

# Impact des **Big Data** sur la recherche en sciences du vivant et leurs implication / applications



Antoine Cornuéjols

*AgroParisTech* – INRA MIA 518

[antoine.cornuejols@agroparistech.fr](mailto:antoine.cornuejols@agroparistech.fr)

# Plan

---

1. La révolution des données
2. Qu'en fait-on ?
3. Les défis
4. Les compétences requises, les acteurs, les outils
5. L'avenir

# Où l'on parle d'« **avalanche** » de données

---

- Des données capturées à foison quand nous allons **sur Internet**
  - Sur quels sites
  - Combien de temps, les clics, les durées, les achats, ...
- **Smartphones**
  - Localisation même si on a dit non
  - Des tas d'applications pleines de curiosité
- **Bracelets** connectés
- Moyens de **paiement** (banques)
- Capteurs dans les **véhicules** (assurances)
- Compteurs Linky

# Où l'on parle d'« **avalanche** » de données

---

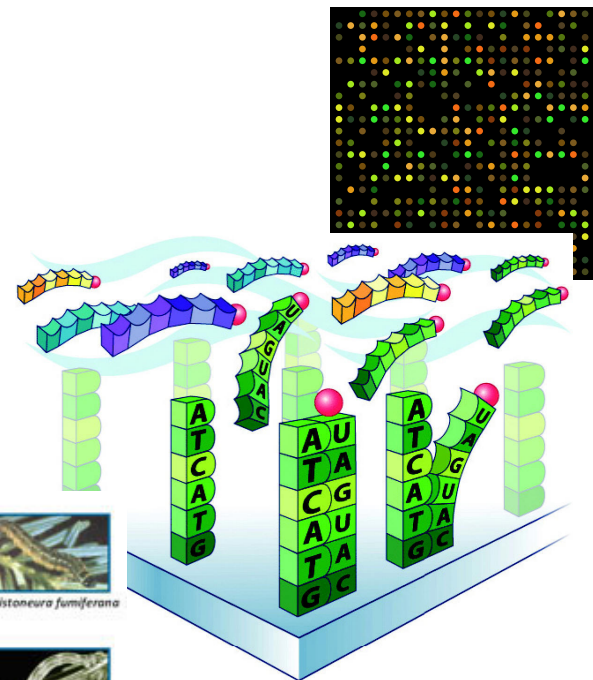
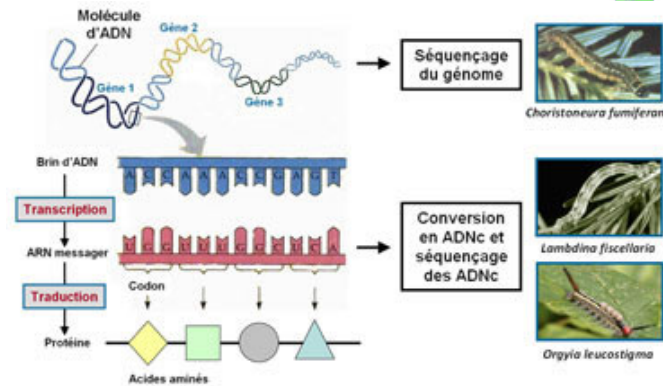
- Des **caméras** dans les panneaux de publicité dans les rues
- Bientôt dans les **vitrines** des magasins
- **Factures** dans les supermarchés (carte de fidélité)
- **Smart cities**

# Exemples de domaines renouvelés

- La bio-informatique

- Gros volume de données

- ◆ Puces à ADN
- ◆ Machines à séquencer



# Exemples de domaines renouvelés

---

- La sociologie
  - Gros volume de données
    - ◆ Réseaux sociaux
    - ◆ Smartphones
    - ◆ Websites consultations

# Exemples de domaines renouvelés

---

- La e-medecine (le me-data)

- Gros volume de données

- ◆ Smartphones
- ◆ Objets connectés
- ◆ Forums
- ◆ WATSON
- ◆ Google Flu

# Exemples de domaines renouvelés

- L'agriculture numérique

- Gros volume de données

- ◆ Capteurs
- ◆ Drones
- ◆ Réseaux sociaux et pro





# Exemples de domaines renouvelés

- Le domaine juridique

- Gros volume de données

- ◆ Archives numérisées
- ◆ Réseaux sociaux et professionnels



# Des **contre-exemples**

---

- **L'alimentation**

- Enquête **Nutrinet**

- ~ 277 000 internautes théoriquement sur des années
    - **Mais**
      - ◆ à 80% des femmes
      - ◆ Milieux socio-professionnels élevés
      - ◆ Abandonnent après quelques jours

**Manque de données  
représentatives**

- **L'éducation**

- Peu de données sur ce qui se passe en classe ou devant un écran

# Les sciences du vivant de l'environnement et de l'agronomie

# Spécificités

---

## 1. Des systèmes **naturels**

- Les **modèles** sous-jacents, les liens de **causalité** sont **inconnus**

## 2. Très **complexes**

- très **multi-échelles** spatiales et temporelles

## 3. Très **multi-factoriels**

« Défricher »

## 4. **Adaptatifs** donc difficiles à prévoir

- E.g. émergence de **bio-résistance**

Importance des  
liens de **causalité**

# Spécificités

---

1. La disposition de données permet
  - D'**explorer** : recherche de « patterns »
2. Possibilité d'une **approche heuristique**
  - On ne cherche pas nécessairement un **modèle explicatif** ou causal
  - On peut se contenter de **modèles prédictifs** (dans un 1er temps)

# Changement de paradigme

---

## 1. Ancien paradigme

- Construire une hypothèse (e.g. tel traitement devrait avoir tel effet)
- Construire un plan d'expérience pour **tester la validité de l'hypothèse**
- Le dispositif expérimental et les données récoltées **ne servent qu'à tester cette hypothèse**

---

## 1. Ancien paradigme

- Construire une hypothèse (e.g. tel traitement devrait avoir tel effet)
- Construire un plan d'expérience pour tester la validité de l'hypothèse
- Le dispositif expérimental et les données récoltées ne servent qu'à tester cette hypothèse

## 2. Nouveau paradigme

- Esprit « ouvert » : on cherche des patterns (inattendus) dans la masse de données disponibles
- Ré-utilisation possible à l'infini des données (non récoltées pour un but précis)

C'est le « data mining »



# Plan

---

1. La révolution des données
2. Qu'en fait-on ?
3. Les défis
4. Les compétences requises, les acteurs, les outils
5. L'avenir

# Apprentissage **descriptif** non supervisé

# Apprentissage **descriptif**

---

À propos d'un *échantillon d'apprentissage*  $s = \{(x_i)\}_{1,m}$

identifier des **régularités** rendant compte de  $S$

- E.g. sous la forme de **clusters** (e.g. *mélange de Gaussiennes*)
  - **CLUSTERING**
- E.g. sous la forme de **motifs fréquents** (fouille de données)

