



Universität des Saarlandes
Max-Planck-Institut für Informatik
AG5



Learning Rules With Categorical Attributes from Liked Data Sources

Masterarbeit im Fach Informatik
Master's Thesis in Computer Science
von / by

André de Oliveira Melo

angefertigt unter der Leitung von / supervised by

Prof. Dr. Gerhard Weikum

betreut von / advised by

Dr. Martin Theobald

begutachtet von / reviewers

Dr. Max Mustermann

Prof. Dr. Gerhard Weikum

November / November 2012

Hilfsmittelerklärung

Hiermit versichere ich, die vorliegende Arbeit selbständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt zu haben.

Non-plagiarism Statement

Hereby I confirm that this thesis is my own work and that I have documented all sources used.

Saarbrücken, den 22. November 2012,

(André de Oliveira Melo)

Einverständniserklärung

Ich bin damit einverstanden, dass meine (bestandene) Arbeit in beiden Versionen in die Bibliothek der Informatik aufgenommen und damit veröffentlicht wird.

Declaration of Consent

Herewith I agree that my thesis will be made available through the library of the Computer Science Department, Saarland University.

Saarbrücken, den 22. November 2012,

(André de Oliveira Melo)

To my father Cicero, my mother Marlene and my sister Carolina

- Andre

Abstract

With millions of articles in multiple languages, Wikipedia has become the de-facto source of reference on the Internet today. Each article on Wikipedia contains encyclopedic information about various topics (people, events, inventions, etc.) and implicitly represents an entity. Extracting the most important facts about such entity will help users to find desired information more quickly and effectively. However, this task is challenging due to the incomplete and noisy nature of Wikipedia articles. This calls for a mechanism to detect and summarize the most important information about an entity on Wikipedia.

This thesis proposes and implements CATE (**C**ontext-**A**ware **T**imeline for **E**ntity Exploration), a framework that utilizes Wikipedia to summarize and visualize the important aspects of entities in a timeline fashion. Such a system will help users to draw quickly an informative picture of an entity (e.g. life of a person, or evolution of a research topic, etc.). The novelty of CATE lies in seeing the entity in different contexts, synchronous with contemporaneous events. In addition, CATE puts the entity in a relationship with other entities, and thus offers a broader portrait about it. In order to efficiently query and visualize the events related to the entity, a number of techniques have been developed, combining information extraction and information retrieval with a novel ranking model. The thesis also discusses several experiments and evaluation results to show the effectiveness of the methods proposed.

Acknowledgements

Firstly, I would like to thank my advisor Dr. Martin Theobald, for his invaluable guidance. I feel deeply grateful for his technical assistance and motivational encouragement.

A special note of thanks to Prof. Gerhard Weikum for giving me the opportunity to pursue this thesis at Information and Database Systems department under his supervision. It was a very enriching and pleasant experience to write my Master Thesis here.

Contents

Abstract	vi
Acknowledgements	viii
1 Introduction	1
1.1 Motivation	2
1.2 Contributions	4
1.3 Outline	4
2 Related Work	7
2.1 Logic Programming	7
2.2 Inductive Logic Programming	7
2.3 Mining Optimized Rules for Numeric Attributes	7
2.4 Minimum Description Length	7
2.5 Semantic Web	7
2.6 Linked Open Data	7
3 Correlation Lattice	9
3.1 Categorical Relation Definition	10
3.2 Support	11
3.2.1 Independence between Nodes	11
3.3 Heuristics	12
4 Algorithmic Framework	13
4.1 Knowledge Base Backend	13
4.2 Preprocessing	13
4.2.1 Relation Preprocessing	13
4.2.2 Correlation Lattice	14
4.3 ILP Core Algorithm	15
List of Figures	15
List of Tables	19

Bibliography

21

Chapter 1

Introduction

In the last years, the volume of semantic data available, in particular RDF, has dramatically increased. Initiatives like the W3C Semantic Web, which provides a common standard that allows data to be shared and reused across different applications, and the Linked Open Data, which provides linkages between different datasets that were not originally interconnected, have great contribution in such development. Moreover, advances in information extraction have also made strong contribution, by crawling multiple non-structured resources in the Web and extracting RDF facts.

Nevertheless, information extraction still has its limitations and many of sources might contain contradictory or uncertain information. Therefore, many of the extracted datasets suffer from incompleteness, noise and uncertainty.

