

# Apprentissage automatique et méga données

Stéphane Canu

[asi.insa-rouen.fr/enseignants/~scanu](http://asi.insa-rouen.fr/enseignants/~scanu)

Séminaire ICube 2015, Strasbourg

June 11, 2015

# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



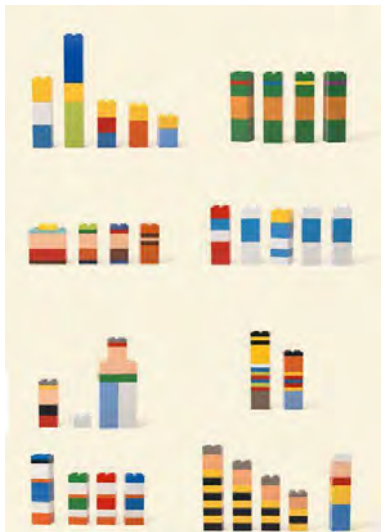
# Apprentissage : humain vs. machine

## Les apprentissages d'un enfant

- marcher : un an
- parler : deux ans
- raisonner : le reste



# Apprendre à raisonner



# Deux définitions de l'apprentissage

## Arthur Samuel (1959)

Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.



## Tom Mitchell (The Discipline of Machine Learning, 2006)

“How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?”

A computer program  $CP$  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  $E$

# Une tâche $T$ « apprenable »

## Machine Learning

A computer program CP 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**

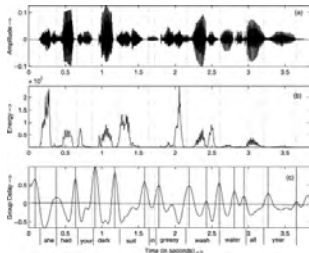


## Construire un programme CP

- **experience E** : statistiques
- **performance measure P** : optimisation
- **tasks T** : utilité
  - ▶ traduction automatique
  - ▶ jouer aux échec
  - ▶ conduire, ... et faire ce que l'on fait



# Les tâches T





## L'expérience E : les données

capteurs	→	variables qualitatives / ordinales / quantitatives
texte	→	chaîne de caractères
parole	→	temps - série temporelle
images/vidéos	→	dépendances 2/3 d
réseaux	→	graphes
jeux	→	séquences d'interactions
flots	→	tickets de caisse, web logs, trafic. . .
étiquettes	→	information d'évaluation



### Les mégadonnées (volume, vitesse, variété, véracité et valeur)

Des données qui arrivent sans qu'on ait décidé de les collecter

→ l'importance des pré-traitements (nettoyage, normalisation, codage. . .)

→ et de la représentation : de la donnée au vecteur

# Objectifs, tâches et mesures de performance $P$

non supervisé

▶ clustering

homogénéité des groupes

supervisé

▶ classification

- ★ bi classe
- ★ multi classe
- ★ détection

nombre d'erreur (0/1)

▶ ordonnancement

nombre de non détection

▶ regression

erreur de rang

erreur quadratique

renforcement

probabilité de gagner

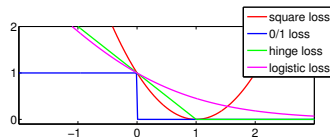
expliquer

qualité du résumé

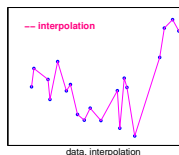
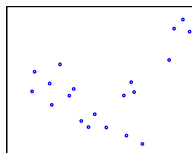
... et semi supervisé, transductif, séquentiel, actif...

## Généraliser

→ bien faire (minimiser  $P$ ) sur les futures données qu'on ne connaît pas



# L'apprentissage comme un problème d'optimisation bi critère

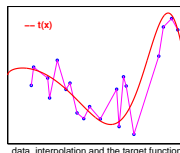
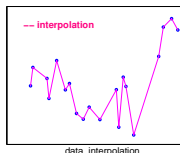
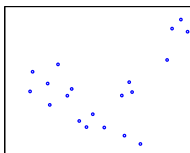


$\left\{ \begin{array}{l} \min \\ f \in H \end{array} \right. \quad \text{s'ajuster aux données}$

Apprendre c'est sélectionner une bonne hypothèse

- un bon model  $H$  (ensemble d'hypothèses universelles)
  - ▶ cout d'ajustement aux données

# L'apprentissage comme un problème d'optimisation bi critère



$$\left\{ \begin{array}{l} \min_{f \in H} \quad \text{s'ajuster aux données} \\ \min_{f \in H} \quad \text{bien généraliser} \end{array} \right.$$

## Apprendre c'est sélectionner une bonne hypothèse

- un bon model  $H$  (ensemble d'hypothèses universelles)
- $P =$  deux critères contradictoires
  - ▶ cout d'ajustement aux données
  - ▶ pénalité favorisant une certaine régularité
- une bonne gestion de cette contradiction

## Exemples de model $H$ (ensemble d'hypothèses universelles)

- dictionnaire fixe

$$f(x) = \sum_{j=1}^p \alpha_j \phi_j(x)$$

- dictionnaire adapté aux observations

$$f(x) = \sum_{j=1}^p \alpha_j \phi_j(x, x_j)$$

- ▶ les machines à noyaux  $f(x) = \sum_{i=1}^n \alpha_i k(x, x_i)$

- dictionnaire adapté aux observations et aux étiquettes

$$f(x) = \sum_{i=1}^p \alpha_i \phi_i(x, x_i, y_i)$$

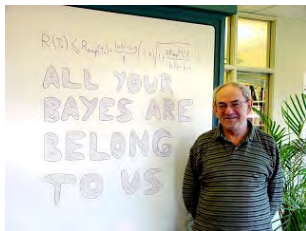
- ▶ les réseaux de neurones de type perceptron multicouche

$$f(x) = \sum_{k=1}^p \alpha_k \Phi \left( \sum_{j=1}^p \beta_j \mathbf{x} \right)$$

- ▶ méthodes d'ensemble

# Les objectifs de l'apprentissage automatique

- Algorithmes d'apprentissage
  - ▶ bonne généralisation
  - ▶ généricité
  - ▶ passage à l'échelle
- Questions théoriques (Vapnik's Book, 1982 - Valiant, 1984)
  - ▶ apprenabilité, sous quelles conditions ?
  - ▶ complexité (en temps, en échantillon)

