# Using Redescriptions and Association Rules for Mining Definitions in Linked Data

Justine Reynaud, Yannick Toussaint and Amedeo Napoli
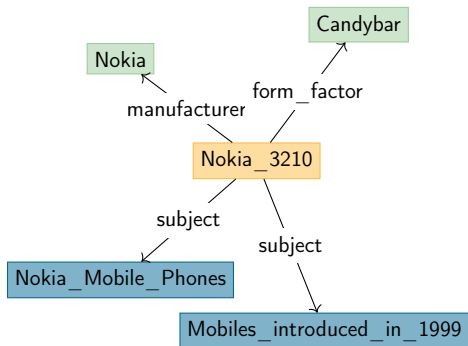
LORIA (Université de Lorraine, INRIA, CNRS), Vandœuvre-les-Nancy, France

« *The DBpedia Ontology is a shallow, cross-domain ontology, which has been **manually created** based on the most commonly used **infoboxes within Wikipedia**.* »

## Introduction — Problem statement

```
Nancy in France      Paris in France
Nancy in Europe      Paris in Europe
Nancy a City         Paris a City

Rome in Italy        Le_Louvre in France
Rome in Europe       Le_Louvre in Europe
Rome a City          Le_Louvre a Museum
```

$$French\_Cities = \{Paris, Nancy\}$$

**How to infer *definitions* in order to *complete* the web of data ?**

```
French_Cities ≡ (a, City) ⊓ (in, France)
```

# Data representation

# Data representation in FCA

```
Nancy in France    Paris in France
Nancy a City       Paris a City
Rome in Italy      Le_Louvre in France
Rome a City        Le_Louvre a Museum
```

|            | (in, France) 🇫🇷 | (in, Italy) 🇮🇹 | (in, Europe) 🇪🇺 | (a, City) 🏙️ | (a, Museum) 🏺 |
|------------|:---:|:---:|:---:|:---:|:---:|
| Nancy      | × |   | × | × |   |
| Rome       |   | × | × | × |   |
| Paris      | × |   | × | × |   |
| Le_Louvre  | × |   | × |   | × |

# Data representation in FCA

Nancy in France      Paris in France        French_Cities = {Paris, Nancy}
Nancy a City         Paris a City           Museums_in_Paris = {Le_Louvre}
Rome in Italy        Le_Louvre in France    European_Capital = {Paris, Rome}
Rome a City          Le_Louvre a Museum

|  | (in, France) 🇫🇷 | (in, Italy) 🇮🇹 | (in, Europe) 🇪🇺 | (a, City) 🏰 | (a, Museum) 🏺 | French_Cities FC | Museums_in_Paris MP | European_Capital EC |
|---|---|---|---|---|---|---|---|---|
| Nancy | × |  | × | × |  | × |  |  |
| Rome |  | × | × | × |  |  |  | × |
| Paris | × |  | × | × |  | × |  | × |
| Le_Louvre | × |  | × |  | × |  | × |  |

# Derivation operators and concepts

| | 🇫🇷 | 🇮🇹 | 🇪🇺 | 🏙 | 🏺 | FC | MP | EC |
|---|---|---|---|---|---|---|---|---|
| Nancy | × | | × | × | | × | | |
| Rome | | × | × | × | | | | × |
| Paris | × | | × | × | | × | | × |
| Le_Louvre | × | | × | | × | | × | |

$$\{\texttt{Nancy}\}' = \{🇫🇷, 🇪🇺, 🏙, \text{FC}\}$$
$$\{🇫🇷, 🏙\}' = \{\texttt{Nancy, Paris}\}$$

## Concept

Given $A \subseteq G$ and $B \subseteq M$, the pair $(A, B)$ is a concept if $A' = B$ and $B' = A$.

**({Nancy, Paris}, { 🇫🇷, 🇪🇺, 🏙, FC }) is a concept.**

# Concept lattice



**Concepts are partially ordered.**

- Implication:

  $$(in, France) \Rightarrow (in, Europe)$$

- Definition:

  $$(a, Museum) \Leftrightarrow MP$$

- Association rule:

  $$(in, France) \rightarrow (a, City)$$

# Association rules and redescriptions

# Association rules

- Searching for dependencies between sets of attributes
- Quality metrics based on confidence

$$conf(X \to Y) = \frac{|\, X' \cap Y' \,|}{|\, X' \,|}$$

# Association rules

- Searching for dependencies between sets of attributes
- Quality metrics based on confidence

$$conf(X \to Y) = \frac{\mid X' \cap Y' \mid}{\mid X' \mid}$$

## Example

| | 🇫🇷 | 🇮🇹 | 🇪🇺 | 🏰 | 🏺 | FC | MP | EC |
|---|---|---|---|---|---|---|---|---|
| Nancy | × | | × | × | | × | | |
| Rome | | × | × | × | | | | × |
| Paris | × | | × | × | | × | | × |
| Le_Louvre | × | | | × | | × | × | |

$$conf(\{🇫🇷\} \to \{🏰\}) = \frac{\mid 🇫🇷' \cap 🏰' \mid}{\mid 🇫🇷' \mid} = \frac{\mid \{Nancy, Paris\} \mid}{\mid Nancy, Paris, Le\_Louvre \mid} = \frac{2}{3}$$