In order to reduce such problems, one can apply to the knowledge base a set of inference rules that describes its domain. With that, it's possible to resolve contradictions or strengthen or weaken their confidence values. It's also possible to derive new facts that are originally not existent due to incompleteness. Such inference rules can be of two types:

1. *Hard Rules*: Consistency constraints which might represent functional dependencies, functional or inverse-functional properties of predicates or Mutual exclusion. For example:

- $marriedTo(x, y) \leftarrow marriedTo(y, x), (x \neq y)$
- $grandChildOf(x, y) \leftarrow childOf(x, z), childOf(z, y)$
- $parentOf(x, y) \leftarrow childOf(y, x)$
- $(z = y) \leftarrow wasBornIn(x, z), wasBornIn(x, y)$

2. *Soft Rules*: Weighted rules that frequently, but not always hold in the real world. As they might also produce incorrect information, each rule itself must have a confidence value which should be applied to derived facts, for example married people live in the same place as their partner has confidence 0.8:

$$livesIn(x, y) \leftarrow marriedTo(x, z)livesIn(z, y) [0.8]$$

So, if we have an incomplete knowledge base, which lacks information about where *Michelle Obama* lives, but we know that she's married to *Barack Obama* and he lives in *Washington, D.C.*, both with confidence 1, we could then apply this soft rule to derive the fact *livesIn(MichelleObama, WashingtonDC)* with confidence 0.9.

Such rules are rarely known beforehand, or are too expensive to be manually extracted. Nevertheless, the data itself can be used to mine these rules using *Inductive Logic Programming (ILP)*.

ILP is a well-established framework for inductively learning relational descriptions (in the form of logic programs) from examples and background knowledge. Given a logical database of facts, an ILP system will generate hypothesis in a pre-determined order and test them against the examples. However in a large knowledge base, ILP becomes too expensive as the search space grows combinatorially with the knowledge base size and the larger the number of examples, the more expensive it is to test each of the hypothesis.

rules with constants might be really interesting.

Moreover, testing hypothesis with constants increases the search space dramatically, making it unfeasible to test all possible hypothesis with constants in a large knowledge base. In such case, it's necessary to arbitrarily reduce the search by restricting the set of constants to be included in ...

data mining, rule mining, datalog rules

talk a bit about ilp

1.1 Motivation

Given the huge size of search space and the great interestingness of rules with constants, we need to smartly prune constants or combinations of constants that , learning datalog rules can be already extremely costly.

Numerical constants are a special case, and they need to be treated differently. Depending on the numerical attribute domain, in case of a continuous real number domain for

example, setting a numerical constant as an individual value will very likely have a very low support and also that would result and extremely large number of possible constants. Therefore, we split the attribute's domain into k buckets and then check if any of the buckets present any gain in comparison with its correspondent numerical constant-free hypothesis.

For example if we test the hypothesis and we find support=100 and confidence=0.4:

$$isMarriedTo(x, y) \leftarrow hasAge(x, z)$$

and then we split z into three buckets:

- $k = 1 : z \in [0, 20]$
- $k = 2 : z \in (20, 40]$
- $k = 3 : z \in (40, \infty]$

we then test the hypothesis for each of the three buckets and we obtain

- $isMarriedTo(x, y) \leftarrow hasAge(x, z), z \in [0, 20]$
support=40, confidence=0.1
- $isMarriedTo(x, y) \leftarrow hasAge(x, z), z \in (20, 40]$
support=40, confidence=0.5
- $isMarriedTo(x, y) \leftarrow hasAge(x, z), z \in (40, \infty]$
support=20, confidence=0.8

as we see, for $k=2$ and $k=3$, the hypothesis we have significant gain by specifying numerical constants. Adding a relation to the body might produce totally different confidence support and confidence distributions along the buckets. For example, if we add the relation $hasChild(x, a)$, we could obtain other interesting rules:

- $isMarriedTo(x, y) \leftarrow hasAge(x, z)hasChild(x, a)$
support=50, confidence=0.625
- $isMarriedTo(x, y) \leftarrow hasAge(x, z)hasChild(x, a), z \in [0, 20]$
support=2, confidence=0.5
- $isMarriedTo(x, y) \leftarrow hasAge(x, z)hasChild(x, a), z \in (20, 40]$
support=30, confidence=0.7

