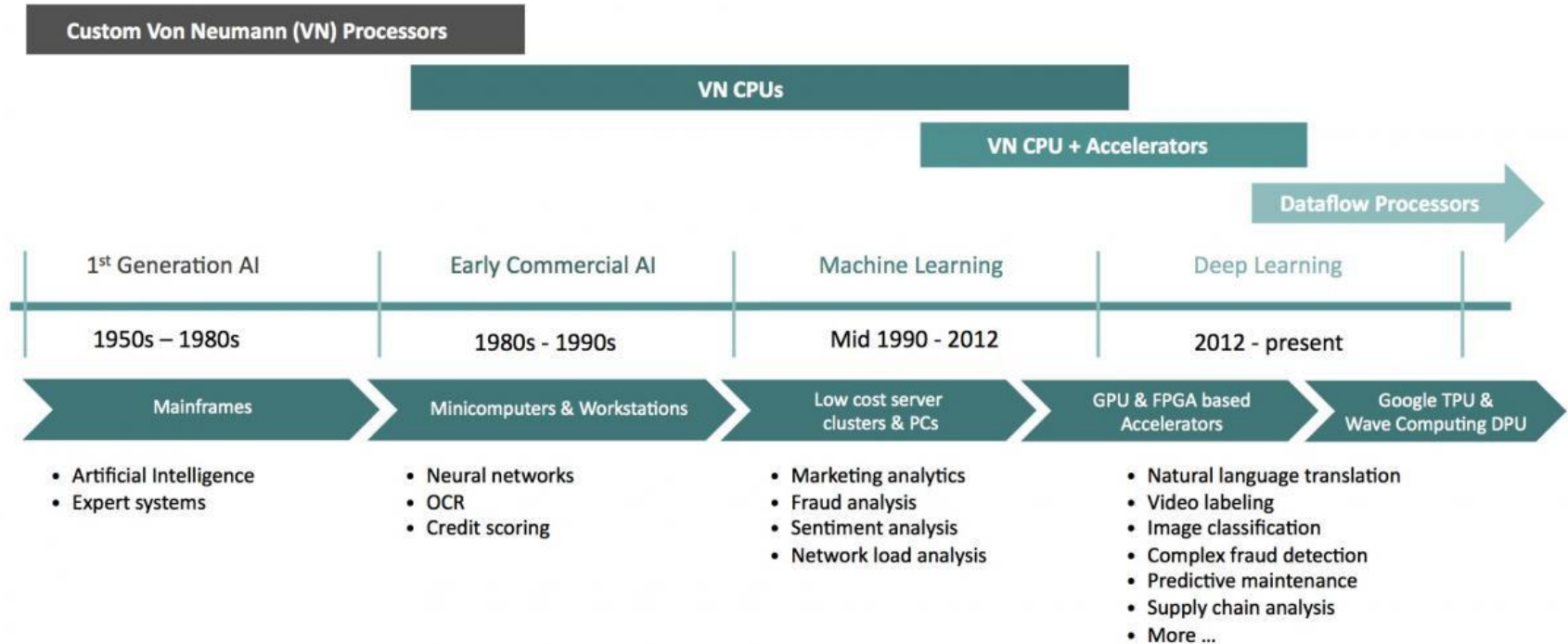
## Machine Learning

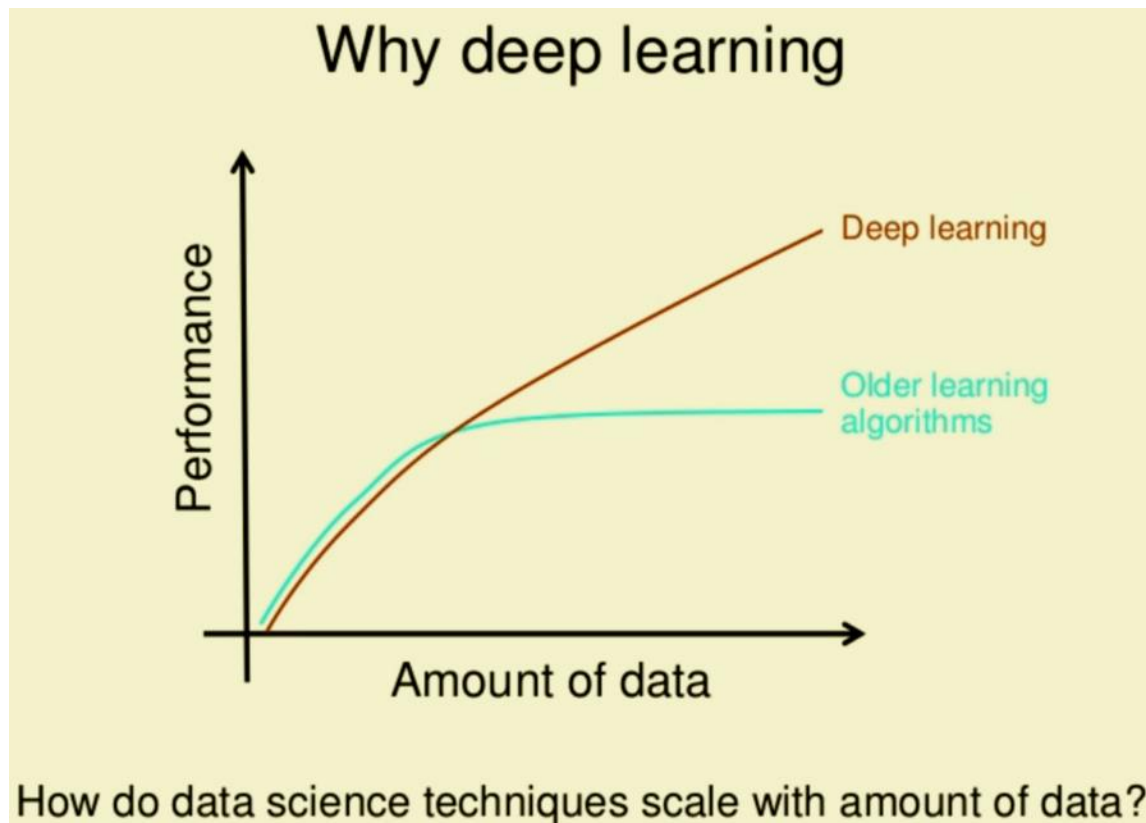


## Deep Learning



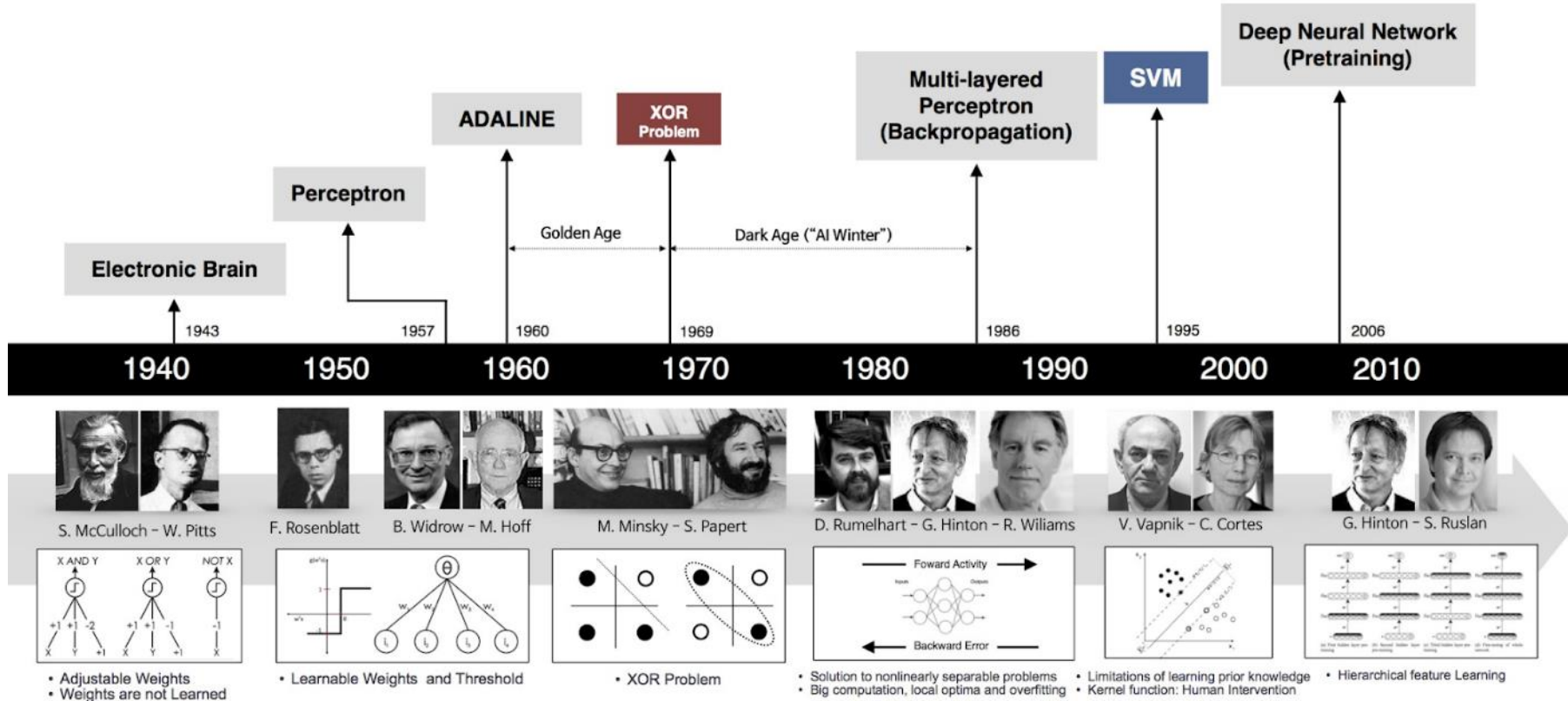
# HISTORIQUE DES MÉTHODES







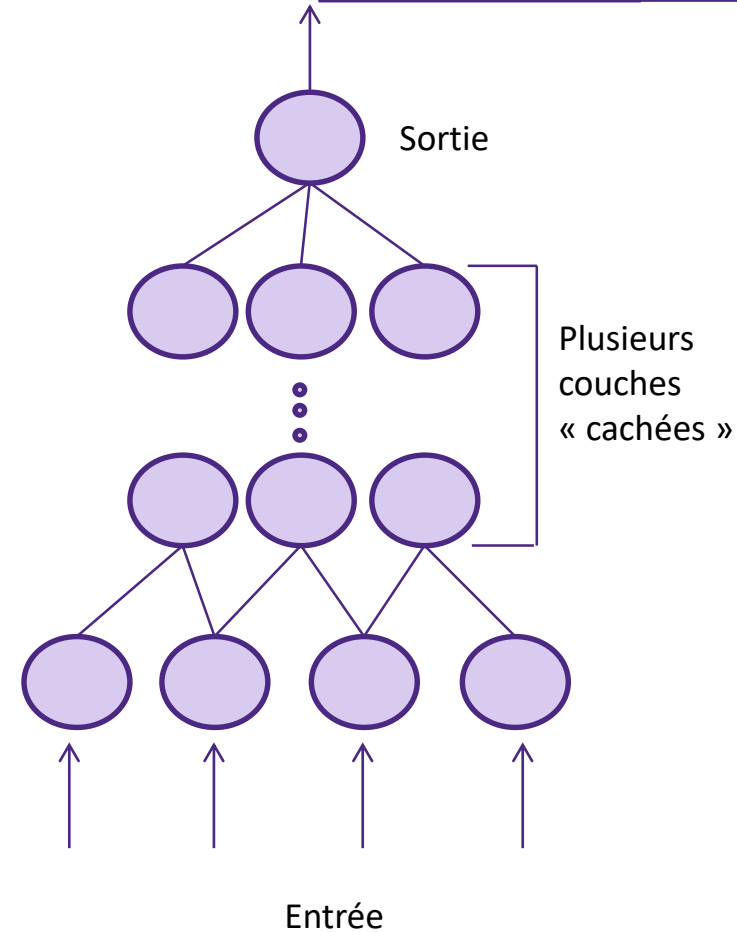
# HISTORIQUE DES RÉSEAUX DE NEURONES



# HISTORIQUE DES RÉSEAUX DE NEURONES (2)

## Une histoire ancienne :

- **Premier modèle de Neurone formel en 1943 :**
  - Pas d'apprentissage
- **Règle de Hebb en 1949**
  - Hypothèse sur l'apprentissage des connexions
- **Le Perceptron (1958)**
  - Introduction de l'apprentissage
  - Problème du « OU » Exclusif et abandon (1969)