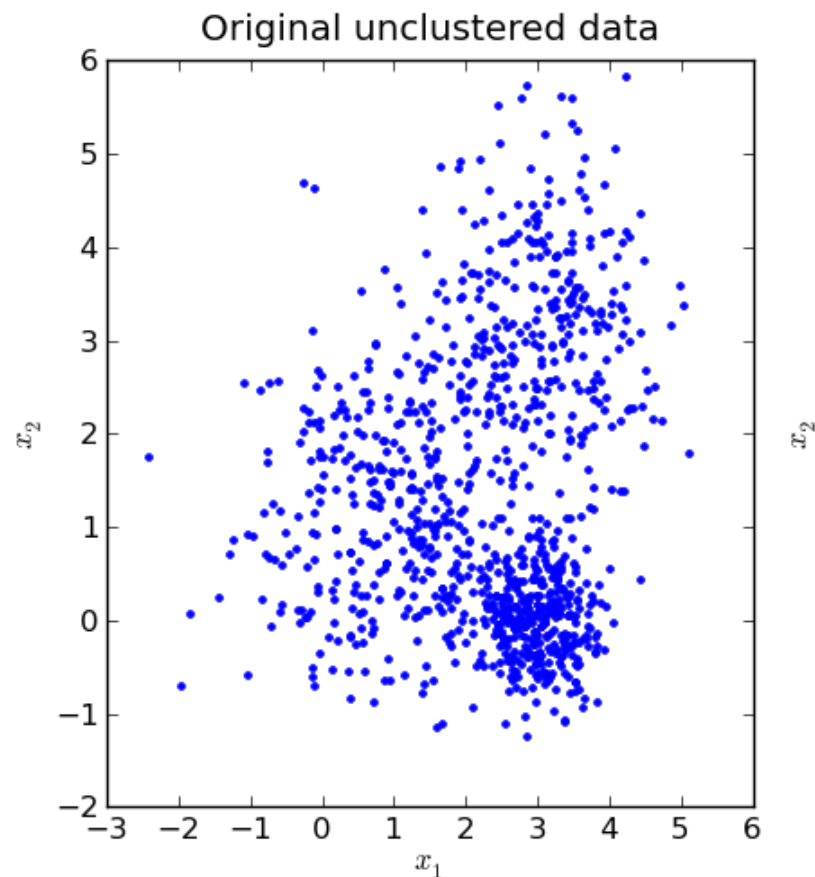
**pour résumer, suggérer des régularités, comprendre ...**



# Clustering / Catégorisation

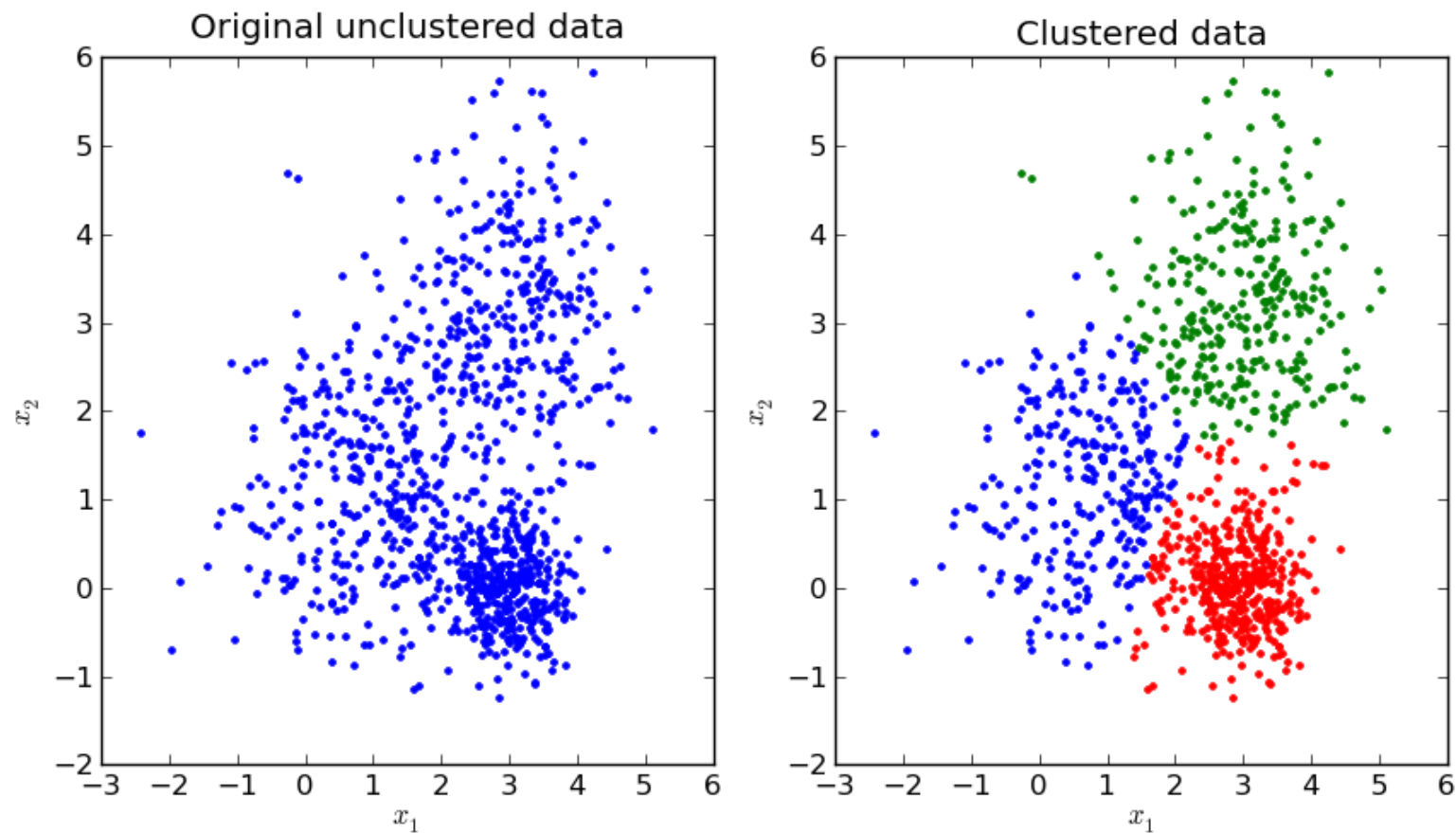
# Les grands types d'apprentissage

- Apprentissage « **descriptif** » (non supervisé)



# Les grands types d'apprentissage

- Apprentissage « **descriptif** » (non supervisé)

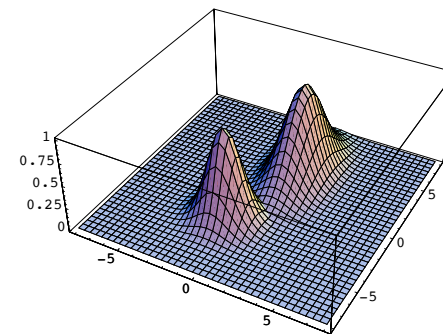


- **Catégorisation** de consommateurs

- Base de données sur les répondants de la base Nutrinet

- ~ 280 000
- Données sur *âge, nb de personnes dans la famille, catégorie socio-professionnelle, ...*
- Données sur consommations alimentaires sur une certaine durée

