

Intelligence Artificielle

MESR - Plan de Formation National

Jamal Atif

Professeur des Universités
PSL, Université Paris-Dauphine, LAMSADE, CNRS



30 mai 2017

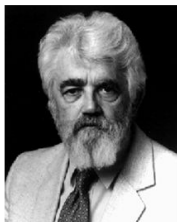
Définir l'intelligence artificielle

Une entreprise périlleuse!

Acte de naissance

Conférence de Dartmouth en 1956...4 ans après le décès tragique de A. Turing

Dartmouth Conference: The Founding Fathers of AI



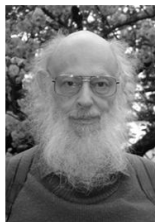
John McCarthy



Marvin Minsky



Claude Shannon



Ray Solomonoff

Alan Newell



Herbert Simon



Arthur Samuel



And three others...

Oliver Selfridge
(Pandemonium theory)

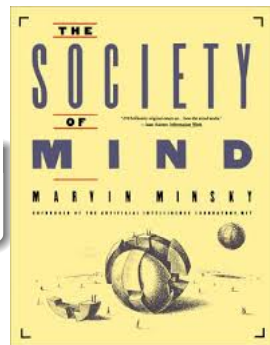
Nathaniel Rochester
(IBM, designed 701)

Trenchard More
(Natural Deduction)

Une définition ... discutable !

Marvin Minsky

Science qui consiste à faire faire aux machines ce que l'homme ferait moyennant une certaine intelligence.



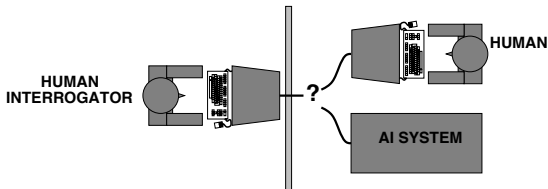
Ecueils

- ▶ Récursivité de la définition : dépend de la définition de l'intelligence humaine
- ▶ Portée de "faire faire" ?

L'héritage de Turing

Turing (1950) "Computing machinery and intelligence":

- ▶ "Can machines think?" → "Can machines behave intelligently?"
- ▶ Test opérationnel : le jeu de l'imitation



- ▶ Turing a prédit qu'en 2000, une machine aurait une chance de 30% de tromper un humain pendant 5 minutes.
- ▶ A anticipé tous les arguments contre l'IA dans les 50 ans suivants.
- ▶ A suggéré les composantes principales de l'IA : connaissances, raisonnement, TAL, apprentissage

Problèmes : le test n'est ni *reproductible*, ni *constructif*, ni apte à une *analyse mathématique*

Autres définitions... tout aussi discutables

- ▶ “Tout problème pour lequel il n'existe pas d'algorithme connu, ou de coût raisonnable, relève de l'I.A.”
- ▶ “L'I.A. doit permettre de proposer des solutions logicielles permettant aux programmes de raisonner logiquement”
- ▶ “L'IA est le domaine de l'informatique qui étudie comment faire faire à l'ordinateur des tâches pour lesquelles l'homme est aujourd'hui encore le meilleur”
- ▶ “Le but de l'Intelligence Artificielle est de construire un objet pouvant réussir avec fiabilité le Test de Turing”
- ▶ “L'IA est ce qui est publié dans les conférences et journaux de l'IA”

Autres définitions... tout aussi discutables

- ▶ “Tout problème pour lequel il n'existe pas d'algorithme connu, ou de coût raisonnable, relève de l'I.A.”
- ▶ “L'I.A. doit permettre de proposer des solutions logicielles permettant aux programmes de raisonner logiquement”
- ▶ “L'IA est le domaine de l'informatique qui étudie comment faire faire à l'ordinateur des tâches pour lesquelles l'homme est aujourd'hui encore le meilleur”
- ▶ “Le but de l'Intelligence Artificielle est de construire un objet pouvant réussir avec fiabilité le Test de Turing”
- ▶ “L'IA est ce qui est publié dans les conférences et journaux de l'IA”

Autres définitions... tout aussi discutables

- ▶ “Tout problème pour lequel il n'existe pas d'algorithme connu, ou de coût raisonnable, relève de l'I.A.”
- ▶ “L'I.A. doit permettre de proposer des solutions logicielles permettant aux programmes de raisonner logiquement”
- ▶ “L'IA est le domaine de l'informatique qui étudie comment faire faire à l'ordinateur des tâches pour lesquelles l'homme est aujourd'hui encore le meilleur”
- ▶ “Le but de l'Intelligence Artificielle est de construire un objet pouvant réussir avec fiabilité le Test de Turing”
- ▶ “L'IA est ce qui est publié dans les conférences et journaux de l'IA”

Autres définitions... tout aussi discutables

- ▶ “Tout problème pour lequel il n'existe pas d'algorithme connu, ou de coût raisonnable, relève de l'I.A.”
- ▶ “L'I.A. doit permettre de proposer des solutions logicielles permettant aux programmes de raisonner logiquement”
- ▶ “L'IA est le domaine de l'informatique qui étudie comment faire faire à l'ordinateur des tâches pour lesquelles l'homme est aujourd'hui encore le meilleur”
- ▶ “Le but de l'Intelligence Artificielle est de construire un objet pouvant réussir avec fiabilité le Test de Turing”
- ▶ “L'IA est ce qui est publié dans les conférences et journaux de l'IA”

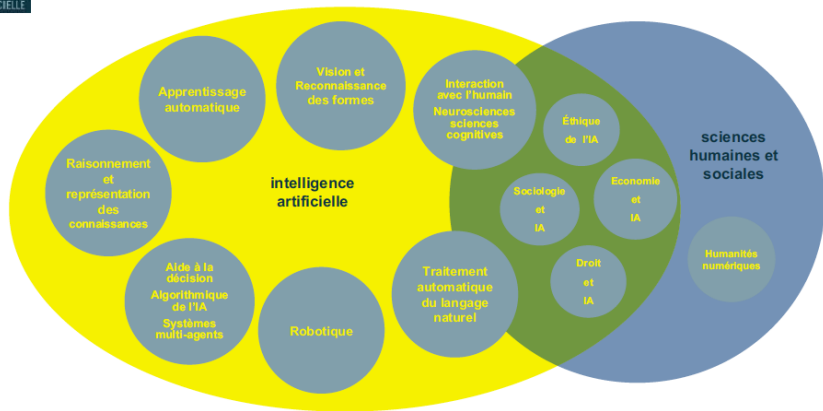
Autres définitions... tout aussi discutables