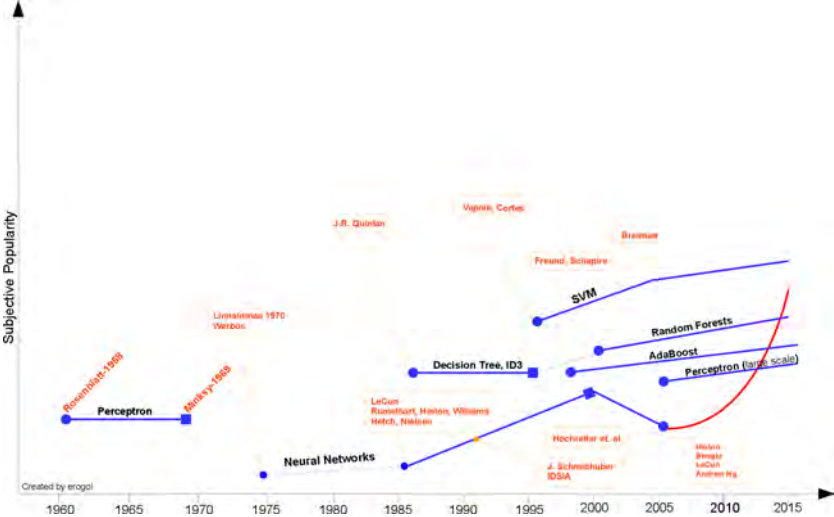


# Une brève histoire de l'apprentissage artificiel

- 1940 Etude de la logique en temps qu'objet mathématique (Shannon, Godel et al).
- 1950 Test de Turing,  
1956 Dartmouth conférence : *artificial intelligence*.
- 1960 "*Within a generation ... the problem of creating 'artificial intelligence' will substantially be solved.*"
- 1960 *Le perceptron*
- 1970 *AI symbolique : les systèmes experts*
- 1980 *AI Winter & Réseaux de neurones (perceptron non linéaire)*
- 1990 *SVM & COLT*
- 2000 *les applications : Google, Amazon (recommandation), spam filtering, OCR, parole, traduction...*
- 2010 *Big data, Go player*



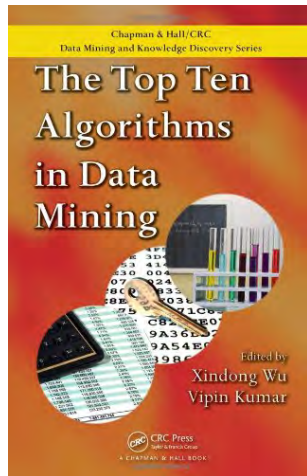
# Une brève histoire des modes en apprentissage





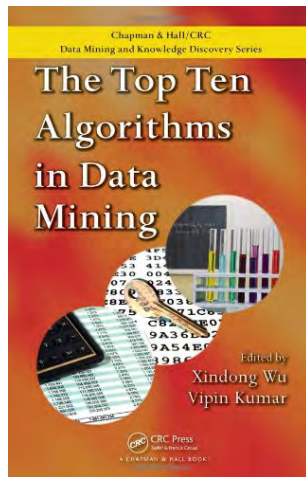
# Top 10 algorithmes en fouille de données

- Arbres de décision (et les forêts)
  - ▶ C4.5,
  - ▶ CART,
- SVM,
- AdaBoost,
- kNN,
- Bayésien
  - ▶ k-Means,
  - ▶ Apriori,
  - ▶ EM,
  - ▶ Naive Bayes,
- PageRank,
- 



# Top 10 algorithmes en fouille de données

- Arbres de décision (et les forêts)
  - ▶ C4.5,
  - ▶ CART,
- SVM,
- AdaBoost,
- kNN,
- Bayésien
  - ▶ k-Means,
  - ▶ Apriori,
  - ▶ EM,
  - ▶ Naive Bayes,
- PageRank,
  
- en 2014 un 11 ème : deep networks



Meilleure compréhension du monde à partir d'observations  
en vue de prédiction



C'est beau mais

→ Comment ça marche ?

# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



# Etude de cas

- web : google, facebook. . .
- Marketing : Walmart & big data
  - ▶ volume
  - ▶ variété
  - ▶ vitesse
- Quelques challenges + image + la traduction automatique : *deep learning strikes back*



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



# Une brève histoire des moteurs de recherche

1994



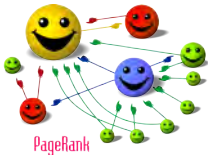
solution manuelle

1995



solution logique

1996



modèle statistique  
+ valeur propre

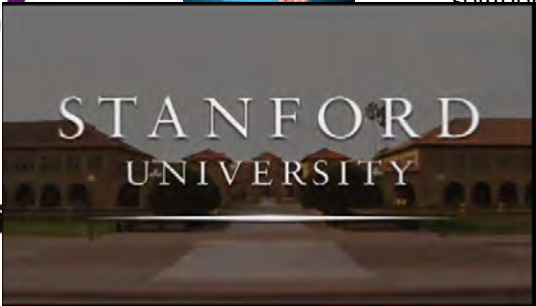
# Une brève histoire des moteurs de recherche

1994



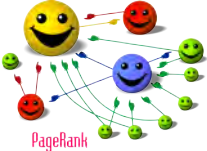
solution manuelle

1995



logique

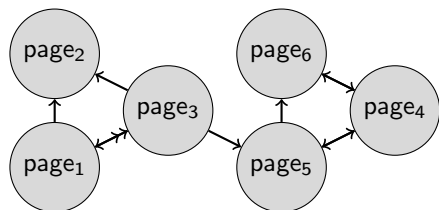
1996



modèle statistique  
+ valeur propre



## La recherche sur le web comme un problème de matrice

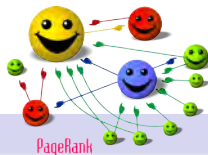


$$Y = \begin{pmatrix} 0 & 0 & 1/3 & 0 & 0 & 0 \\ 1/3 & 0 & 1/2 & 0 & 0 & 0 \\ 2/3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 1/3 & 1/2 & 0 & 0 \\ 0 & 0 & 0 & 1/2 & 1/2 & 0 \end{pmatrix}$$

$\mathbb{P}(\text{visiter la page}) =$  une mesure de l'intérêt de la page

$Y : \mathbb{P}(j \rightarrow i) =$  matrice des transitions

# Formalisation du problème



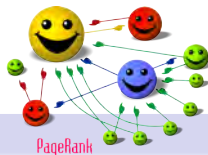
## Modèle

$$\begin{aligned} \mathbb{P}(\text{visiter } p_i) &= \sum_j \mathbb{P}(\text{visiter } p_j) \mathbb{P}(j \rightarrow i) + \mathbb{P}(\text{hasard}) \\ \mathbb{P} &= (1 - \delta) Y \mathbb{P} + \delta \mathbf{1} \mathbb{P} \\ \textit{observation} &= \textit{information} + \textit{bruit} \end{aligned}$$

avec  $Y$  la matrice des transitions

et  $\mathbb{P}$  le vecteur des probabilités de visiter une page

# Formalisation du problème



## Modèle

$$\begin{aligned} \mathbb{P}(\text{visiter } p_i) &= \sum_j \mathbb{P}(\text{visiter } p_j) \mathbb{P}(j \rightarrow i) + \mathbb{P}(\text{hasard}) \\ \mathbb{P} &= (1 - \delta) Y \mathbb{P} + \delta \mathbf{I} \mathbb{P} \\ \textit{observation} &= \textit{information} + \textit{bruit} \end{aligned}$$

avec  $Y$  la matrice des transitions

et  $\mathbb{P}$  le vecteur des probabilités de visiter une page

## Problème : retrouver $\mathbb{P}(\text{visiter } p_i)$

$$M \mathbb{P} = \mathbb{P} \quad \text{avec} \quad M = (1 - \delta) Y + \delta \mathbf{I}$$

