

Approches prédictives dans le digital

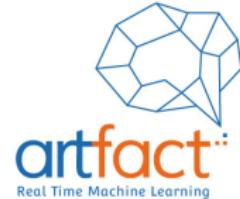
Online and Unsupervised Learning

Sébastien Loustau, CEO

Chaire DAMI, le 28 Mars 2017



Data Analytics & Models for insurance



Contents

Machine Learning : Gentle start

Cas d'application : segmentation clients temps réel

Cas d'application : détection de communautés temps réel

Contents

Machine Learning : Gentle start

Cas d'application : segmentation clients temps réel

Cas d'application : détection de communautés temps réel

Machine learning

$x \longrightarrow$ **nature** $\longrightarrow y$



artfact
Real Time Machine Learning

Machine learning

$x \longrightarrow$ **nature** $\longrightarrow y$

- ▶ prédire la réponse y à partir de x ,



artfact
Real Time Machine Learning

Machine learning

$x \longrightarrow$ nature $\longrightarrow y$

- ▶ prédire la réponse y à partir de x ,
- ▶ comprendre le lien entre x et y .

Machine learning

$$x \longrightarrow \boxed{\text{nature}} \longrightarrow y$$

- ▶ prédire la réponse y à partir de x ,
- ▶ comprendre le lien entre x et y .

$$x \longrightarrow \boxed{\text{algorithm}} \longrightarrow \hat{y}$$


artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$.



artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$. New x arrives. We build a model/algorithm thanks to \mathcal{D}_n and predict \hat{y} .



artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$. New x arrives. We build a model/algorithm thanks to \mathcal{D}_n and predict \hat{y} .

Online Learning

Data arrives sequentially.



artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$. New x arrives. We build a model/algorithm thanks to \mathcal{D}_n and predict \hat{y} .

Online Learning

Data arrives sequentially. At each time t , we want to make a decision based on past observations.



artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$. New x arrives. We build a model/algorithm thanks to \mathcal{D}_n and predict \hat{y} .

Online Learning

Data arrives sequentially. At each time t , we want to make a decision based on past observations. **No assumption over the data mechanism.**



artfact
Real Time Machine Learning

Statistical Learning vs Online Learning

Statistical Learning

We observe a training set $\mathcal{D}_n = \{(X_i, Y_i), i = 1, \dots, n\}$. New x arrives. We build a model/algorithm thanks to \mathcal{D}_n and predict \hat{y} .

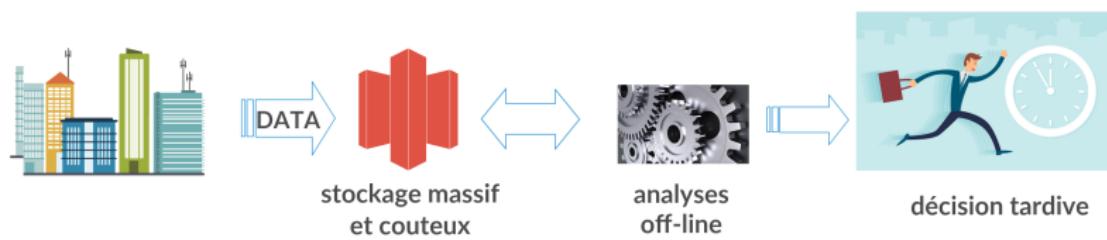
Online Learning

Data arrives sequentially. At each time t , we want to make a decision based on past observations. **No assumption over the data mechanism.**



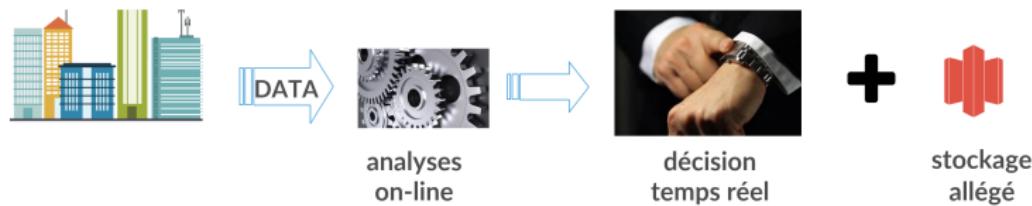
artfact
Real Time Machine Learning

Real time ML : un enjeu de taille



artfact
Real Time Machine Learning

Real time ML : un enjeu de taille



artfact
Real Time Machine Learning

Supervised learning

$$x \longrightarrow \boxed{\text{nature}} \longrightarrow y$$

- ▶ prédire la réponse y à partir de x ,
- ▶ comprendre le lien entre x et y .

$$x \longrightarrow \boxed{\text{algorithm}} \longrightarrow \hat{y}$$


artfact
Real Time Machine Learning

Today : unsupervised learning

nature $\longrightarrow x$



artfact
Real Time Machine Learning

Today : unsupervised learning

nature $\longrightarrow x$

- ▶ décrire les observations x ,



artfact
Real Time Machine Learning

Today : unsupervised learning

nature $\longrightarrow x$

- ▶ décrire les observations x ,
- ▶ réduire la dimension de x ,

Today : unsupervised learning

nature $\longrightarrow x$

- ▶ décrire les observations x ,
- ▶ réduire la dimension de x ,
- ▶ grouper les observations x .

Today : unsupervised learning

nature $\longrightarrow x$

- ▶ décrire les observations x ,
- ▶ réduire la dimension de x ,
- ▶ grouper les observations x .

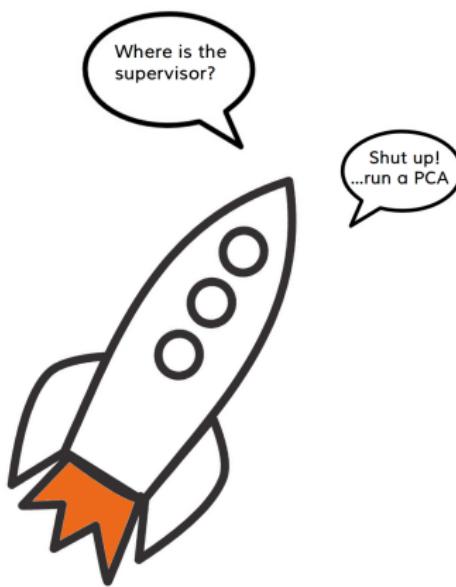
Vanilla algorithms and recent hot topics

PCA, k -means, gaussian mixtures, change points or more recently words embeddings, community detection, generative models (GAN)...



artfact
Real Time Machine Learning

Qu'est-ce que ça change ?



UNSUPERVISED LEARNING

Supervised routine

Supervised routine

$x \longrightarrow$ **algorithm** $\longrightarrow \hat{y}$



artfact
Real Time Machine Learning

Supervised routine

$$x \longrightarrow \boxed{\text{algorithm}} \longrightarrow \hat{y}$$

- ▶ Build a model from a training sample $\{(X_i, Y_i), i = 1, \dots, n\}$.



artfact
Real Time Machine Learning

Supervised routine

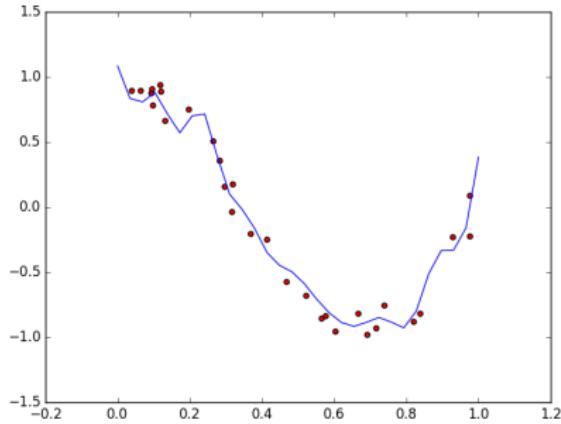
$$x \longrightarrow \boxed{\text{algorithm}} \longrightarrow \hat{y}$$

- ▶ Build a model from a training sample $\{(X_i, Y_i), i = 1, \dots, n\}$.
- ▶ Observe new x and predict \hat{y} as above.