- ▶ “Tout problème pour lequel il n'existe pas d'algorithme connu, ou de coût raisonnable, relève de l'I.A.”
- ▶ “L'I.A. doit permettre de proposer des solutions logicielles permettant aux programmes de raisonner logiquement”
- ▶ “L'IA est le domaine de l'informatique qui étudie comment faire faire à l'ordinateur des tâches pour lesquelles l'homme est aujourd'hui encore le meilleur”
- ▶ “Le but de l'Intelligence Artificielle est de construire un objet pouvant réussir avec fiabilité le Test de Turing”
- ▶ “L'IA est ce qui est publié dans les conférences et journaux de l'IA”

Un objet protéiforme



Plusieurs domaines de recherche & domaines connexes SHS



Une histoire mouvementée

- 1943 McCulloch & Pitts: Boolean circuit model of brain
- 1950 Turing's "Computing Machinery and Intelligence"
- 1952–69 Look, Ma, no hands!
- 1950s Early AI programs, including Samuel's checkers program, Newell & Simon's Logic Theorist, Gelernter's Geometry Engine
- 1956 Dartmouth meeting: "Artificial Intelligence" adopted
- 1965 Robinson's complete algorithm for logical reasoning
- 1966–74 AI discovers computational complexity
- 1969 Minsky and Papert's "Group Invariance Theorem"
Neural network research almost disappears
- 1972 Prolog by Alain Colmerauer and Philippe Roussel
- 1969–79 Early development of knowledge-based systems
- 1980–88 Expert systems industry booms
- 1988–93 Expert systems industry busts: "AI Winter"
- 1985–95 Neural networks return to popularity (Geff Hinton is there!)
- 1988– Resurgence of probability; general increase in technical depth (machine learning)
"Nouvelle AI": ALife, GAs, soft computing
- 1995– Agents, agents, everywhere ...
- 2006– Human-level AI and neural networks (deep learning)
back on the agenda (Geff Hinton is there!)

Que sait on faire aujourd'hui ?

Jouer au Go - AlphaGo, le désormais retraité



Figure: Mars 2016. AlphaGo : 4 Lee Sedol (9 d) : 1

Que sait on faire aujourd'hui ?

Jouer au Go - AlphaGo, le désormais retraité



Figure: 25 mai 2017. AlphaGo : 3 Kee Jie (9 d) : 0

Que sait on faire aujourd'hui ?

Jouer au Poker ?



Figure: Libratus: heads-up, no-limit Texas Hold'em. Counterfactual Regret

Minimization (CFR)

Que sait on faire aujourd'hui ?

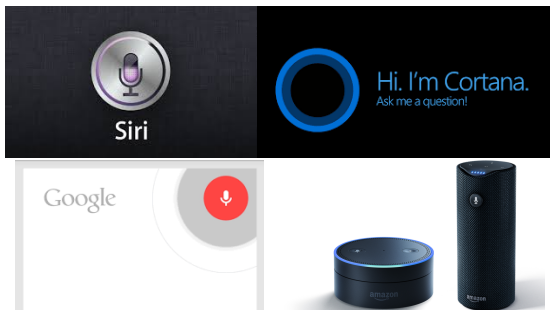
Conduire de façon autonome ?



Figure: Waymo Google Car

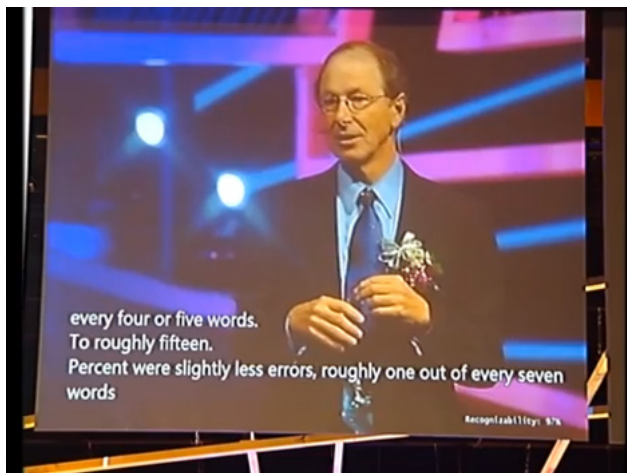
Que sait on faire aujourd'hui ?

Converser ?



Que sait on faire aujourd'hui ?

Traduire en temps réel ?



Que sait on faire aujourd'hui ?

Dépasser les capacités humaines en reconnaissance de formes ?

IMAGENET 14,197,122 images, 21841 synsets indexed [SEARCH](#) [Home](#) [Explore](#)
[About](#) [Download](#)

Not logged in. [Login](#) | [Signup](#)

Fish

Any of various mostly cold-blooded aquatic vertebrates usually having scales and breathing through gills; "the shark is a large fish"; "in the living room there was a tank of colorful fish"

1307 pictures 91.33% Popularity Percentile Wordnet IDs

Numbers in brackets: (the number of synsets in the subtree).

- ImageNet 2011 Fall Release (32321)
 - plant, flora, plant life (4486)
 - geological formation, formation
 - natural object (1112)
 - sport, athletics (176)
 - artifact, artefact (10504)
 - fungus (308)
 - person, individual, someone, son
 - animal, animate being, beast, bn
 - invertebrate (766)
 - homeotherm, homiotherm, t
 - work animal (4)
 - darter (0)
 - survivor (0)
 - range animal (0)
 - creepy-crawly (0)
 - domestic animal, domestic
 - molter, moulter (0)
 - varmint, varment (0)
 - mutant (0)

Treemap Visualization **Images of the Synset** **Downloads**

ImageNet 2011 Fall Release > Aquatic vertebrate > Fish

Bony **Cartilaginous** **Food** **Climbing** **Spawner** **Bot fec** **Rough**

Que sait on faire aujourd'hui ?

- ▶ Composer de la musique ?
- ▶ Jouer à la bourse comme des traders professionnels ?
- ▶ Détecter des tumeurs ?
- ▶ Traduire l'activité cérébrale en signal moteur ?
- ▶ ...

Raisons d'un engouement !

- ▶ Avancées dans l'ensemble des champs disciplinaires : Rep. des connaissances (knowledge graph), raisonnement (complexité), SAT, théorie de l'apprentissage, etc.
- ▶ Disponibilité de grandes masses de données
- ▶ Disponibilité des moyens de calcul
- ▶ Percée de l'apprentissage automatique et des modèles neuronaux profond
- ▶ Représentation des connaissances et raisonnement sur des données à grande échelle

Ouvrons la boîte !