- C'est un problème de **valeur propre** de très grande taille
- avec  $\delta = 0.15$

# Les leçons de l'histoire des moteurs de recherche

- l'arrivée des big data  
mégadonnées
- innover avec la recherche  
fondamentale & académique  
créer des partenariats
- poser les problèmes : faire des  
maths et des statistiques  
programmation à base  
d'exemples
- poser les vrai usages : interroger  
un moteur de recherche  
organiser l'information  
au niveau mondial
- business model : la valeur c'est  
la requête  
vers la recommandation



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



# La recommandation de tous les jours

The screenshot shows the Amazon.com homepage. At the top, there's a navigation bar with the Amazon logo, a search bar containing "MP3 Downloads", and a shopping cart icon. Below the search bar, there are several promotional banners, including one for "amazonmp3" with the text "Play Anywhere, DRM-Free Music Downloads". A "Hello, Sign in to get personalized recommendations. New customer? Start here" message is visible. The main content area features a "New & Notable" section with four album covers: E-MC<sup>2</sup> (Maiah Carey), Hannah Montana (Hannah Montana), Worlds Collide DeL... (Apologetiqa), and Lady Antel... (Lady Antel).

The screenshot shows the docstoc website. The main content area displays a document titled "How to use PLS path modeling for analyzing multiblock data sets". The document is available for download and printing. The website also features a sidebar with various documents, including "Next Generation PCR", "Easy Decision Trees", "SPSS Data Mining Secrets", "Data Mining Automation", and "Vacation Rental Software". The document preview shows a slide with the title "How to use PLS path modeling for analyzing multiblock data sets" and a background image of a landscape.

## Systèmes existants

Libre: Myrrix (facto, Web Mahout), Easyrec (item-based, Web), Duine (multi approches), Propriétaires: KXpro...

# Quelle méthode choisir ?

## Netflix Challenge (2007-2009)

- Recommandation de films
- Moteur maison : CineMatch
- Améliorer les performances de 10%
- critère : erreur quadratique



## Les données

- 450 000 spectateurs
- 17 000 films
- 100 millions d'évaluations (de 1 à 5)
- taux de remplissage : 1,3 %
- test sur les 3 millions les plus récents
- 48 000 téléchargements des données

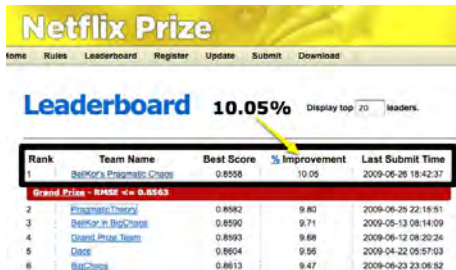
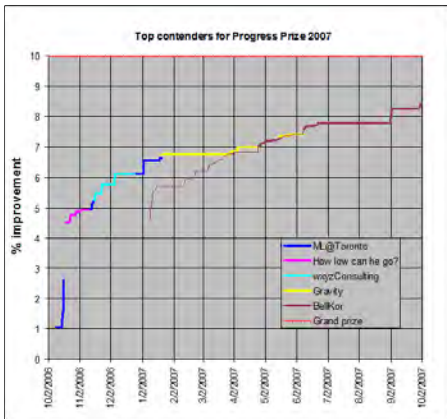


17,700 Movies  
 in the  
**Netflix Competition**

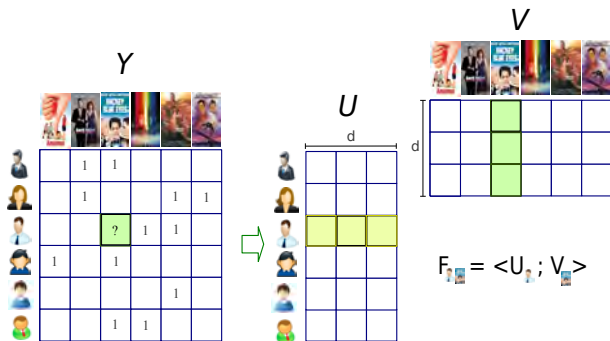
Todd Holloway@gmail.com 03/25/2007



# Faire de la science un sport !



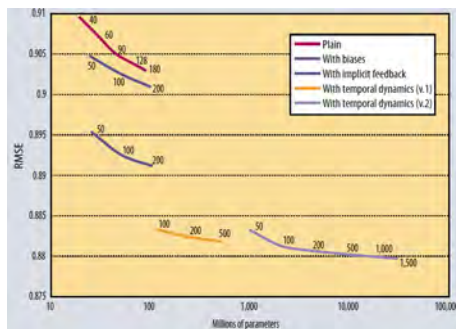
# La recommandation comme un problème de matrice



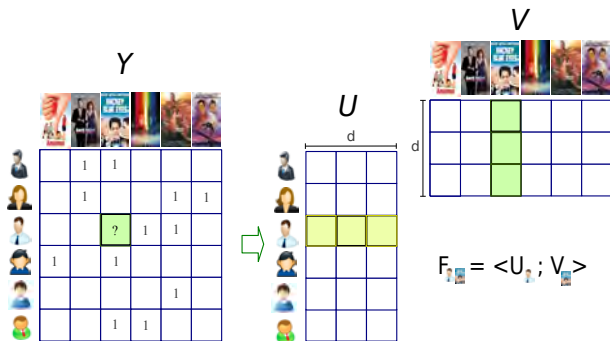
$\langle U_i, V_j \rangle$  une mesure d'appétence de l'utilisateur  $i$  pour le produit  $j$

# Les résultats sur Netflix

- erreur initiale : 0.9514
- factorisation - 4 %
- nombres de facteurs - .5 %
- améliorations - 2.5 %
  - ▶ normalisation
  - ▶ poids
  - ▶ temps
- bagging : - 3 % = 0.8563  
mélange 100 méthodes



# La recommandation comme un problème de matrice



$\langle U_i, V_j \rangle$  une mesure d'appétence de l'utilisateur  $i$  pour le produit  $j$

# Formalisation du problème

## Modèle

$$\begin{array}{rcccl} Y & = & F & + & R \\ \textit{observation} & = & \textit{information} & + & \textit{bruit} \end{array}$$

avec  $F$  régulière et  $R$  une **matrice** de bruit.

# Formalisation du problème

## Modèle

$$\begin{array}{rclcl} Y & = & F & + & R \\ \textit{observation} & = & \textit{information} & + & \textit{bruit} \end{array}$$

avec  $F$  régulière et  $R$  une **matrice** de bruit.