Supervised routine

$x \longrightarrow$ **algorithm** $\longrightarrow \hat{y}$

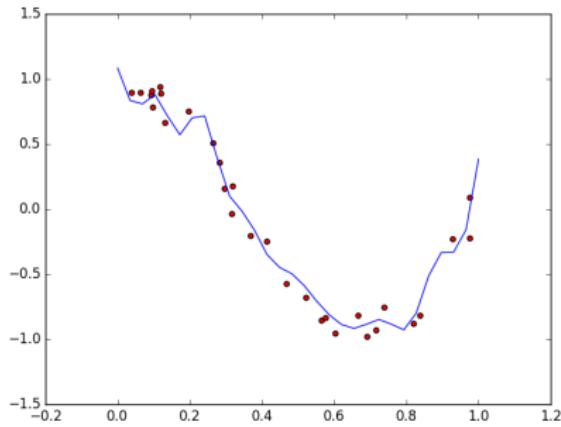
- ▶ Build a model from a training sample $\{(X_i, Y_i), i = 1, \dots, n\}$.
- ▶ Observe new x and predict \hat{y} as above.
- ▶ Problem : overfitting !



Supervised routine

$x \longrightarrow$ **algorithm** $\longrightarrow \hat{y}$

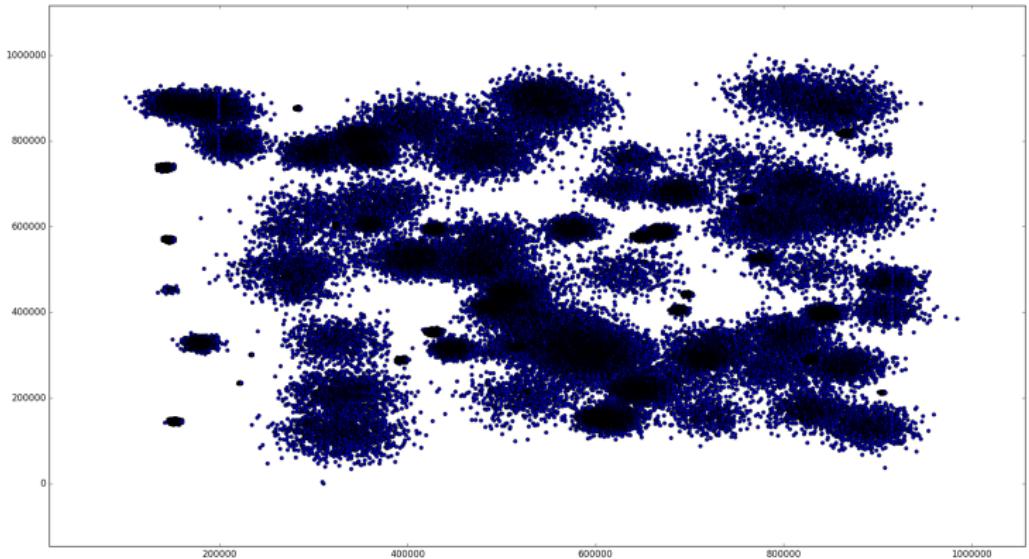
- ▶ Build a model from a training sample $\{(X_i, Y_i), i = 1, \dots, n\}$.
- ▶ Observe new x and predict \hat{y} as above.
- ▶ Problem : overfitting !



artfact
Real Time Machine Learning

- ▶ Solution : Training set + test set, Leave-One-Out, V-fold Cross validation.

Unsupervised : science or art ?



