

EDA 2015
Bruxelles
2 – 3 avril 2015

TLabel: Text clustering and labelling in OLAP environment

TLabel

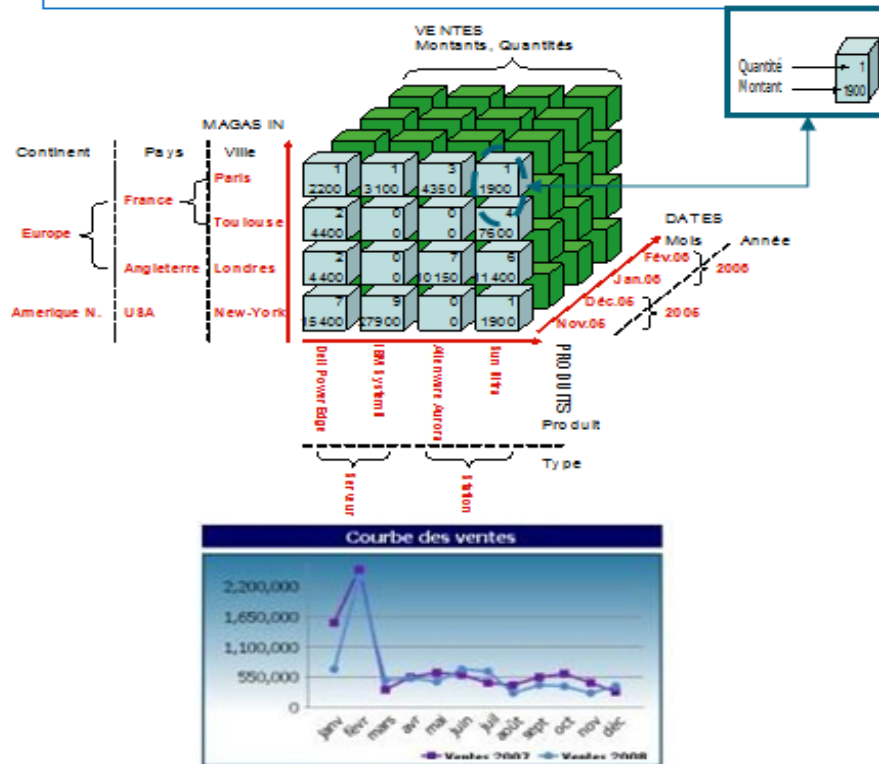
Nouvel opérateur d'agrégation
par catégorisation
dans les cubes de textes