# Association rules – Eclat [Zaki, 2000]

- Exhaustive enumeration
- Rules are unidirectional
- Post-processing in order to select rules satisfying criteria

## Quasi-definition

A quasi-definition $X \leftrightarrow Y$ holds with a confidence $\theta$ iff

$$min(conf(X \to Y), conf(Y \to X)) = \theta$$

# Redescriptions – ReReMi [Galbrun and Miettinen, 2012]

- Searching for two sets of attributes that occurs in the same objects
- Rules are bidirectional and more expressive than association rules
- Quality metrics based on Jaccard coefficient

$$Jacc(X \leftrightarrow Y) = \frac{|\, X' \cap Y' \,|}{|\, X' \cup Y' \,|}$$

# Redescriptions – ReReMi [Galbrun and Miettinen, 2012]

- Searching for two sets of attributes that occurs in the same objects
- Rules are bidirectional and more expressive than association rules
- Quality metrics based on Jaccard coefficient

$$Jacc(X \leftrightarrow Y) = \frac{\mid X' \cap Y' \mid}{\mid X' \cup Y' \mid}$$

## Example

| | 🇫🇷 | 🇮🇹 | 🇪🇺 | 🏙️ | 🏺 | FC | MP | EC |
|---|---|---|---|---|---|---|---|---|
| `Nancy` | × | | × | × | | × | | |
| `Rome` | | × | × | × | | | | × |
| `Paris` | × | | × | × | | × | | × |
| `Le_Louvre` | × | | × | | × | | × | |

$$\{🇫🇷\} \leftrightarrow \{FC\} \qquad Jacc(\{🇫🇷\} \leftrightarrow \{FC\}) = \frac{2}{3}$$

- Searching for two sets of attributes that occurs in the same objects
- Rules are bidirectional and more expressive than association rules
- Quality metrics based on Jaccard coefficient

$$Jacc(X \leftrightarrow Y) = \frac{|\ X' \cap Y'\ |}{|\ X' \cup Y'\ |}$$

## Example

|  | 🇫🇷 | 🇮🇹 | 🇪🇺 | 🏙️ | 🏺 | FC | MP | EC |
|---|---|---|---|---|---|---|---|---|
| Nancy | × |  | × | × | × |  |  |
| Rome |  | × | × | × |  |  |  | × |
| Paris | × |  | × | × | × |  |  | × |
| Le_Louvre | × |  | × |  | × |  | × |  |

$\{🇫🇷, 🏙️\} \leftrightarrow \{FC\}$     $Jacc(\{🇫🇷, 🏙️\} \leftrightarrow \{FC\}) = \frac{2}{2} = 1$

# Experiments

- Datasets extracted from DBpedia, thanks to a SPARQL query
- Various sizes and domains

|        | Person              | Object                      | Film           |
|--------|---------------------|-----------------------------|----------------|
| Small  | Turing_Award        | Samsung_Galaxy              | Hospital_films |
| Medium | Women_Mathematicians | Smartphones<br>Sports_cars | Road_movies    |
| Large  | Mathematicians      | —                           | French_films   |

# Experiments : Datasets (statistics)

| Dataset | Triples | $|G|$ | $|M|$ | $M_{subj}$ | $M_{descr}$ | $|P|$ | $\delta$ |
|---|---|---|---|---|---|---|---|
| Samsung_Galaxy | 940 | 59 | 277 | 30 | 247 | 33 | 5.2e−2 |
| Turing_Award_laureates | 2642 | 65 | 1360 | 503 | 857 | 35 | 2.2e−2 |
| Hospital_films | 1984 | 71 | 1265 | 490 | 775 | 46 | 1.6e−2 |
| Women_mathematicians | 9652 | 552 | 4243 | 1776 | 2467 | 98 | 2.9e−3 |
| Smartphones | 8418 | 598 | 2089 | 359 | 1730 | 98 | 5.8e−3 |
| Sports_cars | 9047 | 604 | 2730 | 435 | 2295 | 61 | 4.7e−3 |
| Road_movies | 20056 | 689 | 9314 | 2652 | 6662 | 103 | 2.4e−3 |
| Mathematicians | 32536 | 1660 | 12279 | 3848 | 8431 | 202 | 1.2e−3 |
| French_films | 121496 | 6039 | 25487 | 6028 | 19459 | 111 | 6.4e−4 |

