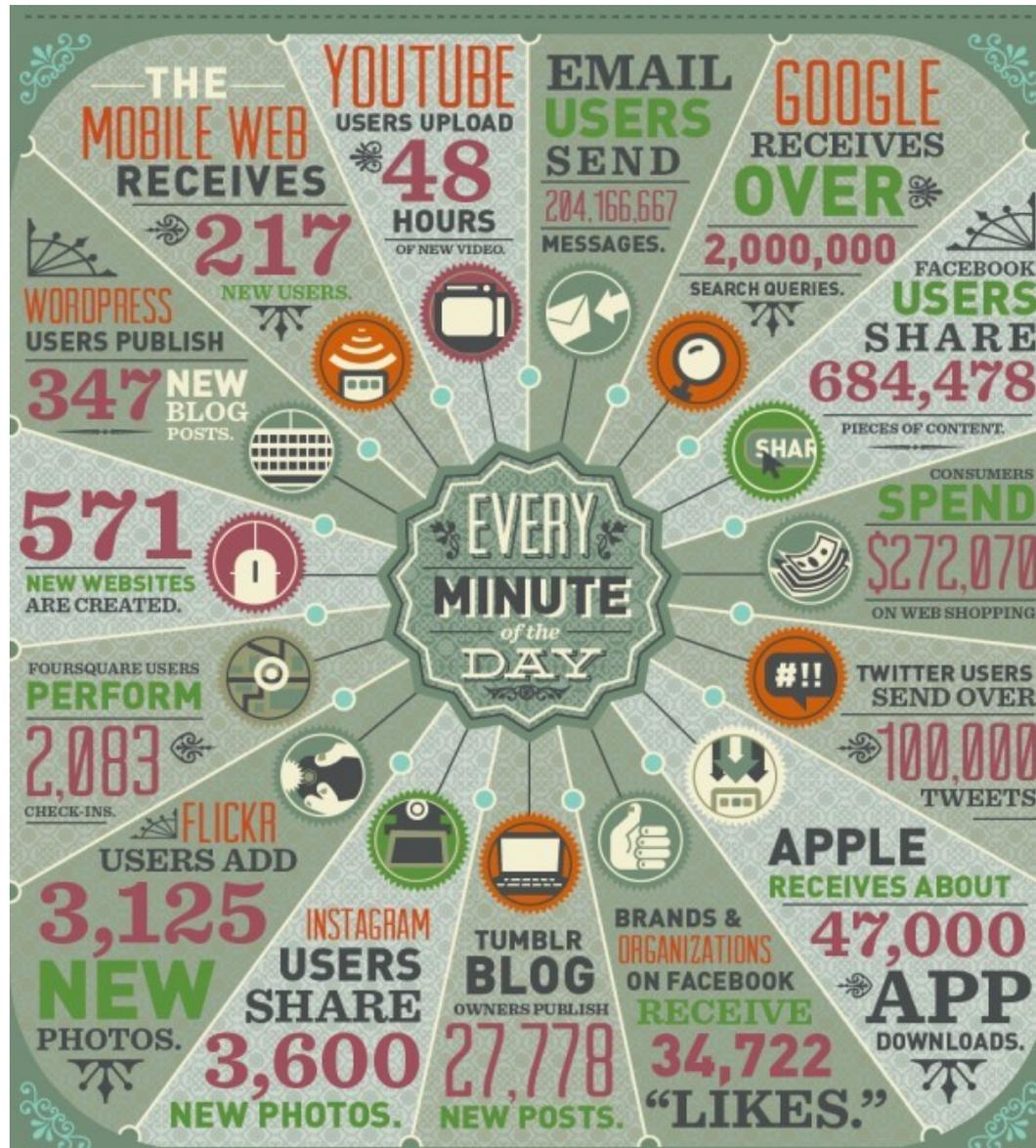


Data Mining

Exposés logiciels, systèmes et réseaux.