Lamia Oukid, Omar Boussaid,
Nadjia Benblidia et Fadila Bentayeb

Context

Classical OLAP

20% structured data of information system



Text OLAP

80% non-structured data of information system



Need new OLAP operators for text data



Towards text OLAP

■ Classical OLAP

- Exploring and navigating through data

■ OLAP limitation

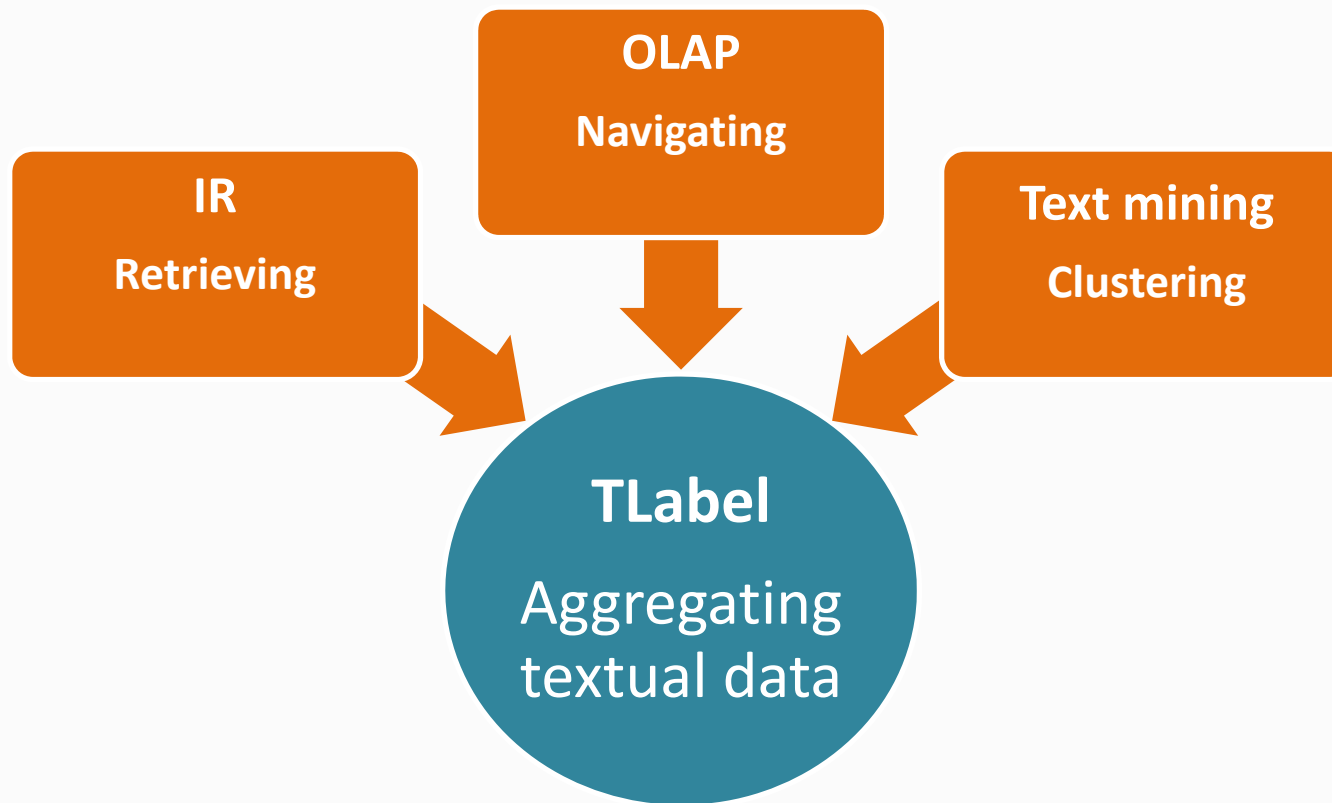
- Not adapted to textual data

■ Objectives

- Text cubes
- Text OLAP operators
- Integrating data semantics in OLAP analysis



OLAP / IR/Text mining





Outline

- Text cube
- Textual query analysis
- TLabel : Clustering operator in text OLAP environment
- Experiments and results
- Conclusion and future work