- **Rétro-propagation du gradient (1985)**
  - Perceptron multi-couches
  - Ancêtre du deep learning



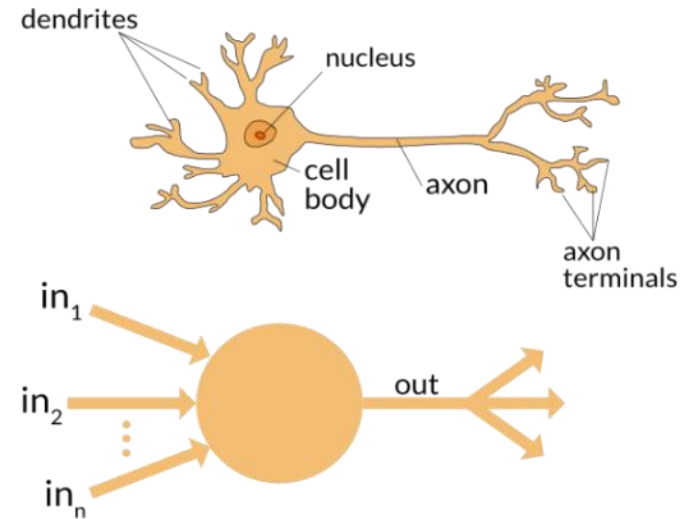
**Des succès dans de nombreux domaines à partir des années 2000 :**

- **Premiers succès d'un Long short-term Memory (LSTM) pour la reconnaissance du dialogue (1998)**
- **Entraînement effectif de réseaux avec de très nombreuses couches (2006)**
  - **Début du deep learning avec des réseaux de plus en plus profond**
- **Les réseaux de neurones convolutifs deviennent state-of-the-art en classification d'images (2012).**
- **Alphago bat le champion du monde de Go en mai 2017**

