

A System For Large Scale Data Exploration and Organization

Luis Rei, Blaž Fortuna, Marko Grobelnik, and Dunja Mladenić

Artificial Intelligence Laboratory, Jožef Stefan Institute
and
Jožef Stefan International Postgraduate School

<http://github.com/lrei/elycite>

Objectives

- Understand and organize large datasets
- Visual
- Semi-Automatic: Man + Machine
- Make Machine Learning easy to use
- Collaborative
- Interactive
- Easy to integrate with other software & services
- Powerful – give expert users fine-grained control

Elycite

- *Elycites* information from unlabelled textual data
- Organizes information as an ontology...
 - Displays it using multiple visualizations
- ... Using Machine Learning interactively
- Multiple clients, one server
- Web Interface (HTML + JavaScript)
- Everything (machine learning, data access, organization, ...) that can be done interactively can also be done using a REST API (Web Services)

Built with QMiner

- A data analytics platform for processing large-scale real-time streams containing structured and unstructured data
- In-memory database with indexing, search, aggregations and Machine Learning built-in
- Text Mining, Social Networks, Time Series, ...
- JavaScript API, C++ library
- JSON-like data: Strings, Numbers, Booleans, Dates, Vectors
- Open Source
- Runs the server component of Elycite
- **<http://qminer.ijs.si>**

Elycite Basics

- **Documents** (Qminer Records)
- Organized into **Concepts** (sets of documents)
- Organized with relationships - **Ontology**
- Semi-automatic or manual
- E.g.:
 - **Documents**: News articles
 - **Concepts**: Politics, Economy, Entertainment, Sports, ...
 - **Hierarchy**: News -> Sports -> Football -> World Cup, Injuries, Commentary, FC Porto, ...

Unsupervised Methods

TF/IDF keyword extraction is used throughout the application to provide a summary of a set of documents (e.g. a concept).

KMeans++ based unsupervised concept suggestion.

The screenshot shows the OntoGen-QMiner interface. A dialog box titled "Suggest Sub Concepts for news" is open, displaying a table of suggested sub-concepts. The dialog has a "Number of sub concepts" dropdown set to 4 and a "Suggest" button. The table lists four suggestions, each with a name, keywords, document count, and an "Add" button. The background shows a concept map with nodes like "bus, children, crash", "gm, toyota, barra", "palestinian, talks, kerry", "pollard, release, israel", "app, facebook, company", "asia_economy_china_japan", "company, stocks, quarter", "colbert, letterman, cbs", "bank, banks, credit", and "film".

Name	Keywords	Docs	
police, court, state	police, court, state, people, government, year, party, mr, man, president	6476	Add
people, year, school	people, year, school, time, day, show, years, life, search, world	9140	Add
points, game, season	points, game, season, league, team, year, games, players, win, back	4231	Add
percent, company, year	percent, company, year, million, people, market, health, business, government, billion	9278	Add

Semi-Supervised Methods

SVM based Active Learning with user supervision provided initially by a query followed by a sequence of Yes or No questions. The concept is refined after each answer, the user can decide when to stop depending on how satisfied he is with the current suggestion (based on the keywords an number of documents).

The screenshot shows the OntoGen-QMiner web interface. A dialog box is open with the title "Does this document belong to the query whatsapp". The dialog contains a text block with a paragraph about the Kitestring app. Below the text is a section titled "Current Concept" which contains a table with columns for Name, Keywords, and Docs. The table has one row with the following data:

Name	Keywords	Docs	
MESSAGE, OBJECTIONABLE, RESPONSIBLE	MESSAGE, OBJECTIONABLE, RESPONSIBLE, APP, FACEBOOK, MESSAGING, WHATSAPP, USERS, MESSAGES, SERVICE	31	<input type="button" value="Add"/>

At the bottom of the dialog, there are three buttons: "Yes" (blue), "No" (white), and "Cancel" (red). The background interface shows a search bar with "news" entered, a "Build" button, and a horizontal tree view on the right side with various nodes like "colbert, letterman", "film", "facebook, whatsapp", etc.

Supervised Methods

A **SVM** classifier can be built from a concept and used to classify other sub-concepts in the same ontology, create concepts in another ontology, or as service to classify new documents for another application.

