

Big Data et Intelligence Artificielle: panorama et quelques zooms

Marie-Christine Rousset

LIG (Laboratoire d'Informatique de Grenoble)

Responsable du Labex PERSYVAL-lab

Université de Grenoble-Alpes et Institut Universitaire de France



1

Qu'est-ce que l'IA ?

- Une branche de l'Informatique, dont le but est :
la **représentation**, l'**extraction** et l'**exploitation** de
connaissances interprétables et **manipulables**
dans des **formalismes à dominante symbolique**,
en vue d'**aider** à la résolution de **problèmes de**
décision, planification, conception, diagnostic

2

Ambition de l'IA

- Rendre la machine capable :
 - de **raisonner** sur une situation *statique* ou *dynamique*
 - pour faire un *diagnostic*, proposer une *décision*, un *plan d'action*, résoudre une *tâche complexe*
 - d'**expliquer** et de communiquer ses conclusions
 - d'**abstraire**, d'**apprendre**, de **découvrir** à partir de données
- Par des méthodes **génériques** susceptibles de s'adapter à de larges classes de situations.

3

Problématiques fondamentales

- Représentation de connaissances et de raisonnements
 - **Systemes-experts: les premiers systemes d'IA (70-80)**
XCON: premier SE industriel, DEC, configuration d'ordinateurs
 - Mise à jour et fusion d'informations
 - Problèmes de décision
 - Fouille de données
 - Apprentissage automatique
 - Algorithmes génériques de résolution de problèmes
- ⇒ **Nouvel essor avec l'émergence du Big Data**
(disponibilité de données massives et hétérogènes)

4

Les thèmes d'applications

- Les jeux
- Planification d'actions
- Ordonnancement
- Diagnostic
- Les systèmes multi-agents
- La robotique
- Le web sémantique
- Le e-learning

5

Quelques zooms

- Web sémantique
 - **Intégration de données via des ontologies**
- Raisonnement par contraintes
 - **Ordonnancement/affectation de ressources sous contraintes**
- Raisonnement incertain
 - **Réseaux Bayésiens**
- Apprentissage automatique
 - Fouille de données (approche symbolique)
 - **Recherche de patterns et de règles d'association**
 - Classification automatique (approches numériques)
 - **Réseaux de Neurones**

6

Web sémantique

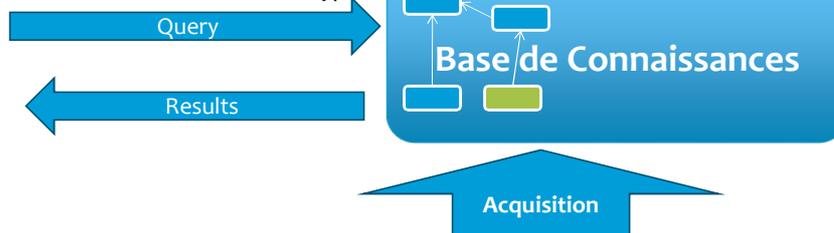
- Un ensemble de technologies et de normes recommandées par le W3C
 - RDF, RDFS, OWL, SPARQL
- Raisonnement à base d'ontologies sur des méta-données (RDF) de ressources du Web
- Application à d'autres domaines que le web
 - Intégration sémantique de données hétérogènes
 - Ontology-based Data Access

7

Projet ANR CONTINUUM (2008-2012)