Text definition

- Set of terms
 - Textuel content

- Set of metadata
 - Information on textual data

- Set of concepts
 - Extracted from domain ontology
 - Enrich text contents

Text cube modelling

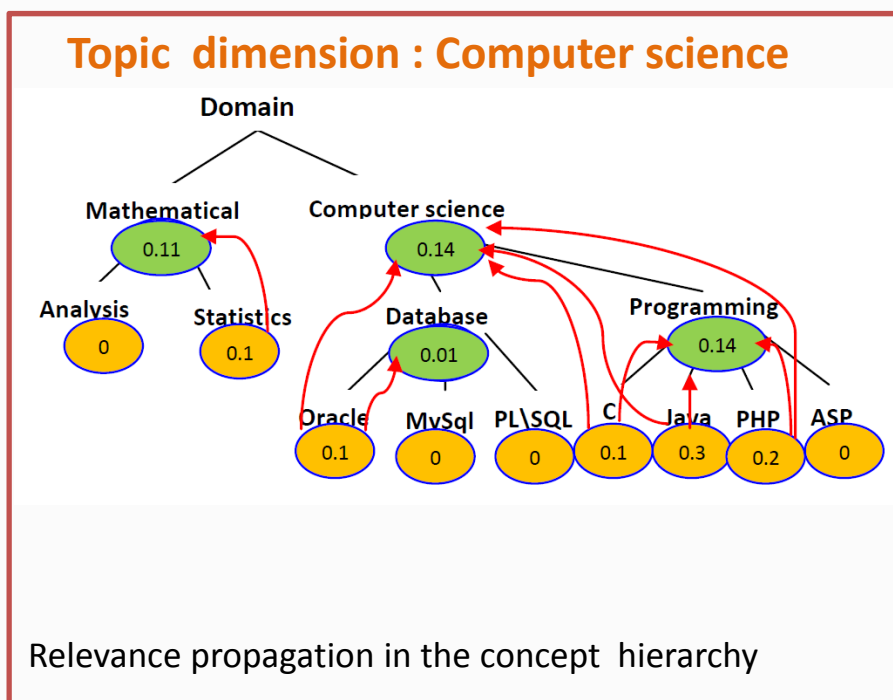
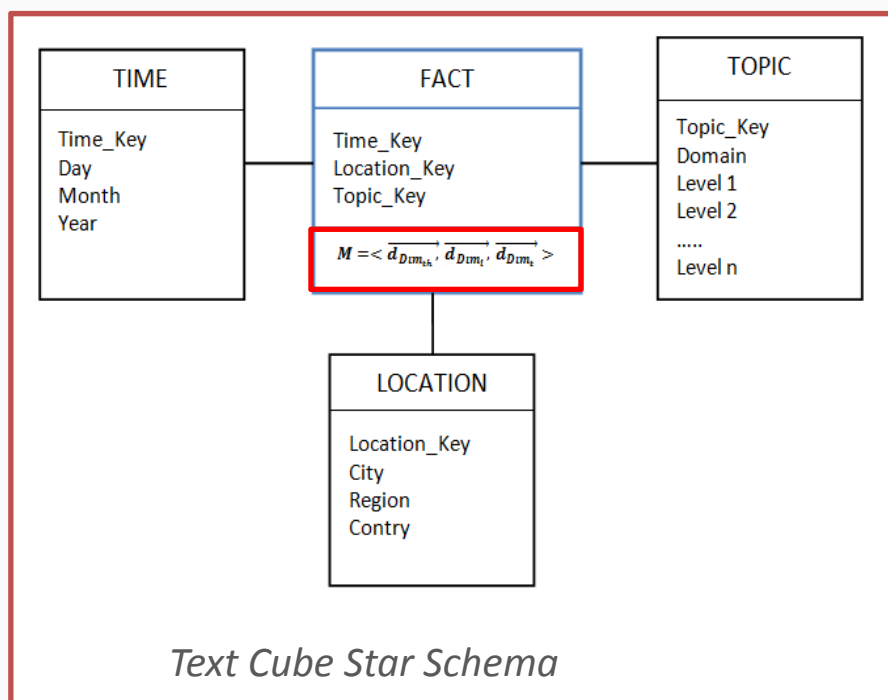
- Dimensions
 - Semantic dimensions
 - Metadata dimensions

- Textual measure M
 - Vector of weighted concepts
 - One vector per dimension

- Vector of concepts

$$M = \langle \overrightarrow{d_{Dim_1}}, \overrightarrow{d_{Dim_2}}, \dots, \overrightarrow{d_{Dim_*}} \rangle$$

Text cube: Example



Text cube for CV collection



Query modeling

- Given text cube with n dimensions: Simple Query
 - $Q = \langle V_1, V_2, \dots, V_n \rangle$
- Decision query with user preferences
 - The user can assign weight to each dimension
- Applied method
 - Generalized Cosinus Similarity : between query-document
- Result: Relevant text documents

More than extracting relevant text documents...

- Information Retrieval
 - Searching for relevant text documents

- Extracting knowledge from text documents
 - Ranking
 - Clustering
 - Resume
 - ...

- Text mining
 - Supervised methods
 - No supervised methods



TLabel: Text Labelling

- Combining OLAP/IR/Text mining

- Aggregating by clustering
 - Clusters of documents
 - Adapted K-means

- Assigning labels to clusters of documents
 - Domain ontology



TLabel: Text Labelling

- **Clustering step:** OCluster - OLAP-Cluster
 - Adapted K-means
 - Clusters of documents

- **Labelling step**
 - For each cluster, compute its **DResume**: resume document
 - DResume is a vector of weighted terms
 - Assign to each cluster one label obtained from DResume

Documents Clustering

■ OCluster : OLAP-Cluster

- Input: set of documents obtained from decision query
- Output: set of documents clusters
- Method: K-means with similarity function *ORank*

■ *ORank*: Computes the similarity between documents

$$ORank(d, ct) = \frac{\sum_{i=1}^n (\alpha_i \times Sim(\overrightarrow{d_{Dim_i}}, \overrightarrow{ct_{Dim_i}}))}{n}$$

- α_i : user preferences
- n: number of dimensions

DResume

- For each cluster of documents
 - Computes its **DResume** document

$$DResume = \langle \overrightarrow{DResume_{Dim_1}}, \overrightarrow{DResume_{Dim_2}}, \dots, \overrightarrow{DResume_{Dim_*}} \rangle$$

$$\overrightarrow{DResume_{Dim_i}} = \frac{\sum_{i=1}^N \overrightarrow{d_{Dim_i}}}{N}$$