Damien Jubeau IR3

Lundi 19 novembre 2012



Plan

Data mining : définition, utilisations et concepts

Wolfram|Alpha : extraction de données d'un compte Facebook

Le problème du « badge » avec Weka

I – Data mining : définition, utilisations et concepts

I – Data mining : définition, utilisations et concepts

Définition

se traduit par

Exploration de données

ou encore

Extraction de connaissances à partir de données

« Processus non-trivial d'identification de structures inconnues, valides et potentiellement exploitables dans les bases de données »

(Fayyad, 1996)

I – Data mining : définition, utilisations et concepts

Définition

Analyse de données et statistiques exploratoires existent depuis plus de 30 ans.

On peut voir le Data mining comme un prolongement :

- ajout de techniques d'IA
- données non structurées
- caractère business marqué

Big Data : 4ème révolution industrielle ?

I – Data mining : définition, utilisations et concepts

Domaines d'application

Scoring (réduction du coût d'acquisition ou de conservation d'un client – assurances, banques, opérateurs téléphoniques...)

Prévention du crime

Reconnaissance vocale

Recherche médicale

Détection de fraudes

Poker !

...etc

I – Data mining : définition, utilisations et concepts

Domaines d'application

Les recherches Google :

- Google spell checker
- Autocomplétion
- Recherche locale

« Storing and analyzing logs of user searches is how Google's algorithm learns to give you more useful results. Just as data availability has driven progress of search in the past, **the data in our search logs will certainly be a critical component of future breakthroughs.** »

Log

IP – Cookie – Recherche – Date & heure

<http://www.google.org/flutrends>

I – Data mining : définition, utilisations et concepts

Démarche

Mode projet

- + importance de l'objectif
- + Préparation des données en vue du traitement structurées (base de données, fichiers tabulaires...) ou non structurées (textes, images...)
- + Élaboration, choix des modèles à appliquer (monde des statistiques ou de l'apprentissage automatique)
- + Évaluation et validation des connaissances extraites

I – Data mining : définition, utilisations et concepts

Des outils issus des statistiques et de l'IA

2 familles d'algorithmes :

Les **méthodes descriptives** permettent d'organiser, de simplifier et d'aider à comprendre l'information à partir des sources de données.

Recherche d'associations / Recherche de séquences similaires ...etc

Les **méthodes prédictives** visent à expliquer ou prévoir plusieurs phénomènes observables et effectivement mesurés. On cherche à prédire la valeur d'une **variable cible** à partir des valeurs de prédicteurs.

Régression linéaire multiple / Réseaux de neurones / Arbres de régression...

I – Data mining : définition, utilisations et concepts

Problèmes et limites

Hétérogénéité des données

→ importance de la préparation des données

Volume des données

→ calcul distribué (Hadoop)

Risque d'apporter une réponse hors scope

→ problème éthique

Danger pour la vie privée

→ explosion des données personnelles sur Internet

II – Extraction de données d'un compte Facebook

II – Extraction de données d'un compte Facebook

Présentation de Wolfram|Alpha

Outil de calcul en langage universel – lancé 2009

Analyse des données (10 milliards d'informations) issues de disciplines très variées (mathématiques, physique, chimie, nouvelles technologies, données socio-économiques...)

Utilisé par Bing, DuckDuckGo et Siri.

II – Extraction de données d'un compte Facebook

Présentation de Wolfram|Alpha

Outil de calcul en langage universel – lancé 2009

Analyse des données (10 milliards d'informations) issues de disciplines très variées (mathématiques, physique, chimie, nouvelles technologies, données socio-économiques...)