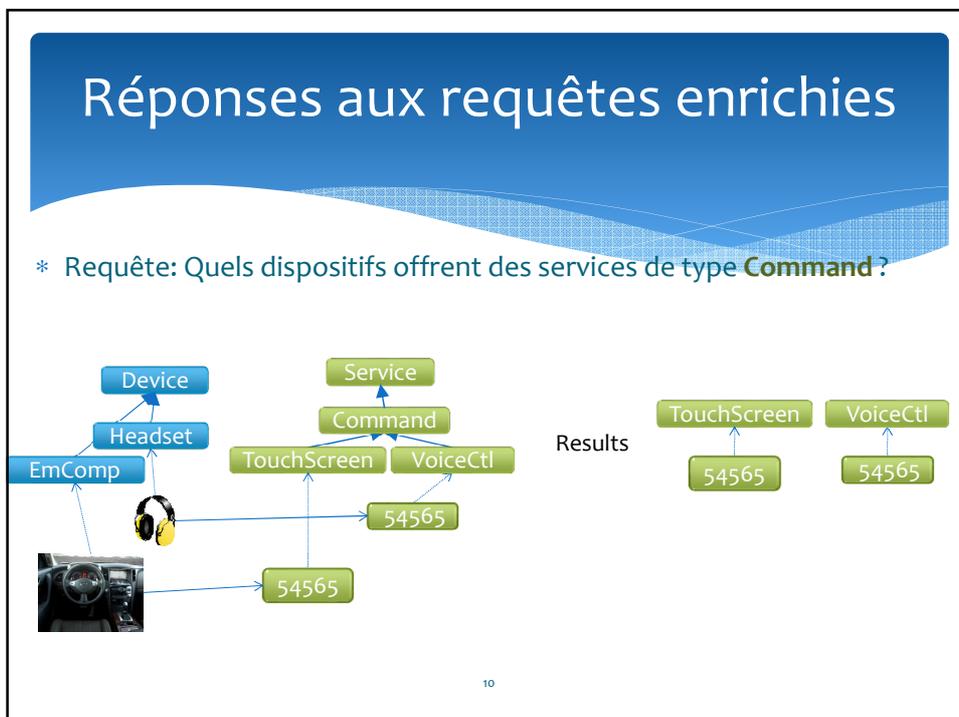
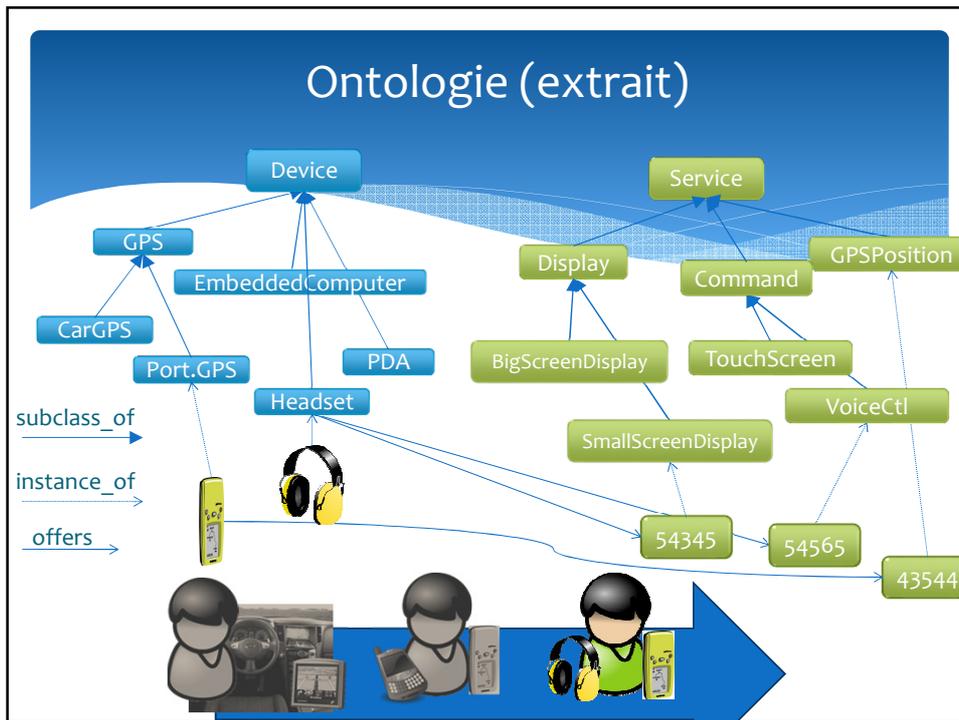
CONTinuité de service en Informatique UbiqUitaire et Mobile

Trouver des dispositifs de l'environnement qui offrent des services d'un certain type