- $isMarriedTo(x, y) \leftarrow hasAge(x, z)hasChild(x, a), z \in (40, \infty]$
support=18, confidence=0.9

adding some relations might not bring any gain or even loss in confidence, but when bucketing per age, present a different distribution,

nevertheless, adding some relations might not generate any interesting rules, like

d

1.2 Contributions

In this Thesis, we propose a pre-processing step to build a graph we call Correlation Lattice for each numerical property. In each graph, that has a numerical property as root, we first query the frequency distribution on the numerical attribute, then split them in k buckets. Then we pick a set of c categorical properties that can be joined with the root, and analyze how the distribution of sub-population created by joining them with the root is affected. Afterwards we try to combine each of the categories and see if they still produce interesting sub-populations, like in frequent set mining.

We also evaluate different heuristics and interestingness measures.

In a hypothesis containing a numerical attribute in the body, we can obtain a support and confidence value for each of the buckets, and

With that, during the ILP algorithm, once we add one of the root properties, we can then search for the most interesting categorical properties that could result in different accuracy distributions. For every categorical property we can also suggest the most interesting constants and other categorical properties to be combined in a subcategory of both.

1.3 Outline

The remainder of this thesis is structured as follows. In Chapter ??, we provide technical background on MapReduce and BigTable. In Chapter ??, we present a summary of previous work in the areas of duplicate and near-duplicate detection, information retrieval on web archives, and MapReduce applications in graph processing. Following that, we state our problem and describe solutions in Chapter ??. In Chapter ??, we describe an implementation of our solution using the MapReduce framework. In Chapter ??, we

present our experimental results. We conclude this thesis and outline directions of future research in Chapter ??.

Chapter 2

Related Work

2.1 Logic Programming

2.2 Inductive Logic Programming

2.3 Mining Optimized Rules for Numeric Attributes

[\[1\]](#)

2.4 Minimum Description Length

[\[? \]](#)

2.5 Semantic Web

2.6 Linked Open Data

Chapter 3

Correlation Lattice

The idea is to build during preprocessing a graph inspired in the Itemset Lattice that describes the influence of different categorical relations on a given numerical attribute's distribution. We call such graph a Correlation Lattice. To illustrate the idea, let's analyze a simple real-world example with the *hasIncome* relation. If we have two categorical relations, one strongly correlated to income, e.g. *hasEducation*, and one uncorrelated (or very weakly correlated), e.g. *wasBornInMonth*.

Let's assume that for the relation *wasBornInMonth*(x,y) we have the 12 months from the Gregorian Calendar as constants and for *hasEducation*(x,y) we can have 10 different categorical constants for y : "Preschool", "Kindergarten", "ElementarySchool", "MiddleSchool", "Highschool", "Professional School", "Associate's degree", "Bachelor's degree", "Master's degree" and "Doctorate degree".

It's expected that the income distribution will be roughly the same for people born in any of the months, whereas for different education levels, e.g. Elementary School and Doctoral Degree, their income distribution are expected to be different between them and different from the overall income distribution.

In a further step, we try to join every possible pair of categorical relations and including the constants. For the given example with the relations *hasEducation* and *wasBornInMonth* we would then create the nodes:

hasIncome(x,y)*wasBornInMonth*($x, "January"$),*hasEducation*($x, "Preschool"$)
hasIncome(x,y)*wasBornInMonth*($x, "January"$),*hasEducation*($x, "Kindergarten"$)
...
hasIncome(x,y)*wasBornInMonth*($x, "January"$),*hasEducation*($x, "Doctorate Degree"$)

$hasIncome(x,y)wasBornInMonth(x, "February"),hasEducation(x, "Preschool")$
 $hasIncome(x,y)wasBornInMonth(x, "February"),hasEducation(x, "Kindergarten")$
 \dots
 $hasIncome(x,y)wasBornInMonth(x, "February"),hasEducation(x, "Doctorate Degree")$
 \dots
 $hasIncome(x,y)wasBornInMonth(x, "December"),hasEducation(x, "Preschool")$
 $hasIncome(x,y)wasBornInMonth(x, "December"),hasEducation(x, "Kindergarten")$
 \dots
 $hasIncome(x,y)wasBornInMonth(x, "December"),hasEducation(x, "Doctorate Degree")$

Based on this idea, we basically check how different categorical relations affect a numerical distribution. Such information, together with other measures like support, provides valuable cues on what categorical attributes and what categorical constants might be the most interesting to be added to the hypothesis in the core ILP algorithm.