# QU'EST-CE QU'UN RÉSEAU DE NEURONES ?

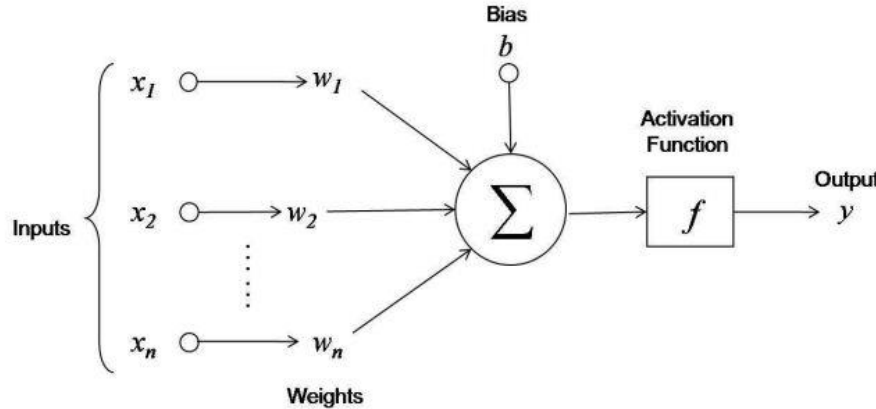
## UN PEU DE VOCABULAIRE

- **Algorithme de Machine Learning inspiré du fonctionnement du cerveau**
- **Composé d'un ensemble d'unités appelées « neurones » fonctionnant en parallèle**
- **Ces neurones sont reliés entre eux par des connexions (« synapses »)**



# QU'EST-CE QU'UN RÉSEAU DE NEURONES ?

## UN PEU DE VOCABULAIRE



Pré-activation :

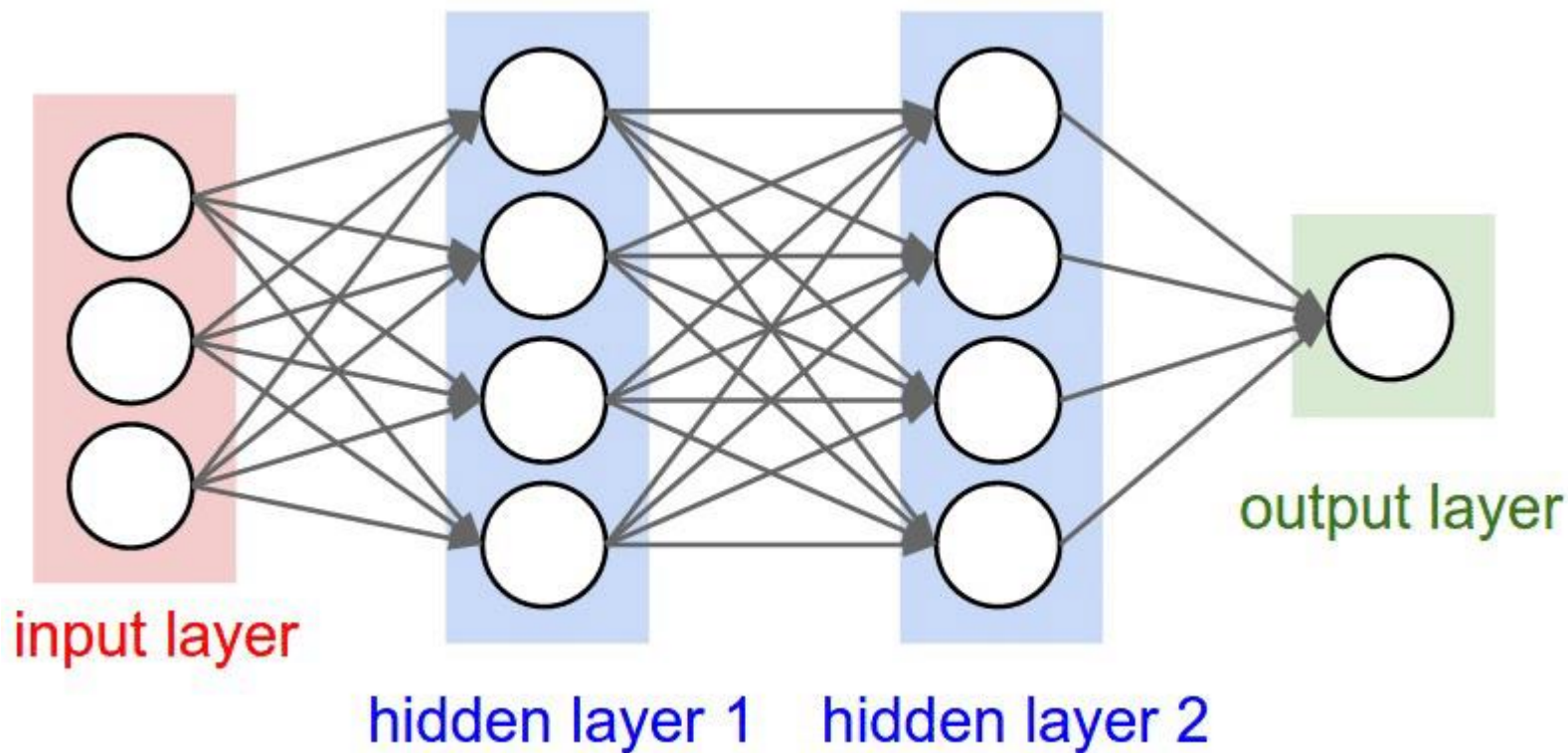
$$a(x) = b + \sum_i w_i x_i$$

Activation (output) :

$$h(x) = g(a(x)) = g(b + \sum_i w_i x_i)$$

- Le réseau apprend en faisant varier les valeurs des connexions, ou « poids »
- Le terme « Deep learning » désigne les réseaux de neurone avec un nombre important de couches

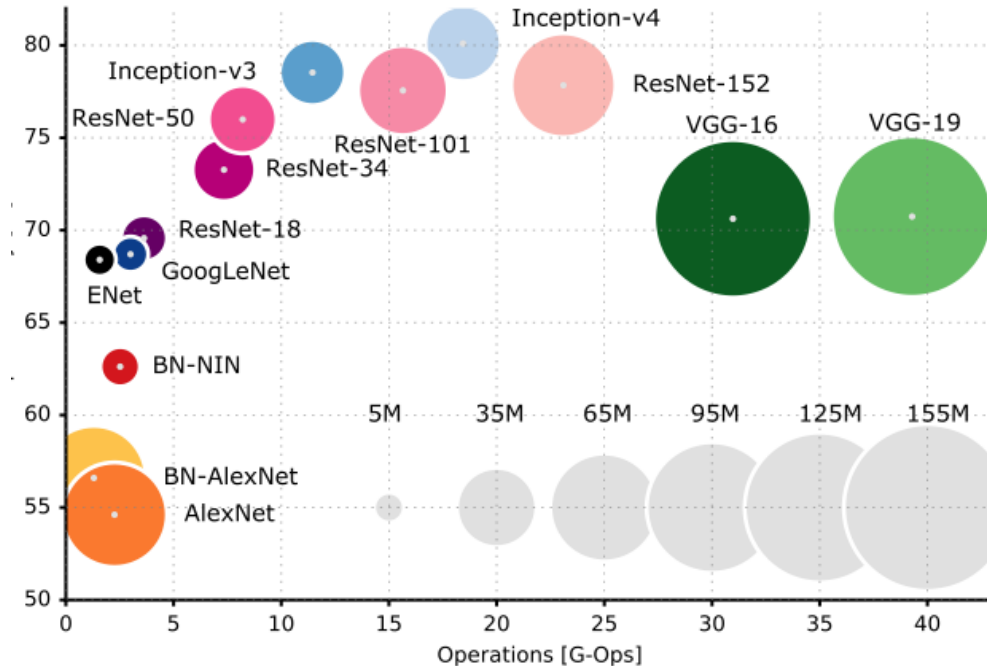
- $\mathbf{w} \Rightarrow$  poids
- $\mathbf{b} \Rightarrow$  biais
- $\mathbf{g}(\cdot) \Rightarrow$  fonction d'activation