Environnement

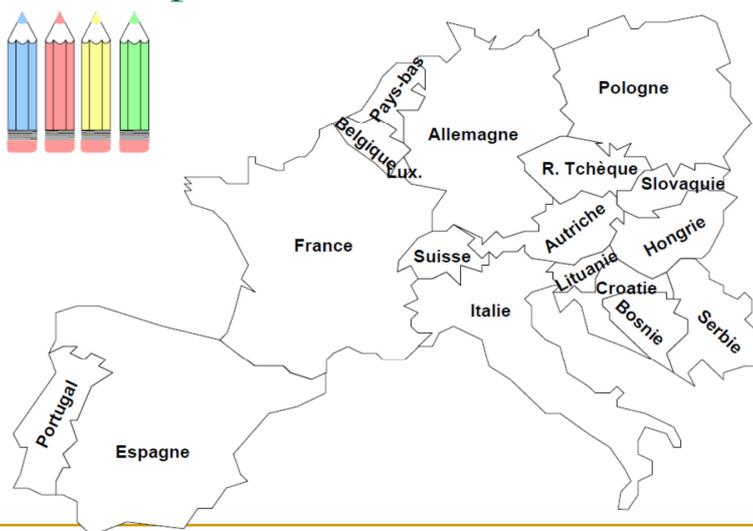




Raisonner par contraintes

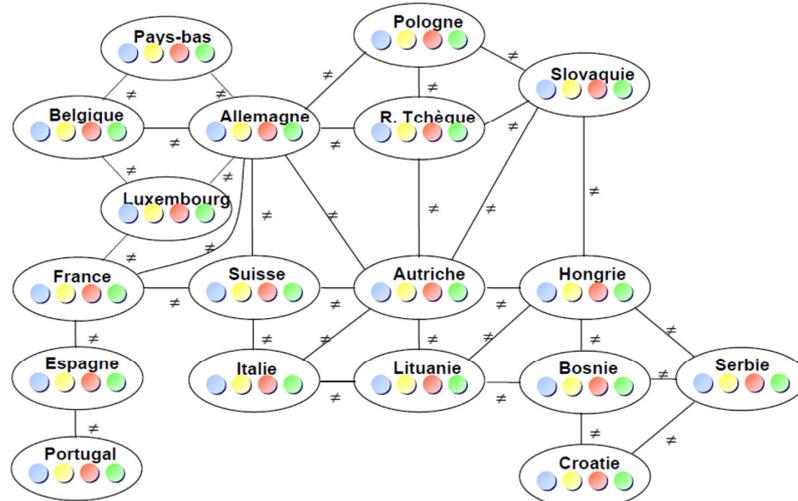
- Tâche centrale
 - Résoudre des problèmes combinatoires modélisés par un réseau de contraintes (booléennes ou numériques)
 - L'utilisateur modélise ses contraintes dans un formalisme déclaratif
 - Un ensemble fini de variables entières et leurs domaines
 - Exemple en ordonnancement: dates de démarrage des tâches
 - Un ensemble de contraintes sur ces variables
 - dates de disponibilité, de livraison, précédences et durées des tâches
 - Un algorithme CSP prend en entrée le réseau de contraintes et trouve (si elle existe) une instanciation des variables qui satisfait les contraintes
 - de coût minimal si une fonction de coût est définie
- Depuis 10 ans: énormes progrès des solveurs SAT et CSP (heuristiques, parallélisme, machine learning)¹¹
 - ⇒ passage à l'échelle de problèmes industriels

Un exemple : coloration d'une carte I



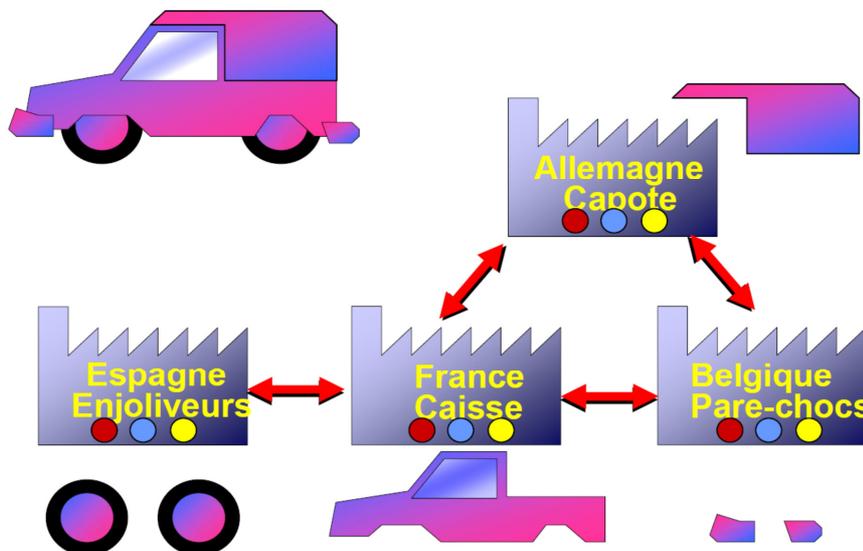
Exemple tiré d'un cours de Ruslan Sadykov ¹²