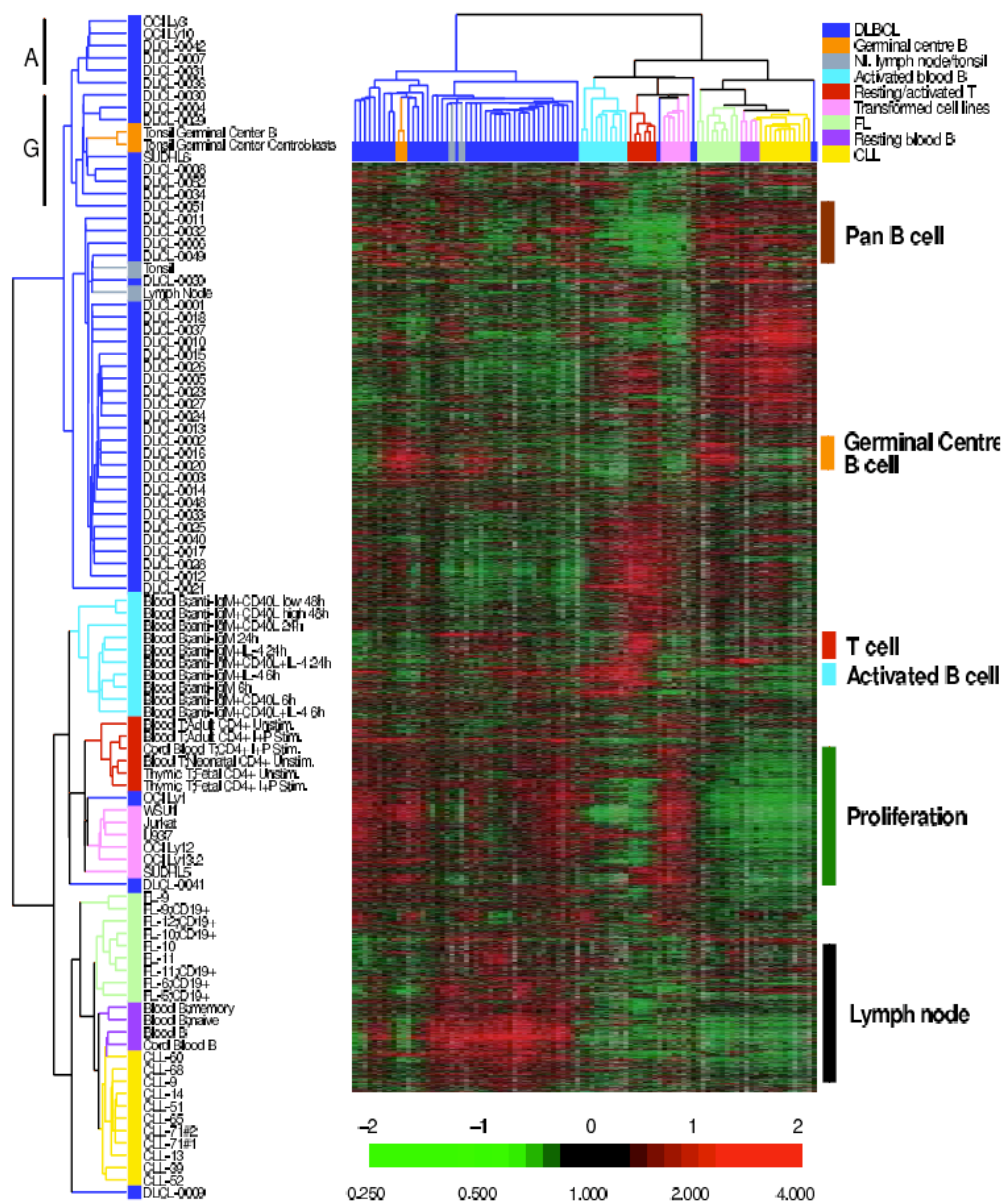
- Y a-t-il émergence de **groupes** distincts ?



Apprentissage  
Non supervisé

# Clustering

Bi-clustering  
gènes - patients







Recherche de motifs fréquents

*Frequent Item Sets*



Recherche de règles d'association

- Extraire des **régularités**
  - Base de données sur les **consommations alimentaires**
  - Peut-on identifier des « **patterns** » de consommation ?

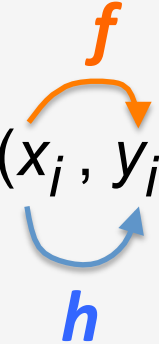
- Extraire des **régularités**
  - Base de données sur les **consommations alimentaires**
  - Peut-on identifier des « **patterns** » de consommation ?

- **Reconnaissance** d'animaux malades ou en chaleur
  - Mesures en continu sur leur comportement
    - Vidéos
    - Capteurs « embarqués »
      - Mobilité (nb de pas / minute ; distance parcourue à l'heure)
      - Lieux visités
      - ...
  - *Reconnaissance de **comportements types***

# Apprentissage **prédictif** supervisé

# Apprentissage prédictif (*supervisé*)

- Un *échantillon d'apprentissage*

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_j, y_j), \dots, (x_m, y_m)\}$$


The diagram illustrates the relationship between the input  $x_j$  and the output  $y_j$  in the training set  $S$ . An orange arrow labeled  $f$  points from  $x_j$  to  $y_j$ , representing the true function. A blue arrow labeled  $h$  points from  $x_j$  to  $y_j$ , representing the hypothesis function.

Prédiction pour de **nouveaux** exemples  $x \xrightarrow{h} y ?$

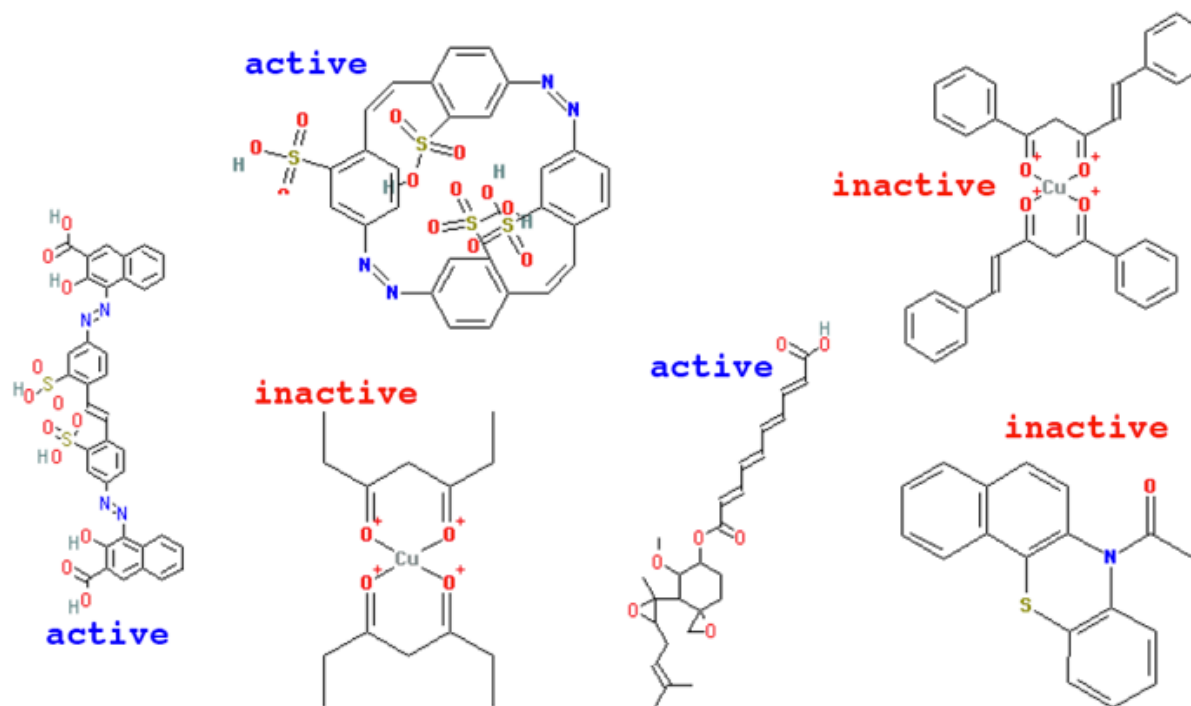
- **Reconnaissance** d'insectes ravageurs
  - Base d'images d'insectes dans des cuvettes
    - *Reconnaissance du type d'insectes*
    - *Comptage*



# Association / Prédiction

Apprentissage  
supervisé

- Prédire si une molécule est bio-active ou pas






NCI AIDS screen results (from <http://cactus.nci.nih.gov>).



# Analyse de textes

- Reconnaissance de **sentiments** exprimés dans des textes

	Electronics	Video games
	(1) <u>Compact</u> ; easy to operate; very good picture quality; looks <u>sharp</u> !	(2) A very <u>good</u> game! It is action packed and <u>full of excitement</u> . I am very much <u>hooked</u> on this game.
	(3) I purchased this unit from Circuit City and I was very <u>excited</u> about the quality of the picture. It is really <u>nice</u> and <u>sharp</u> .	(4) Very <u>realistic</u> shooting action and good plots. We played this and were <u>hooked</u> .
	(5) It is also quite <u>blurry</u> in very dark settings. I will <u>never_buy</u> HP again.	(6) It is so boring. I am extremely <u>unhappy</u> and will probably <u>never_buy</u> UbiSoft again.

