

Three approaches to find definitions in RDF data

Justine Reynaud, Yannick Toussaint and Amedeo Napoli

LORIA (Université de Lorraine, INRIA, CNRS), Vandœuvre-les-Nancy, France

FCA4AI — July 14, 2018



Introduction

From RDF assertions, such as

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

French_Cities = {Paris, Nancy}

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

French_Cities \equiv (a, City) \sqcap (in, France)

We compare three algorithms with different approaches.

Data representation

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

	<i>(in, France)</i>	<i>(in, Italy)</i>	<i>(a, City)</i>	<i>(a, Museum)</i>
Nancy	×		×	
Rome		×	×	
Paris	×		×	
Le_Louvre	×			×

Data representation

Nancy in France
Nancy a City
Rome in Italy
Rome a City

Paris in France
Paris a City
Le_Louvre in France
Le_Louvre a Museum

French_Cities = {Paris, Nancy}
Museums_in_Paris = {Le_Louvre}
European_Capital = {Paris, Rome}

	<i>(in, France)</i>	<i>(in, Italy)</i>	<i>(a, City)</i>	<i>(a, Museum)</i>	<i>French_Cities</i>	<i>Museums_in_Paris</i>	<i>European_Capital</i>
Nancy	x		x		x		
Rome		x	x				x
Paris	x		x		x		x
Le_Louvre	x			x		x	

Data representation

Nancy in France
Nancy a City
Rome in Italy
Rome a City

Paris in France
Paris a City
Le_Louvre in France
Le_Louvre a Museum

French_Cities = {Paris, Nancy}
Museums_in_Paris = {Le_Louvre}
European_Capital = {Paris, Rome}

	<i>(in, France)</i>	<i>(in, Italy)</i>	<i>(a, City)</i>	<i>(a, Museum)</i>	<i>French_Cities</i>	<i>Museums_in_Paris</i>	<i>European_Capital</i>
Nancy	×		×		×		
Rome		×	×				×
Paris	×		×		×		×
Le_Louvre	×			×		×	

$\{Nancy\}' = \{(in, France), (a, City), French_Cities\}$

$\{(in, France), (a, City)\}' = \{Nancy, Paris\}$

Association rules – Eclat [Zaki, 2000]

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

$$\text{conf}(X \rightarrow Y) = \frac{|X' \cap Y'|}{|X'|}$$

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

Quasi-definition

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

$$\min(\text{conf}(X \rightarrow Y), \text{conf}(Y \rightarrow X)) = \theta$$

- Searching for two sets of attributes that occurs in the same objects
- Quality measure based on Jaccard coefficient

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

- Rules are bidirectional

- Searching for two sets of attributes that occurs in the same objects
- Quality measure based on Jaccard coefficient

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

- Rules are bidirectional

	FR	IT	City	Museum	FC	MP	EC
Nancy	×		×		×		
Rome		×	×				×
Paris	×		×		×		×
Le_Louvre	×			×		×	

(in, France) \leftrightarrow French_Cities

- Searching for two sets of attributes that occurs in the same objects
- Quality measure based on Jaccard coefficient

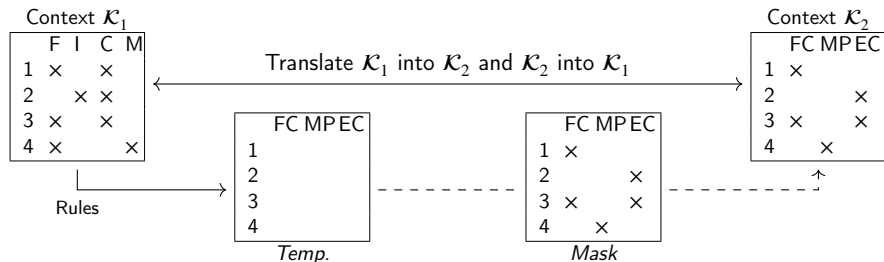
$$Jacc(X \leftrightarrow Y) = \frac{|X' \cap Y'|}{|X' \cup Y'|}$$

- Rules are bidirectional

	FR	IT	City	Museum	FC	MP	EC
Nancy	×		×		×		
Rome		×	×				×
Paris	×		×		×		×
Le_Louvre	×			×		×	