# UN EXEMPLE D'ÉVOLUTION DES RÉSEAUX DE NEURONES

## LE CHALLENGE IMAGENET

- En 2011, 25% d'erreurs
- Envahi par les réseaux de neurones depuis 2012
- 16% d'erreur par un réseau de convolution en 2012
- Meilleur que l'humain depuis 2015
- ... mais sur moins de catégories



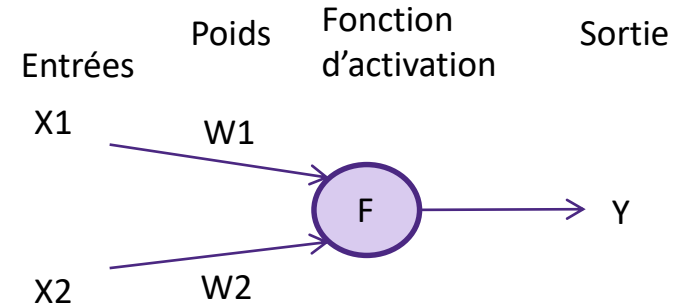
## LE PERCEPTRON SIMPLE

- **Perceptron simple = Un neurone**
- **La sortie est une fonction des poids et des entrées**

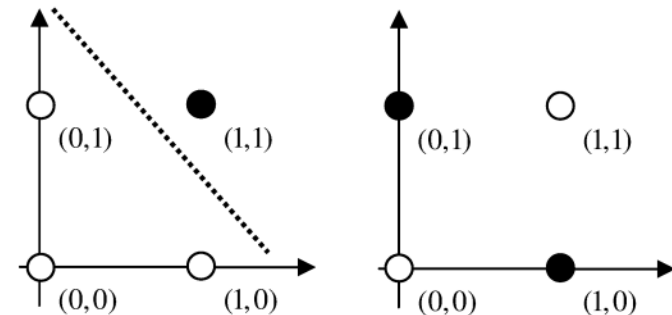
$$Y = f(W1 \cdot X1 + W2 \cdot X2)$$

- **Ne peut pas résoudre le « Ou exclusif »**
- **En pratique, les calculs sont fait couche par couche avec des opérations matricielles**
  - + efficace, notamment par GPU

### Le perceptron simple



### Le « Ou exclusif »



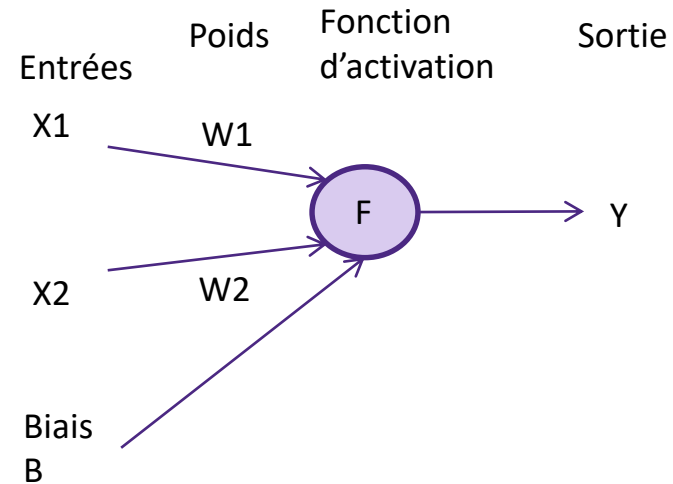
Réponse d'un neurone



- Ajout d'un biais sur les neurones

$$Y = f(W1 \ X1 + W2 \ X2 + \textcolor{violet}{B})$$

- Permet de décaler la fonction d'activation
- Les poids changent la pente de la fonction
- Permet une réponse du neurone même si les entrées sont vides

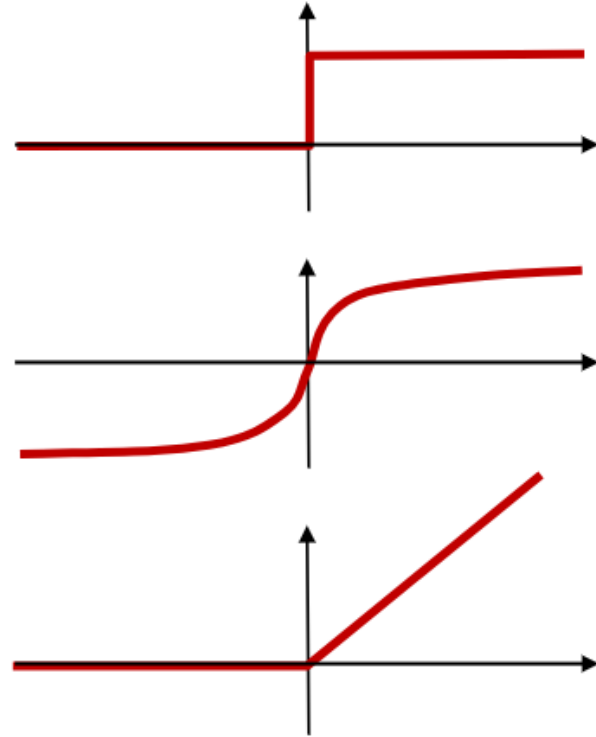


## LES FONCTIONS D'ACTIVATION

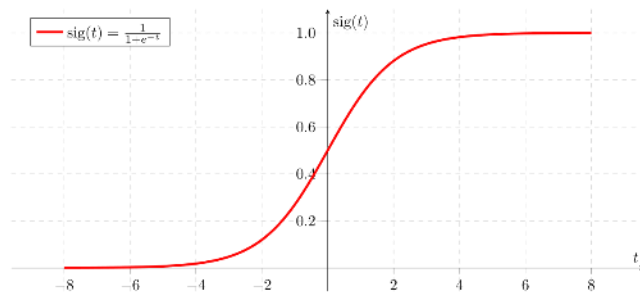
Différentes fonctions d'activation.

Comment les choisir ?