# GIEC : filtrage de documents

---

- Estimation de l'**émission de gaz à effet de serre par les sols agricoles**
  - En particulier N<sub>2</sub>O (influence des engrais azotés)
- Par une méta-analyse des **articles scientifiques pertinents**
  - **Plus de 10<sup>6</sup> articles** scientifiques publiés / an
  - (plus ou moins) disponibles sur Internet

**Filtrage** nécessaire de ces articles

En optimisant **précision** et **rappel**

(et **interprétabilité** du filtre)

# Apprentissage prescriptif pour « intervenir »

# Apprentissage **prescriptif**

---

- Apprentissage « **prescriptif** » (recherche de *causalités*)

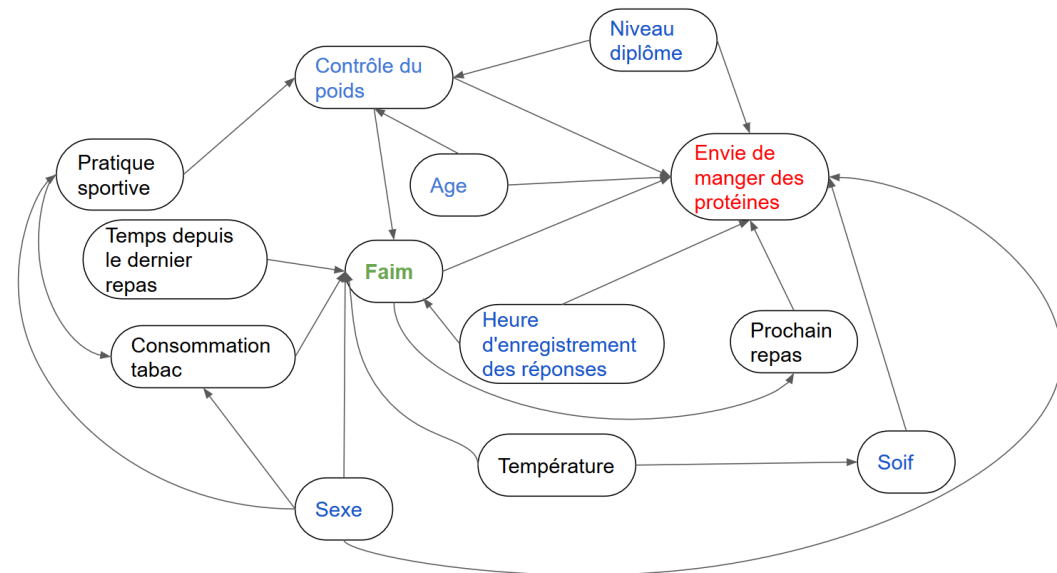
1. J'observe que les gens qui mangent des glaces sont souvent en maillot de bain
2. Je voudrais vendre davantage de glaces

→ Je demande aux gens de se mettre en maillot de bain

# La recherche de relations causales

Qu'est-ce qui **cause** l'appétence pour des plats protéinés ?

- La **faim** ?
- L'**heure** dans la journée ?
- Le **genre** ?
- L'**aspect visuel** ?
- L'**aspect olfactif** ?
- La richesse en **protéines** des **repas précédents** ?
- ...



- Quelles **recommandations** faire à un consommateur pour qu'il baisse sa consommation d'aliments carnés ?
- Quel impact **si on double le prix** de ... ?
- Quel rendement aurais-je eu l'année dernière **si j'avais** planté du ... au lieu de ...

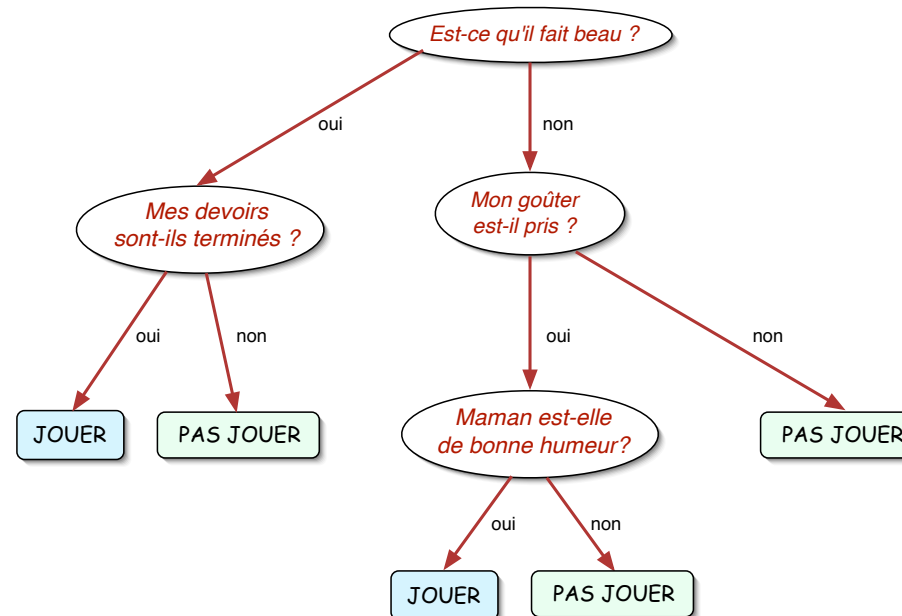
# Quels modèles ou hypothèses ?