Un exemple : coloration d'une carte II

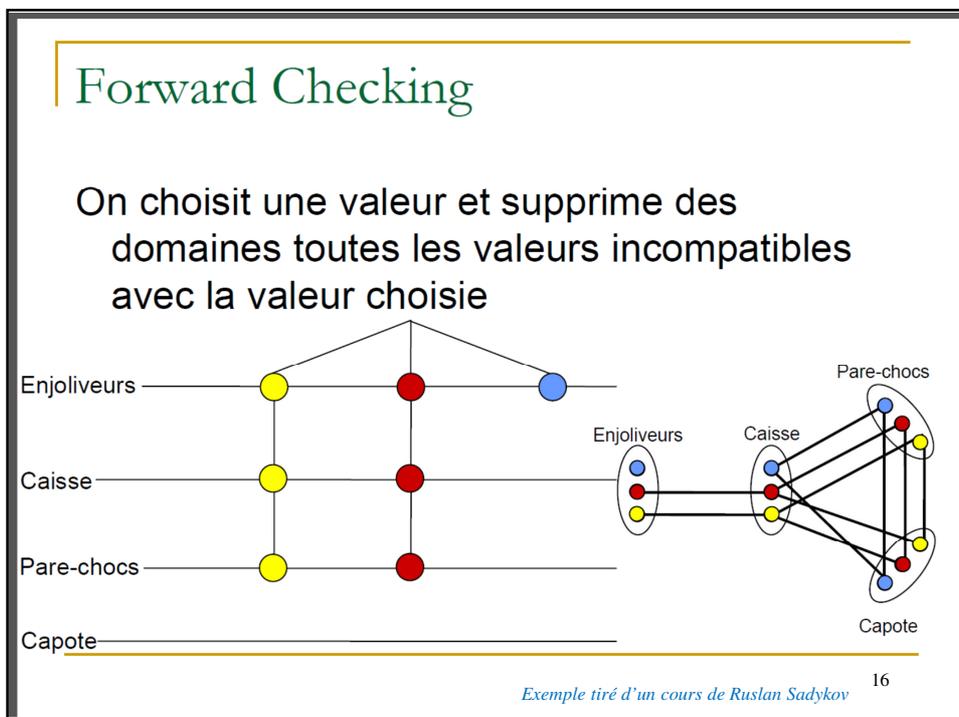
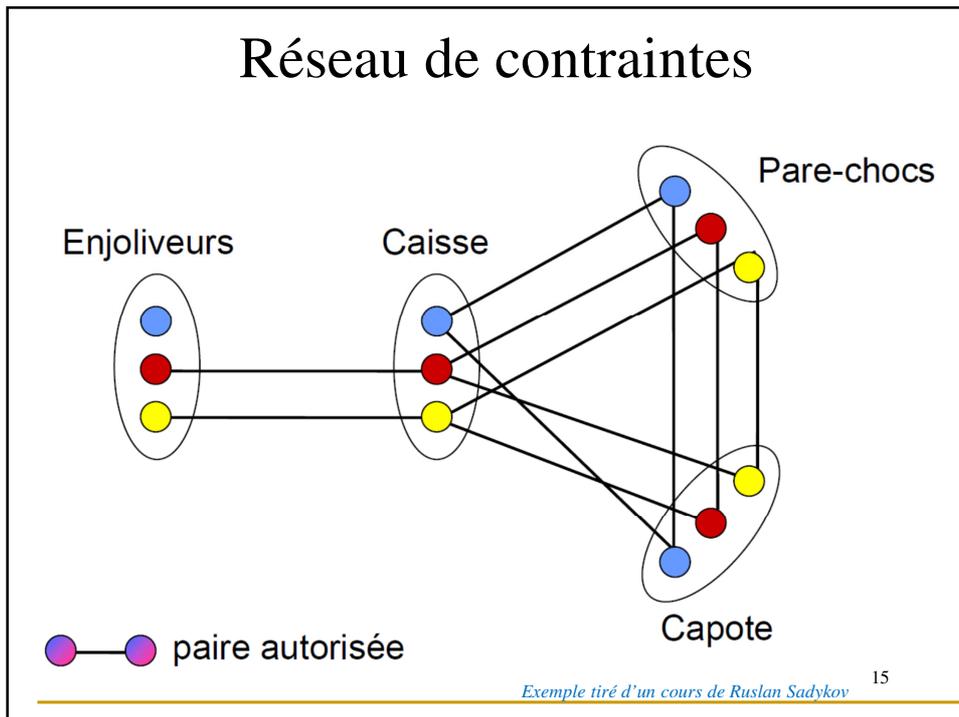