3.1 Categorical Relation Definition

In this section, we formally define a categorical relation as used in the Correlation Lattice.

First of all, a candidate relation must be joined with root relation's 1st argument (assuming that the numerical attribute is in the 2nd argument).

A candidate categorical relation $r(x, y)$, should be equivalent a non-injective function:

$$r(x, y) \equiv f : X \rightarrow Y, s.t. |Y| < |X| \text{ and } |Y| = n, n > 1$$

$$\nexists g : Y \rightarrow X, s.t. f(g(x)) = x, \forall x \in X$$

We can define subsets of $X_i \in X$, with which of them belonging to one category $y_i \in Y$:

$$X_i \subset X, s.t. X_i = \{x \in X \mid f(x) = y_i, y_i \in Y\}$$

$$X = \bigcup_{i=1}^n X_i \text{ and } X_i \cap X_j = \emptyset, \forall i, j \in [1, n], i \neq j$$

We can also broaden this definition by composing functional relations to a categorical or multiple categorical relations:

If we have:

$r_1(x, y) \equiv f_1 : X \rightarrow Y$ (categorical or not)

$r_2(y, z) \equiv f_2 : Y \rightarrow Z$ (categorical relation)

Then,

$r'(x, z) \equiv f : X \rightarrow Z$, where $r'(x, z) = r_2(f_1(x), z)$ is also categorical

Numerical relations can also be turned into a categorical, by simply applying a bucketing function that maps a numerical domain into a finite set of k buckets:

$b : \mathbb{N} \rightarrow B$, where $B = \{b_1, b_2, \dots, b_k\}$

So a numerical relation:

$r(x, y) \equiv f : X \rightarrow \mathbb{N}$

combined with a bucketing function b , $r'(x, b(y))$ would be categorical

(Then talk about non-categorical relations as categorical by considering its presence/absence)

3.2 Support

As described in ([2]), in top-down ILP every refinement causes the support to decrease, therefore we know that for every node in the Correlation Lattice, its support will be greater or equal than any of its children, so support is a monotonically decreasing measure so we can safely prune a node that doesn't reach the minimum support threshold.

3.2.1 Independence between Nodes

By simplicity, we assume that every possible pair of categorical relations are independent and we search for evidence to prove the contrary.

For 2 nodes to be joined, they must have a common parent, i.e. two nodes at level l (with $l + 1$ literals) are joinable if they share l literals. Therefore, it's straightforward to calculate the conditional probabilities of each of the joining nodes given the common parent, and estimate the frequency distribution for the conditional independence case.

If we are joining `hasEducation()`...

For every bucket b_i in the frequency histogram, we can calculate the conditional probability $p_i(n_1|p)$ and $p_i(n_2|p)$ assuming conditional independence given p in order to estimate $\hat{h}_i(n_1, n_2)$:

$$\begin{aligned}
p_i(n_1|n_2, p) &= p_i(n_1|p) \\
&= \frac{h_i(n_1)}{h_i(p)} \\
p_i(n_2|n_1, p) &= p_i(n_2|p) \\
&= \frac{h_i(n_2)}{h_i(p)}
\end{aligned} \tag{3.1}$$

$$\begin{aligned}
\hat{h}_i(n_1, n_2) &= p_i(n_1|n_2, p) * p_i(n_2|p) * h_i(p) \\
&= p_i(n_1|p) * h_i(n_2) \\
\hat{h}_i(n_1, n_2) &= p_i(n_2|n_1, p) * p_i(n_1|p) * h_i(p) \\
&= p_i(n_2|p) * h_i(n_1)
\end{aligned}$$

After that, we query the actual frequency distribution on the Knowledge Base and do an Pearson's chi-squared independence test. As null hypothesis and alternative hypothesis we have:

$H_0 = n_1$ and n_2 are conditionally independent given their common parent p $H_1 = n_1$
and n_2 are conditionally dependent given their common parent p

Number of degrees of freedom is the number of buckets minus one:

$$df = k - 1$$

We calculate the χ^2 value:

$$\chi^2 = \sum_{i=1}^k \frac{(h_i - \hat{h}_i)^2}{\hat{h}_i} \tag{3.2}$$

[3]

3.3 Heuristics

Then it's possible to obtain the p-value and check whether there's enough confidence to reject the null hypothesis H_0 .