Unsupervised : the hype



artfact
Real Time Machine Learning

Contents

Machine Learning : Gentle start

Cas d'application : segmentation clients temps réel

Cas d'application : détection de communautés temps réel



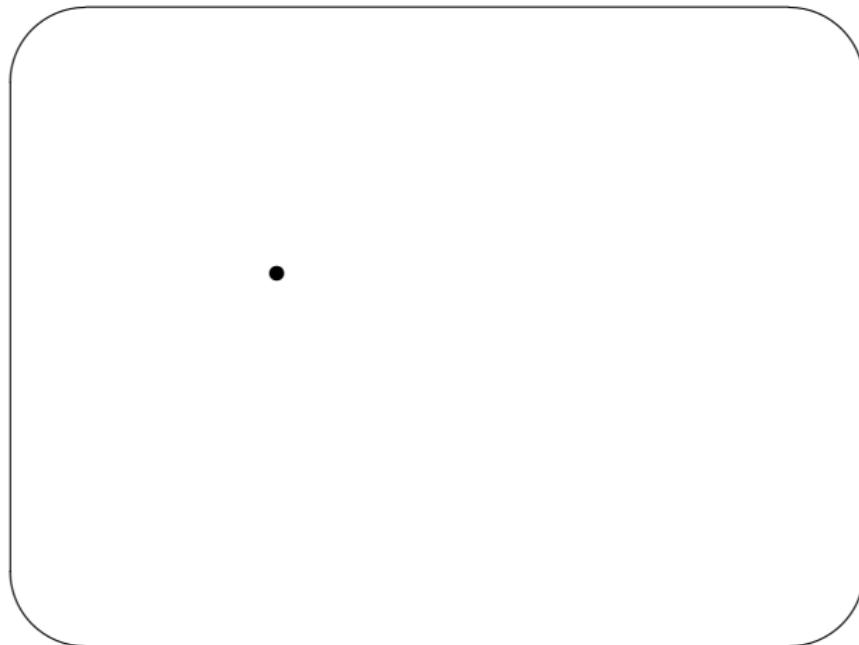
artfact
Real Time Machine Learning

The problem of Online Clustering



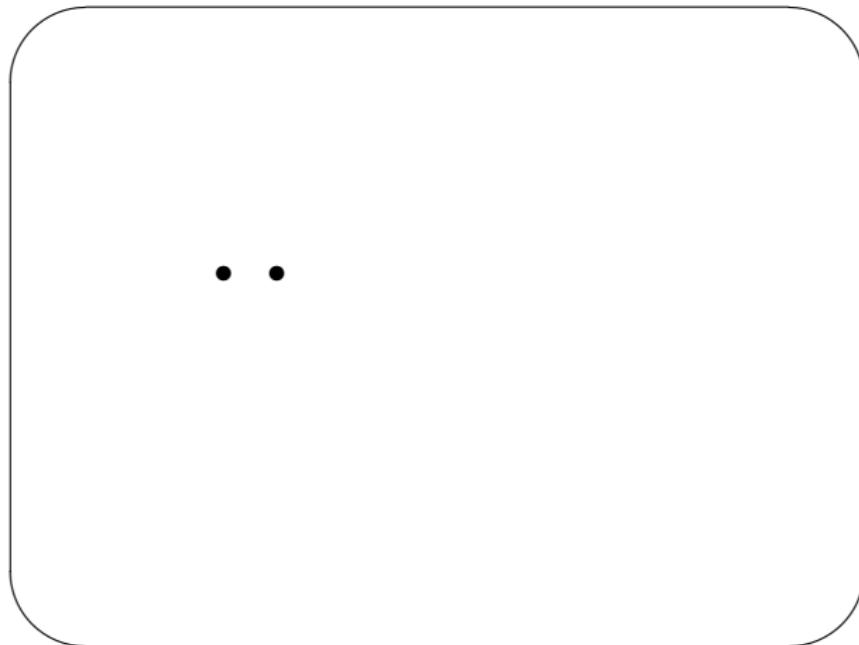
artfact
Real Time Machine Learning

The problem of Online Clustering



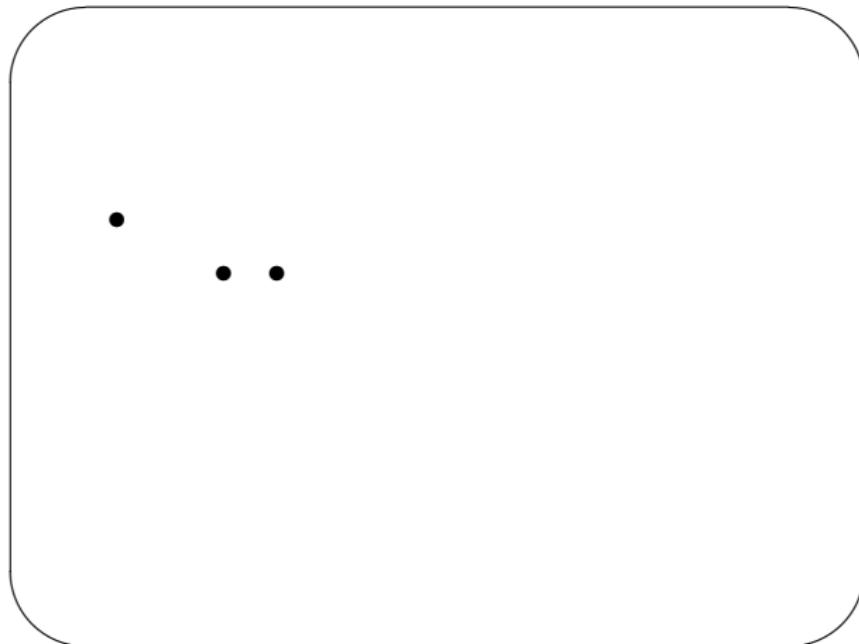
artfact
Real Time Machine Learning

The problem of Online Clustering



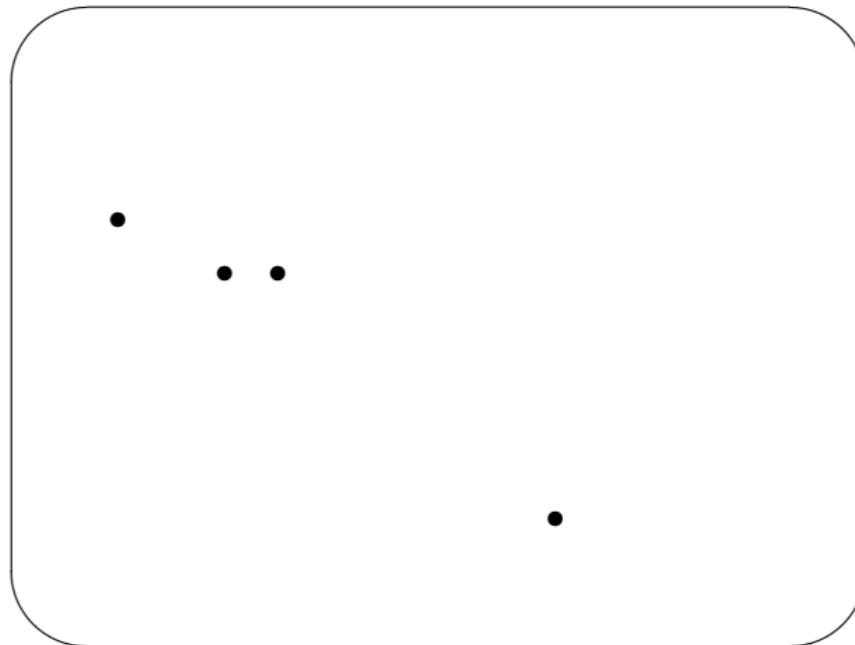
artfact
Real Time Machine Learning

The problem of Online Clustering



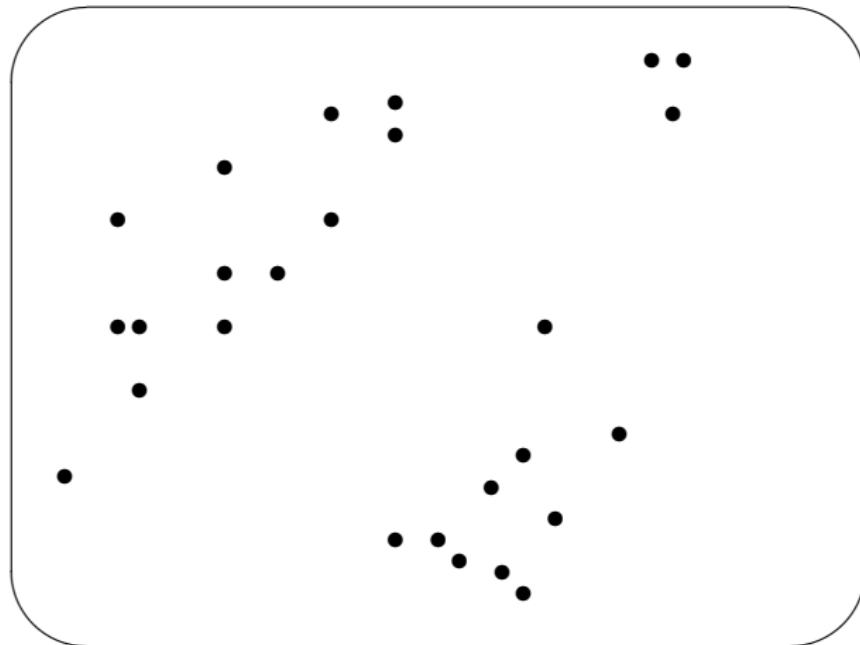
artfact
Real Time Machine Learning

The problem of Online Clustering



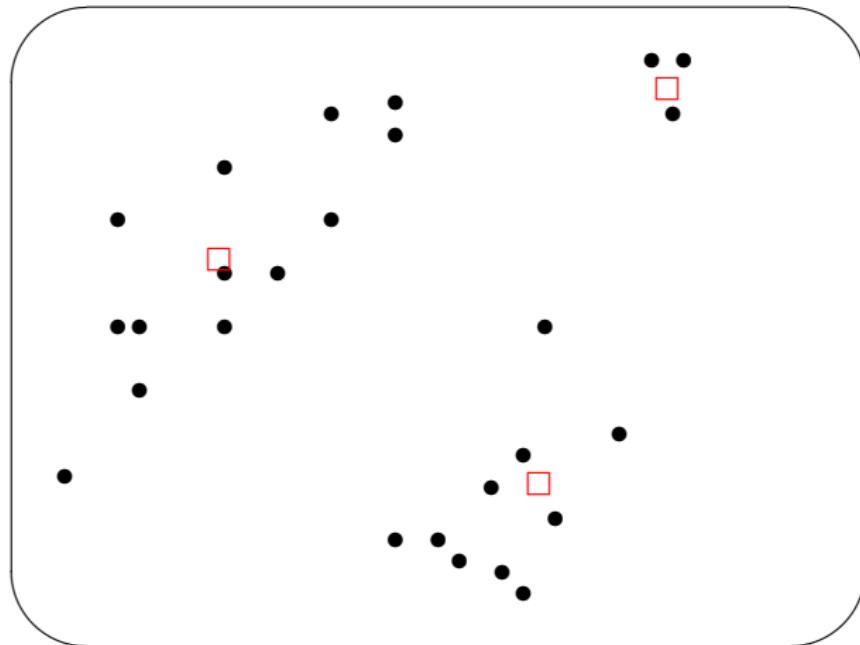
artfact
Real Time Machine Learning

The problem of Online Clustering



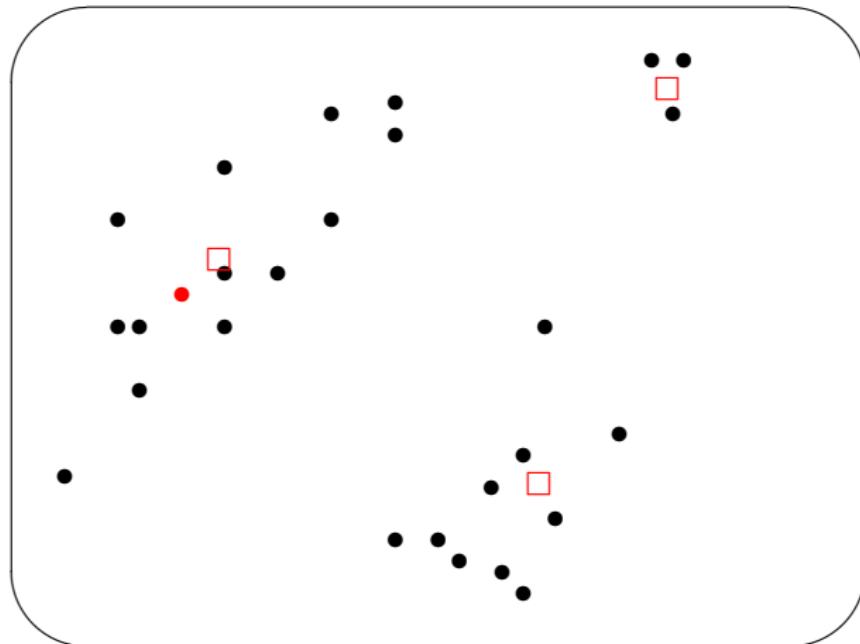
artfact
Real Time Machine Learning

The problem of Online Clustering



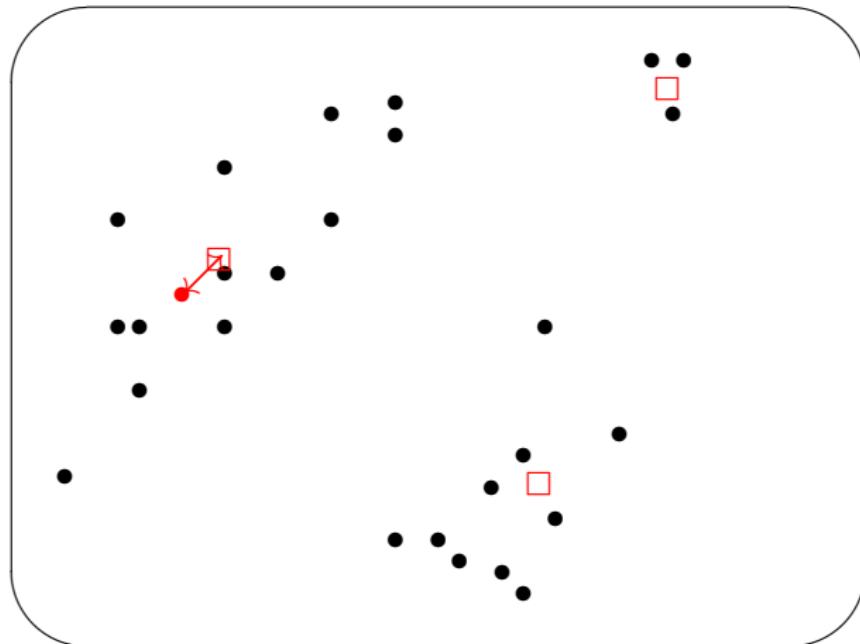
artfact
Real Time Machine Learning

The problem of Online Clustering



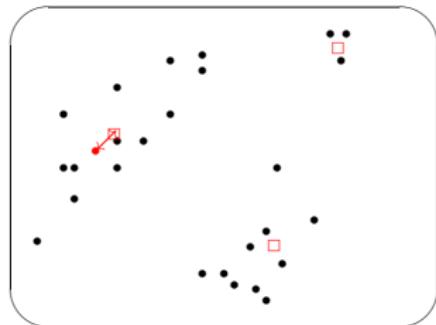
artfact
Real Time Machine Learning

The problem of Online Clustering



artfact
Real Time Machine Learning

The problem of Online Clustering

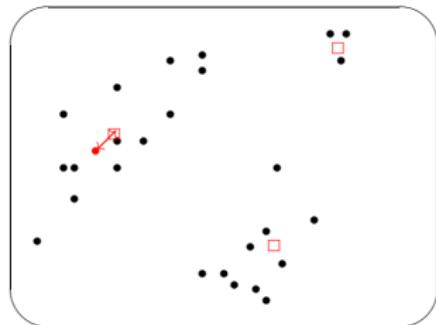


- ▶ Principle : clustering as prediction !



artfact
Real Time Machine Learning

The problem of Online Clustering

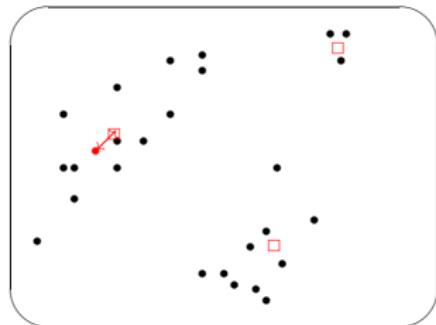


- ▶ Principle : clustering as prediction !
- ▶ Sparsity assumption : points are grouped into s clusters.



artfact
Real Time Machine Learning

The problem of Online Clustering

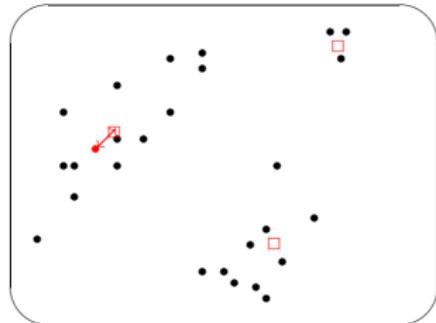


- ▶ Principle : clustering as prediction !
- ▶ Sparsity assumption : points are grouped into s clusters.
- ▶ Use PAC-Bayesian regularization to choose the number of clusters.



artfact
Real Time Machine Learning

The problem of Online Clustering