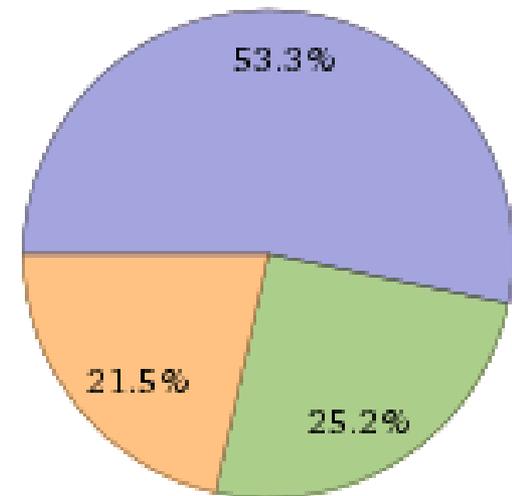
Utilisé par Bing, DuckDuckGo et Siri.

II – Extraction de données d'un compte Facebook

Publications de l'utilisateur

Types:

type	count	ratio	
statuses	245		
posted links	116		
uploaded photos	99		



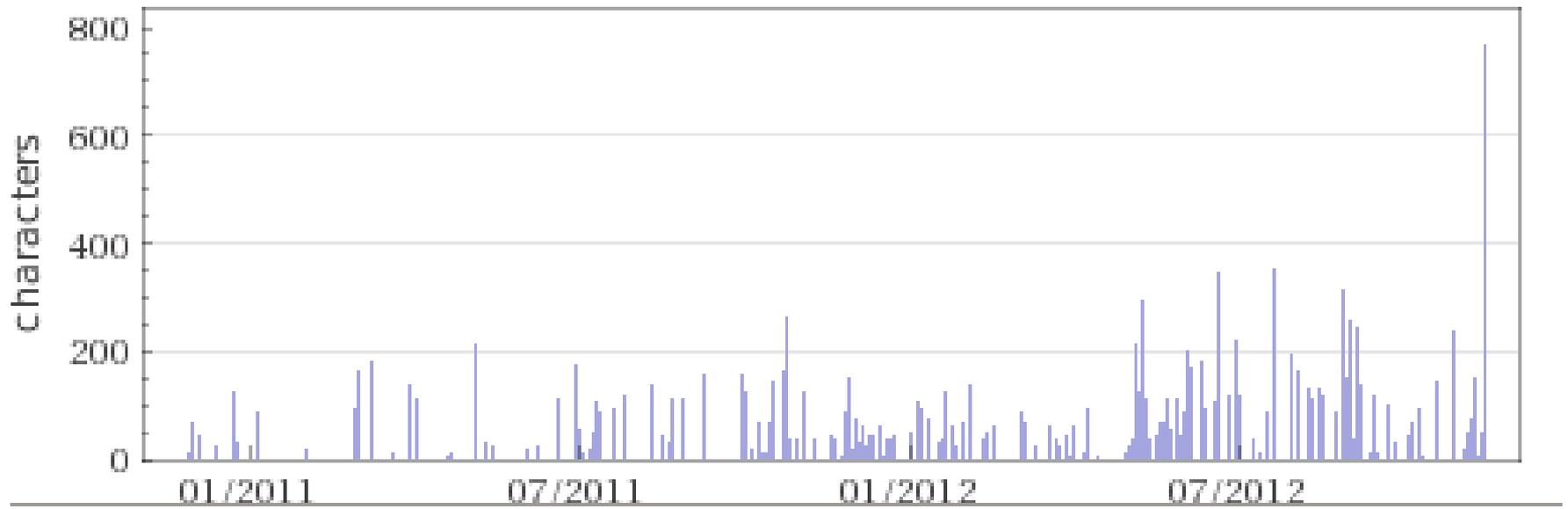
(based on 460 total activities)