### Association Rules

| | | |
|---|---|---|
| **Harvard_University_alumni** | ≡ | ∃almaMater.Harvard_University ⊓ Agent ⊓ Person ⊓ Scientist |
| **Harvard_University_alumni** | ≡ | ∃almaMater.Harvard_University ⊓ ∃award.Turing_Award ⊓ Agent ⊓ Person ⊓ Scientist |
| **National_Medal_of_Science_l** | ≡ | ∃award.National_Medal_of_Science ⊓ Agent ⊓ Person ⊓ Scientist |
| **M._I._T._faculty** | ≢ | ∃award.Turing_Award ⊓ Agent ⊓ Person ⊓ ∃birthPlace.New_York_City |

### Redescriptions

| | | |
|---|---|---|
| **Harvard_University_alumni** | ≡ | ∃almaMater.Harvard_University |
| **Stanford_University_alumni** | ≡ | ∃almaMater.Stanford_University |
| **National_Medal_of_Science_l.** | ≡ | ∃award.National_Medal_of_Science |
| **British_computer_scientists** | ≢ | ∃award.Fellow_of_the_Royal_Society |

# Experiments: extracted rules — Smartphones

## Association Rules

| | | |
|---|---|---|
| **Nokia_mobile_phones** | ≡ | ∃manufacturer.Nokia ⊓ Device |
| **Samsung_Galaxy** | ≡ | ∃manufacturer.Samsung_Electronics ⊓ Smartphones ⊓ Device |
| **Mobile_operating_systems** | ≡ | Software ⊓ Work |
| **Sony_mobile_phones** | ≢ | ∃input.Capacitive_sensing ⊓ ∃input.Proximity_sensor ⊓ ∃input.Touchscreen |

## Redescriptions

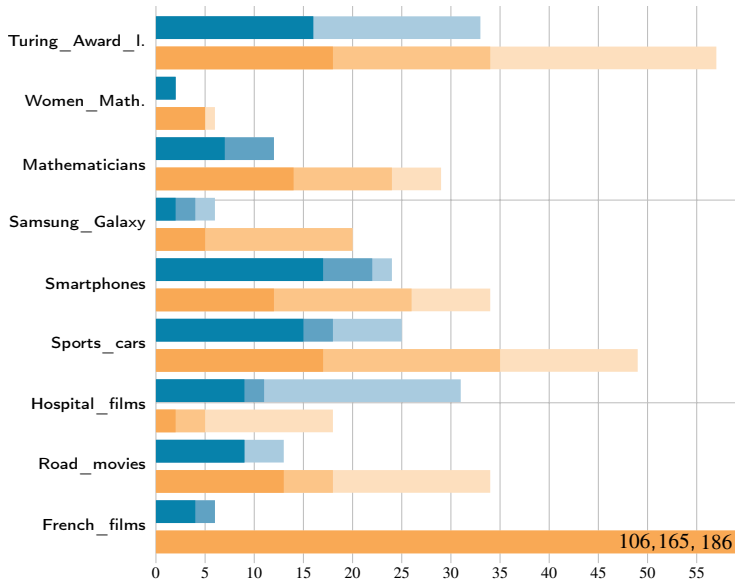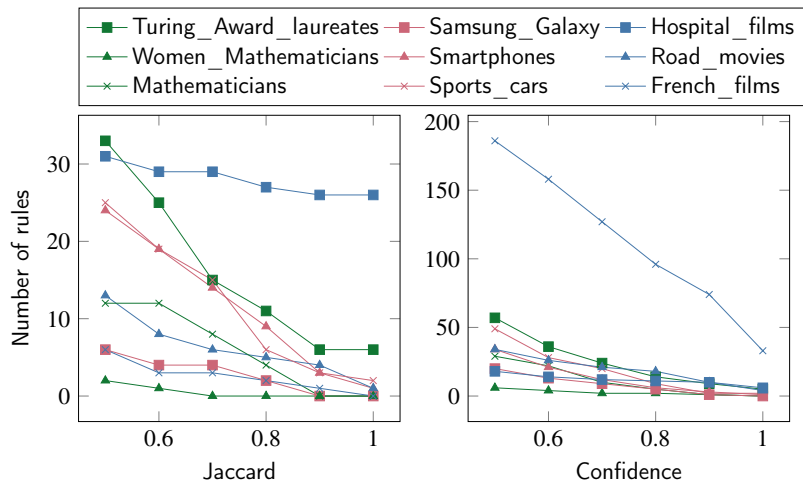| | | |
|---|---|---|
| **Nokia_mobile_phones** | ≡ | ∃manufacturer.Nokia |
| **Samsung_Galaxy** | ≡ | ∃manufacturer.Samsung_Electronics ⊓ ∃operatingSystem.Android_OS |
| **Mobile_operating_systems** | ≡ | Software ⊓ Work |
| **MeeGo_Devices** | ≢ | ∃operatingSystem.Sailfish_OS |