- ▶ Principle : clustering as prediction !
- ▶ Sparsity assumption : points are grouped into s clusters.
- ▶ Use PAC-Bayesian regularization to choose the number of clusters.

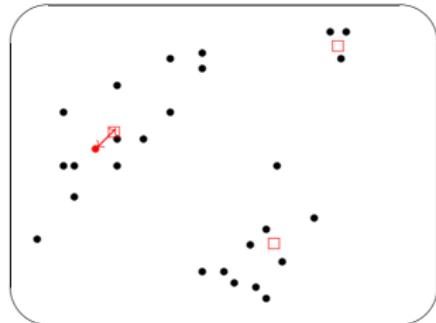
We prove new kind of sparsity regret bounds:

$$\sum_{t=1}^T \ell(\hat{\mathbf{c}}_t, x_t) - \inf_{\mathbf{c} \in \mathbb{R}^{dp}} \left\{ \sum_{t=1}^T \ell(\mathbf{c}, x_t) + \lambda |\mathbf{c}|_0 \right\},$$



artfact
Real Time Machine Learning

The problem of Online Clustering



- ▶ Principle : clustering as prediction !
- ▶ Sparsity assumption : points are grouped into s clusters.
- ▶ Use PAC-Bayesian regularization to choose the number of clusters.

We prove new kind of sparsity regret bounds:

$$\sum_{t=1}^T \ell(\hat{\mathbf{c}}_t, x_t) - \inf_{\mathbf{c} \in \mathbb{R}^{dp}} \left\{ \sum_{t=1}^T \ell(\mathbf{c}, x_t) + \lambda |\mathbf{c}|_0 \right\},$$

where $|\mathbf{c}|_0 = \text{card}\{j = 1, \dots, p : c_j \neq 0_{\mathbb{R}^d}\}$ and

$$\ell(\mathbf{c}, x) = \min_{j=1, \dots, p} \|c_j - x\|_2^2.$$

Business case

- ▶ Client : éditeur de logiciel dans le commerce conversationnel



artfact
Real Time Machine Learning

Business case

- ▶ **Client** : éditeur de logiciel dans le commerce conversationnel
- ▶ **Inputs (X)**: comportement de navigations

Business case

- ▶ **Client** : éditeur de logiciel dans le commerce conversationnel
- ▶ **Inputs (X)**: comportement de navigations
- ▶ **Objectifs** : satisfaction clients, instantanéité, conversion sans canibalisme

Business case

- ▶ **Client** : éditeur de logiciel dans le commerce conversationnel
- ▶ **Inputs (X)**: comportement de navigations
- ▶ **Objectifs** : satisfaction clients, instantanéité, conversion sans canibalisme



L'existant

- ▶ Absence de stockage

L'existant

- ▶ Absence de stockage
- ▶ Absence d'architecture Big Data

L'existant

- ▶ Absence de stockage
- ▶ Absence d'architecture Big Data
- ▶ Réglage manuel des règles de ciblage