Clip 'n Share

II – Extraction de données d'un compte Facebook

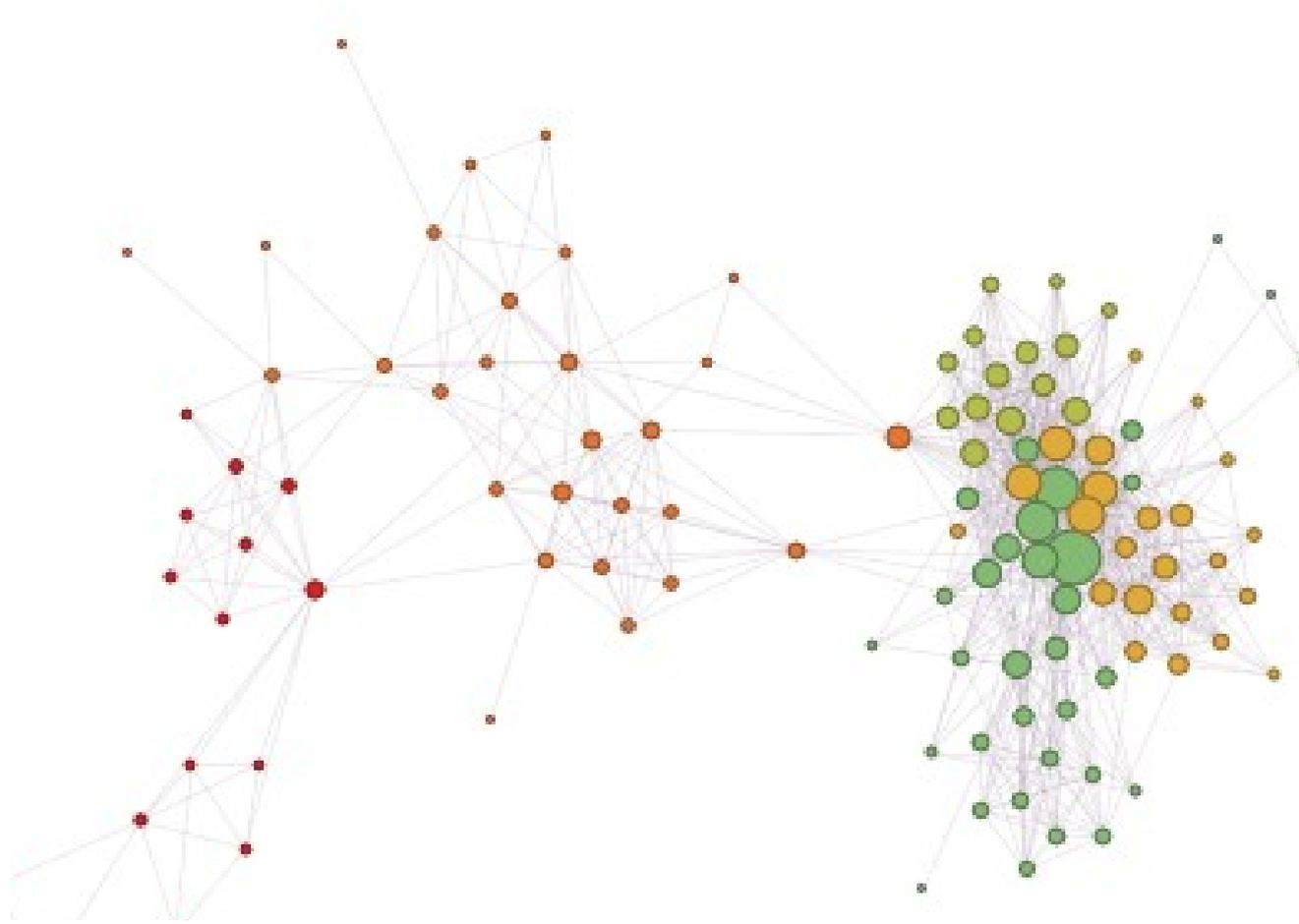
Taille des statuts

Post lengths:



II – Extraction de données d'un compte Facebook

Réseau de connaissances



II – Extraction de données d'un compte Facebook

Analyse d'un compte Facebook

Et bien d'autres : opinions politiques, mots les plus utilisés...