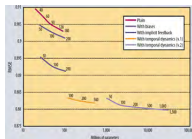
Problème : construire  $\hat{F}$  un estimateur de  $F$

$$\min_{\hat{F}} \underbrace{\mathcal{L}(Y, \hat{F})}_{\text{attache aux données}} + \lambda \underbrace{\Omega(\hat{F})}_{\text{pénalisation}}$$

- optimiser un cout bien choisi
- avec  $\hat{F} = UV$  factorisation (...à l'aide de valeurs propres)

# Les leçons de Netflix

- trop d'efforts spécifiques pour être généralisée
- Modèle statistique : **factorisation** *is the solution*
- prendre en compte **les spécificités** (ici l'effet temps et les feedback implicites)
- une optimisation de type **gradient stochastique**
  - ▶ flexible
  - ▶ qui passe à l'échelle
- quels **usages** pour la recommandation ?



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

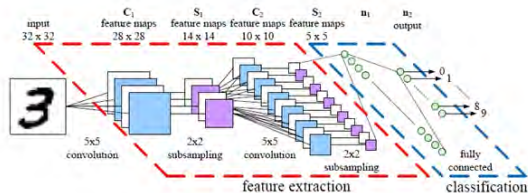
## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

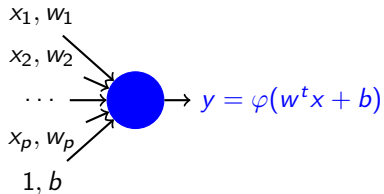
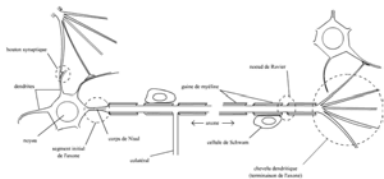
## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion





# The formal neuron



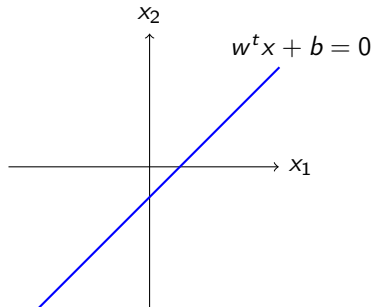
$x$  input

$w$  weight,  $b$  bias

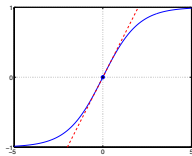
$\varphi$  activation function

$y$  output

Defines a hyperplane



$$\varphi(t) = \tanh(t)$$



# The formal neuron as a learning machine

## Fitting the data for the 2 classes classification problem

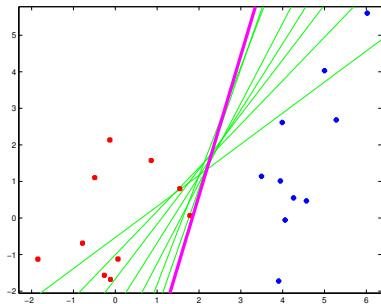
given  $n$  pairs of input–output data  $\mathbf{x}_i, y_i, i = 1, n$

find  $w$  such that

$$\underbrace{\varphi(\mathbf{w}^t \mathbf{x}_i)}_{\text{prediction of the model}} = \underbrace{y_i}_{\text{ground truth}}$$

make some initial guess for  $w$

- pick some data
- extract information
- improve  $w$



what about new data – generalization?

# Fitting the data with an energy-based model

Fitting the data

$$\min_{\mathbf{w} \in \mathbb{R}^{p+1}} \sum_{i=1}^n \left( \underbrace{\varphi(\mathbf{w}^t \mathbf{x}_i)}_{\text{prediction}} - \underbrace{y_i}_{\text{truth}} \right)^2$$

Improve  $\mathbf{w}$

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \rho \underbrace{(\varphi(\mathbf{w}^t \mathbf{x}_i) - y_i) \nabla_{\mathbf{w}} \varphi(\mathbf{w}^t \mathbf{x}_i)}_{\text{extracted information}}$$

---

## Algorithm 1 Adaline

---

**Data:**  $\mathbf{w}$  initialization,  $\rho$  fixed stepsize

**Result:**  $\mathbf{w}$

**while** *not converged* **do**

$\mathbf{x}_i, y_i \leftarrow$  pick a point  $i$  at random

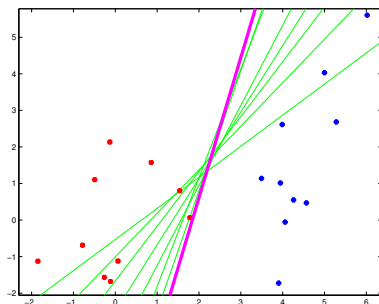
$\mathbf{d} \leftarrow (\mathbf{w}^t \mathbf{x}_i - y_i) \mathbf{x}_i$

$\mathbf{w} \leftarrow \mathbf{w} - \rho \mathbf{d}$

    check that the cost decrease

**end**

---



# Fitting the data with an energy-based model (part of $P$ )

## Coding the output (truth?)

$$y_i = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{class 1} \\ -1 & \text{else} \end{cases}$$

Fitting the data

$$\min_{\mathbf{w} \in \mathbb{R}^{p+1}} \sum_{i=1}^n \left( \underbrace{\varphi(\mathbf{w}^t \mathbf{x}_i)}_{\text{prediction } p} - \underbrace{y_i}_{\text{truth } t} \right)^2$$
$$(\varphi(\mathbf{w}^t \mathbf{x}_i) - y_i)^2 = y_i^2 \underbrace{(y_i \varphi(\mathbf{w}^t \mathbf{x}_i) - 1)}_{\text{score } z_i}$$

As a loss

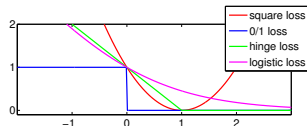
$$\min_{\mathbf{w} \in \mathbb{R}^{p+1}} \sum_{i=1}^n \text{cost}(z_i) \quad z_i = \varphi(\mathbf{w}^t \mathbf{x}_i) y_i$$

0/1 loss

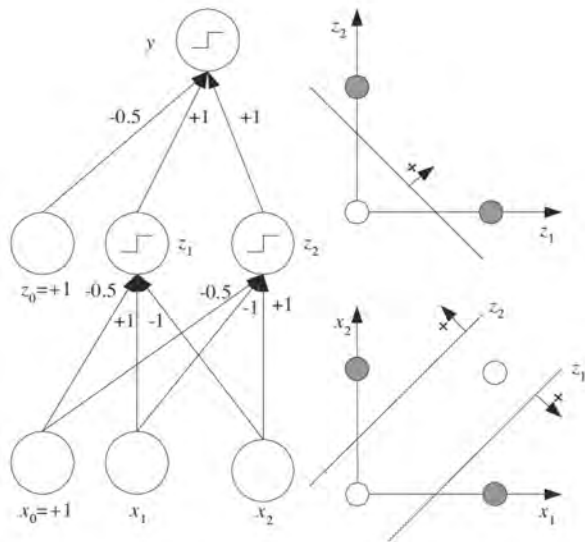
$$\text{cost}(z_i) = \begin{cases} 0 & \text{if } \text{sign}(\mathbf{w}^t \mathbf{x}_i) = \text{sign}(y_i) \\ 1 & \text{else} \end{cases}$$