# Modèles interprétables

- Régression linéaire

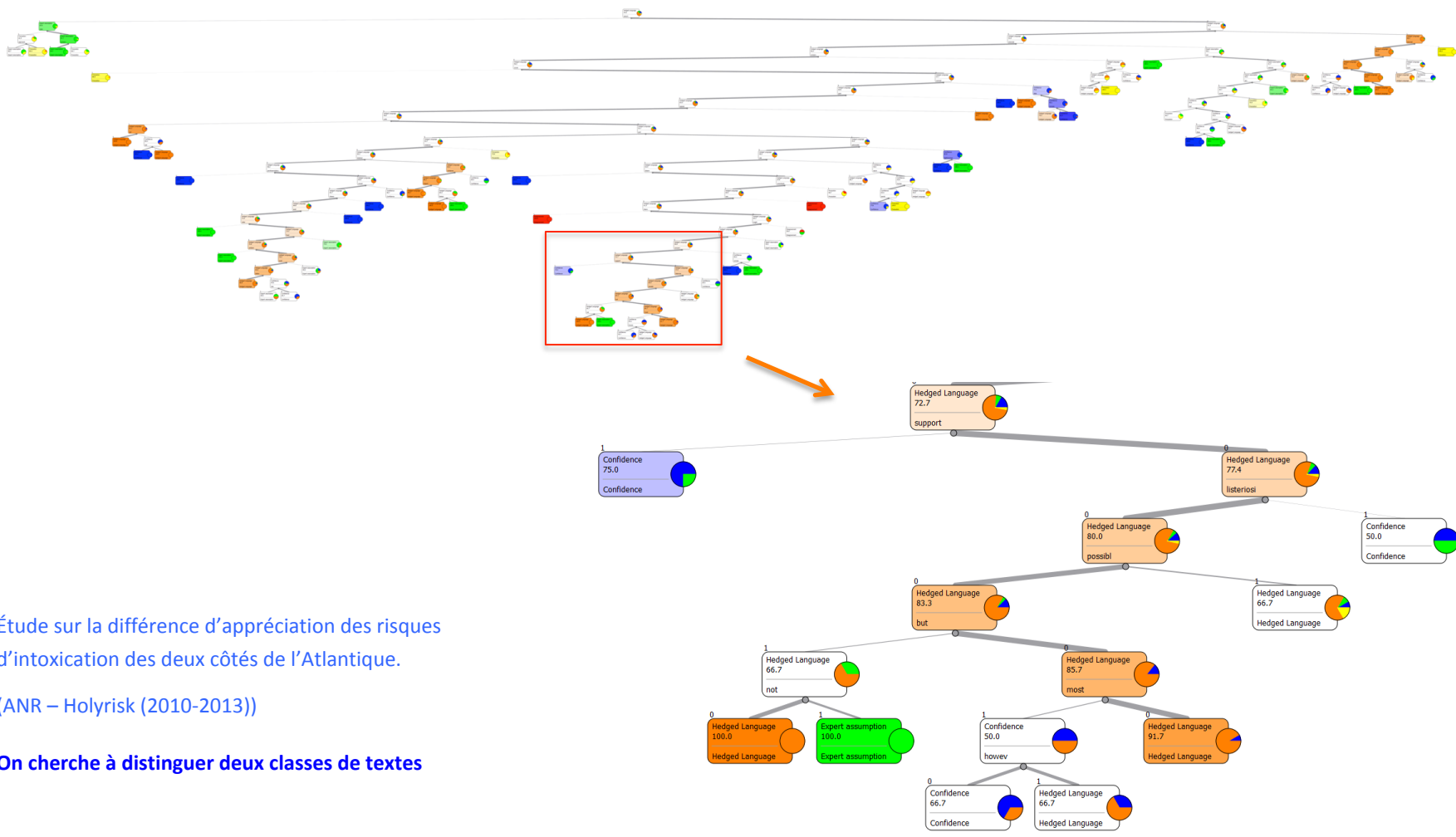
$$y = \sum_{i=1}^N \alpha_i x_i$$

- Arbre de décision





# Exemple : arbre de décision



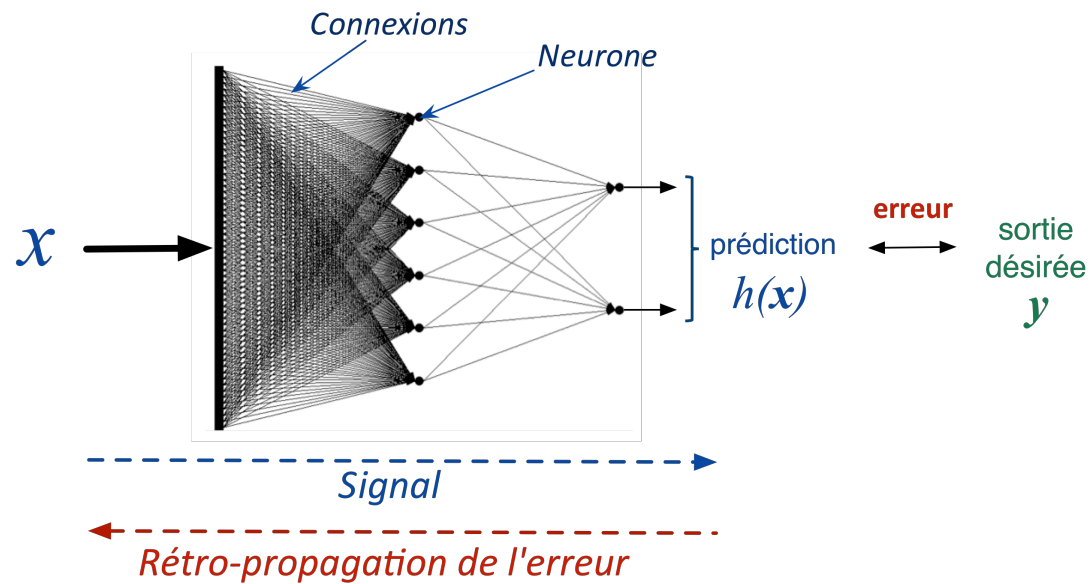
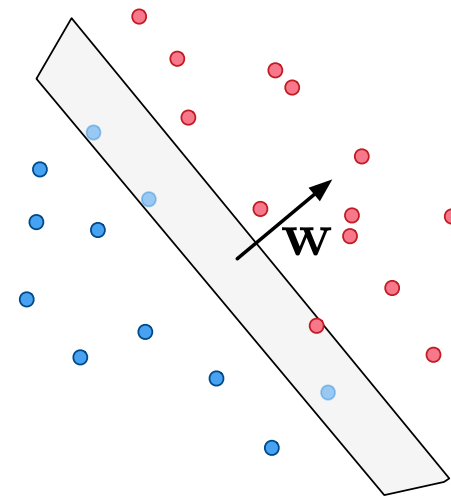
Étude sur la différence d'appréciation des risques  
d'intoxication des deux côtés de l'Atlantique.

(ANR – Holyrisk (2010-2013))

On cherche à distinguer deux classes de textes

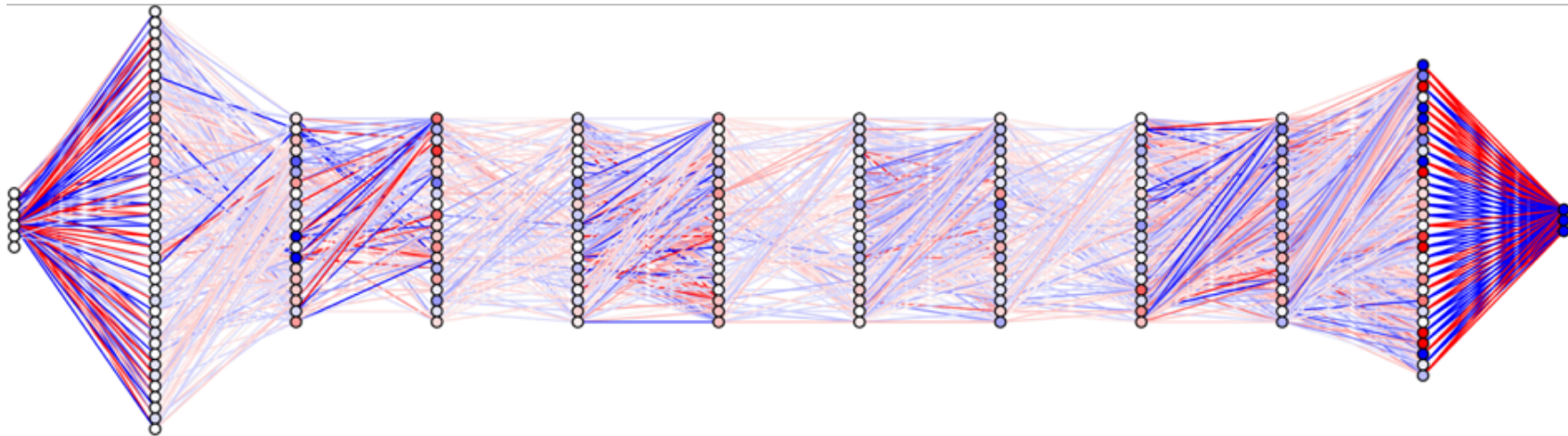
# Modèles opaques

- Régression linéaire généralisé
- Réseaux de neurones



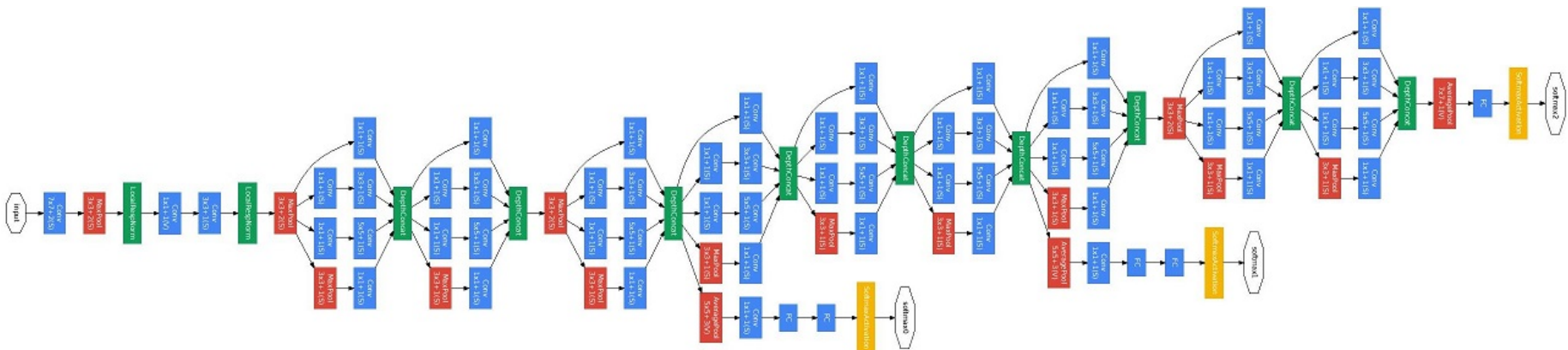
# Les « réseaux de neurones **profonds** »