Exemple tiré d'un cours de Ruslan Sadykov 13

Exemple simpliste de conception/production



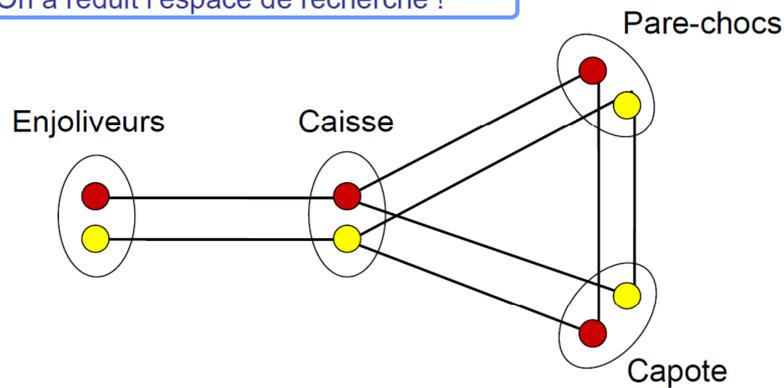
Exemple tiré d'un cours de Ruslan Sadykov 14



Consistance locale I

On n'a pas changé l'ensemble des solutions :
On a un réseau de contraintes **équivalent**

On a réduit l'espace de recherche !



Exemple tiré d'un cours de Ruslan Sadykov 17

Raisonnement dans l'incertain

- Tâche centrale
 - Elaborer des croyances sur une situation à partir de connaissances génériques sur le monde et d'observations.
- Beaucoup d'informations sont incertaines et/ou non Booléennes
 - Les connaissances génériques ont des exceptions
 - Les informations sur la situation courante ne sont pas forcément fiables
 - Les mots de la langue naturelle sont peu précis
- La logique classique ne suffit pas.

18

Réseaux Bayésiens

- Représentation compacte de distributions de probabilités
 - toute probabilité jointe positive est factorisable par un produit de probabilités conditionnelles formant un **graphe acyclique orienté** (DAG) visualisant des relations d'indépendance conditionnelle

Robot modélisé par 4 variables booléennes:

B: batterie ok, **M:** le bras bouge

L: le bloc peut être soulevé

G: la jauge indique que la batterie est chargée

Structure du réseau (modélisation qualitative)

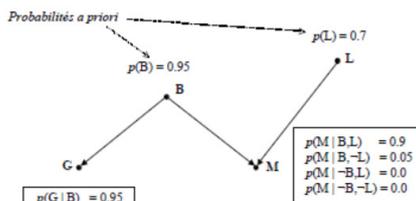
Arcs : influence/causalité directe

Une variable est **conditionnellement indépendante** de ses **non-descendants** étant donné ses **parents**

- L et B sont conditionnellement indépendants
- L et G sont conditionnellement indépendants

+ Modèle probabiliste local:

Distribution de probabilités conditionnelles pour chaque nœud sachant ses parents