L'existant

- ▶ Absence de stockage
- ▶ Absence d'architecture Big Data
- ▶ Réglage manuel des règles de ciblage



artifact
Real Time Machine Learning

Proof of concept (POC) : périmètre

- ▶ 11 sites
- ▶ 20M pvues/mois
- ▶ 160k events/heure

Proof of concept (POC) : objectifs

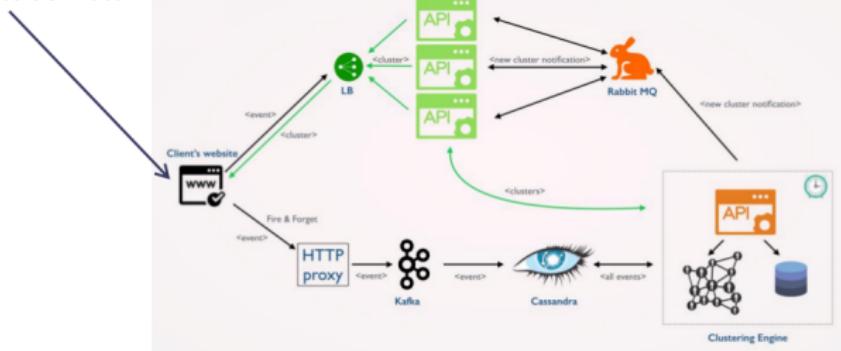
- ▶ Automatiser le ciblage
- ▶ Simplifier l'administration et la mise à jour
- ▶ Améliorer les performances