Apprentissage automatique

Apprentissage automatique

Définition imprécise

Doter les machines de capacités

- ▶ d'extraction automatique de "connaissances" à partir de masses de données
- ▶ et d'auto-amélioration à partir d'expérience

Définition moins imprécise (Tom Mitchell)

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

Jeu de dame : T jouer au dame, P % de parties gagnées, E mouvements connus ou pratique du jeu

Apprentissage automatique

Définition imprécise

Doter les machines de capacités

- ▶ d'extraction automatique de "connaissances" à partir de masses de données
- ▶ et d'auto-amélioration à partir d'expérience

Définition moins imprécise (Tom Mitchell)

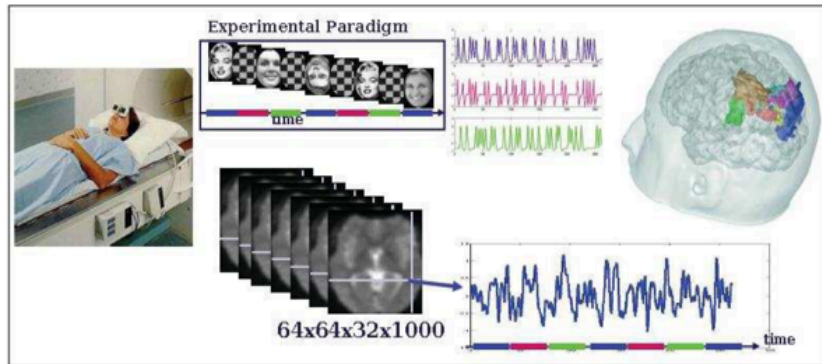
A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

Jeu de dame : T jouer au dame, P % de parties gagnées, E mouvements connus ou pratique du jeu

Grandes figures de l'apprentissage automatique

- ▶ Apprentissage supervisé
- ▶ Apprentissage non-supervisé
- ▶ Apprentissage par renforcement
- ▶ Apprentissage actif
- ▶ Transfert d'apprentissage, par analogie, de préférences, etc.

- **Apprentissage supervisé : interprétation d'IRMf**



Trouble de la reconnaissance de visage ou non

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Principe: étant donné un échantillon de données étiquetées $\mathcal{S} = \{\langle x_i, y_i \rangle\}_{1 \dots n}$, apprendre une fonction/densité de prob. de prédiction qui lie les données aux étiquettes.

$$\mathcal{X} \ni x \xrightarrow[p(\cdot, \cdot)]{h \in \mathcal{H}} y \in \mathcal{Y}$$

- ▶ $\mathcal{Y} \equiv \mathbb{R}$: problème de régression
- ▶ $\mathcal{Y} \equiv$ ensemble discret (e.g. $\{0, 1\}$): problème de classification
- ▶ \mathcal{H} peut être un espace fonctionnel ou de densités de probabilités
- ▶ Choix de la fonction de perte et du risque à minimiser (erreur en généralisation) :

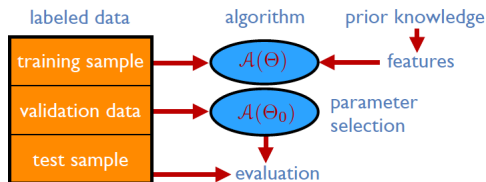
$$R(h) = \mathbb{E}_{(x,y) \sim D}[\ell(h(x), y)] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(h(x), y) p_{XY}(x, y) dx dy$$

- ▶ Minimisation du risque empirique

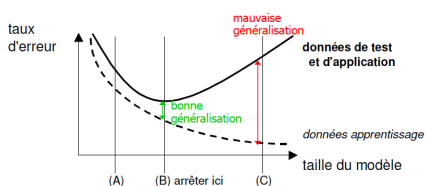
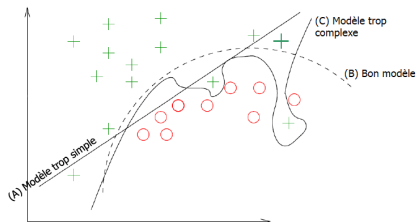
$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)$$

Apprentissage supervisé

Procédure



Apprentissage \neq mémorisation : généralisation vs spécialisation



Apprentissage supervisé

Risque structurel

Soit une séquence infinie d'ensembles d'hypothèses ordonnés par inclusion,
 $\mathcal{H}_1 \subset \mathcal{H}_2 \subset \dots \subset \mathcal{H}_m \subset \dots$

$$\hat{h} = \arg \min_{h \in \mathcal{H}_m, m \in \mathbb{N}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) + \text{penalty}(H_m, m)$$

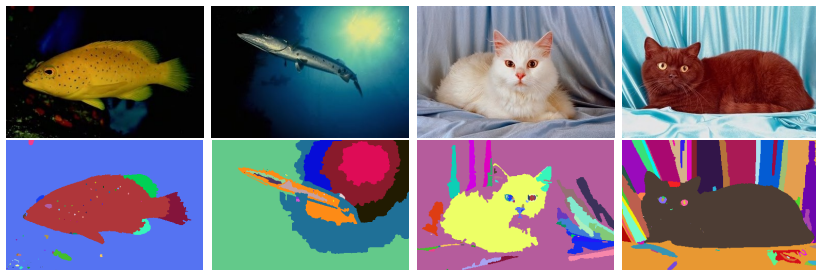
- ▶ Garanties théoriques fortes
- ▶ Complexité de calcul
- ▶ Cadre pour la régularisation :

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) + \lambda \text{Reg}(h) \quad (1)$$

$$\hat{h} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i) + \lambda \|h\|_0 \quad (2)$$

Apprentissage non-supervisé

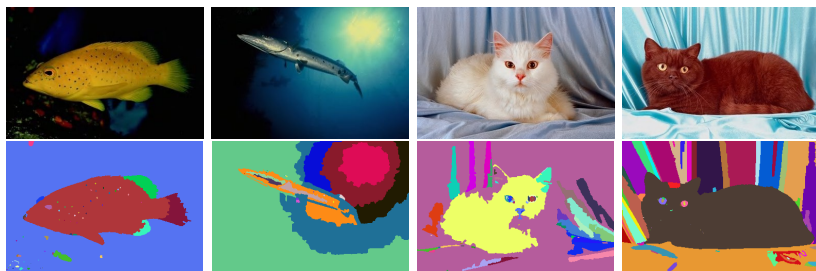
Principe: étant donné un échantillon de données non-étiquetées $\mathcal{S} = \{x_i, i = 1, \dots, n\}$, découvrir des régularités en créant des groupes homogènes.



- ▶ Cadre théorique mal maîtrisé
- ▶ Challenge pour les années à venir

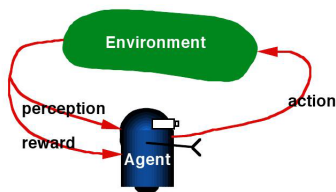
Apprentissage non-supervisé

Principe: étant donné un échantillon de données non-étiquetées $\mathcal{S} = \{x_i, i = 1, \dots, n\}$, découvrir des régularités en créant des groupes homogènes.