N: Number of the documents in the cluster



Cluster Labelling

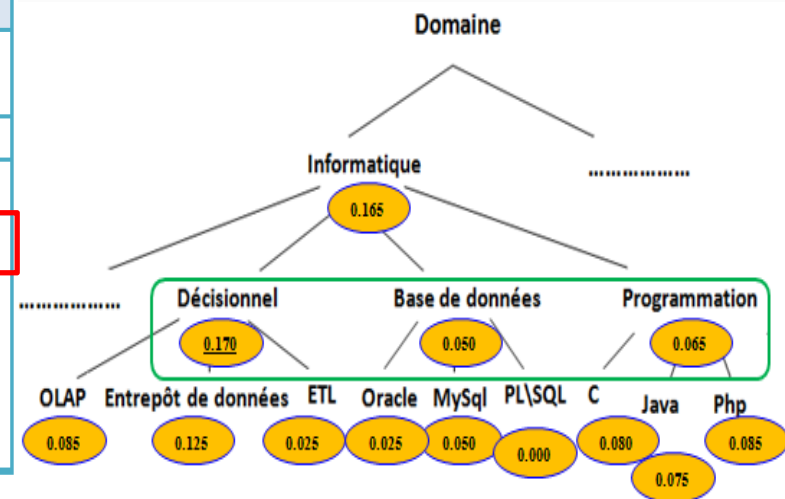
- Input
 - DResume*
 - One dimension
 - Domain ontology

- Method
 - DResume* Projection on the domain ontology

- Output
 - Documents clusters labelled
 - One label for one cluster

TLabel: Example

Opérateur d'agrégation <i>TLabel</i>			THEMATIQUE						
			Domaine	Informatique					
LOCALISATION	Ville	TEMPS	Année	Classe	Documents	Label			
	France		2014						
							Cl 1	{d1, d5}	Décisionnel
							Cl 2	{d3, d4}	Base de données
			Cl 3	{d2}	Programmation				



OLAP analysis with TLabel

Labelling documents cluster Cl 1

Experiments

- Data sources
 - 2000 CVs of candidates
 - Topic dimension : Computer science

- Ontology: hierarchy of concepts
 - Wikipedia

Experiments

- Preparing data sources
 - Text Tokenisation
 - Drop stop words
 - Term Lemmatisation: Tree tagger

- Loading data into text cube
 - Semantic dimensions: Topic and Location
 - For each dimension, load the concept hierarchy from the corresponding domain ontology
 - Time Dimension

Experiments & results

Results

■ Query: <Topic= Computer Science, Location= France , Time= 2014>

■ *Ocluster*

<i>OCluster</i>	CI 1	CI 2	CI 3	CI 4	CI 5	CI 6	CI 7	CI 8
Documents Number	98	181	84	179	215	178	1	216
Total	1152							