As a constraint

$$\text{sign}(\mathbf{w}^t \mathbf{x}_i) \text{sign}(y_i) \geq 0$$

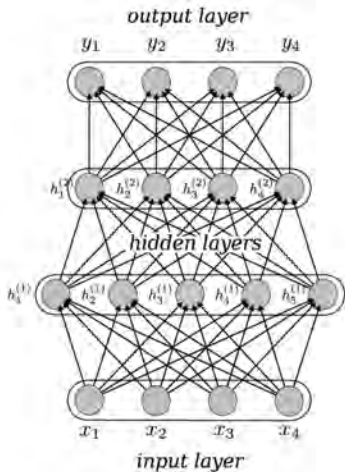


## Non linearity combining linear neurons



# Neural networks as universal approximator

Running several neurons at the same time



$$y = \varphi(W_3 h^{(2)})$$

$$h^{(2)} = \varphi(W_2 h^{(1)})$$

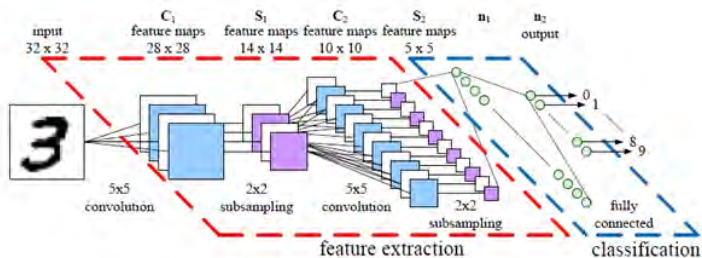
$$h^{(1)} = \varphi(W_1 x)$$

$x$

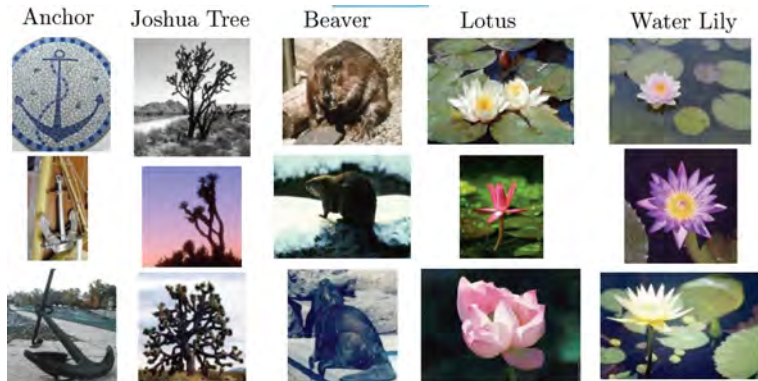
Multilayered neural network: learning internal representation  $W_1, W_2, W_3$

# OCR: the MNIST database (Y. LeCun, 1989)

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



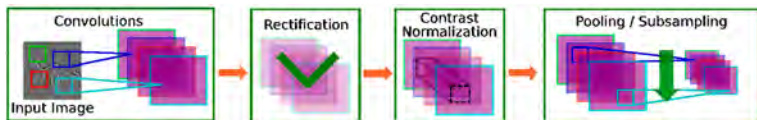
# The caltech 101 database (2004)



- 101 classes
- 30 training images per category
- ...and the winner is NOT a deep network
  - ▶ dataset is too small



# Deep architecture lessons from caltech 101 (2009)



Single Stage System: $[64.F_{CSG}^{U \times U} - R/N/P^{5 \times 5}] - \log_{reg}$					
	$R_{abs} - N - P_A$	$R_{abs} - P_A$	$N - P_M$	$N - P_A$	$P_A$
U <sup>+</sup>	54.2%	50.0%	44.3%	18.5%	14.5%
R <sup>+</sup>	54.8%	47.0%	38.0%	16.3%	14.3%
U	52.2%	43.3%(±1.6)	44.0%	17.2%	13.4%
R	53.3%	31.7%	32.1%	15.3%	12.1%(±2.2)
G	52.3%				
Two Stage System: $[64.F_{CSG}^{U \times U} - R/N/P^{5 \times 5}] - [256.F_{CSG}^{U \times U} - R/N/P^{4 \times 4}] - \log_{reg}$					
	$R_{abs} - N - P_A$	$R_{abs} - P_A$	$N - P_M$	$N - P_A$	$P_A$
U <sup>+</sup> U <sup>+</sup>	65.5%	60.5%	61.0%	34.0%	32.0%
R <sup>+</sup> R <sup>+</sup>	64.7%	59.5%	60.0%	31.0%	29.7%
UU	63.7%	46.7%	56.0%	23.1%	9.1%
RR	62.9%	33.7%(±1.5)	37.6%(±1.9)	19.6%	8.8%
GT	55.8%				
Single Stage: $[64.F_{CSG}^{U \times U} - R_{abs}/N/P_A^{5 \times 5}] - PMK-SVM$					
U	64.0%				
Two Stages: $[64.F_{CSG}^{U \times U} - R_{abs}/N/P_A^{5 \times 5}] - [256.F_{CSG}^{U \times U} - R_{abs}/N] - PMK-SVM$					
UU	52.8%				

$R_{abs}$  Rectification Layer  
 $N$  Local Contrast Normalization Layer  
 $P_A$  Average Pooling and Subsampling Layer

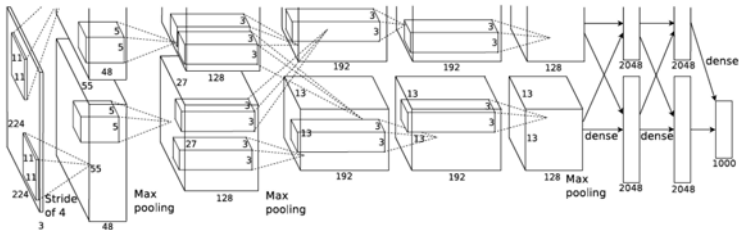
# The image net database (Deng et al., 2012)



ImageNet = 15 million labeled high-resolution images of 22,000 categories.  
Large-Scale Visual Recognition Challenge (a subset of ImageNet)

- 1000 categories.
- 1.2 million training images,
- 50,000 validation images,
- 150,000 testing images.