- ▶ Cadre théorique mal maîtrisé
- ▶ Challenge pour les années à venir

Apprentissage par renforcement



Généralités

- ▶ Un agent, situé dans le temps et l'espace
- ▶ Evoluant dans un environnement incertain (stochastique)
- ▶ But : sélectionner une action à chaque pas de temps,
- ▶ ... afin de maximiser une espérance du gain cumulé à horizon temporel fini ou infini

Qu'apprend-on

une politique = stratégie = {état \rightarrow action}

Allons plus profondément dans la boîte !

Réseaux de neurones ··· profonds

Anatomie (basique) d'une neurone

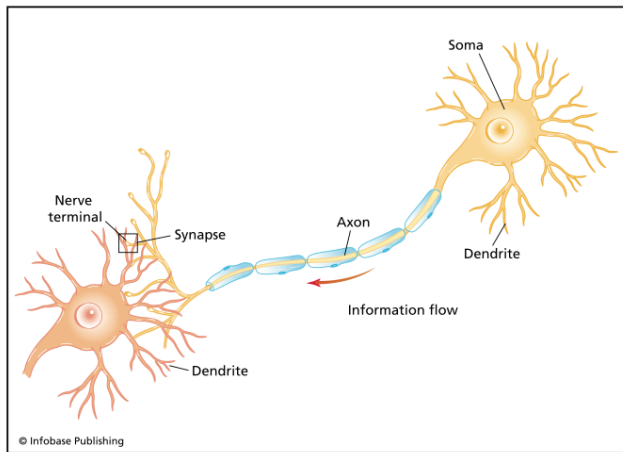
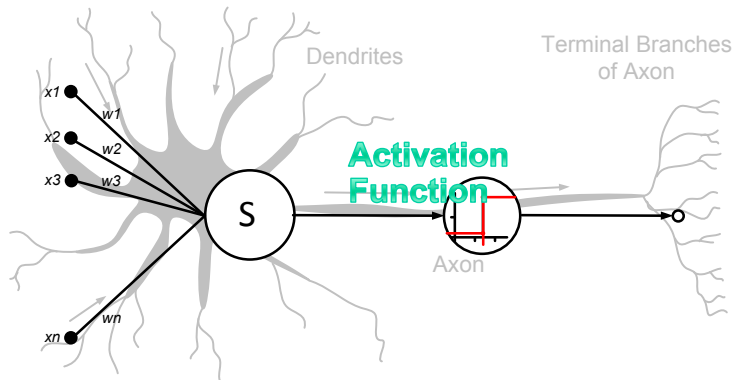


Figure: A neuron's basic anatomy consists of four parts: a **soma** (cell body), **dendrites**, an **axon**, and **nerve terminals**. Information is received by dendrites, gets collected in the cell body, and flows down the axon.

Neurone artificiel



Perceptron

Rosenblatt 1957

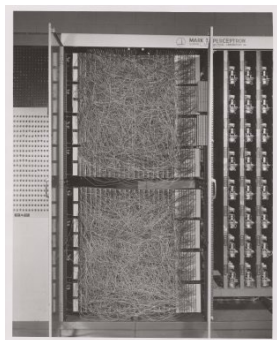
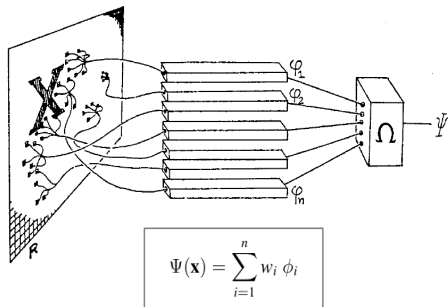


Figure: Mark I Perceptron machine

Perceptron

Apprentissage des poids w_i

Règle de Hebb

En cas de succès, ajouter à chaque connexion quelque chose de proportionnel à l'entrée et à la sortie.

Règle du perceptron : apprendre seulement en cas d'échec

Algorithme 1 : Algorithme d'apprentissage du perceptron

tant que *non convergence* **faire**

si *la forme d'entrée est correctement classée* **alors**

 ne rien faire

sinon

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \mathbf{x}_i y_i$$

fin

 Passer à la forme d'apprentissage suivante

fin

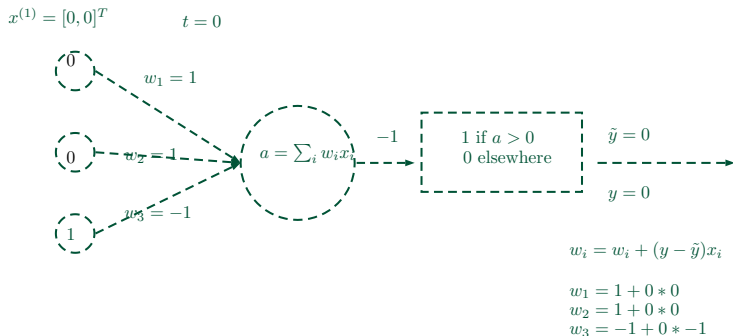
Apprentissage du Perceptron

Exemple : la fonction OR

Initialization: $w_1(0) = w_2(0) = 1, w_3(0) = -1$

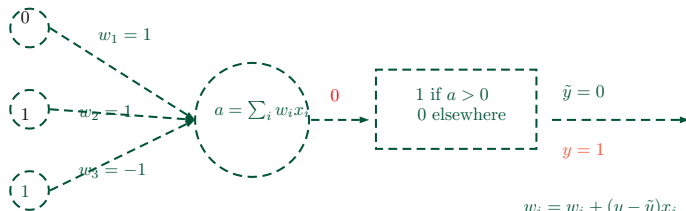
| t | $w_1(t)$ | $w_2(t)$ | $w_3(t)$ | $\mathbf{x}^{(k)}$ | $\sum w_i x_i^k$ | $\tilde{y}^{(k)}$ | $y^{(k)}$ | $\Delta w_1(t)$ | $\Delta w_2(t)$ | $\Delta w_3(t)$ |
|---|----------|----------|----------|--------------------|------------------|-------------------|-----------|-----------------|-----------------|-----------------|
| 0 | 1 | 1 | -1 | 001 | -1 | 0 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | -1 | 011 | 0 | 0 | 1 | 0 | 1 | 1 |
| 2 | 1 | 2 | 0 | 101 | 1 | 1 | 1 | 0 | 0 | 0 |
| 3 | 1 | 2 | 0 | 111 | 3 | 1 | 1 | 0 | 0 | 0 |
| 4 | 1 | 2 | 0 | 001 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 1 | 2 | 0 | 011 | 2 | 1 | 1 | 0 | 0 | 0 |

