rapidminer

DS4DM
Data Search for Data Mining

UNIVERSITÄT
MANNHEIM

# Automated Mechanisms to Discover and Integrate Data from Web-based Tabular Collections

Edwin Yaqub, David Arnu, Ralf Klinkenberg – RapidMiner, Germany
Annalisa Gentile, Chris Bizer, Heiko Paulheim – Universität Mannheim, Germany

## Relevance & Research Question

*Keywords*: *data search, integration, enrichment*

Growth in data is phenomenal in recent years. Data exists on public sources such as the web or corporate intranet sources. However, in order to effectively utilize the information contained in this data, supporting tools are needed to discover new data and integrate it to existing datasets. This raises several research questions:

- How to access heterogeneous data that resides on the web.
- How to discover data that strongly matches an entity-query (comprising of an existing table and attributes to be extended).
- How to incorporate integration and fusion algorithms to control the discovered data (tabular collections) within data mining processes.

## Methods & Data

*Keywords*: *schema & instance matching, text corpus analysis, clustering*
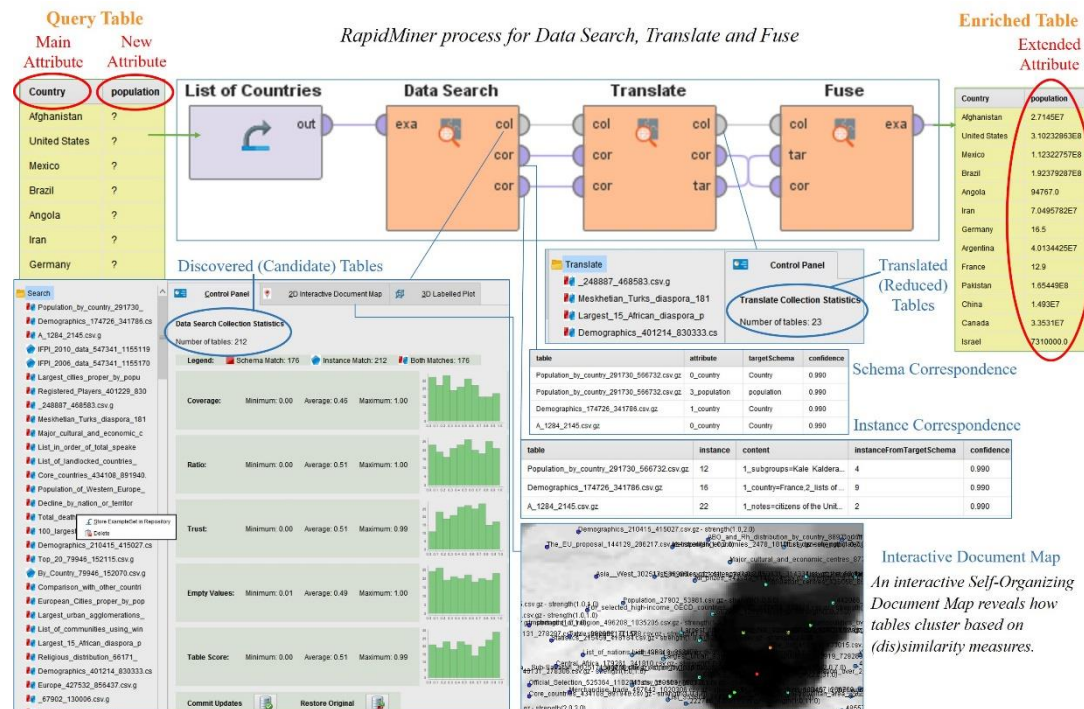
In project DS4DM [1], we developed Data Search extension for RapidMiner (an open source and free data mining tool). The extension provides Search-Join [2] as a data mining process for automated discovery and integration of tabular data.

A collection of half a million Wikipedia data tables is extracted and indexed in Lucene search engine for efficient retrieval. On this, the discovery algorithms compute schema and instance level matches for given query [3]. The result set is a space of candidate tables. The extension enables analyst to:

- Translate and fuse new data to existing tables using provided algorithms, or
- Perform manual refinements through i) exploratory visualizations (expose patterns among the vast corpus of tabular collections) and ii) graphical controls to manipulate intermediate outcomes of Search-Join process in real time (e.g. remove noisy tables) to prevent loss of valuable information.



### Results:

- Domain-independent data enrichment solution to harness web tables.
- Quality of search results is evaluated on the basis of various statistical metrics such as coverage, trust, ratio, empty values, table score.

### Added Value:

- RapidMiner platform has been extended to incorporate data discovery and integration methods as first class citizens of data mining processes.

References
[1] http:// ds4dm.de
[2] Extending RapidMiner with data search and integration capabilities. Gentile, Anna Lisa, et al., International Semantic Web Conference, 2016
[3] Entity Matching on Web Tables: a Table Embeddings approach for Blocking. Gentile, Anna Lisa, et al., (accepted submission), International Conference on Extending Database Technology, 2017