# Deep architecture and the image net



The architecture of the CNN [Krizhevsky, Sutskever, Hinton, 2012]

- 60 million parameters
- using 2 GPU
- regularization
  - ▶ data augmentation
  - ▶ dropout
  - ▶ weight decay

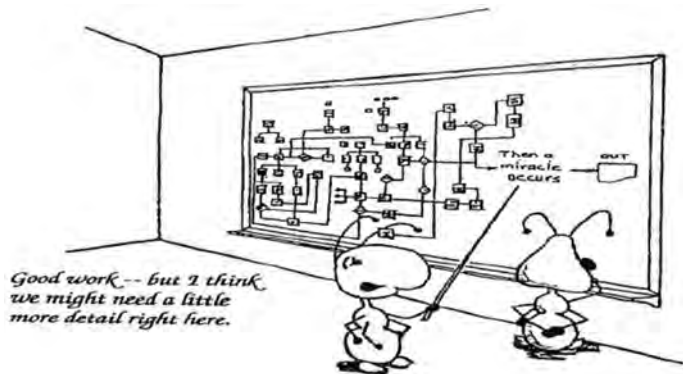


# A new fashion in image processing

2012 Teams	%error	2013 Teams	%error	2014 Teams	%error
Supervision (Toronto)	15.3	Clarifai (NYU spinoff)	11.7	GoogLeNet	6.6
ISI (Tokyo)	26.1	NUS (singapore)	12.9	VGG (Oxford)	7.3
VGG (Oxford)	26.9	Zeiler-Fergus (NYU)	13.5	MSRA	8.0
XRCE/INRIA	27.0	A. Howard	13.5	A. Howard	8.1
UvA (Amsterdam)	29.6	OverFeat (NYU)	14.1	DeeperVision	9.5
INRIA/LEAR	33.4	UvA (Amsterdam)	14.2	NUS-BST	9.7
		Adobe	15.2	TTIC-ECP	10.2
		VGG (Oxford)	15.2	XYZ	11.2
		VGG (Oxford)	23.0	UvA	12.1

Baidu: 5%

# Learning Deep architecture

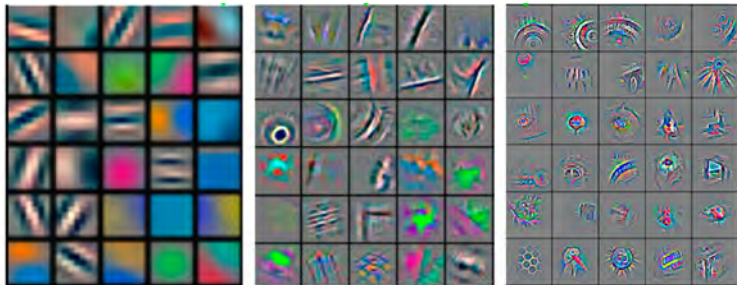
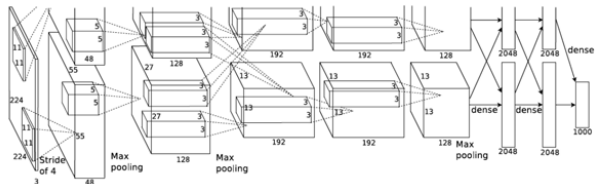


$$\min_{W \in \mathbb{R}^d} \sum_{i=1}^n \|f(x_i, W) - y_i\|^2 + \lambda \|w\|^2$$

- $d = 60 \times 10^6$
- $n = 1,200,000 ++$
- $\lambda = 0.0005$

$f$  is a deep NN

# Then a m. occurs: learning internal representation



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

## Deep architecture: more contests and benchmark records



- speech (phoneme) recognition,
- automatic translation,
- Optical Character Recognition (OCR),
- ICDAR Chinese handwriting recognition benchmark,
- Grand Challenge on Mitosis Detection, Road sign recognition
- Higgs boson challenge

# Machine Learning in High Energy Physics

The screenshot shows the Higgs Boson Machine Learning Challenge website. At the top left is the logo for the Higgs challenge, featuring the text 'Higgs challenge' and a stylized 'H' made of yellow and black squares. To the right of the logo, it says 'Completed • 513,000 • 1,785 teams' and 'Higgs Boson Machine Learning Challenge' in bold. Below that, it indicates the dates 'Mon 12 May 2014 - Mon 15 Sep 2014 (8 months ago)'. A navigation menu on the left includes 'Dashboard', 'Home', 'Data', 'Make a submission', 'Information', 'Forum', and 'Leaderboard'. The main content area has 'Competition Details' and links for 'Get the Data' and 'Make a submission'. The central text reads 'Use the ATLAS experiment to identify the Higgs boson'. Below this is a banner image for the ATLAS EXPERIMENT, showing a particle detector structure with colorful tracks and the text 'ATLAS EXPERIMENT', 'Run: 204153', 'Event: 35369265', and '2012-05-30 20:31:28 UTC'.

The winning model is "brute force"

- a bag of 70 dropout neural networks
- three hidden layers of 600 neurons
- produced by 2-fold stratified cross-validations



# Comment ça marche ?

- Programmation par l'exemple (et beaucoup - des méga données)
- un bon critère :
  - ▶ intègre des connaissances a priori
  - ▶ adjustment aux données (énergie)
  - ▶ un ou des mécanismes de régularisation
- une bonne procédure d'optimisation
- des ressources informatiques suffisantes

## Ce que l'on ne comprend pas

- gérer et protéger les **données** et les contextes ?
- qu'est-ce qu'apprendre ?
  - ▶ pourquoi les architectures profondes fonctionnent ?
  - ▶ apprendre à apprendre (transfert, transport) ?
  - ▶ quelles **garanties** peut on avoir ?
  - ▶ comment gérer les **dynamiques** ?
- comment **optimiser** efficacement des fonctionnelles non convexes et non régulières ?
- comment **automatiser** la chaîne de traitement nécessaire à l'apprentissage ?



Statistiques + optimisation + informatique

vers la science des données

# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion

### Data Science



# New trends in ML fashion: big data

## ICML 2014 workshops

- Designing Machine Learning Platforms for **Big Data**  
Xiangxiang Meng, Wayne Thompson, Xiaodong Lin
- New Learning Frameworks and Models for **Big Data**  
Massih-Reza Amini, Eric Gaussier, James Kwok, Yiming Yang
- Deep Learning Models for Emerging **Big Data** Applications  
Shan Suthaharan, Jinzhu Jia
- Unsupervised Learning for Bioacoustic **Big Data**  
Hervé Glotin, P. Dugan, F. Chamroukhi, C. Clark, Yann LeCun
- Knowledge-Powered **Deep Learning for Text Mining**  
Bin Gao, Scott Yih, Richard Socher, Jiang Bian
- Optimizing Customer Lifetime Value in **Online Marketing**  
Georgios Theodorou, Mohammad Ghavamzadeh, Shie Mannor

Comment traiter ces *big data* : data science

## **this joke makes perfect sense**

4 January 2013 12:28

Big data is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it

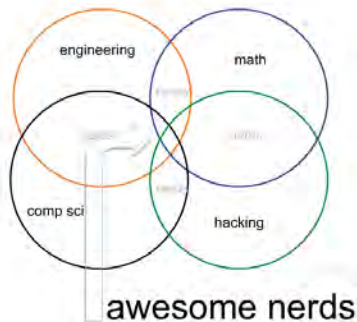
@Behaishiyi

[weibo.com/behaisiyyi](http://weibo.com/behaisiyyi)

# New trends in ML fashion: data science



## Data science?



# Data Scientist Voted Sexiest Job of 21st Century

## The Rise of Data Scientists

BEFORE



nobody cared for a "math geek" in parties.