Perceptron : illustration



Perceptron : illustration

$$x^{(2)} = [0, 1]^T \quad t = 1$$



$$w_i = w_i + (y - \tilde{y})x_i$$

$$w_1 = 1 + 1 * 0 = 1$$

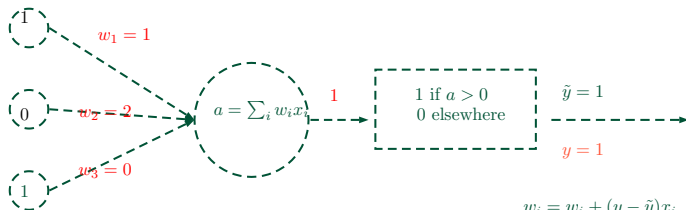
$$w_2 = 1 + 1 * 1 = 2$$

$$w_3 = -1 + 1 * 1 = 0$$

Perceptron : illustration

$$x^{(3)} = [1, 0]^T$$

$$t = 2$$



$$w_i = w_i + (y - \tilde{y})x_i$$

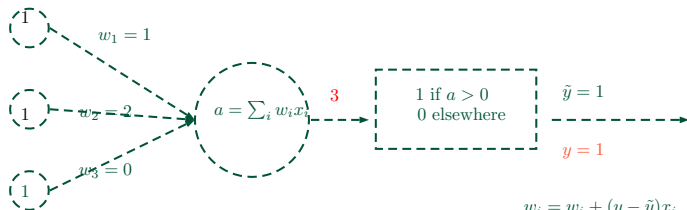
$$w_1 = 1 + 0 * 0 = 1$$

$$w_2 = 2 + 0 * 1 = 2$$

$$w_3 = 0 + 0 * 1 = 0$$

Perceptron : illustration

$$x^{(4)} = [1, 1]^T \quad t = 3$$



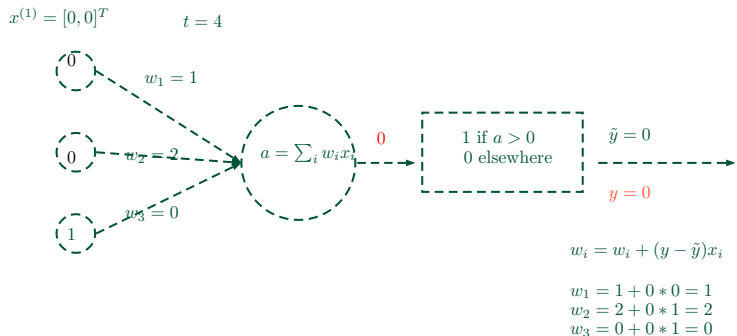
$$w_i = w_i + (y - \tilde{y})x_i$$

$$w_1 = 1 + 0 * 0 = 1$$

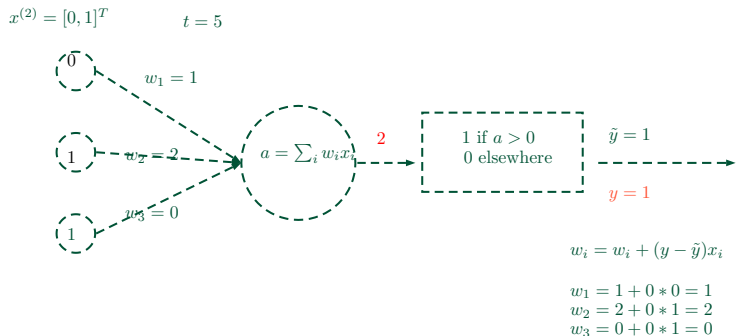
$$w_2 = 2 + 0 * 1 = 2$$

$$w_3 = 0 + 0 * -1 = 0$$

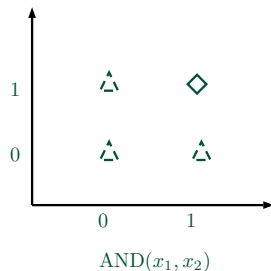
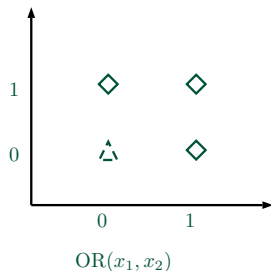
Perceptron : illustration



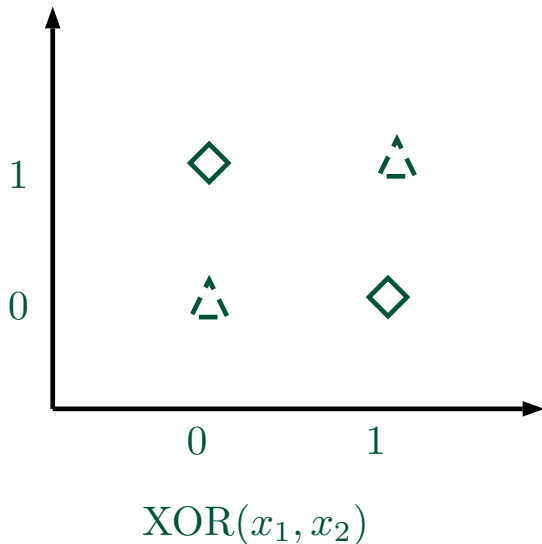
Perceptron : illustration



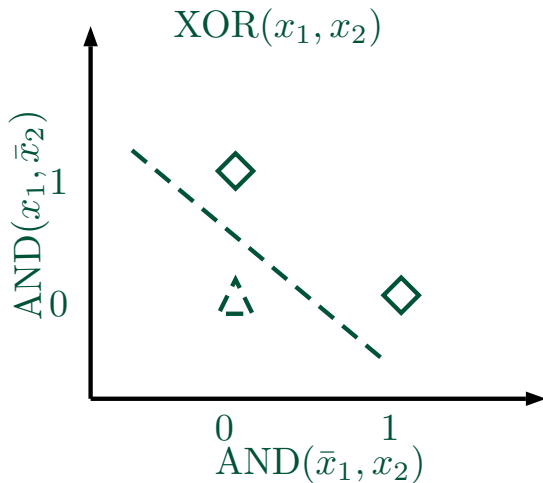
Perceptron : capacité



Perceptron : capacité

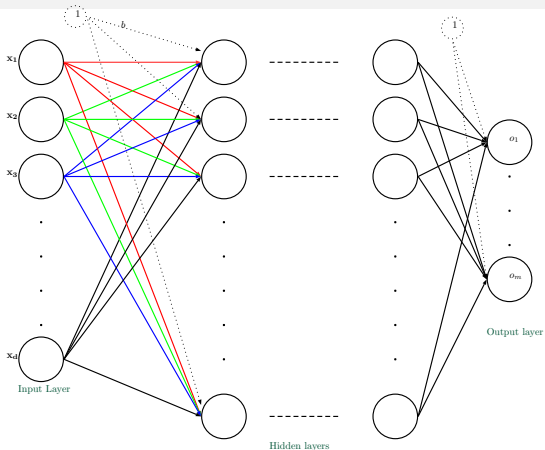


Mais !