- Des réseaux de neurones artificiels
  - à grand nombre de couches (parfois > qqs 100)
  - et **très grand nombre de paramètres** (qqs  $10^7 - 10^8$  paramètres)

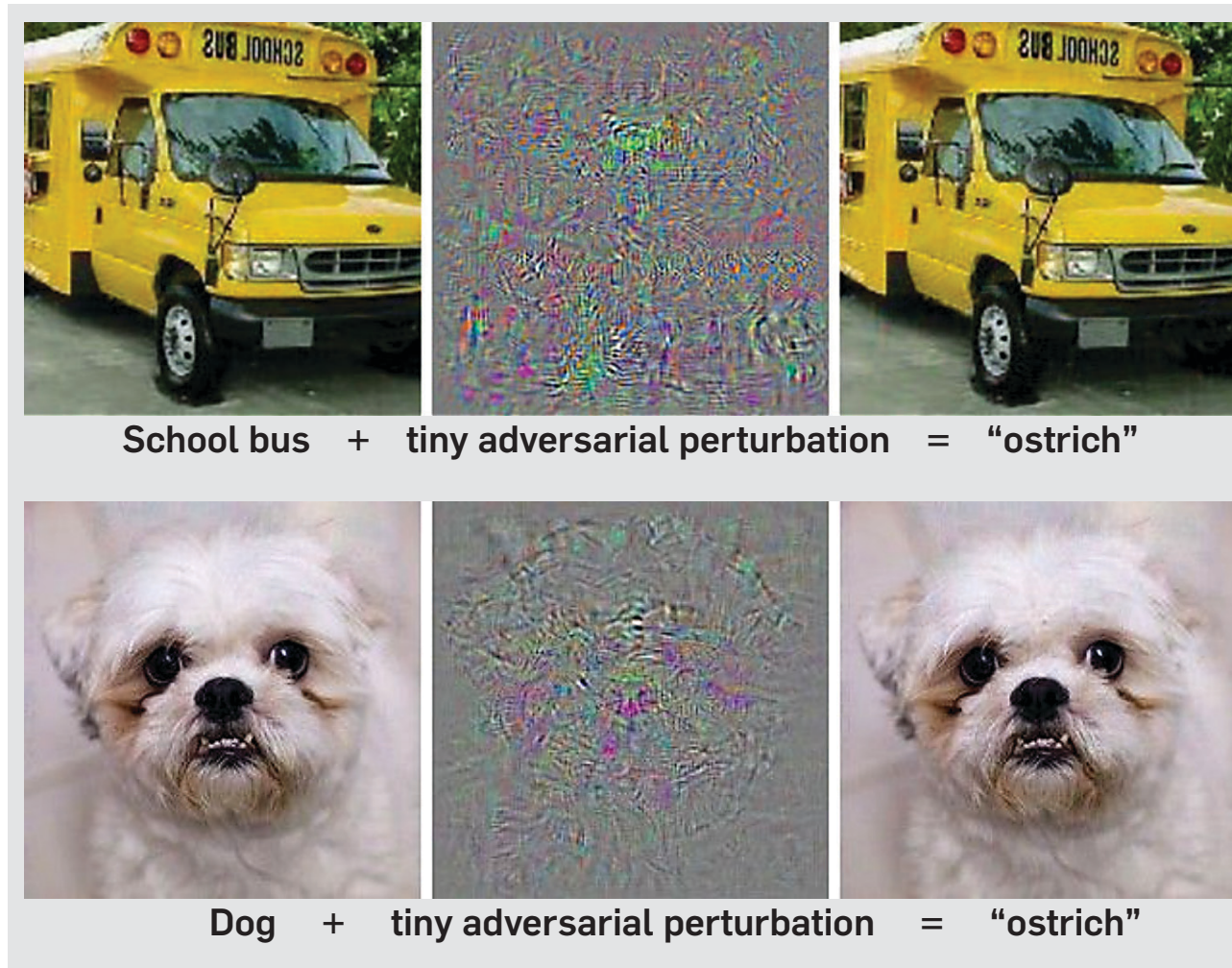


# GoogleNet

- Un **mécano** de réseaux de neurones



# Des **erreurs** difficiles à comprendre



**Adversarial input can fool a machine-learning algorithm into misperceiving images.**

Illustration

# Plan

---

1. La révolution des données
2. Qu'en fait-on ?
3. Les défis
4. Les compétences requises, les acteurs, les outils
5. L'avenir

# Les défis

---

## 1. Le recueil des données

# Obtenir les données

---

## Souvent **difficile** !!!

- Les données ne sont **pas encore disponibles**
- Le donneur d'ordre n'est **pas détenteur des données**
  - Pas le même service / département
- Les données sont **protégées par des droits**
- Une partie des données **reste à recueillir**



# Les questions juridiques

---

## Essentiel !!!

- Données **personnelles**
- **Obtenir l'autorisation**
  - CNIL
  - RGPD
    - Depuis le **25 mai 2018**, le Règlement Général Européen sur la Protection des Données (**RGPD**) affecte toutes les organisations traitant les **données personnelles identifiables (DPI)** de résidents européens.

# Les défis

---

1. Le **recueil** des données
2. Les **prétraitements** des données

# Les prétraitements

---

- **90%** du temps d'un projet
- **Recueil** des données
- Mise dans un **format adéquat**
- **Nettoyage**
  - **Bruit** dans les données
  - Données **manquantes**
  - Données **aberrantes**
  - **Doublons**
  - **Normalisation** des mesures
  - **Discrétisation** de valeurs continues
  - **Rendre continues** des valeurs discrètes
- Élimination des **attributs redondants** / calcul de **nouveaux attributs**
- **Précision / incertitude**
- Intégration de plusieurs **sources de données (hétérogènes)**
- ...

# Le traitement des documents en .pdf

- Structure des pages ... en .pdf

The image shows a screenshot of a PDF document page with several red arrows pointing to specific structural elements. The document is titled "Synthetic fertilizer management for China's cereal crops has reduced N<sub>2</sub>O emissions since the early 2000s" by Wenjuan Sun and Yao Huang. The page includes a header with the journal name "Environmental Pollution" and the Elsevier logo. The main content is divided into sections: "ARTICLE INFO", "ABSTRACT", "1. Introduction", and "2. Materials and methods". The "ARTICLE INFO" section contains metadata such as "Article history", "Received 3 May 2011", "Received in revised form 10 August 2011", and "Accepted 3 September 2011". The "ABSTRACT" section provides a summary of the study. The "1. Introduction" section discusses the impact of synthetic nitrogen fertilizer on N<sub>2</sub>O emissions. The "2. Materials and methods" section describes the data sources and analysis. Red arrows point to the following elements: the journal title, the article title, the authors' names, the "ARTICLE INFO" section, the "ABSTRACT" section, the "1. Introduction" section, the "2. Materials and methods" section, and the footer information.