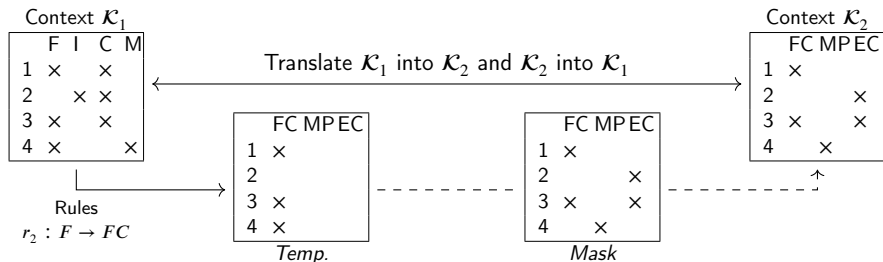
(in, France), (a, City) \leftrightarrow French_Cities

Translation rules – Translator [van Leeuwen and Galbrun, 2015]



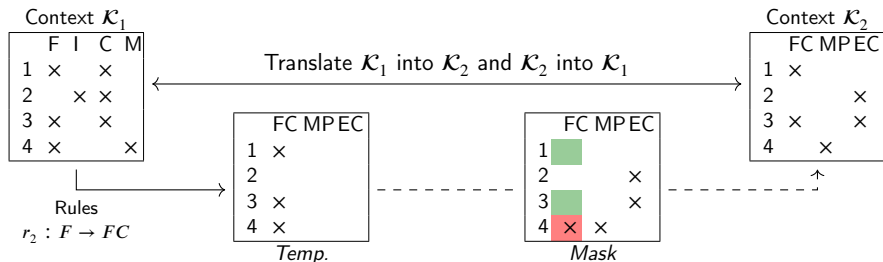
- Searching for rules that allow to construct one context from the other
- Rules may be unidirectional or bidirectional

Translation rules – Translator [van Leeuwen and Galbrun, 2015]



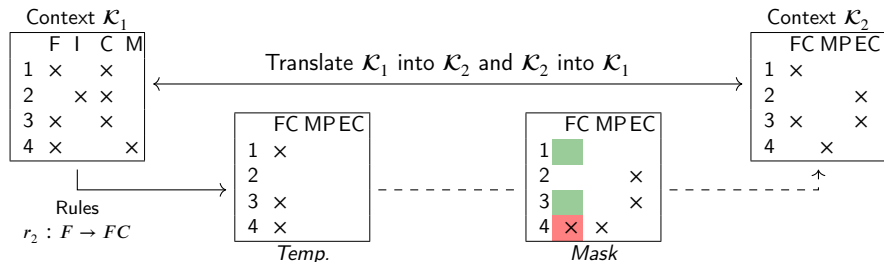
- Searching for rules that allow to construct one context from the other
- Rules may be unidirectional or bidirectional

Translation rules – Translator [van Leeuwen and Galbrun, 2015]



- Searching for rules that allow to construct one context from the other
- Rules may be unidirectional or bidirectional

Translation rules – Translator [van Leeuwen and Galbrun, 2015]



- Searching for rules that allow to construct one context from the other
- Rules may be unidirectional or bidirectional
- Adds the best rule at each step
- Quality metric inspired from *minimum description length* (MDL)

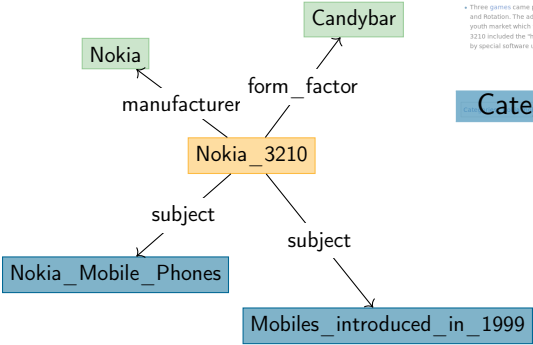
$$\Delta(X \rightarrow Y) = \underbrace{L(\text{Mask}^-) - L(\text{Mask}^+)}_{\text{Information gain}} - \underbrace{L(X \cup Y)}_{\text{Rule length}} \quad L(X) = - \sum_{x \in X} \log_2 P(x | \mathcal{K})$$

Algorithms: Comparison

	Eclat	ReReMi	Translator
Data	Bool.	Bool., Num., Cat.	Bool.
Quality measure	Confidence $\frac{ X' \cap Y' }{ X' }$	Jaccard $\frac{ X' \cap Y' }{ X' \cup Y' }$	Compression based on MDL
Symmetric rule	No	Yes	Both

- Eclat needs a post-process to build bi-directional rules
- ReReMi and Eclat compute confidence in a very similar way
ReReMi should return a subset of the rules found by Eclat
- Translator aims to mine a good set of rules instead of a set of good rules

From Wikipedia to DBpedia



Resource name

From Wikipedia, the free encyclopedia

The **Nokia 3210** is a *OSM cellular phone*, announced by Nokia on March 18, 1999.^[1] With 160 million units sold,^[2] the 3210 is one of the most popular and successful phones in history.