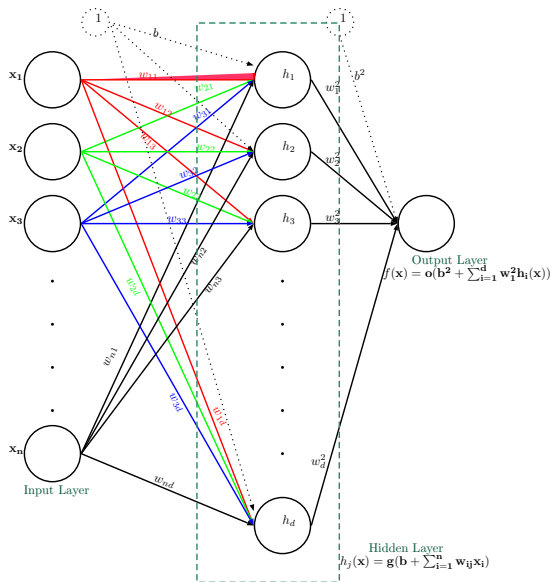
Perceptron Multi-Couches

Paul Werbos, 84. Rumelhart, Hinton et al, 86



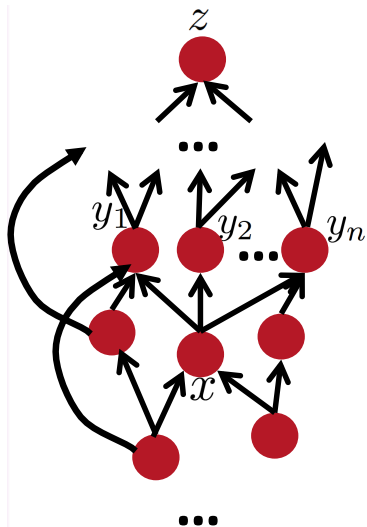
Universal approximation theorem (Cybenko, 89; Hornik 91) :
Sous certaines conditions sur les fonctions d'activation, le PMC avec une seule couche cachée composée d'un nombre fini de neurones, peut approcher avec une erreur arbitraire toute fonction dans \mathbb{R}^n

PMC avec une couche cachée



PMC: entraînement par retropropagation

Chain rule généralisée



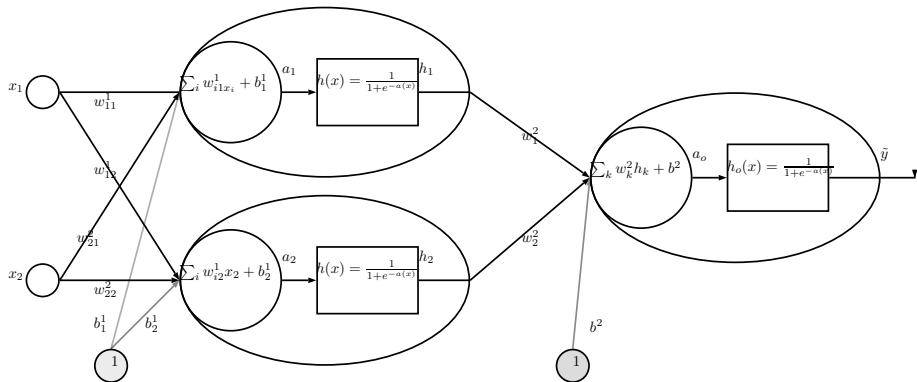
Flow graph: any directed acyclic graph
node = computation result
arc = computation dependency