- **Fonction Heaviside :**
  - Pour prendre une décision
  - Ex: pour la couche de sortie du réseau
- **Sigmoïde (et Linéaire) :**
  - Evite que l'activité dans les réseaux n'explose
  - Problème de disparition du gradient de l'erreur
  - Adaptée si peu de couches
- **Rectified Linear Unit (ReLU) :**
  - Pas de problème de disparition du gradient
  - « Sparse »
  - State-of-the-art dans des réseaux à nombreuses couches.

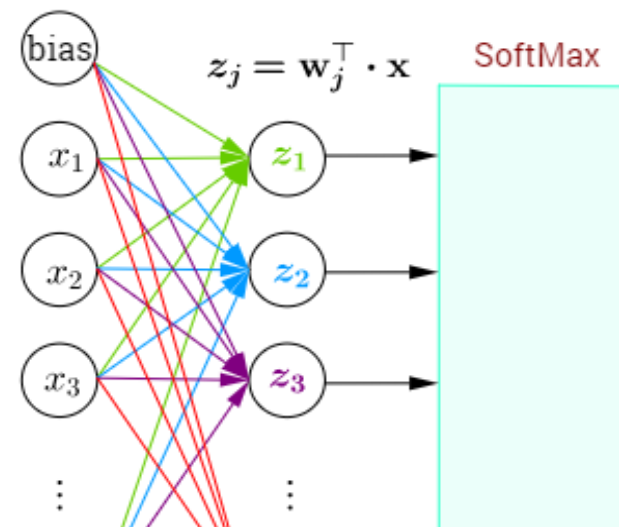


## Binaire



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

## Multi-classes

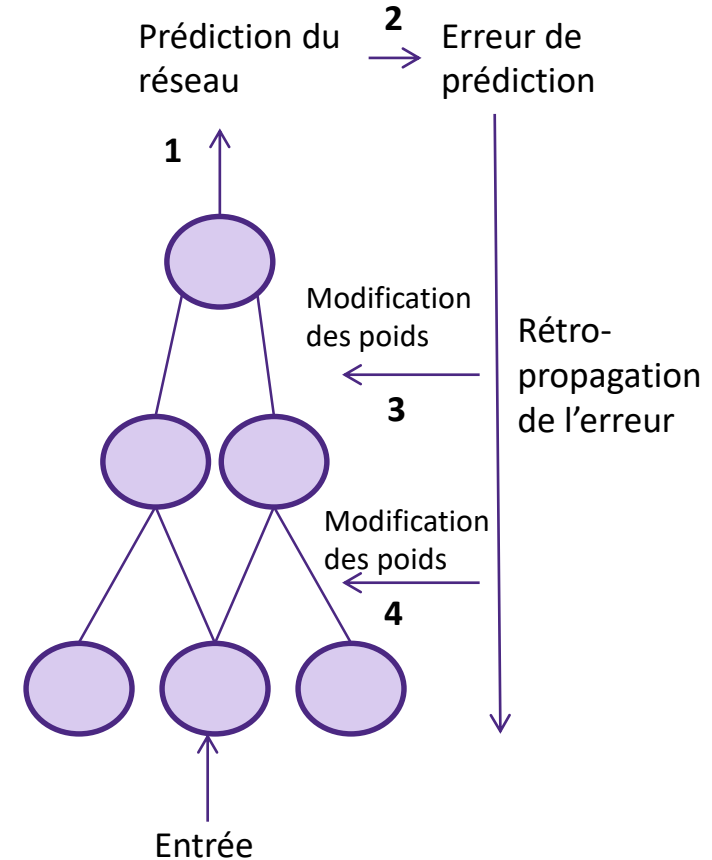


Softmax :  $p(y = j | \mathbf{x}) = \frac{e^{(\mathbf{w}_j^T \mathbf{x} + b_j)}}{\sum_{k \in K} e^{(\mathbf{w}_k^T \mathbf{x} + b_k)}}$

**Objectif : minimiser l'écart entre la réponse du réseau et la réponse voulue en faisant varier les poids du réseau.**

**Comment savoir quel poids bouger ?**

- **Algorithme de rétro-propagation du gradient de l'erreur**
  - Permet de déterminer les nouveaux poids par calcul de dérivée
  - Intuitivement : la dérivée permet de voir quel effet aurait la variation de chaque poids sur la réponse finale.
- **Différents algorithmes de descente du gradient**
  - Appelés « Optimizer »
  - Différences de performance pour éviter les minima locaux.