artfact
Real Time Machine Learning

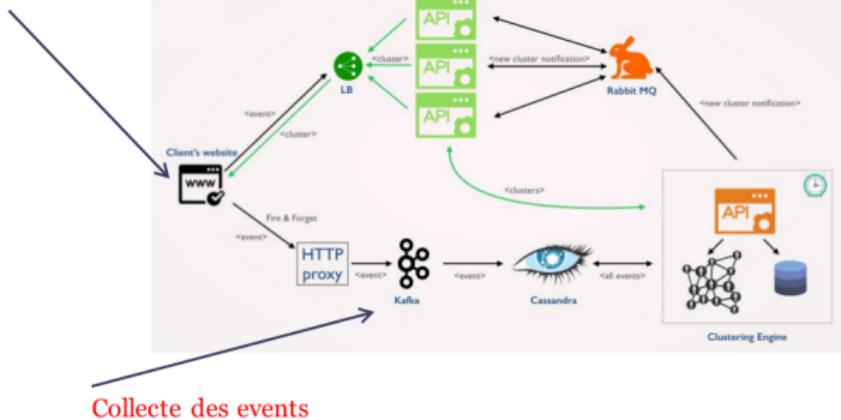
Industrialisation : architecture

Site e-commerce:
Source des données



Industrialisation : architecture

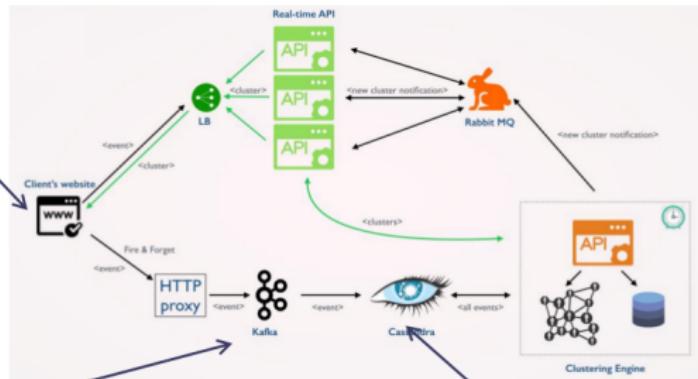
Site e-commerce:
Source des données



artfact
Real Time Machine Learning

Industrialisation : architecture

Site e-commerce:
Source des données



Collecte des events

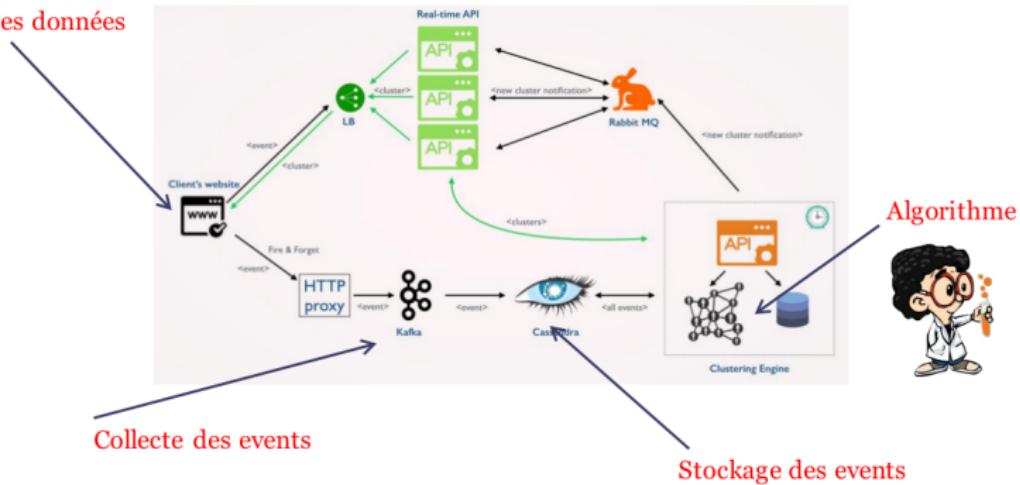
Stockage des events



artfact
Real Time Machine Learning

Industrialisation : architecture

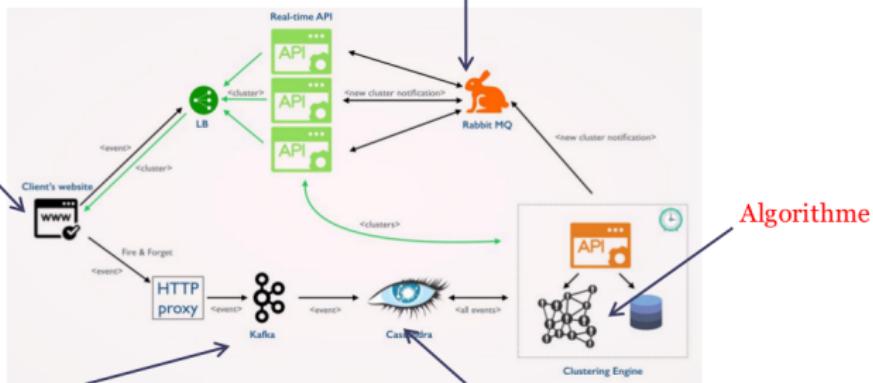
-commerce:
des données



artfact
Real Time Machine Learning

Industrialisation : architecture

:rce:
inées

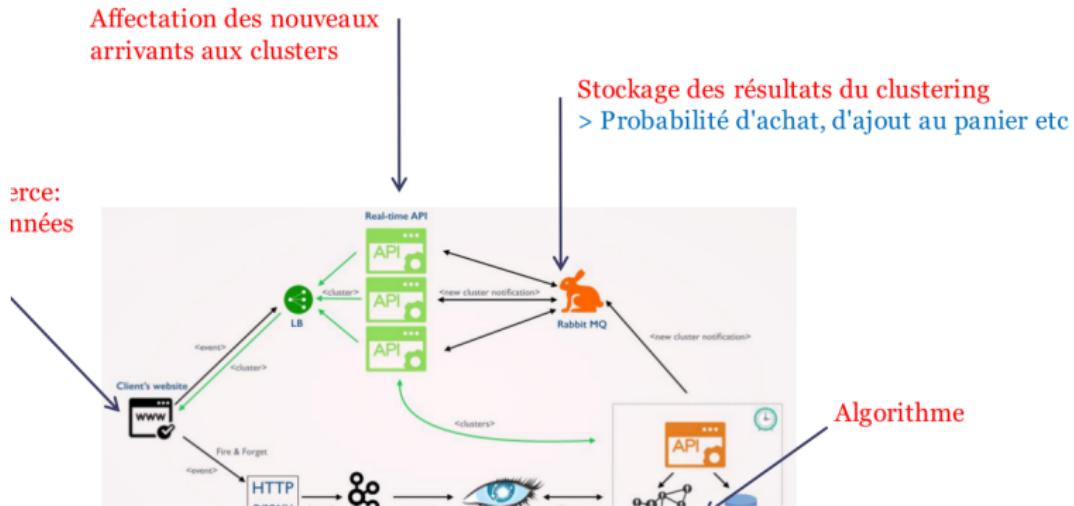


Algorithmme



artfact
Real Time Machine Learning

Industrialisation : architecture

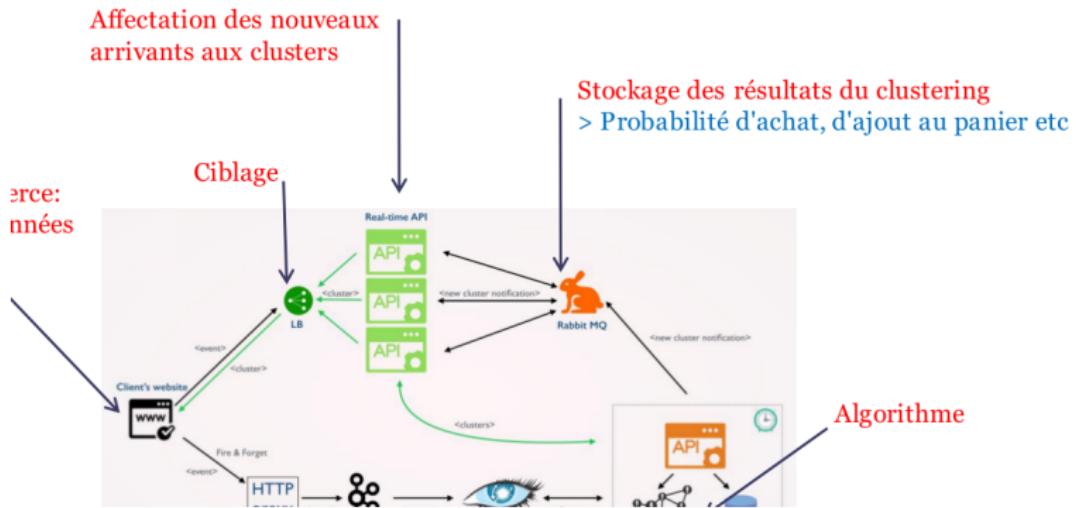


Algorithme



artfact
Real Time Machine Learning

Industrialisation : architecture



artfact
Real Time Machine Learning

Contents

Machine Learning : Gentle start

Cas d'application : segmentation clients temps réel

Cas d'application : détection de communautés temps réel

Graph clustering

- ▶ Proposer l'analyse précédente sur un graphe



artfact
Real Time Machine Learning

Graph clustering

- ▶ Proposer l'analyse précédente sur un graphe
- ▶ Nouveaux challenges



artfact
Real Time Machine Learning

Graph clustering

- ▶ Proposer l'analyse précédente sur un graphe
- ▶ Nouveaux challenges
- ▶ Have a look !



artfact
Real Time Machine Learning

Ingrédients

- ▶ Théorie des graphes
- ▶ Optimisation par Monte Carlo Markov Chain (MCMC)
- ▶ Calculs parallèles



artfact
Real Time Machine Learning

Qu'est-ce qu'un graphe ?

Définition ([Wikipédia](#)) :

Un graphe est un ensemble de points nommés nœuds (parfois sommets ou cellules) reliés par des traits (segments) ou flèches nommées arêtes (ou liens ou arcs). L'ensemble des arêtes entre nœuds forme une figure similaire à un réseau.



artfact
Real Time Machine Learning

Qu'est-ce qu'un graphe ?

Définition ([Wikipédia](#)) :

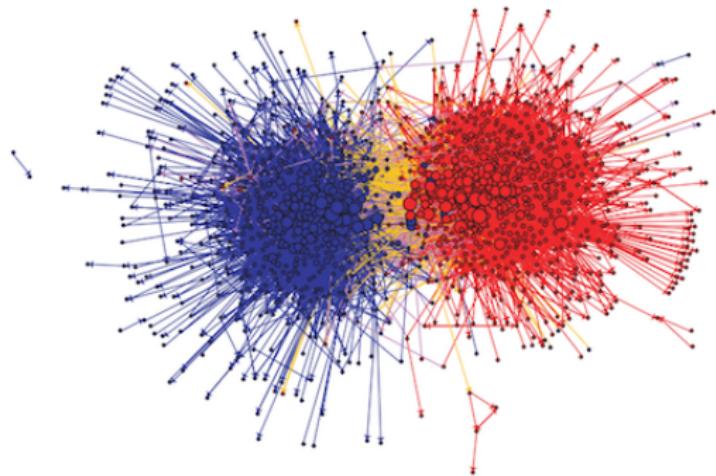
Un graphe est un ensemble de points nommés nœuds (parfois sommets ou cellules) reliés par des traits (segments) ou flèches nommées arêtes (ou liens ou arcs). L'ensemble des arêtes entre nœuds forme une figure similaire à un réseau.

Key point : **relational structure**.



artfact
Real Time Machine Learning

Illustrations : political blog

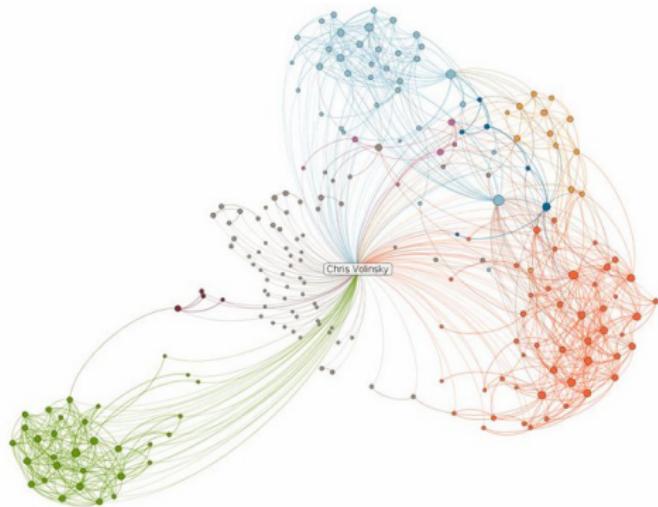


Network of U.S.
political blogs by
Adamic and Glance
(2004)



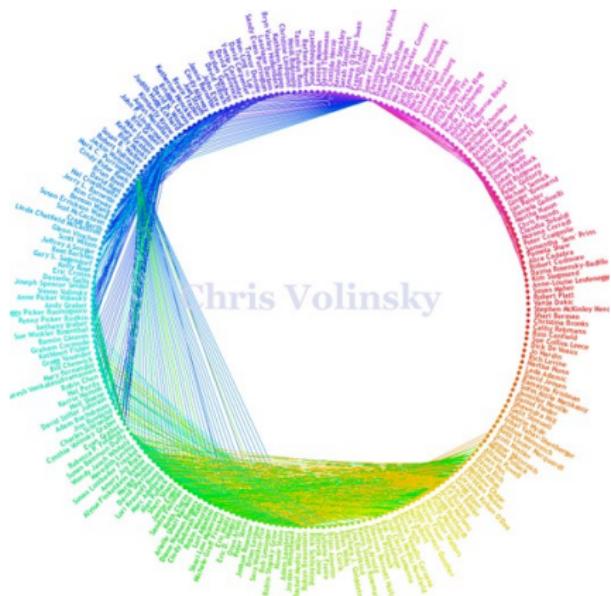
artfact
Real Time Machine Learning

Illustrations : LinkedIn community



LinkedIn community
of Chris Volinsky
(2011, Columbia
University), see
www.socilab.com

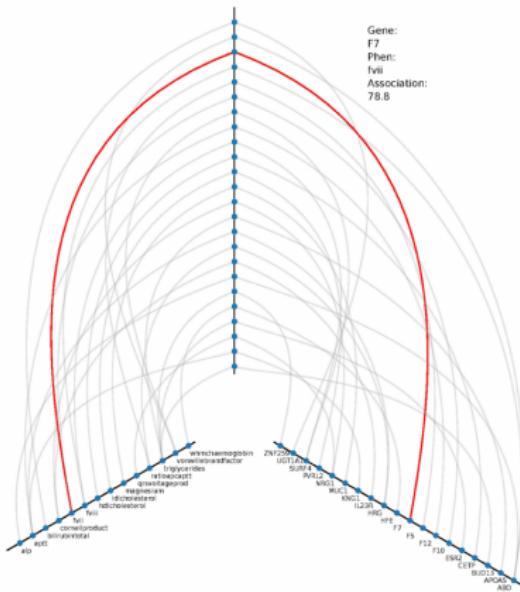
Illustrations : Facebook friends



Facebook of Chris
Volinsky (2011,
Columbia University)

Illustrations : Gene/phenotype interactions

Hive diagram test



Canonical Correlation Analysis for Gene-Based Pleiotropy Discovery (Plos Computational Biology, 2016)



artfact
Real Time Machine Learning

Formalisme mathématique

Definition

Un graphe G est une structure constituée d'un couple (E, V) où :



artfact
Real Time Machine Learning

Formalisme mathématique

Definition

Un graphe G est une structure constituée d'un couple (E, V) où :

- ▶ $V = \{v_1, \dots, v_N\}$ est un ensemble de N nœuds,



artfact
Real Time Machine Learning

Formalisme mathématique

Definition

Un graphe G est une structure constituée d'un couple (E, V) où :

- ▶ $V = \{v_1, \dots, v_N\}$ est un ensemble de N nœuds,
- ▶ E est un ensemble d'arêtes ou liens.



artfact
Real Time Machine Learning

Formalisme mathématique

Definition

Un graphe G est une structure constituée d'un couple (E, V) où :

- ▶ $V = \{v_1, \dots, v_N\}$ est un ensemble de N nœuds,
- ▶ E est un ensemble d'arêtes ou liens.