NOW



People love ~~math geeks~~ data scientists!  
RK

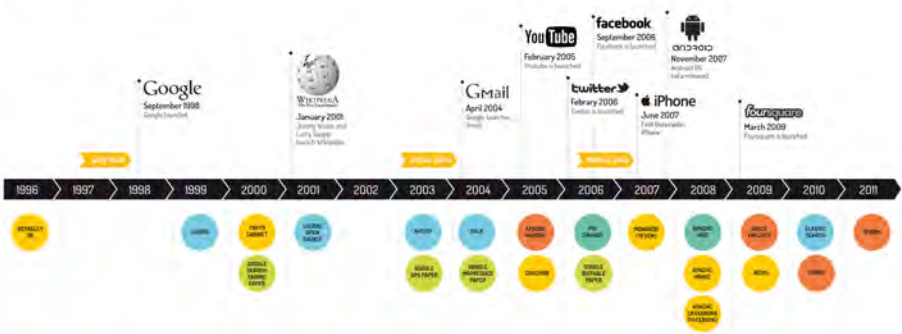
## New trends in ML fashion: big data hiring boom



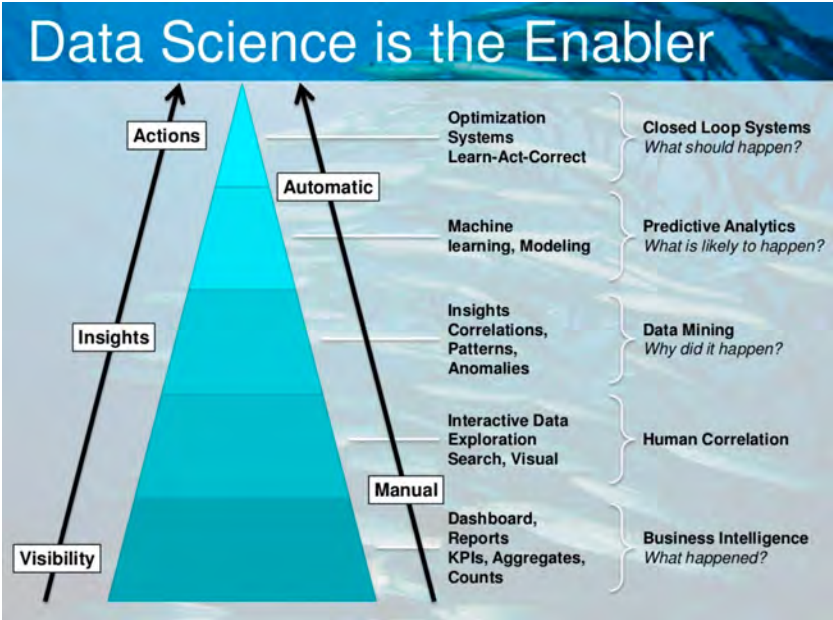


# Big data brief history

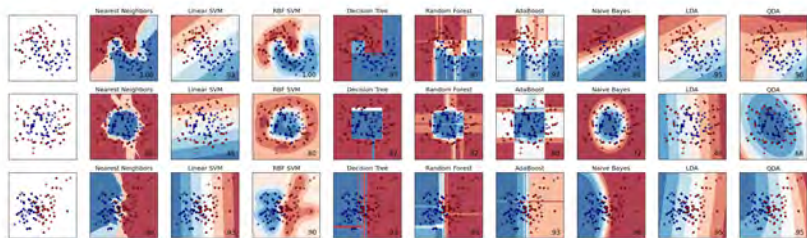
# BIG DATA A BRIEF HISTORY



Original source at <http://invent.ge/Uah9YO>



# M.L. fashion as classifier comparison (from scikit-learn)

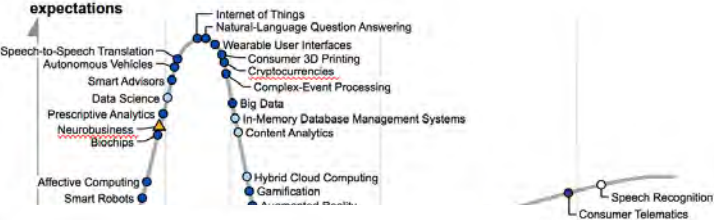


# Une brève histoire de la mode en technology (by Gartner)

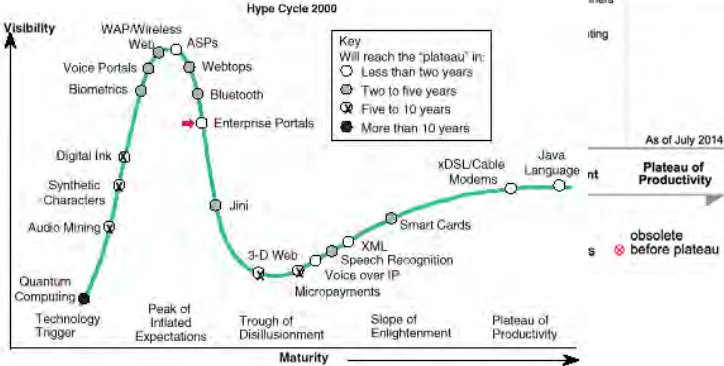


Gartner's 2014 Hype Cycle for Emerging Technologies Maps the Journey to Digital Business  
<http://www.gartner.com/newsroom/id/2819918>

# Une brève histoire de la mode en technology (by Gartner)



Volu

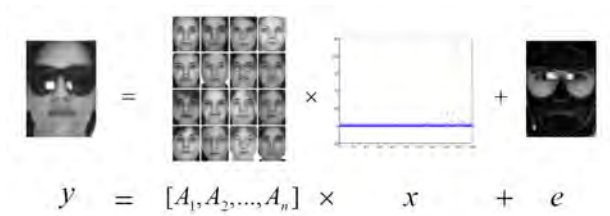


Gartner's 2000 Hype Cycle for Emerging Technologies Maps the Journey to Digital Business

<http://www.gartner.com/newsroom/id/2819918>

## Quelles propriétés

passage à l'échelle → parcimonie : sélectionner les exemples et les variables  
randomisation → traiter la masse par le hasard


$$y = [A_1, A_2, \dots, A_n] \times x + e$$



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

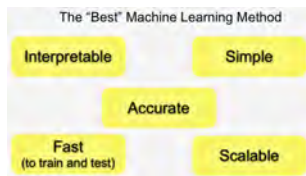
## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