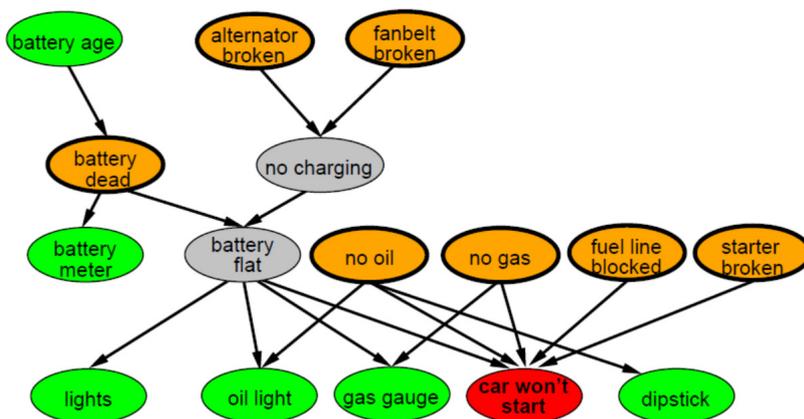
$$p(G,B,M,L) = p(G|B)p(M|B,L)p(B)p(L)$$

Inférences de prédiction: $p(M|L) = p(M,B|L) + p(M,-B|L) = p(M|B,L)p(B|L) + p(M|-B,L)p(-B|L)$
 $= p(M|B,L)p(B) + p(M|-B,L)p(-B)$ (L et B indépendants) = 0.1855

Inférences de diagnostic: $p(-L | \neg M)$

Exemple tiré d'un cours d'Antoine Cornuejols

Autre exemple



Variable observée: en rouge

Variables testables: en vert

Variables de diagnostic: en orange

Variables cachées: en gris

⇒ calcul de $p(\text{Diagnostic} | \text{Observations})$

20

Exemple tiré d'un cours de Stuart Russel

Apprentissage automatique

- Problématique du « Machine Learning »
 - Découvrir des régularités ou des corrélations dans un ensemble de données et les utiliser pour mieux prédire sur de nouveaux cas
- Problématique reliée : Fouille de données
 - Recherche d'informations remarquables dans une base de données
 - Résumé (de masses) d'informations
- Deux visions
 - Apprentissage *supervisé* de connaissances génériques à partir d'ensembles d'exemples et de contre-exemples, ou plus généralement d'exemples étiquetés.
 - Apprentissage *non supervisé*: regroupement de données semblables en « clusters »

21

Apprentissage automatique

- Spécificités de l'approche IA du Machine Learning
 - approches *symboliques*: pour leur interprétabilité
 - apprentissage *relationnel*: programmation logique inductive
 - apprentissage de *structures* : réseaux bayésiens...
- Évolution récente vers des approches numériques
 - Apprentissage statistique: apprendre la distribution de probabilités des données à partir d'un échantillon
 - Renouveau des approches de type réseaux de neurones avec le Deep Learning

⇒ Voir les 2 exposés suivants

22

Fouille de données

Recherche d'itemsets fréquents (pour l'analyse de tickets de caisse)

TID	Pain	Lait	Couches	Bière	Oeufs	Coca
1	1	1	0	0	0	0
2	1	0	1	1	1	0
3	0	1	1	1	0	1
4	1	1	1	1	0	0
5	1	1	1	0	0	1

Itemsets fréquents: {Pain, Lait}, {Bière, Couches, Lait}, ...

Règles d'association: Couches, Lait => Bière (+ calcul du support et de la confiance)

=> Aide au marketing et à la gestion de magasins

23

Réseaux de Neurones

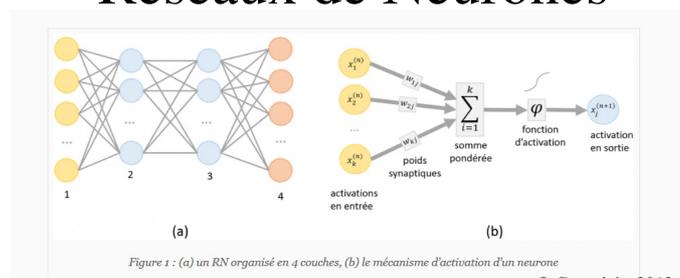


Figure 1 : (a) un RN organisé en 4 couches, (b) le mécanisme d'activation d'un neurone

© Copyright 2012 - 2015 [SQLI](#)

Applications phares:

- Classification/reconnaissance d'images
- Reconnaissance de chiffres manuscrits

Principes:

- Entraînement du réseau sur un ensemble d'images d'entraînement
- ⇒ Trouver la meilleure approximation possibles des poids w_{ij} pour calculer en sortie: $p(\text{labellimage})$

Trade-off:

- La qualité et la robustesse augmentent avec le nombre de couches
- Le temps d'entraînement croit rapidement avec le nombre de couches

24