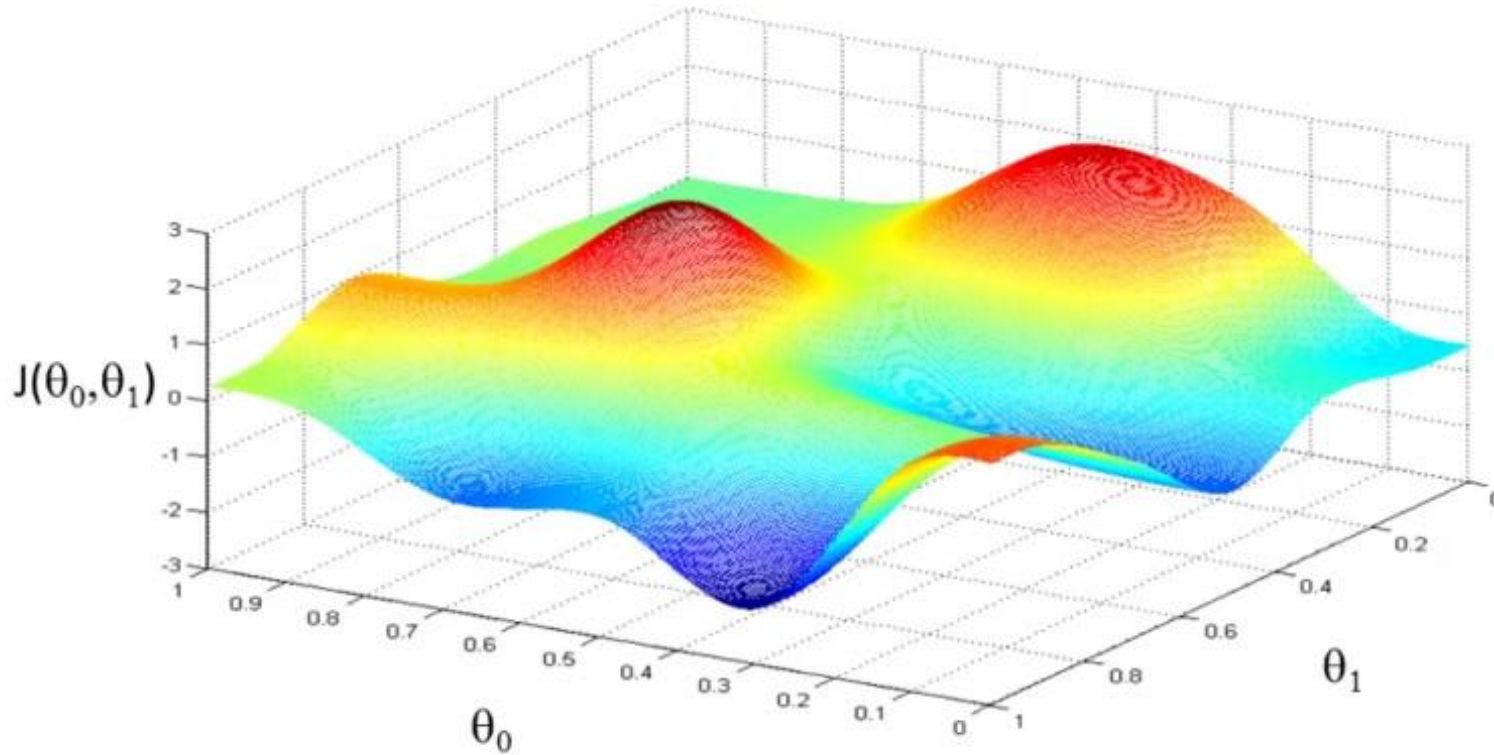


## tf playground

$$J(\theta) = 1/N \sum_{n=1}^N (\hat{y}(x_n, \theta) - y(x_n))^2$$

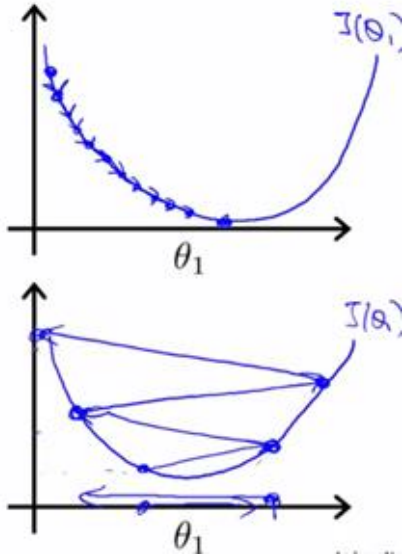
**=> Objectif : Minimiser la Cost function**

# COST FUNCTION



$$J(\theta) = 1/N \sum_{n=1}^N (\hat{y}(x_n, \theta) - y(x_n))^2$$

Dérivée partielle : 
$$\frac{\partial J(\theta)}{\partial \theta_1} = 1/N \sum_{n=1}^N (\hat{y}(x_n, \theta) - y(x_n)) x_n$$



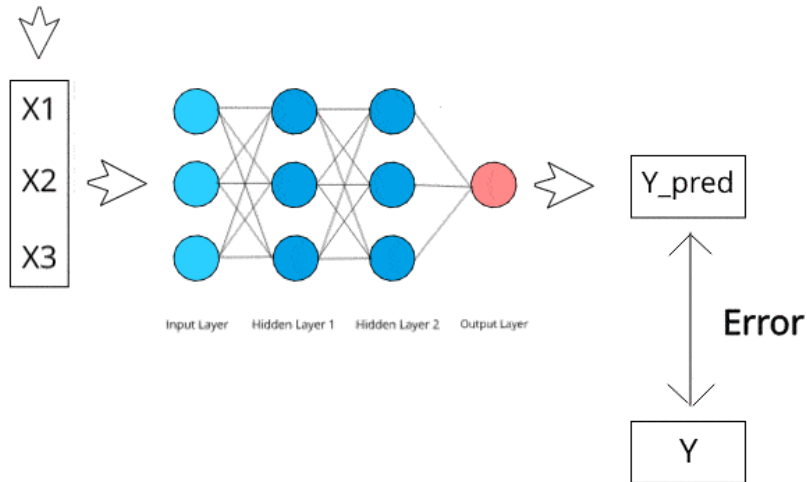
Pour  $i$  allant de 1 à  
nombre\_choisi :

$$\theta_1 = \theta_1 - \alpha \frac{\partial J(\theta)}{\partial \theta_1}$$

Avec  $\alpha > 0$ , le pas d'avancement



Feed new data



## Idée :

- Grâce au théorème de la dérivée en chaîne, on peut calculer plus facilement les dérivées de fonctions composées
- Un réseau de neurones est constituées de fonctions composées
- On calcule l'erreur de chaque couche de la dernière à la première
- On utilise l'erreur pour calculer les dérivées partielles
- On met à jour les poids du réseau (Learning rate)

Initialiser les poids du réseau

**Initialiser**  $\Delta w_{kj} \leftarrow 0, \Delta w_{ji} \leftarrow 0$

Répéter jusqu'à terminaison

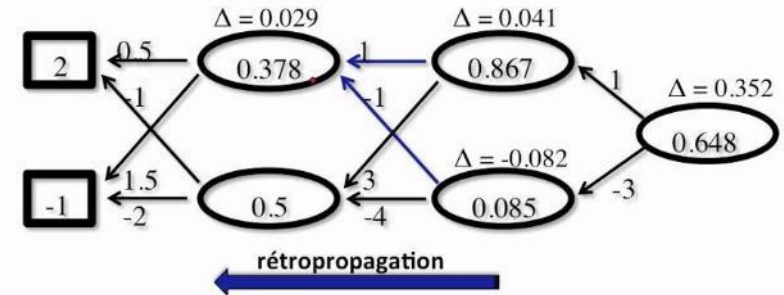
Pour chaque exemple d'apprentissage faire

- 1 Appliquer le réseau et calculer les sorties
- 2 Calculer et cumuler les deltas
  - Pour chaque unité de sortie  $k$   
Calculer  $\delta_k$ ;  $\Delta w_{kj} \leftarrow \Delta w_{kj} - \eta \delta_k z_j$
  - Pour chaque unité cachée  $j$   
Calculer  $\delta_j$ ;  $\Delta w_{ji} \leftarrow \Delta w_{ji} - \eta \delta_j x_i$

**Ajuster les poids**

- $w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$
- $w_{kj} \leftarrow w_{kj} + \Delta w_{kj}$

