# WP 2 : Automatic data integration in multidimensional data warehouses

Yuzhao YANG

Franck RAVAT

Olivier TESTE

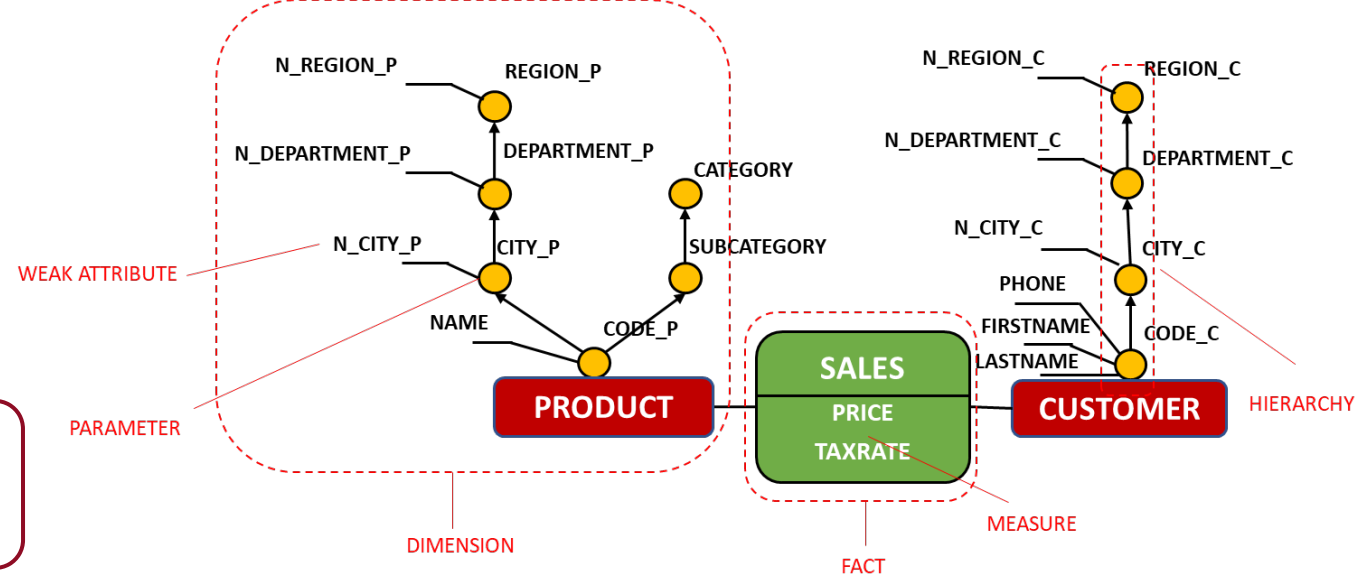Jérôme DARMONT

# Contents

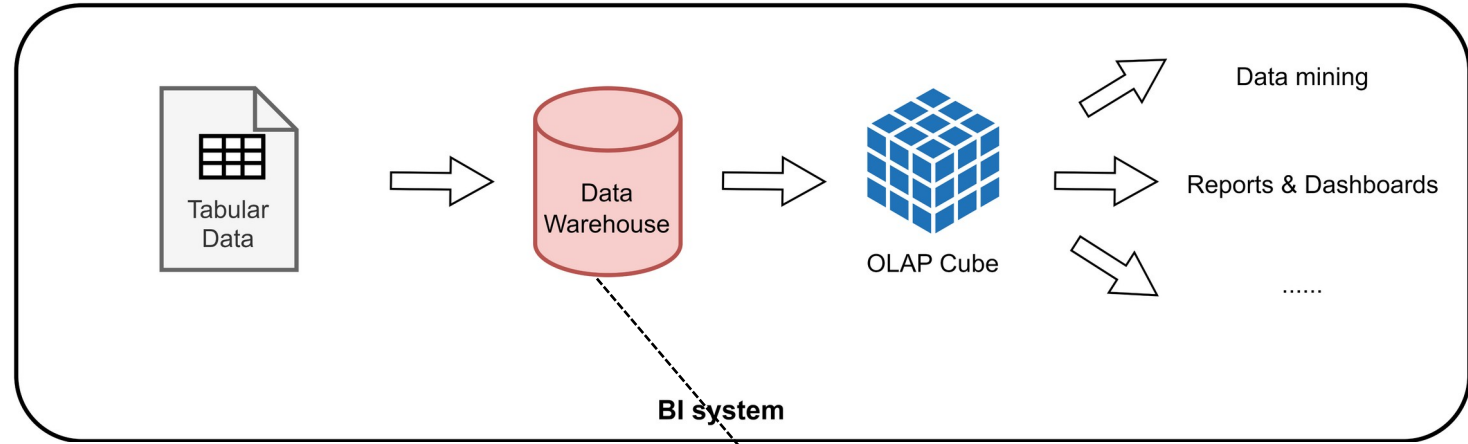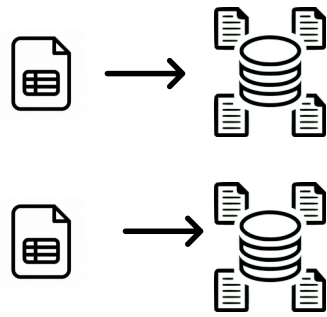**Large Companies** ✓
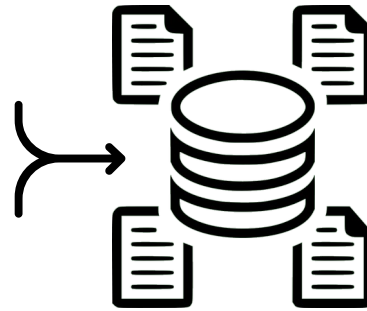
**Small Entities** ✗ Lack of budget and experts
✓

BI4people (Business intelligence for the people)

How can we automatically integrate tabular data integration in multidimensional data warehouses?

Automatic DW design and implementation

DW merging

Data imputation

Semi-automatic process
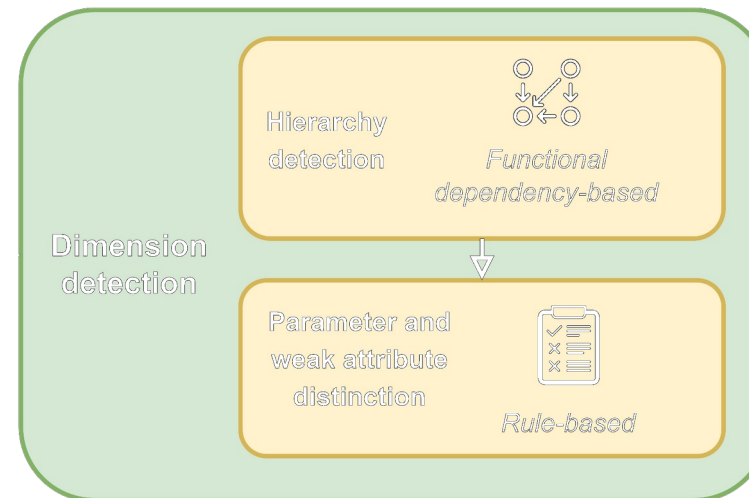
# Contents

**How to automatically generate a DW from tabular data?**

- **Lack of schema**
- **Complex DW structure**



Measure detection — *Machine learning-based*

Dimension detection

Hierarchy detection — *Functional dependency-based*

Parameter and weak attribute distinction — *Rule-based*

Data warehouse

# Automatic DW Design
## Measure Detection



| Feature Category | Feature |
|---|---|
| General feature | Data Type |
| | Positive/Negative/Zero value ratio |
| | Unique value ratio |
| Statistical feature | Same digital number |
| | Average/Minimum/Maximum/Median/Upper quartile/Lower quartile values |
| | Coefficient of variation |
| | Range ratio |
| Inter-column feature | Location ratio |
| | Numerical column ratio |
| | Multiple functional dependencies |
| | Numerical neighbor |

- **Support vector machine (SVM)**
- **Decision tree (DT)**
- **Random forest (RF)**
- **K-nearest neighbors (KNN)**

**Random forest**

- Best F-score
- Stable distribution

7

# Contents

How to merge two DWs having commun elements?