# Les défis

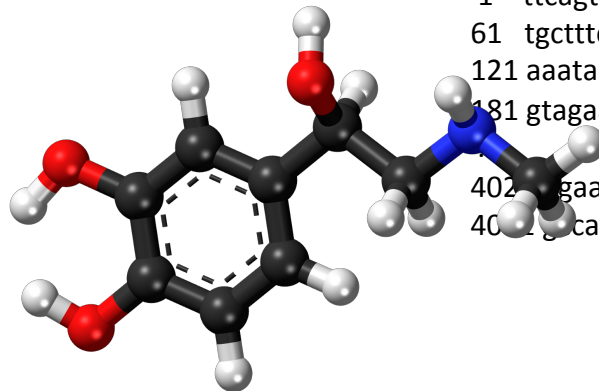
---

1. Le **recueil** des données
2. Les **prétraitements** des données
3. Les **sources multiples** et hétérogènes

# Intégration de **multiple sources** de données

- Annotation de protéines

Protéine « sp|P00004|CYC\_HORSE » is activated by ...



```
1 ttcagttgtg aatgaatgga cgtgccaaat agacgtgccg ccgccgctcg attcgactt
61 tgctttcggg ttgcccgtcg tttcacgcgt ttagttccgt tcggttcatt cccagttctt
121 aaataccgga cgtaaaaata cactctaacg gtcccgcgaa gaaaaagata aagacatctc
181 gtagaatat taaaataat tcctaaagtc gttggtttct cgttcacttt cgctgcctgc
402 ggaacagcc gaggtccat tcatagcacc acttcgctgt ctaatcccc tcctcatcc
403 gcatggcgg tgcaaaaaat aaaaagaact c
```

# Intégration de **multiple sources** de données

- **GIEC**

- Documents scientifiques multiples
- Tableaux
- mesures

Moore's Law has, for nigh half a century, reliably predicted the growth in efficiency of processors: Moore's Law states that the number of transistors that can be placed on a given surface area doubles every two years [Intel Corporation, 2003]. As a consequence, the number of transistors – and consequently, the computing power – of processors has grown exponentially until recently. However, this growth can no longer be sustained due to a combination of several factors. The most important cause are quantum mechanical effects which raise the electrical resistance of the transistors and thus cause heat dissipation problems which result in energy loss [Freyman, 1985; Tanenbaum, 1990].

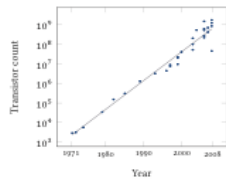


Figure 1: Moore's Law illustrated by the number of transistors of typical processors for each year. Note that the y axis is logarithmic.

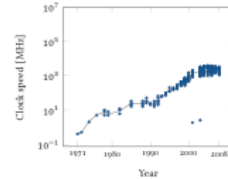


Figure 2: Clock speed (in MHz) of intel processors over the years and their mean values for each year.

On the other hand, we're dealing with ever increasing amounts of data that our grams have to process. Figure 3 illustrates this using the example of the number o

	MaxEnt			MaxEnt + GE			Unsup GE		
	P	R	F	P	R	F	P	R	F
BKG	.38	.19	.25	.49	.48	<b>.48</b>	.49	.44	.46
PROB	0	0	0	.38	.23	.29	.28	.38	<b>.32</b>
METH	0	0	0	.29	.50	<b>.37</b>	.08	.56	.14
RES	0	0	0	.68	.51	<b>.58</b>	.08	.51	.14
CON	.69	.96	.80	.81	.84	<b>.82</b>	.74	.69	.71
CN	.35	.06	.10	.39	.29	<b>.33</b>	.40	.13	.20
DIFF	0	0	0	.21	.30	<b>.25</b>	.12	.13	.12
FUT	0	0	0	.24	.44	.31	.26	.61	<b>.36</b>

## Document Ranking using Customizes Vector Method

**Priyanka Mesariya**  
Computer Engineering, Gujarat Technological University, India

**Nishi Madia**  
Computer Engineering, Gujarat Technological University, India

### ABSTRACT

Information retrieval (IR) system is about positioning reports utilizing client's question and get the important records from extensive dataset. Archive positioning is fundamentally looking the pertinent record as per their rank. Document ranking is basically search the relevant document according to their rank. Vector space model is traditional and widely applied information retrieval models to based on similarity values. Term are the significant of an inform and it is query used in docu ranked calculates the term using query on basis of term who documents. When user enter q documents in which the query is it will count the term calculate th the highest weight of value it v documents.

### KEYWORD

Information retrieval, term fi frequency, vector space model, C

### 1. INTRODUCTION

In the information retrieval (IR) are ranked optimally by using the relevant documents from lar dataset [1]. When the user gives consulted to archives the most the relevant documents are then of their degree of relevance. May rely on search engines for extra providing a query from any i queries are processed by the a certain information retrieval or applied to obtain the cluster of the query. After the retrieval of important task is to present this where documents at the top are more relevant for the user. This

©IJTSRD | May-Jun-2017  
Available Online @www.ijtsrd.com

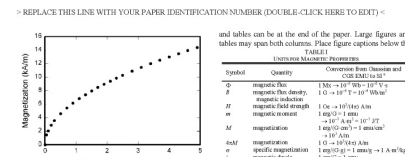


Fig. 3. Illustration of a document. This is a n x m matrix. n is good practice option.

E. Copyright Form as IEEE copyright submission. You can find the details of the IEEE-1318 are responsible for ob

If you are using # for equations in your Manuscript Equations or should not be selected.

Use either SI (MKS strongly encouraged) units (in parentheses storage. For example exception is when Eng such as "3% in disk units, such as carrier overach. This often is not "business dimension clearly state the units if the SI unit for mag

of documents [15] Information retrieval system is a set of documents to discover convenient information equivalent to a user's query. In information retrieval basically data can be fetching from web structure information that can be type of content, pictures, graph etc. Several components make this task challenging: (i) normally unstructured information is in document database, (ii) reports are typically composed in unstructured characteristics (dialect, mix)

REPLACE THIS LINE WITH YOUR PAPER IDENTIFICATION NUMBER (DOUBLE-CLICK HERE TO EDIT) <

and tables can be at the end of the paper. Large figures and tables may span both columns. Place figure captions below the TABLE

Applied Mechanics and Materials  
ISSN: 1662-7482, E-ISSN: 1688-4184  
doi:10.4028/www.scientific.net/AMM.543-547.4180  
© 2014 Trans Tech Publications, Switzerland

**Research and Improvement Strategies on Disaster Education for Primary and Secondary School**  
Yingqian Hu<sup>1,a</sup>, Man Zhang<sup>2,b</sup>  
<sup>1</sup> Jiangxi Science and Technology Normal University, Nanchang, Jiangxi, P.R.China.  
<sup>2</sup> School of Information Engineering, Nanchang University, Nanchang, Jiangxi, P.R.China.  
\*Email: 1328675451@qq.com; \*Email: manzhang201010@163.com

**Keywords:** Disaster Education; Primary and Secondary School; Strategies

**Abstract.** The frequent occurrence of disasters make people pay more attention on disaster education, but the situation of primary and secondary school on disaster education in China is not ideal. The paper verified the viewpoint from the analysis of documents on the theme retrieved through CNKI. The paper proposed the point above and proposed an improvement strategies model to improve the situation according to the analysis of the data collected for the paper.

### Introduction

China is one of the countries most affected by the natural disasters in the world. The frequently occurred disasters affect economic development and social stability of the country, causing a great economic losses and casualties. Table 1 is part of economic losses and casualties caused by disasters choose from China Statistical yearbook , 2011. Especially after the Wenchuan earthquake, experts and scholars in China begin to focus more attention on disaster education research, and have achieved some success. However, researches on primary and secondary school are in a low level contrast to disaster education to other groups.

Table 1. The economic losses and casualties caused by disasters

Year	Direct economic losses caused by earthquake (million)	Direct economic losses caused by natural and Oceanic disaster (billion)	Casualties caused by earthquake (frequency)	Casualties caused by disaster (frequency)
2000	1467.92	12.08	2855	79
2001	1484.49	10.01		401
2002	147.74	6.59	362	124
2003	4660.40	8.05	7465	128
2004	949.59	5.42	696	140
2005	2628.11	33.24	882	371
2006	799.62	21.85	229	492
2007	2019.22	8.84	422	161
2008	859495.94	20.61	446293	152
2009	2737.82	10.02	407	95
2010	23610.77	13.28	13795	137

Source: China Statistical yearbook, 2011  
Disaster education first introduced to the public of China was by two professors Wang Hong and Zongqun in the year 1996, but they were failed to give a definition of its concept. Even near 20 years past, scholars still haven't given a unified and standard definition of disaster education in China, but we can get a understanding of it by reading papers on disaster education of scholars from home and abroad. A definition widely accepted but not standard on Disaster Education by many researchers in China is defined as education on improving citizens' awareness and ability to cope

All rights reserved. No part of contents of this paper may be reproduced or transmitted in any form or by any means without the written permission of Trans Tech Publications, www.scientific.net (2017) 17376-17380.