- Exemple:  $x = [2, -1]$ ,  $y = 1$



$$\Delta = 0.378 * (1 - 0.378) * (1 * 0.041 + -1 * -0.082) = 0.029$$

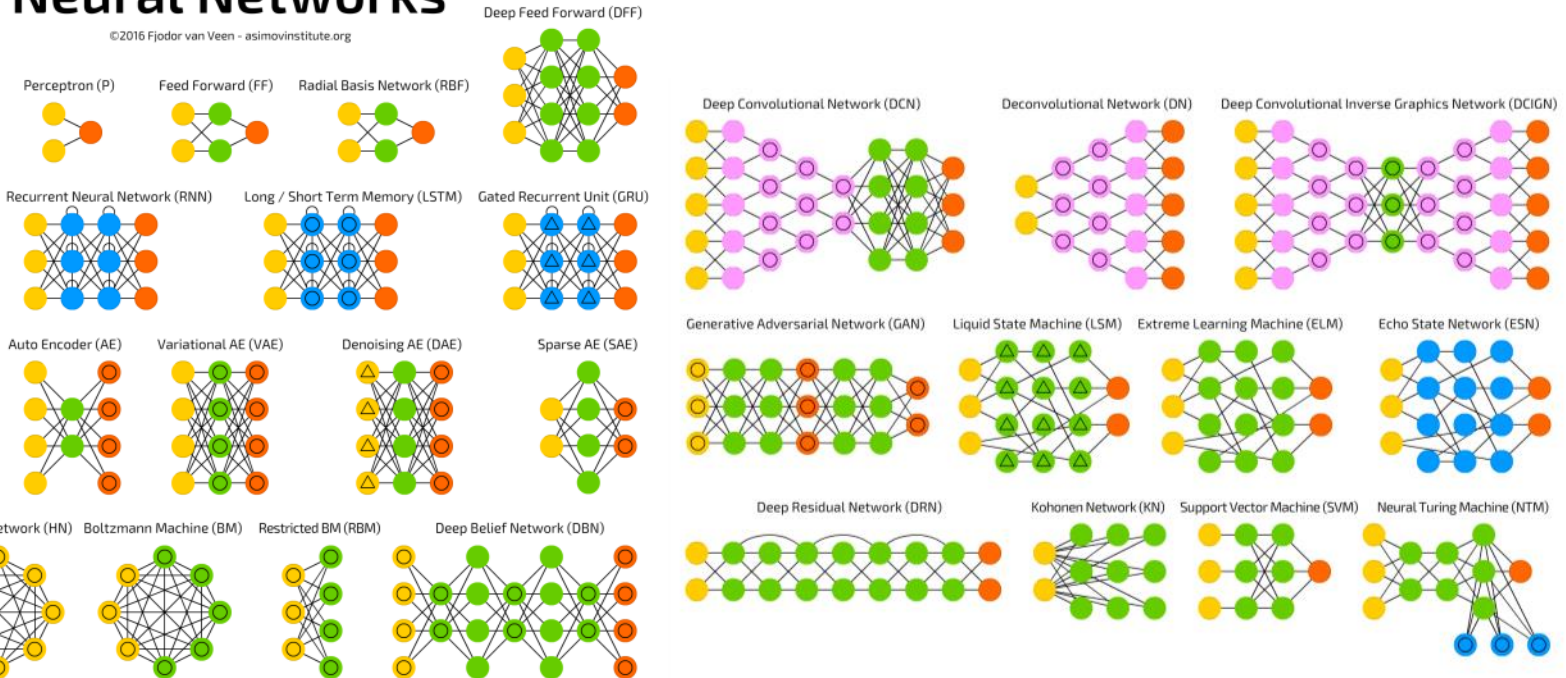
# LES DIFFÉRENTES ARCHITECTURES DE RÉSEAUX DE NEURONES

A mostly complete chart of

## Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool



Source : <http://www.asimovinstitute.org/neural-network-zoo/>

# LES DIFFÉRENTES ARCHITECTURES DE RÉSEAUX DE NEURONES

**La liste précédente est très loin d'être exhaustive !**

**Quelques grandes familles que nous explorerons plus en détail demain :**

- **Les réseaux de neurones convolutifs (CNN)**
- **Les réseaux de neurones récurrents (RNN)**
- **Les auto-encodeurs**
- **Les GANs**



# **ENTRAINER ET ÉVALUER SON RÉSEAU**

**Les bonnes pratiques d'un projet de Machine Learning :**

- **Isoler le jeu de test des création du jeu de donnée, si possible sur une machine séparée.**
- **Ne pas prendre de décision sur le jeu de test, uniquement pour confirmation !**
- **Les jeux de validation et de tests doivent coller le plus près aux exemples réels**
  - **Proportions respectées avec la réalité**
  - **Exemples les plus proches des cas réels**
- **Equilibrer les exemples positifs/négatifs dans le jeu d'entraînement mène généralement à de meilleurs performances.**
  - **Undersampling/oversampling/smote**
  - **La meilleure méthode dépend souvent du cas d'usage**
- **Sauvegarder les données après pré-processing si le pré-processing est couteux en calcul**

Quelques particularités liés aux réseaux de neurones :

- Si peu de données, du feature engineering donnera de meilleurs résultats qu'un algorithme de type deep learning.
- Beaucoup de paramètres (architecture/apprentissage/...)
  - On ne peut généralement pas faire de méthode automatique sur tous ces paramètres
  - On se limite souvent au taux d'apprentissage (voir pas du tout).
  - Ajouter un « decay » (décroissance) au taux d'apprentissage améliore presque toujours les performances
- L'apprentissage peut être très long
  - Certains réseaux mettent des semaines à apprendre sur plusieurs GPU
  - Souvent trop long pour faire de la validation croisée.
- Il est difficile de « comprendre » ce qui se passe à l'intérieur du modèle
  - Intérêt de la visualisation
  - Intérêt d'afficher de nombreuses informations sur l'état du réseau en temps réel au cours de l'entraînement
- Mettre en concurrence différentes architectures.

**Utiliser les métriques classiques de Machine learning :**

- **Afficher la performance sur le jeu de validation.**
- **L'accuracy (pourcentage de bonnes réponse) ne suffit pas !**
- **La matrice de confusion, la précision, le rappel et l'aire sous la courbe roc sont de bons indicateurs.**
- **La courbe ROC permet de définir le seuil de décision comportant le meilleur compromis par rapport au besoin.**

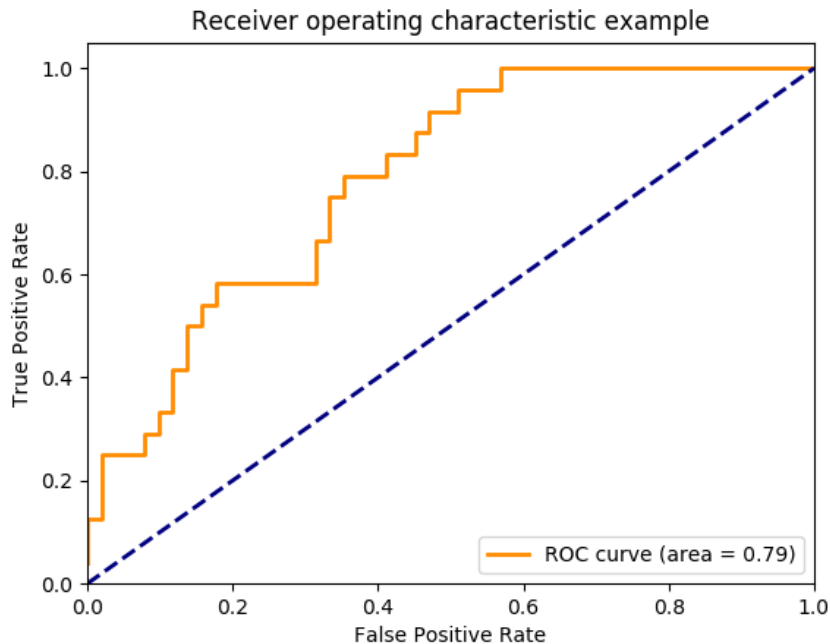


# LES MÉTRIQUES

## LA MATRICE DE CONFUSION

		TP	FN
Bonne réponse	oui	Nombre de Type A bien identifiés	Nombre de Type A mal identifiés
	non	Nombre de Type B mal identifiés	Nombre de Type B bien identifiés
		oui	non
		Réponse du réseau	

- Permet de trouver où se situent les erreurs
- Adaptable aux problèmes avec de multiples classes
- Permet de calculer la précision  $TP/(TP+FP)$  et le rappel  $TP/(TP+FN)$



## Exemple

- L'aire sous la courbe est une bonne mesure de la performance du réseau
- Adaptable aux problèmes avec de multiples classes
- Chaque point de la courbe peut être un compromis atteint par le réseau
- Fixer le seuil de décision du réseau permet de choisir quel point de la courbe prendre
- Seuil souvent ajouté par l'ajout d'une couche/fonction d'activation spécifique au réseau