- Merging at both schema and instance levels
- Generation of different types of schema (star, constellation)

# DW Merging
## Contribution



Level merging (Schema) — Algo. 5 mergeHierarchies, p68

Hierarchy merging (Schema) — Algo. 5 mergeHierarchies, p68

Dimension merging (Schema and instance) — Algo. 6 mergeDimensions, p73

Fact merging (Schema and instance) — Algo. 8 mergeStars, p78

# Content

**How to carry out data imputation to ensure consistent analysis?**

- **Dimension**
- **Categorical**

- **Deducted values**
- **Predicted values**

**Hie - OLAPKNN**

**Hierarchical imputation** → **OLAPKNN imputation**

Functional dependencies in hierarchies

Non-parametric and instance-based

Suitable for different types of data

Relatively high accuracy

✓ Deducted values

✗ Limited

✓ Specific distance metric

✓ Consideration of dependencies

# Contributions on Automatic DW Design from Tabular Data

- **Mesure Detection**
  - *Machine learning classification*
  - *Random forest : +17%*
  - *Relevant features*
  - *Generic model*

- **Dimension Detection**
  - *Hierarchy: functional dependency*
  - *Parameter and weak attribute: rules*
  - *Dimension: 100%*
  - *Hierarchy: 67% - 100%*

# Contributions on Automatic DW Merging

- **DW merging**
  - *Schema and instance*
  - *Generation of star or constellation schema*

# Contributions on Data Imputation

- **Hie-OLAPKNN**
  - *Hierarchical imputation*
  - *OLAPKNN: specific distance*
  - *Effective : + 45%*
  - *Efficient : -+19 times*

# Contributions on Tabular Data Integration Application

- **Application**
  - *3 fonctionnalities*
  - *User-friendly interface*
  - *Non-expert and expert version*

# Publications

- **Automatic DW Design from Tabular Data**
  - Yuzhao Yang, Fatma Abdelhédi, Jérôme Darmont, Franck Ravat, Olivier Teste: Automatic Machine Learning-Based OLAP Measure Detection for Tabular Data. **DaWaK** 2022: 173-188
  - Yuzhao Yang, Jérôme Darmont, Franck Ravat, Olivier Teste: Automatic Integration Issues of Tabular Data for On-Line Analysis Processing. **EDA** 2020: 5-18
- *Automatic DW Merging*
  - Yuzhao Yang, Jérôme Darmont, Franck Ravat, Olivier Teste: An Automatic Schema-Instance Approach for Merging Multidimensional Data Warehouses**. IDEAS** 2021: 232-241
- *Data Imputation*
  - *Yuzhao Yang, Jérôme Darmont, Franck Ravat, Olivier Teste: Dimensional Data KNN-Based Imputation.* **ADBIS** *2022: 315-329*
  - *Yuzhao Yang, Fatma Abdelhédi, Jérôme Darmont, Franck Ravat, Olivier Teste: Internal Data Imputation in Data Warehouse Dimensions.* **DEXA** *(1) 2021: 237-244*

# PhD Thesis

- **Tabular data integration for multidimensional data warehouse**
  (https://theses.fr/2022TOU10052)

# Conclusion
Future work

| **Short-term plan** | • Imputation by External Sources |

| **Mid-term plan** | • Schema Evolution of Sources |

| **Long-term plan** | • Imputation Algorithm Generalization for ML algorithms <br> • Data Science with Data Lake <br>  • Merging / Matching Datasets <br>  • Imputing Data <br>  • User-friendly DL |

**Thank you!**