Facebook possède des données similaires
pour plus d'un milliard d'utilisateurs...
...et si le modèle de l'outil intégrait des recoupements ?

Problématique Data mining / vie privée évidente

Utilisation commerciale ? Surveillance abusive des gouvernements ?

III – Le problème du « badge »

III – Le problème du « badge »

Contexte

En 1994, chaque participant à une conférence sur l'apprentissage automatique (Machine Learning) s'est vu remettre un badge « + » ou « - ».

Le but était pour un participant de découvrir la fonction d'attribution des badges, sachant qu'elle était basée uniquement sur les noms des participants. Pour cela la liste des participants et des badges attribués étaient fournies.

Utilisation de Weka, logiciel libre de Data Mining

III – Le problème du « badge »

Préparation des données

Attribute name, and type

name {...}
length numeric
even_odd {0,1}
first_char_vowel {0,1}
second_char_vowel {0,1}
vowels numeric
consonants numeric
vowel_consonant_ratio numeric
spaces numeric
dots numeric
initials
words numeric
class {+,-}

Explanation

all the names (given in the original)
length of name
length of name even or odd?
is first character a vowel?
is second character a vowel?
number of vowels in the name
number of consonants
the ratio of vowels / consonant
number of spaces
number of "." in the name, i.e. name
number of words, i.e number of names
the badge labels (given in the original)

Ex : Ameer Foued, 11 , 0 , 1 , 0 , 6 , 4 , 1.50 , 1 , 0 , 2 , -

III – Le problème du « badge »

Utilisation du logiciel Weka

The screenshot shows the Weka Explorer interface with the 'Classify' tab selected. The 'Current relation' is 'badges' with 294 instances and 12 attributes. The 'Attributes' list includes 'name', 'length', 'even_odd', 'first_char_vowel', 'second_char_vowel', 'vowels', 'consonants', 'vowel_consonant_ratio', 'spaces', 'dots', 'words', and 'badge'. The 'Selected attribute' is 'length', which is numeric with 16 distinct values and 2 unique values (1%). A histogram for the 'length' attribute is displayed, showing the distribution of values from 7 to 25. The histogram bars are colored red and blue, representing two different classes. The status bar at the bottom indicates 'OK' and 'Log'.

Statistic	Value
Minimum	7
Maximum	25
Mean	14.034
StdDev	2.862

Class	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
Class 1 (Red)	5	20	26	87	37	61	21	25	8	3	0	1	0	0	0	0	0	0	0
Class 2 (Blue)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Classification
Simple : + ou -

III – Le problème du « badge »

Utilisation du logiciel Weka

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

Filter: Choose None Apply

Current relation: Relation: badges, Instances: 294, Attributes: 12

Attributes: All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> name
2	<input checked="" type="checkbox"/> length
3	<input checked="" type="checkbox"/> even_odd
4	<input checked="" type="checkbox"/> first_char_vowel
5	<input checked="" type="checkbox"/> second_char_vowel
6	<input checked="" type="checkbox"/> vowels
7	<input checked="" type="checkbox"/> consonants
8	<input checked="" type="checkbox"/> vowel_consonant_ratio
9	<input checked="" type="checkbox"/> spaces
10	<input checked="" type="checkbox"/> dots
11	<input checked="" type="checkbox"/> words
12	<input checked="" type="checkbox"/> badge

Remove

Selected attribute: Name: second_char_vowel, Type: Nominal, Missing: 0 (0%), Distinct: 2, Unique: 0 (0%)