$|V| = N$ est appelé l'ordre de G et $|E|$ la taille de G .



Formalisme mathématique

Definition

Un graphe G est une structure constituée d'un couple (E, V) où :

- ▶ $V = \{v_1, \dots, v_N\}$ est un ensemble de N nœuds,
- ▶ E est un ensemble d'arêtes ou liens.

$|V| = N$ est appelé l'ordre de G et $|E|$ la taille de G .

Definition

Pour un graphe $G = (E, V)$, le degré d'un nœud $v \in V$ est le nombre d'arêtes dans E partant de v .



artfact
Real Time Machine Learning

Propriétés élémentaires

Definition

A **sparse graph** is a graph (E, V) such that $|E|$ is of order $|V|$.



artfact
Real Time Machine Learning

Propriétés élémentaires

Definition

A **sparse graph** is a graph (E, V) such that $|E|$ is of order $|V|$.



artfact
Real Time Machine Learning

Propriétés élémentaires

Definition

A **sparse graph** is a graph (E, V) such that $|E|$ is of order $|V|$.

Definition

A **scale-free network** is a network whose degree distribution follows a power law, that is to say:

$$\mathbb{P}(d_i = k) \sim k^{-\gamma}.$$



artfact
Real Time Machine Learning

Propriétés élémentaires

Definition

A **sparse graph** is a graph (E, V) such that $|E|$ is of order $|V|$.

Definition

A **scale-free network** is a network whose degree distribution follows a power law, that is to say:

$$\mathbb{P}(d_i = k) \sim k^{-\gamma}.$$

Definition

A graph G with N vertices satisfies the **small-world property** if in expectation, 2 randomly chosen vertices i and j have distances proportional to $\log N$.



Détection de communautés (ou graph clustering)

- ▶ Partition des nœuds qui vérifie : beaucoup d'arêtes à l'intérieur des groupes, peu d'arêtes entre les groupes.



artfact
Real Time Machine Learning

Détection de communautés (ou graph clustering)

- ▶ Partition des nœuds qui vérifie : beaucoup d'arêtes à l'intérieur des groupes, peu d'arêtes entre les groupes.
- ▶ Tâche non-supervisée (!)



artfact
Real Time Machine Learning

Détection de communautés (ou graph clustering)

- ▶ Partition des nœuds qui vérifie : beaucoup d'arêtes à l'intérieur des groupes, peu d'arêtes entre les groupes.
- ▶ Tâche non-supervisée (!)
- ▶ Méthodes spectrales consistant à diagonaliser :

$$L = D - A,$$

où A est la matrice d'adjacence et D la matrice des degrés.



artfact
Real Time Machine Learning

Détection de communautés (ou graph clustering)

- ▶ Partition des nœuds qui vérifie : beaucoup d'arêtes à l'intérieur des groupes, peu d'arêtes entre les groupes.
- ▶ Tâche non-supervisée (!)
- ▶ Méthodes spectrales consistant à diagonaliser :

$$L = D - A,$$

où A est la matrice d'adjacence et D la matrice des degrés.

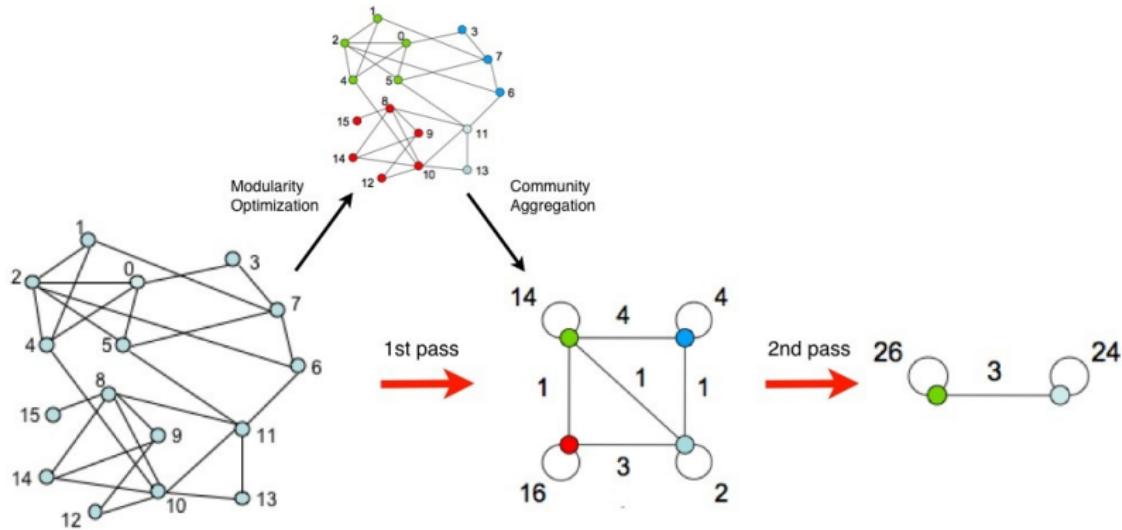
- ▶ Optimisation d'un critère appelé la modularité :

$$\mathcal{M}_G(\mathbf{c}) = \frac{1}{2m} \sum_{(i,j)} \left(A_{ij} - \frac{k_i k_j}{2m} \right) \delta_{\mathbf{c}}(i,j).$$



artfact
Real Time Machine Learning

Algorithme glouton



artfact
Real Time Machine Learning

Algorithme MCMC

- ▶ Select at random a neighborhood C' of the current coloration C
- ▶ Accept the proposal with standard Metropolis accept/reject ratio

Algorithm 1 MH for Community Detection

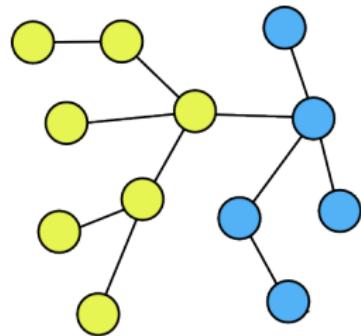
- 1: Initialization $\lambda > 0$, $C^{(0)}$.
- 2: For $k = 1, \dots, N$:
- 3: Draw $C' \sim p(\cdot | C^{(k-1)})$ where $p(\cdot | C^{(k-1)}) \in \mathcal{P}(\mathcal{N}^{C^{(k-1)}})$ is the proposal distribution over $\mathcal{N}^{C^{(k-1)}}$, a neighborhood of $C^{(k-1)}$.
- 4: Update $C^{(k)} = C'$ with acceptance ratio :

$$\rho = 1 \wedge \left(r_{C^{(k-1)} \rightarrow C'} \frac{\exp(\lambda Q^{C'})}{\exp(\lambda Q^{C^{(k-1)}})} \right), \text{ where } r_{C \rightarrow C'} := p(C^{(k-1)} | C') / p(C' | C^{(k-1)}).$$