$\{y_1, y_2, \dots, y_n\}$ = successors of x

$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}$$

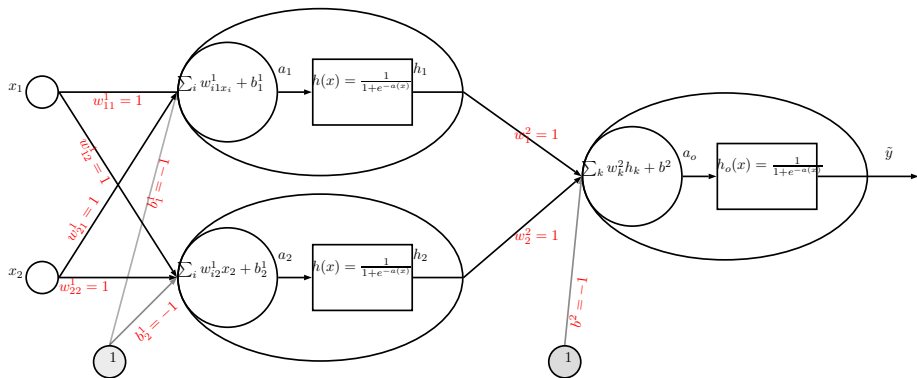
PMC: entraînement par retropropagation

XOR



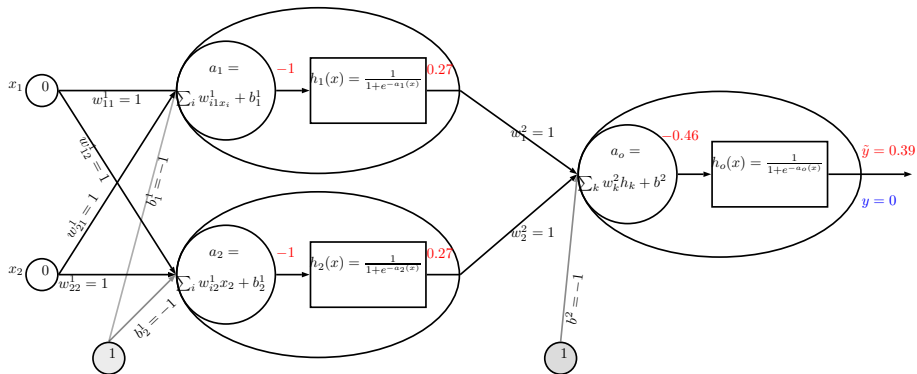
PMC: entraînement par retropropagation

XOR



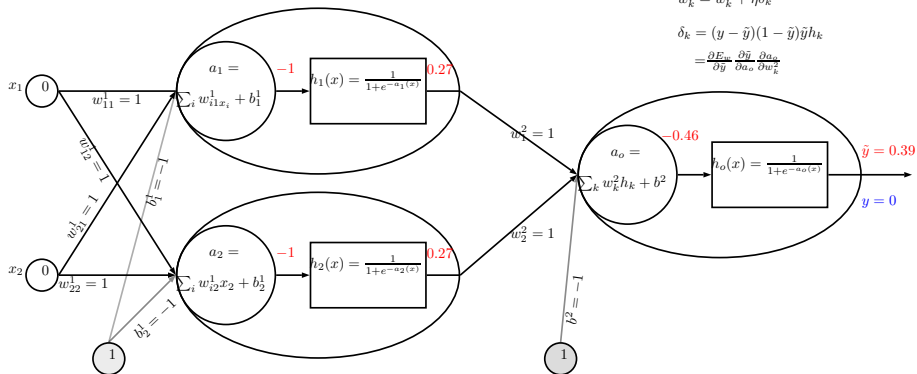
PMC: entraînement par retropropagation

XOR



PMC: entraînement par retropropagation

XOR

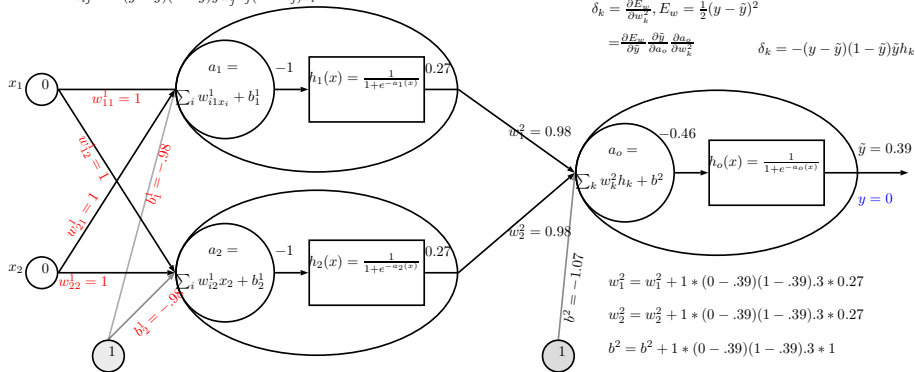


PMC: entraînement par retropropagation

XOR

$$w_{ij}^1 = w_{ij}^1 - \eta \delta_{ij} \quad \delta_{ij} = \frac{\partial E_w}{\partial w_{ij}} = \frac{\partial E_w}{\partial h_o} \frac{\partial h_o}{\partial a_o} \frac{\partial a_o}{\partial h_j} \frac{\partial h_j}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}}$$

$$\delta_{ij} = -(y - \tilde{y})(1 - \tilde{y})\tilde{y}w_j^2 h_j(1 - h_j)x_i$$



$$w_k^2 = w_k^2 - \eta \delta_k$$

$$\delta_k = \frac{\partial E_w}{\partial w_k^2}, E_w = \frac{1}{2}(y - \tilde{y})^2$$

$$= \frac{\partial E_w}{\partial \tilde{y}} \frac{\partial \tilde{y}}{\partial a_o} \frac{\partial a_o}{\partial w_k^2}$$

$$\delta_k = -(y - \tilde{y})(1 - \tilde{y})\tilde{y}h_k$$

$$w_{11}^1 = 1 + 1 * (0 - .39)(1 - .39) * .39 * 1 * .27 * (1 - .27) * 0$$

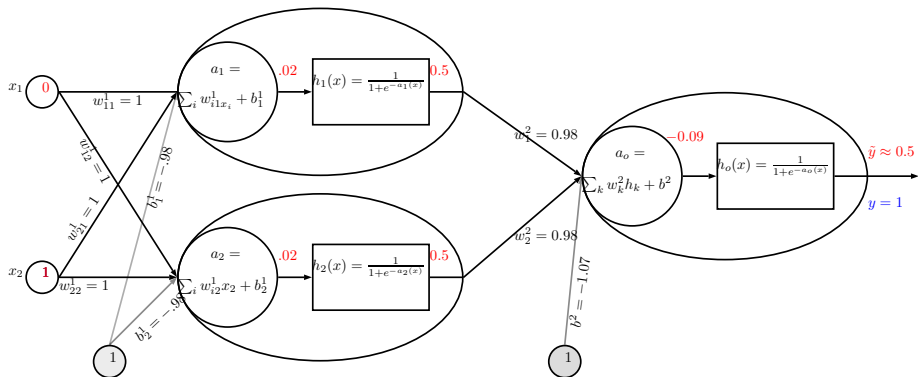
$$w_{22}^1 = 1 + 1 * (0 - .39)(1 - .39) * .39 * 1 * .27 * (1 - .27) * 0$$

⋮

$$b_1^1 = -1 + 1 * (0 - .39)(1 - .39) * .39 * 1 * .27 * (1 - .27) * 0$$

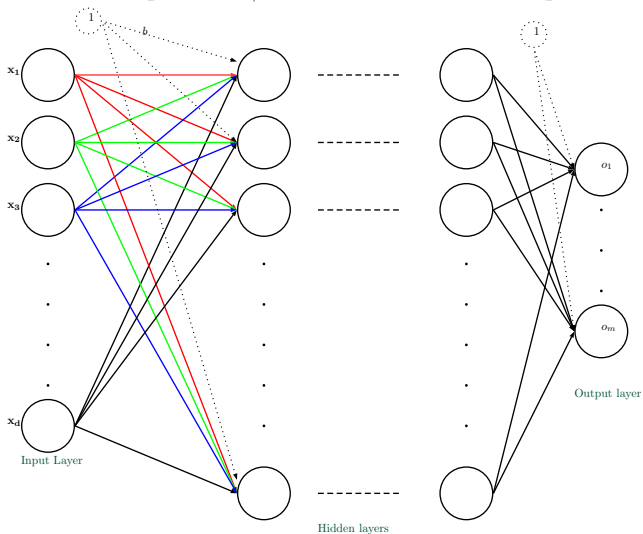
PMC: entraînement par retropropagation

XOR



PMC comme un réseau profond

PMC avec plus de 2/3 couches est un réseau profond



D'où vient la rupture ?

Avant 2006, l'entraînement des architectures profondes était sans succès !