# Experiments : Rules extracted (statistics)

# Experiments : Rules extracted (statistics)
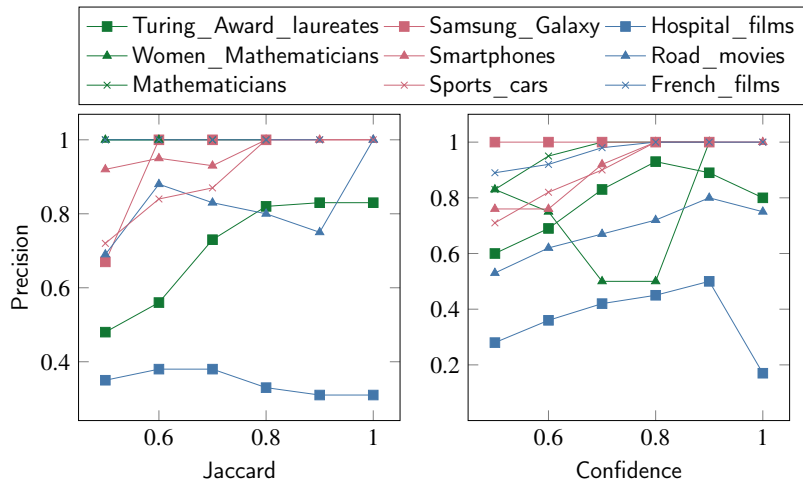
# Experiments : Rules extracted (statistics)
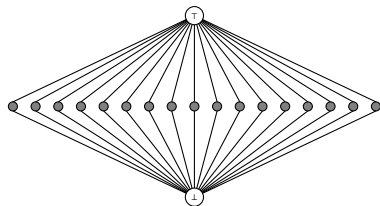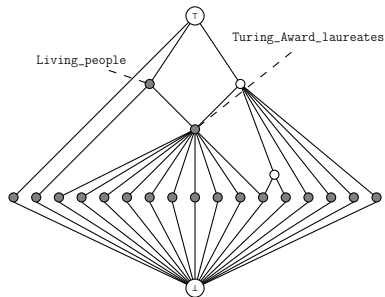
# Discussion – Ordering categories

Concept lattices of the defined categories (`Turing_Award_laureates`).



Redescriptions

Association rules

## Definition of `Harvard_University_alumni`

| | |
|---|---|
| **Red.** | ∃almaMater.Harvard_University |
| **A.R.** | ∃almaMater.Harvard_University ⊓ Agent ⊓ Person ⊓ Scientist |

## Discussion – Predicates

|  | Pred. | Pred ($D_{RD}$) | Pred ($D_{QD}$) |
|---|---|---|---|
| Turing_Award_laureates | 35 | 4 | 5 |
| Women_Mathematicians | 98 | 2 | 3 |
| Mathematicians | 202 | 4 | 5 |
| Samsung_Galaxy | 33 | 2 | 7 |
| Smartphones | 98 | 5 | 8 |
| Sports_cars | 61 | 3 | 5 |
| Hospital_films | 46 | 5 | 2 |
| Road_movies | 103 | 3 | 7 |
| French_films | 111 | 4 | 10 |

**dbo:award (261), *rdf:type (186)***, dbo:knownFor (182), dbo:doctoralStudent (148),
**dbp:workInstitution (123)**, dbo:birthPlace (117), **dbo:almaMater (110), dbo:field (84)**,
dbo:doctoralAdvisor (36), dbo:deathPlace (36), dbp:workplaces (28),
dbp:workInstitutions (26), dbo:influenced (23), dbo:nationality (15)

# Conclusion and Future work

- Redescriptions interesting for defining categories
- Association Rules and Redescriptions complete each other

- *Are definitions operational ?*
  Integration to knowledge base.
- *Can we (should we) use more expressive definitions ?*
  $C \equiv A \sqcup B$ or $C \equiv \neg A$

# Thanks for your attention.
# Questions ?

justine.reynaud@loria.fr

📄 Galbrun, E. and Miettinen, P. (2012).
From Black and White to Full Color: Extending Redescription Mining Outside the Boolean World.
*Statistical Analysis and Data Mining*, 5(4):284–303.

📄 van Leeuwen, M. and Galbrun, E. (2015).
Association Discovery in Two-View Data.
*TKDE*, 27(12):3190–3202.

📄 Zaki, M. J. (2000).
Scalable algorithms for association mining.
*TKDE*, 12(3):372–390.