No.	Label	Count
1	0	84
2	1	210

Class: badge (Nom) Visualize All

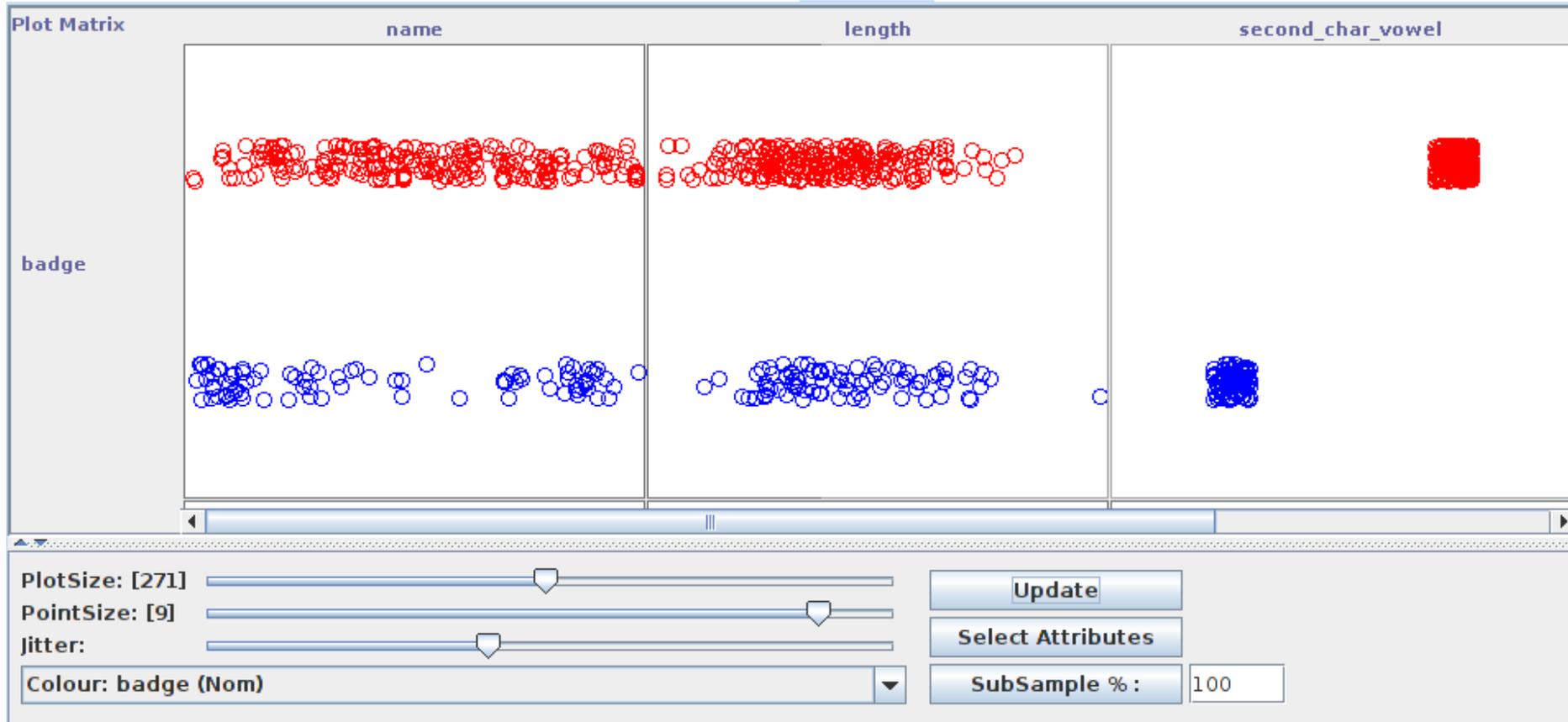
84 210

Status: OK Log x 0

Cas simple : on cherche une association forte : elle est systématique

III – Le problème du « badge »

Visualisation des résultats



III – Le problème du « badge »

D'autres cas d'application

<http://ftp.ics.uci.edu/pub/machine-learning-databases/>

Conclusion

De nombreux logiciels mais nécessitent une expertise

Des techniques qui vont se développer et se démocratiser

De grandes avancées potentielles (médicales...)

Des dangers sans cadre légal

Data Mining

Exposés logiciels, systèmes et réseaux.

Damien Jubeau IR3

Lundi 19 novembre 2012