Contents [show]

Design [edit]

The Nokia 3210 has a total weight of 153g. The handset measures 123.8mm x 50.5mm x 16.7mm (min), 22.5mm (max) and features customizable fascias which clip on. It was the first mass-market phone with an internal antenna, after the feature had been introduced by Nokia on the luxury phone 8810 in 1998. The 3210 was designed by the Nokia Design Center in Nokia's Los Angeles Design Center.^[4]

Notable features [edit]

- Three games came preinstalled: *Snake*, *Memory (pairs-memory game)*, and *Rotation*. The addition of such games encouraged high sales within a youth market which was enlarging at a very fast rate. Some versions of the 3210 included the "hidden" games *React* and *Logic*. They were activated by special software using a data cable.

Nokia 3210



Manufacturer	Nokia
Compatible networks	GSM / GPRS
Availability	1999
Price	\$150
Successor	Nokia 3210
Form factor	Candybar
Dimensions	123.8 x 50.5 x 16.7-22.5 mm
Weight	153g
Memory	Up to 250 names in phonebook
Battery	1350 mAh
Display	Backlit Monochrome
Rear camera	None
Connectivity	None

Categories

Experiment

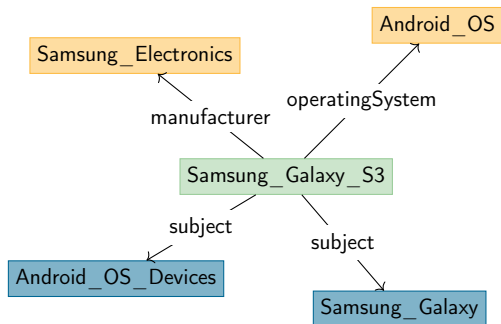
Datasets Triples from four domains of DBpedia

Turing_Award_laureates, Smartphones, Sports_cars, French_Films

Objects Subjects of the triples

Categories Pairs (*subject*, *C*) from the categories

Descriptions Pairs (*p*, *o*) from the other triples



Smartphones

8500 Triples

600 Resources

400 Categories

1800 Descriptions

Results

- R **Samsung_Galaxy**
(manufacturer Samsung_Electronics), (operatingSystem Android_(operating_system))
- ET **Samsung_Galaxy, Samsung_mobile_phones, Smartphones**
(a Device), (manufacturer Samsung_Electronics), (operatingSystem Android_(operating_system))

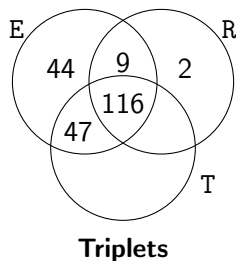
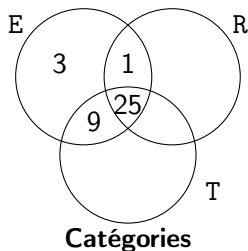
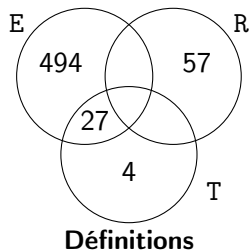
Smartphones				
	X	E	R	T
	$ \mathcal{R}^X $	810	98	41
	$ \mathcal{D}^X $	521	57	31
	Précision	.64	.58	.76
	$\overline{ C_i } - \overline{ D_i }$	4.3	1.6	3.1
	$ C_i - D_i $	7.8	1.8	3.1

Results

R **Samsung_Galaxy**

(manufacturer Samsung_Electronics), (operatingSystem Android_(operating_system))

ET **Samsung_Galaxy, Samsung_mobile_phones, Smartphones**
(a Device), (manufacturer Samsung_Electronics), (operatingSystem Android_(operating_system))



How to include domain knowledge like classes and/or predicates hierarchy?

i.e. dealing with a partial order on the attributes

Can we find class disjointness instead of definitions?

i.e. searching for rules with a very low quality measure

How to deal with scalability?

i.e. evaluating a huge amount of rules

Thanks. Questions ?