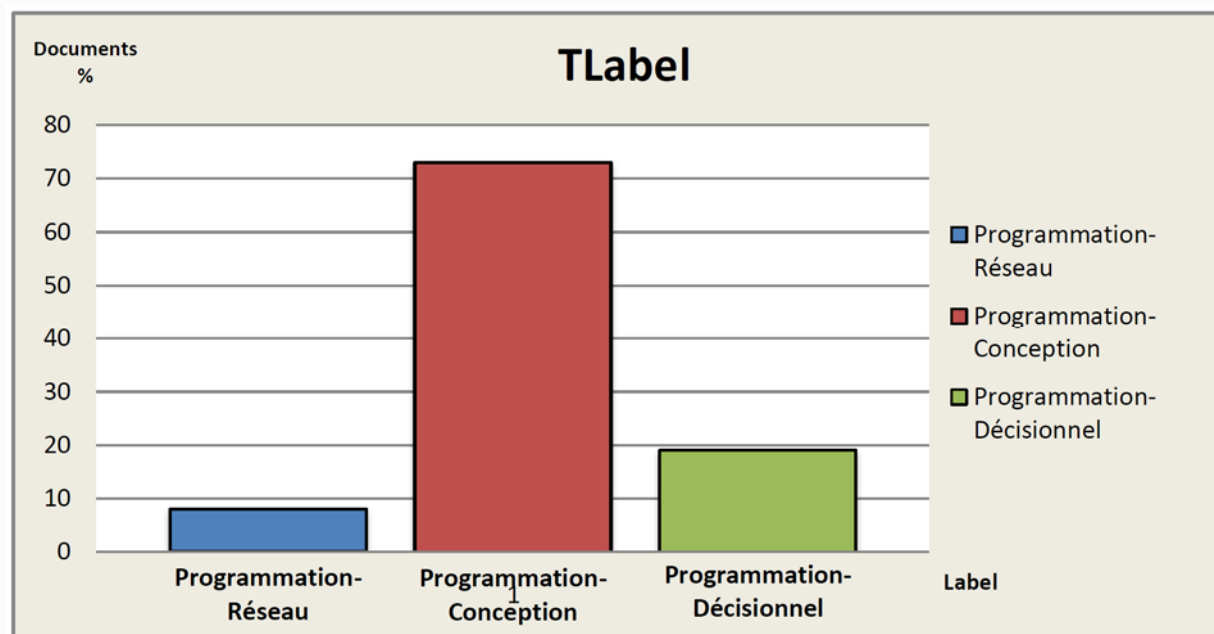
■ Labelling

<i>TLabel</i>	Programmation- Réseau	Programmation- Conception	Programmation- Décisionnel
<i>Ocluster</i>	CI 1, CI 7	CI 2, CI 3, CI 4, CI 6, CI 8	CI 5
Documents Number	99 (8%)	838 (73 %)	215 (19%)

Experiments & results

Results

- Query: <Topic= Computer Science, Location= France, Time= 2014 >

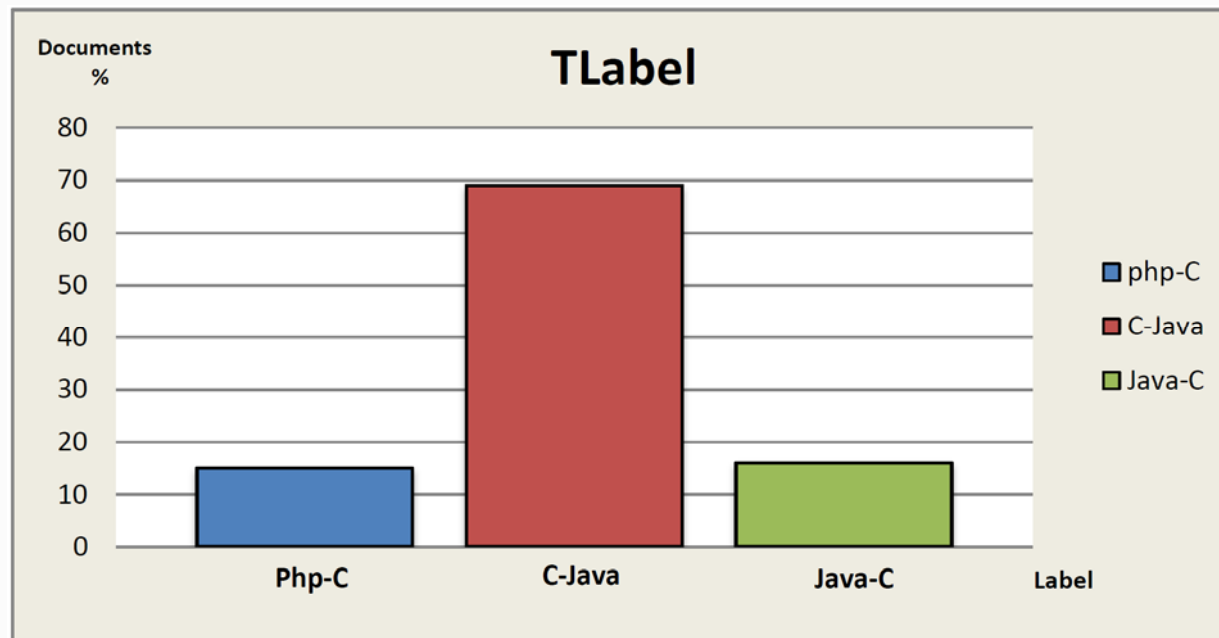


OLAP with TLabel

Experiments & results

Results

- Drill-down on Topic dimension



Drill-Down on Topic dimension with TLabel

Conclusion

- **TextLabel** : Clustering text documents in text OLAP systems

- **Text mining**
 - *Ocluster*: adapting K-means in OLAP environment

- **Documents clusters labelling**
 - Dresume
 - Domain ontology

- **Experiments on CV collections**

Future Work

- Think about other methods to obtain *Dresume document*
- Cluster Labelling according to several dimensions
- Evaluating TLabel with other text collections
- Validating TLabel with known labelled clusters