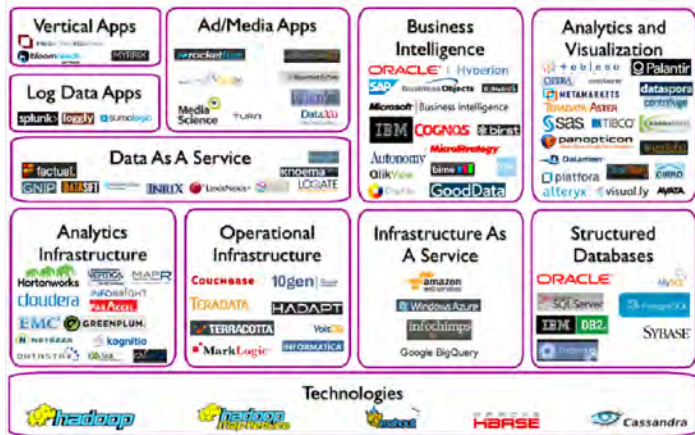
## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion

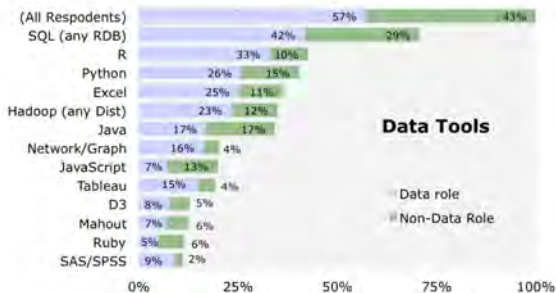


# Les outils de la science des données

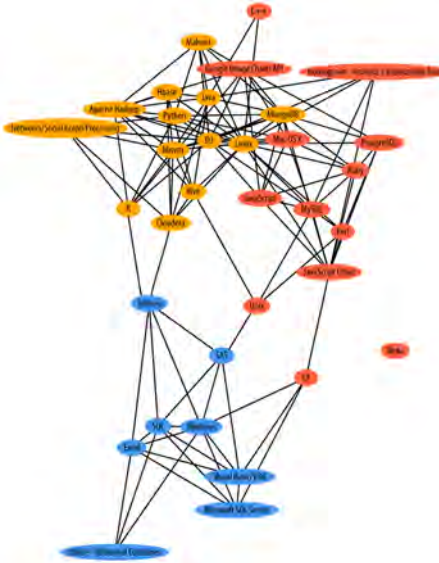
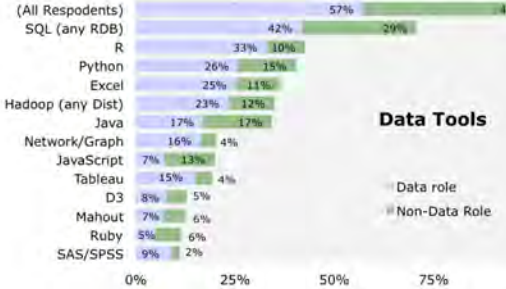




# Les outils de la science des données



# Les outils de la science des données



# Logiciels (mloss.org)

environ. interactifs

- ▶ R (opensource)
- ▶ Matlab,

Java

- ▶ Weka,
- ▶ Rapid miner,
- ▶ MLib (Spark), [spark.apache.org](http://spark.apache.org)

Python

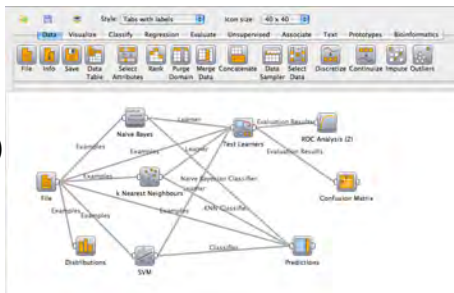
- ▶ Orange, <http://orange.biolab.si>
- ▶ SciPy, <http://www.scipy.org>
- ▶ Shogun, MLpy, ...

C, C++

- ▶ Vowpal Wabbit, [hunch.net/~vw](http://hunch.net/~vw)
- ▶ Torch (Lua), <http://torch.ch/>
- ▶ libLinear, libSVM,

ML as a service

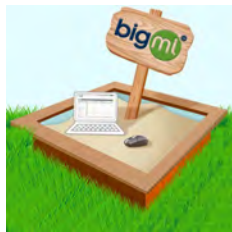
- ▶ Google prediction API,
- ▶ Microsoft Azure Machine Learning API,



# Question méthode



- 1 définir objectifs, les contraintes et les couts
- 2 étude préliminaire
  - ▶ récupérer des données
  - ▶ construire un bac à sable
  - ▶ affiner les objectifs et les couts en interagissant avec les données
  - ▶ proposer un modèle et choisir une solution (un alto)
- 3 mise en oeuvre sur de véritable flots de données
  - ▶ randomisation
  - ▶ outils adaptés



# Plan

## 1 Introduction

- Une brève définition de l'apprentissage statistique
- Les principaux algorithmes d'apprentissage

## 2 Etudes de cas

- Apprendre à rechercher l'information
- Apprendre à recommander
- Apprendre à étiqueter une image avec un réseaux de neurones

## 3 La science des données

## 4 Les outils pour l'apprentissage statistique

## 5 Conclusion



# Quelques prédictions... à propos du futur

- Data science - Applications
  - ▶ multi compétences
  - ▶ chaine de traitements
- Outils pour l'apprentissage
  - ▶ l'apprentissage sans paramètres (off-the-shelf)
  - ▶ passage à l'échelle (4v - big data - mégadonnées)
- Algorithmes d'apprentissage
  - ▶ optimisation (mégadonnées, non convexe)
  - ▶ dynamique (interactions)
  - ▶ apprendre à apprendre (transfert)
- Théorie de l'apprentissage
  - ▶ la nature de l'information





## ● livres

- ▶ Bishop, C. M. 1995. Neural Networks for Pattern Recognition. Oxford: Oxford University Press.
- ▶ Duda, R. O., P. E. Hart, and D. G. Stork. 2001. Pattern Classification, 2nd ed. New York: Wiley.
- ▶ Hastie, T., R. Tibshirani, and J. Friedman. 2001. The Elements of Statistical Learning: Data Mining, Inference, and Prediction. New York: Springer.
- ▶ A. Cornuéjols & L. Miclet, "L'apprentissage artificiel. Concepts et algorithmes". Eyrolles. 2 ème éd. 2010.

## ● conférences

- ▶ xCML, NIPS, COLT, séminaire SMILE, Paris Machine Learning Applications Group

## ● revues

- ▶ JMLR, Machine Learning, Foundations and Trends in Machine Learning, machine learning survey <http://www.mlsurveys.com/>