Eclat

I want to be exhaustive, no matter if they're is a lot of equivalent rules.




ReReMi

I want a few rules easy to interpret and it's important they're valid.

Translator

I want a small set of rules representing the whole set of data, even if it's more difficult to interpret.

References

-  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.

SPARQL Query

```
SELECT DISTINCT ?s ?p ?o WHERE {
  ?s ?p ?o .
  ?s dct:subject dbc:Smartphones .
  ?p a owl:ObjectProperty .
  FILTER (isURI(?o))
  FILTER (!STRSTARTS(STR(?o), "http://www.wikidata.org/"))
  FILTER (!STRSTARTS(STR(?o), "http://dbpedia.org/class/yago/"))
  FILTER (!STRSTARTS(STR(?p), "http://xmlns.com/foaf/0.1/"))
  FILTER (
    (?p != dbp:wordnet_type) AND (?p != dbp:website)
    AND (?p != prov:wasDerivedFrom) AND (?p != dbo:thumbnail)
    AND (?p != rdfs:comment) AND (?p != rdfs:label)
    AND (?p != rdfs:seeAlso) AND (?p != owl:sameAs)
    AND (?p != owl:differentFrom) AND (?p != foaf:depiction)
    AND (?p != dbo:wikiPageExternalLink)
  )
}
```

The query was run on DBpedia 2016-04.

- Searching for a set of rules that enable to construct one context from the other
- Greedy approach : adds the better rule at each step
- Quality measure based on *minimum description length* :

$$\Delta(X \rightarrow Y) = \underbrace{L(\text{Mask}^-) - L(\text{Mask}^+)}_{\text{Information gain}} - \underbrace{L(X \cup Y)}_{\text{Rule length}}$$

$$L(X) = - \sum_{x \in X} \log_2 P(x | \mathcal{K})$$

- Rules may be unidirectional or bidirectional

Statistiques sur les jeux de données extraits

	Triplets	Objets	Attributs	
			Cat.	Descr.
Turing_Award	2 642	65	503	857
Smartphones	8 418	598	359	1 730
Sports_cars	9 047	604	435	2 295
French_films	121 496	6 039	6 028	19 459

Results

Turing_Award_laureates

X	E	R	T
$ \mathcal{R}^X $	47	12	11
$ \mathcal{D}^X $	30	9	9
Précision	.64	.75	.85
$\overline{ C_i } - \overline{ D_i }$	2	1	3 5
$ C_i - D_i $	4	1	5

Smartphones

X	E	R	T
$ \mathcal{R}^X $	810	98	41
$ \mathcal{D}^X $	521	57	31
Précision	.64	.58	.76
$\overline{ C_i } - \overline{ D_i }$	4.3	1.6	3.1
$ C_i - D_i $	7.8	1.8	3.1

Sports_cars

X	E	R	T
$ \mathcal{R}^X $	132	52	31
$ \mathcal{D}^X $	95	30	23
Précision	.72	.68	.74
$\overline{ C_i } - \overline{ D_i }$	2.8	1.3	2.6
$ C_i - D_i $	4.5	1.4	4.1

French_films

X	E	R	T
$ \mathcal{R}^X $	546	36	93
$ \mathcal{D}^X $	371	12	89
Précision	.68	.33	.96
$\overline{ C_i } - \overline{ D_i }$	2.8	1.2	2.3
$ C_i - D_i $	4.4	1.1	4.2

Evaluation

Three expert evaluated each rule as True or False.

From the rules evaluated True, we build a rules base \mathcal{D} of 20 rules.

We say that $X \leftrightarrow Y$ covers $A \leftrightarrow B$ iff $A \subseteq X$ and $Y \subseteq B$.

Given a set of k rules returned by the algorithm X , we can compute the precision and the recall of those rules wrt the rule base :

$$\text{recall}(X) = \frac{|\text{cov}(\mathcal{D}, \mathcal{R}_X)|}{|\mathcal{D}|}$$

$$\text{precision}(X) = \frac{|\{R \in \mathcal{R}_X \mid \exists D \in \mathcal{D}, R \text{ covers } D\}|}{|\mathcal{R}_X|}$$

where $|\text{cov}(\mathcal{D}, \mathcal{R}_X)|$ is the number of rules from \mathcal{D} covered by a rule of \mathcal{R}_X .