Classify news ×

Classifier

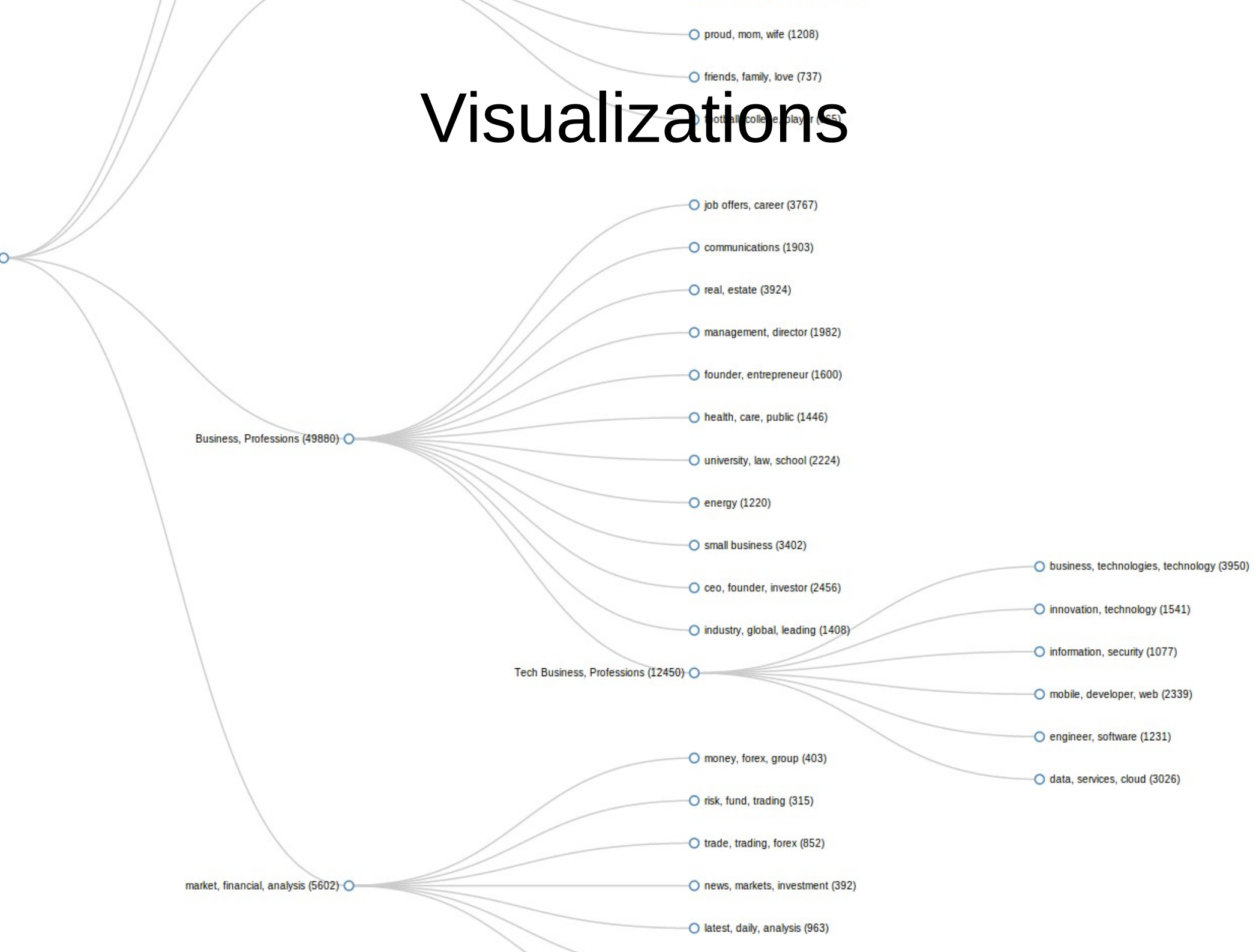
onto_news_sports_classifier ▾

Classify

Name	Keywords	Docs	
points, game, season	points, game, season, team, league, games, year, players, win, day	3893	Add
people, year, police	people, year, police, time, company, state, years, million, government, percent	25232	Add

Close

Visualizations



Future Work

- Interactively handle larger datasets
- More visualizations
- Improve text preprocessing (n-grams, hashing)
- Guided Learning
- Pre-built Classifiers (DMOZ topics, sentiment)
- Data: Numerical, Graph, Image, Time Series
 - Preprocessing
 - Feature extraction/learning
 - Machine Learning (Unsupervised, Semi-supervised, supervised)
 - Visualizations