# LES PROBLÈMES « CLASSIQUES » D'UN RÉSEAU DE NEURONE

## 1. Mon réseau ne converge pas !

**Problème : l'accuracy reste aux alentours du niveau de chance sur les jeux d'entraînement et de validation**

**Causes potentielles :**

- Les données et les labels sont mal alignés et ne correspondent pas.
- Les données sont mal fournies au réseau
- L'apprentissage ne fonctionne pas
  - Les couches de sorties ne sont pas touchées par l'optimizer
  - Les poids des couches ont été initialisés à 0
  - Le taux d'apprentissage est trop faible/le temps d'apprentissage trop court
  - Changer d'optimizer peut aider

# LES PROBLÈMES « CLASSIQUES » D'UN RÉSEAU DE NEURONE

## 1.(bis) Mon réseau ne converge pas !

**Problème : La sortie du réseau explose (Nan)**

**Causes potentielles :**

- Les poids ou les biais initiaux sont trop élevés
- Le taux d'apprentissage est trop fort
  - Le baisser ou augmenter son decay
- Changer de fonction d'activation si le réseau le permet

## 2. Mon réseau a de bonnes performances sur le jeu d'entraînement mais pas sur le jeu de validation

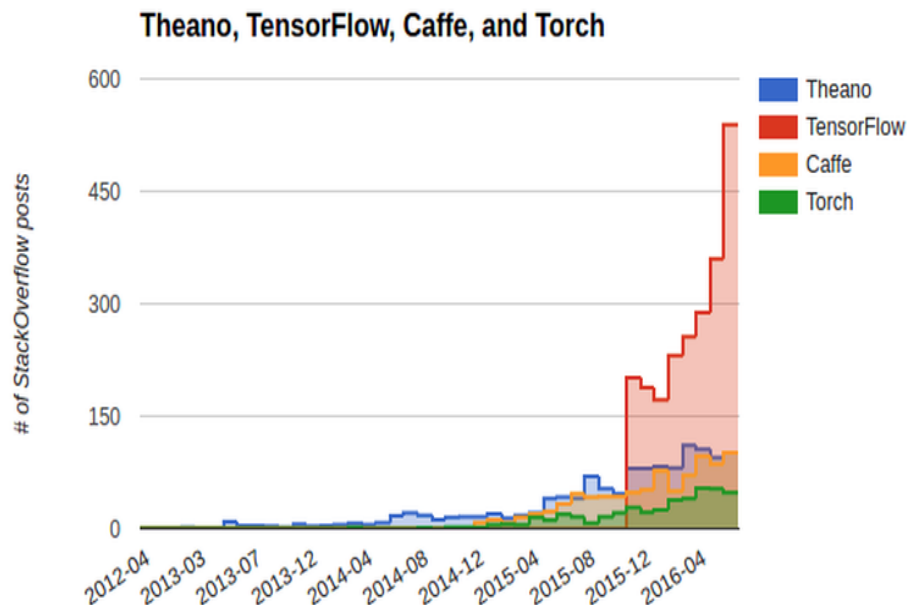
### Causes potentielles :

- Le ratio d'exemples de chaque type est différent dans les jeux d'entraînement et validation
  - Modifier les proportions du jeu de d'entraînement (ajout de donnée ou sampling)
- Le réseau fait du sur-apprentissage (overfitting) du jeu d'entraînement
  - Ajouter de la variation aléatoire dans les données à chaque présentation au réseau (« data augmentation »).
  - Normaliser les batchs
  - Arrêter l'entraînement plus tôt quand les performances divergent (« early stopping »)
  - Régularisation L1 et L2
  - Ajouter du « Dropout »

# LES DIFFÉRENTES TECHNOLOGIES DE DEEP LEARNING

<b>Caffe</b>	<b>Caffe 2</b>	<b>MXNet</b>	<b>Torch</b>	<b>DeepLearning4j</b>
<ul style="list-style-type: none"><li>✓ Plus de communauté</li><li>✓ Migration Caffe 2</li></ul>	<ul style="list-style-type: none"><li>✓ Créer pour la prod</li><li>✓ Utilisé Facebook, MS, Amazon.</li><li>✓ Peu de documentation, d'exemple, de communauté</li></ul>	<ul style="list-style-type: none"><li>✓ Framework de Deep Learning d'Amazon !</li><li>✓ Structure de données propriétaire, pas compatible avec tout numpy</li></ul>	<ul style="list-style-type: none"><li>✓ Ecrit en LUA/C</li><li>✓ Import des modèles « caffe »</li><li>✓ Google migre de Torch vers Tensorflow.</li><li>✓ Twitter est le plus gros utilisateur.</li></ul>	<ul style="list-style-type: none"><li>✓ Java</li><li>✓ Gère les formats des différentes techno « caffe, torch, theano, ... »</li></ul>
<b>Tensorflow</b>	<b>Watson</b>	<b>Azure Cognitive Services</b>	<b>Google Cloud ML Engine</b>	<b>Microsoft Cognitive Toolkit</b>
<ul style="list-style-type: none"><li>✓ La solution technique la plus plebiscitée.</li><li>✓ Soutenue par Google.</li><li>✓ Utilisation en Prod</li><li>✓ API Keras compatible avec Theano</li></ul>	<ul style="list-style-type: none"><li>✓ Impossible de choisir le modèle à utiliser.</li><li>✓ Uniquement solution Cloud</li><li>✓ Coûte très cher</li></ul>	<ul style="list-style-type: none"><li>✓ Impossible de choisir le modèle à utiliser.</li><li>✓ Pas possible d'entraîner un modèle</li></ul>	<ul style="list-style-type: none"><li>✓ 100 % compatible avec Tensorflow</li><li>✓ Coût de 1 server GPU 1,62\$/heure</li></ul>	<ul style="list-style-type: none"><li>✓ Prometteur</li><li>✓ Fonctionne uniquement avec des GPU</li><li>✓ Solution encore jeune</li></ul>

## Nombre de posts sur StackOverflow



## Nombre de stars github

new stars from 2017-04-20 to 2017-07-06		
#1:	7929	tensorflow/tensorflow
#2:	2465	fchollet/keras
#3:	1894	caffe2/caffe2
#4:	1526	BVLC/caffe
#5:	1250	pytorch/pytorch
#6:	1233	Microsoft/CNTK
#7:	979	dmlc/mxnet
#8:	709	deepmind/sonnet
#9:	690	tflearn/tflearn
#10:	485	deeplearning4j/deeplearning4j
#11:	458	Theano/Theano
#12:	452	davisking/dlib
#13:	341	torch/torch7
#14:	303	baidu/paddle
#15:	243	pfnet/chainer

**Les réseaux de neurones sont constitués d'un petit nombre de « briques » et d'architectures.**

**La majorité des réseaux complexes sont des combinaisons de ces éléments de base.**

**Ce sont souvent les mêmes problèmes que l'on rencontre lors de la création et l'entraînement de réseaux.**

**Tensorflow est devenu la référence en bibliothèque de réseaux de neurones.**