Bengio, Hinton, LeCun

Essoufflement du gradient

- ▶ Avancées en optimisation stochastique
- ▶ Pré-entraînement non-supervisé

Sur-apprentissage

- ▶ Techniques de régularisation
- ▶ Stochastic "dropout"

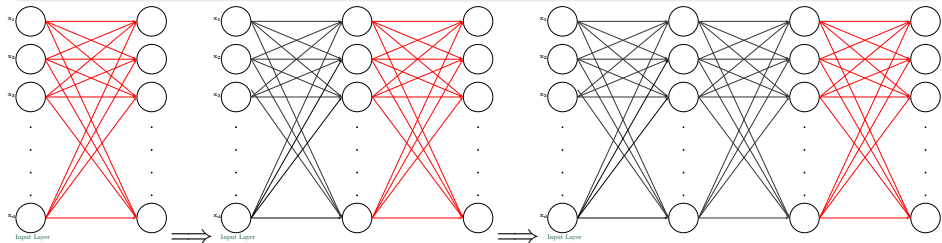
Et surtout

- ▶ Disponibilité de très grandes masses de données
- ▶ Disponibilité de moyens de calcul

Pré-entraînement non-supervisé

Idée principale

Initialiser le réseau de façon non-supervisée pas à pas

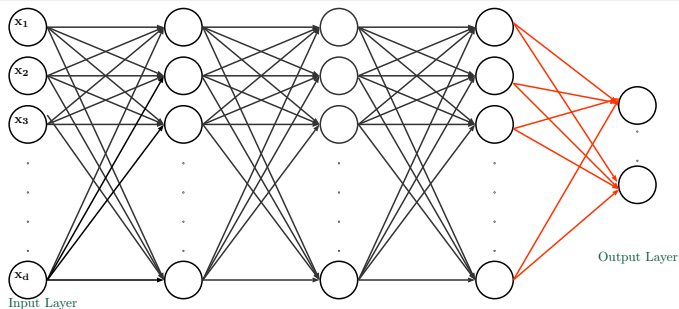


Modèles dédiés : RBM, Auto-encodeurs, et plusieurs variantes

Fine tuning

Comment ?

- ▶ Ajouter la couche de sortie
- ▶ Initialiser ses poids de façon aléatoire
- ▶ Mise à jour par rétropropagation



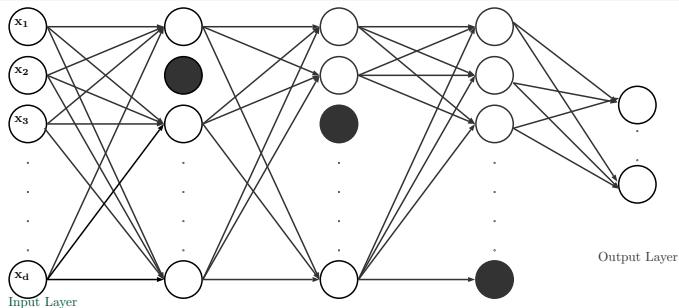
Drop out

Intuition

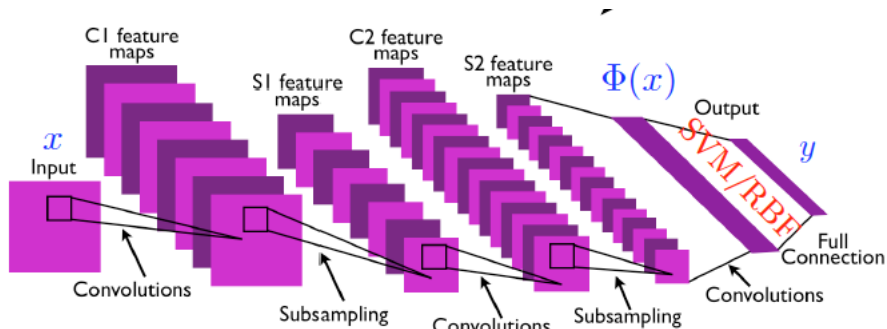
Régulariser le réseau en **annulant** aléatoirement des unités cachées.

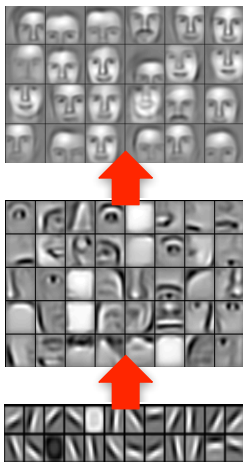
Procédure

Affecter à chaque neurone caché une valeur 0 avec une probabilité p (choix commun : .5)



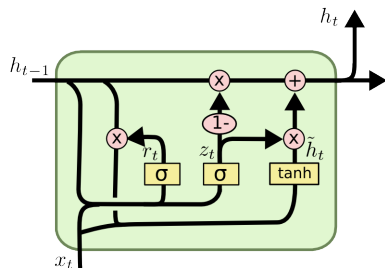
Quelques réseaux d'intérêt





Réseaux récurrents : LSTM

Etat de l'art en TAL



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

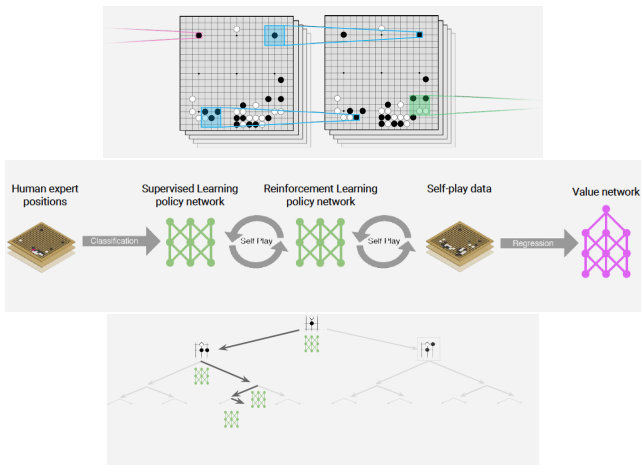
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

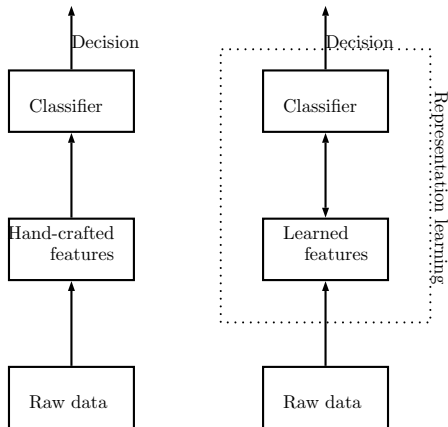
Inside AlphaGo

Combiner CNN, RL et MCTS



Qu'y a-t-il de profond en apprentissage profond

Paradigmes nouveaux vs anciens



Conclusion

- ▶ Avancées considérables en reconnaissance de formes (solution au paradoxe de Moravec ?)
- ▶ Continuum mathématique/informatique en apprentissage automatique
- ▶ Peu de compréhension des réseaux de neurones profonds : problèmes d'interprétation et donc d'acceptabilité
- ▶ Le futur est pour la combinaison des approches
- ▶ Nous sommes loins de l'intelligence artificielle générale
- ▶ ... un (petit) pas au travers du Deep Reinforcement Learning