# Les défis

---

1. Le **recueil** des données
2. Les **prétraitements** des données
3. Les **sources multiples** et hétérogènes
4. La possibilité de l'intervention de **l'expert**



# Disponibilité des experts métier

---

## Essentiel !!!

- **Comprendre** le problème
- Établir un **vocabulaire commun**
- **Évaluer** les résultats
- Orienter / **ré-orienter**
- **S'approprier** les résultats / assurer la suite

# Des algorithmes « transparent »

---

1. Dans lesquels on puisse « injecter » l'expertise humaine
2. Dont les résultats (modèles appris) soient interprétables

# Les défis

---

1. Le **recueil** des données
2. Les **prétraitements** des données
3. Les **sources multiples** et hétérogènes
4. La possibilité de l'intervention de **l'expert**
5. L'identification de **relations causales**
6. Les environnements **non stationnaires**
7. Un **génie logiciel** des systèmes apprenants

# Plan

---

1. La révolution des données
2. Qu'en fait-on ?
3. Les défis
4. Les **compétences** requises, les **acteurs**, les **outils**
5. L'avenir

# Les « data scientists »

---

- **Compétences attendues**

1. Apprentissage artificiel / Statistiques

- Bonne compréhension des questions et des hypothèses sur lesquelles reposent les méthodes

2. Compétences en informatique

- Algorithmique
- Bases de données
- Réseaux

3. Capacités relationnelles

# Les « data scientists »

- **Compétences attendues**

1. Apprentissage artificiel / Statistiques

- Bonne compréhension des questions et des hypothèses sur lesquelles reposent les méthodes

2. Compétences en informatique

- Algorithmique
- Bases de données
- Réseaux

3. Capacités relationnelles

**En très forte  
demande**

100 000 en France  
à l'horizon 2022 !!

- **Formations**

- Quelques dizaines d'heures
- **Master** ou équivalent
- **Doctorat**

**Grand risque** de déconvenue  
si pas les bons recrutements

# Les passages à l'échelle

## 1. Savoir traiter de (très) **gros volumes de données**

### — Méthodes efficaces

- Gradient stochastique
- Apprentissage convexe
- Optimisation du code
  - ✓ Accès mémoire
  - ✓ Complexité computationnelle

### — Distribution des calculs

- Cartes graphiques / cœurs
- Clusters de machines
- Cloud computing
  - ✓ Approches Map Reduce

# Les passages à l'échelle

---

## 2. Savoir traiter de (très) **petits** volumes de données

**Compenser** le manque d'information dans les données

- Par de la **connaissance experte**
- **Enrichissement** des données
  - Ontologies
  - Web sémantique
  - Wikipedia and Co
- Question de la **validation des résultats**
  - Les experts



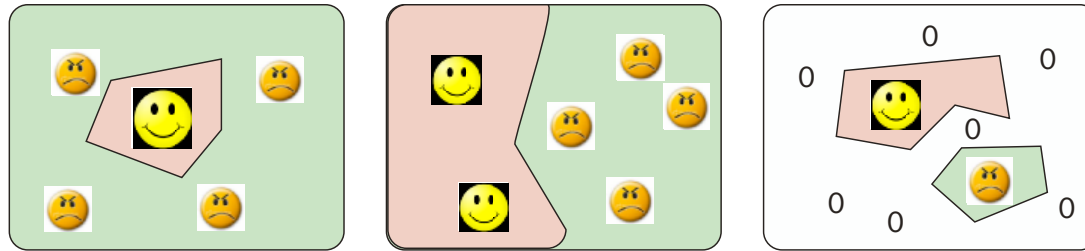
# Les méthodes et algorithmes

---

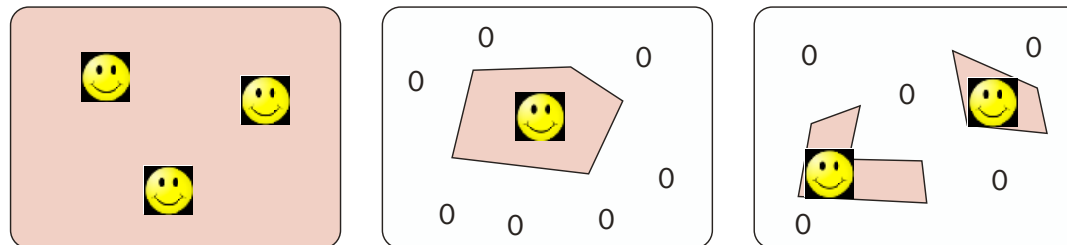
- Bibliothèques / méthodes / algorithmes
  - Sont dans le **domaine public !!!**
    - Publications scientifiques
    - Forums
    - Conférences
    - Bibliothèques (e.g. ScikitLearn)
- Des « **recettes** » privées
  - Réseaux de neurones profonds
  - Traitement d'images / télédétection
  - Connaissances métiers (e.g. alimentation)

# Le no-free-lunch theorem

Possible



Impossible



Il faut **choisir** le **bon** **algorithme** pour la **classe de problèmes** étudiée

# Plan

---

1. La révolution des données
2. Qu'en fait-on ?
3. Les défis
4. Les compétences requises, les acteurs, les outils
5. L'avenir

# Une liste ...

---

## 1. Résoudre les difficultés

- Données **multi-sources** hétérogènes
- **Dialogue** possible avec les **experts** : interprétabilité des modèles produits, compréhension et contrôle raisonné des algorithmes

## 2. Identification de **relations causales**

## 3. Apprendre à partir de (très) **peu d'exemples**

## 4. Apprendre en environnement **non stationnaire**

- Flux de données
- Transfert entre tâches

## 5. **Génie logiciel** pour des **systèmes apprenants**

## Une révolution en cours

1. Tirer profit des données
  - ✓ Numérisation
  - ✓ Capteurs partout
  - ✓ Internet
  - ✓ Des ressources calcul
  - ✓ Des algorithmes
2. Gros progrès en intelligence artificielle

Mais ce n'est pas « magique »

---

Beaucoup d'opportunités

Mais pas de magie



# Conclusions



# 4 approches pour appréhender le monde

---

## 1. Empirique : description et classement



# 4 approches pour appréhender le monde

1. Empirique : description et classement



2. Théorique : Modélisation, construction de théories

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$

# 4 approches pour appréhender le monde

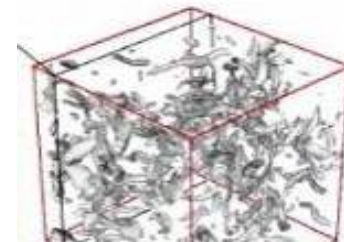
1. **Empirique** : description et classement



2. **Théorique** : Modélisation, construction de théories

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$

3. **Simulation** : systèmes complexes et/ou non reproductibles



# 4 approches pour appréhender le monde

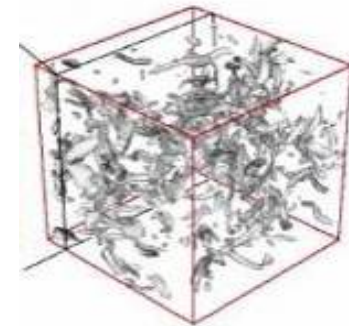
1. Empirique : description et classement



2. Théorique : Modélisation, construction de théories

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K\frac{c^2}{a^2}$$

3. Simulation : systèmes complexes et/ou non reproductibles



4. Exploration de données

- Énormes masses de données numérisées
- Largement disponibles
- Sources et formats très différents



# Le cas AlphaGo

- Un joueur « **extraterrestre** »
- Un jeu **stupéfiant**
- **Révolutionne** la manière de jouer
- **Effervescence** dans les écoles de go



# Le cas AlphaGo : comprendre

Fan Hui, Gu Li, Zhou Ruyang (très forts joueurs de Go) se reconvertissent dans l'analyse des parties jouées par AlphaGo

- Sorte d'exégèse. Explications a posteriori
- Nécessaire pour
  - La communication
  - L'enseignement

Et même AlphaGo peut se tromper