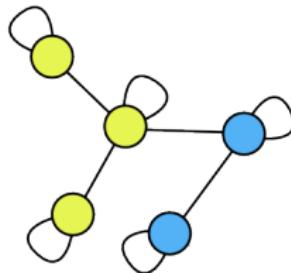


artfact
Real Time Machine Learning

Algorithme final



$L = 0$



$L = 1$

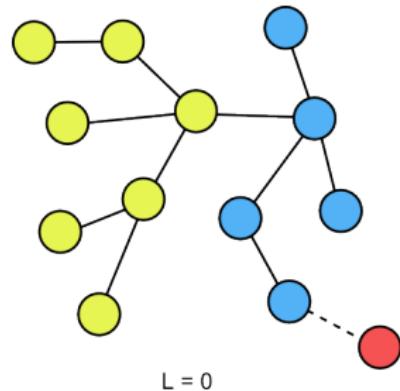


$L = 2$

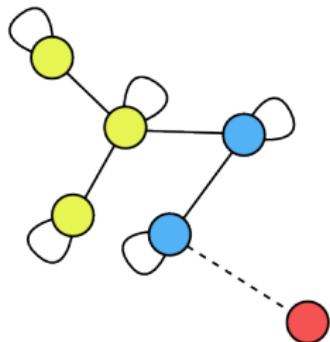


artfact
Real Time Machine Learning

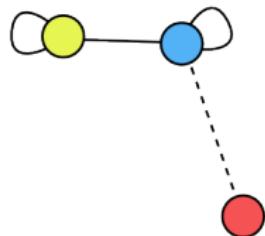
Algorithme final



$L = 0$



$L = 1$

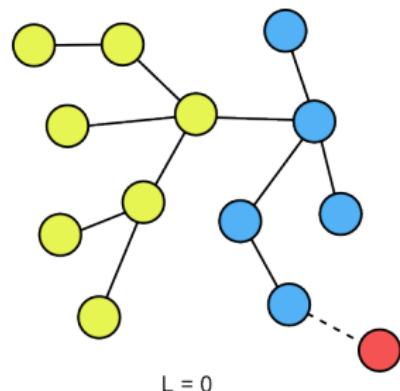


$L = 2$

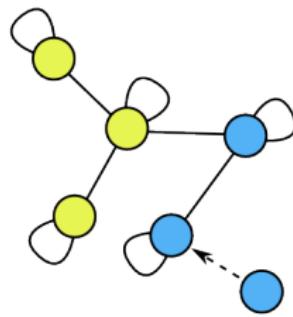


artfact
Real Time Machine Learning

Algorithme final



$L = 0$



$L = 1$

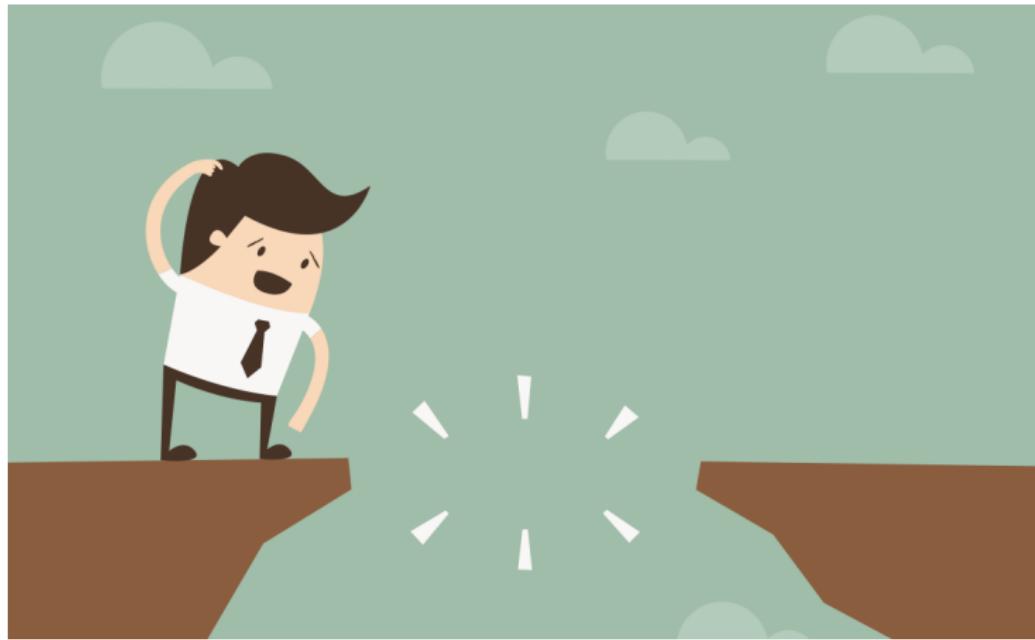


$L = 2$



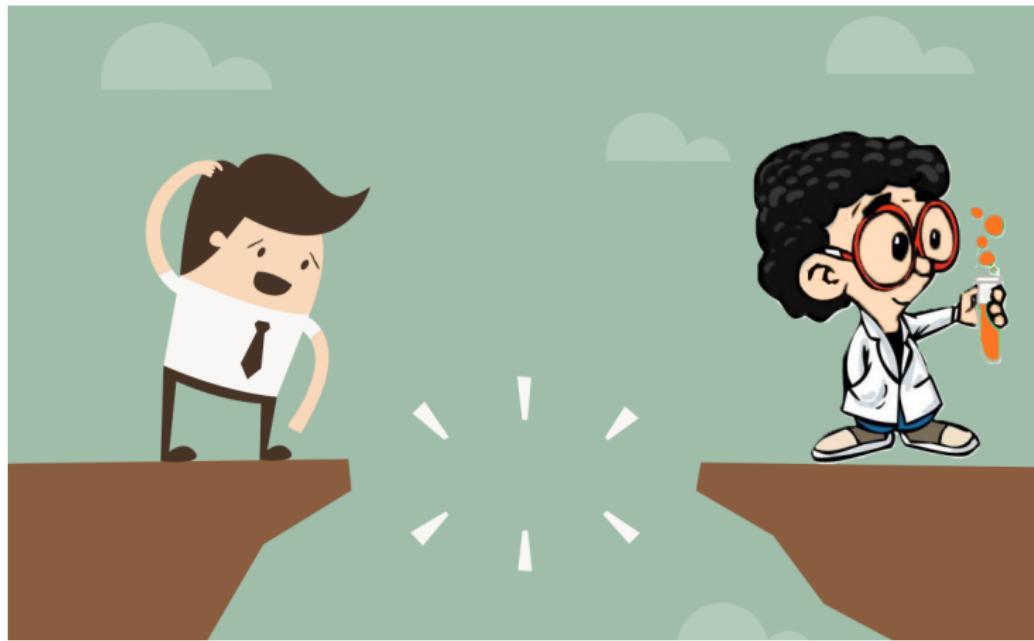
artfact
Real Time Machine Learning

Business case : appel à volontaire



artfact
Real Time Machine Learning

Business case : appel à volontaire



artfact
Real Time Machine Learning

Business case : appel à volontaire



artfact
Real Time Machine Learning

Propaganda :-)

Website : <http://www.artfact-online.fr>

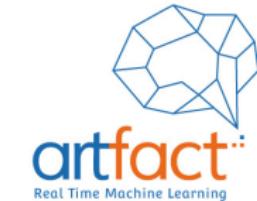


artfact
Real Time Machine Learning

Propaganda :-)

Website : <http://www.artfact-online.fr>

Twitter : @artfact-lab



Propaganda :-)

Website : <http://www.artfact-online.fr>

Twitter : @artfact-lab

Youtube : Le Machine Learning par l'exemple



artfact
Real Time Machine Learning

Propaganda :-)

Website : <http://www.artfact-online.fr>

Twitter : @artfact-lab

Youtube : Le Machine Learning par l'exemple

Meetup Machine Learning : Nantes Rennes Pau



artfact
Real Time Machine Learning

Propaganda :-)

Website : <http://www.artfact-online.fr>

Twitter : [@artfact-lab](#)

Youtube : [Le Machine Learning par l'exemple](#)

Meetup Machine Learning : [Nantes Rennes Pau](#)

Plateforme d'innovation ouverte : [Learnation Square](#)



artfact
Real Time Machine Learning

Propaganda :-)

- Merci de votre attention ! -

Website : <http://www.artfact-online.fr>

Twitter : @artfact-lab

Youtube : [Le Machine Learning par l'exemple](#)

Meetup Machine Learning : Nantes Rennes Pau

Plateforme d'innovation ouverte : [Learnation Square](#)



artfact
Real Time Machine Learning