Chapter 4

Algorithmic Framework

4.1 Knowledge Base Backend

[4]

4.2 Preprocessing

In this section, we will present the preprocessing steps required by our proposed algorithm. It basically consists of building a joinable relations map for each of the four join patterns, according to relations domain and range types as well as support threshold. Afterwards, we search the available categorical properties for each numerical relation that will be used in the Correlation Lattice. At last we describe the Correlation Lattice structure and the algorithm to build it.

4.2.1 Relation Preprocessing

In this step, we focus on creating for each of the four join patterns between two relations:

- Argument 1 on Argument 1: e.g. $hasIncome(\mathbf{x}, y) hasAge(\mathbf{x}, z)$
- Argument 1 on Argument 2: e.g. $hasIncome(\mathbf{x}, y) isMarriedTo(z, \mathbf{x})$
- Argument 2 on Argument 1: e.g. $livesIn(y, \mathbf{x}) isLocatedIn(\mathbf{x}, z)$
- Argument 2 on Argument 2: e.g. $livesIn(y, \mathbf{x}) wasBornIn(z, \mathbf{x})$

4.2.1.1 Exploiting Relation Range and Domain Types

A knowledge base is expected to have an ontology defining the structure of the stored data (the types of entities and their relationships). Additionally, every relation's range (type of 1st argument) and domain (type of 2nd argument) should be defined. These information can help us identify the allowed joining relations for each join pattern.

For every possible pair of relations,

The algorithm is shown in the pseudo-code bellow:

4.2.1.2 Exploiting Support Monotonicity

As seen in (???), support is the only monotonically decreasing measure in top-down ILP. So we know that by adding any literals to the hypothesis, we can only get a smaller or equal support. Therefore, for each pair of joinable relations in each of the join patterns, we can query the knowledge base and check whether they reach the minimum support threshold.

Thus, if any pair of relations doesn't reach the minimum support for a given join pattern, we know that any hypothesis containing such join will therefore fail the support test as well, so we don't need to test such hypothesis in the core ILP algorithm.

The relation preprocessing will result in 4 maps, one for each join pattern. Each map will a relation as key and a set of joinable relations as value. The refinement step at the ILP algorithm, will then access this map when choosing a new literal to be added.

4.2.2 Correlation Lattice

4.2.2.1 Graph Node

Every node essentially contains the following attributes:

- Set of pointers to parent nodes
- Set of pointers to child nodes
- Set of pointers to constant nodes
- Histogram with facts distribution over root numerical property

4.2.2.2 Building the Correlation Lattice

For building the Correlation Lattice, we start with the root node, which has a numerical property as literal and no constants assigned, e.g. *hasIncome(x,y)*. We then query the distribution of positive examples over the property in the whole Knowledge Base.

```
SELECT COUNT ?y WHERE { ?x <hasIncome> ?y } GROUP BY (?y)
```

It's also necessary to specify the bucketing technique and the number of buckets in order to extract the histogram from the obtained query results. These buckets are used to build the histograms of all nodes in the graph.

Afterwards, we select the the categorical properties that will be used in the lattice. For each of the selected properties, we join them with the root numerical property (for simplicity we'll assume all the categorical properties are joined with both 1st arguments) and we query the distribution again. In the first level, it's necessary to extract a histogram for each of the categorical constants in the selected properties. Therefore, it's a good strategy to group the results also by these categorical constants so If we select *hasEducation* for example, we would then fire the following SPARQL query:

```
SELECT COUNT ?z ?y WHERE { ?x <hasIncome> ?y . ?x <hasEducation> ?z }  
GROUP BY (?z,?y)
```

With such query, it's possible to extract a histogram for the node *hasIncome(x,y)hasEducation(x,z)* and its correspondent constants

4.2.2.3 Searching Rules in Correlation Lattice

4.2.2.4

4.3 ILP Core Algorithm

Algorithm 1: Checks whether two relations are joinable for a given join pattern

Relation r_i, r_j , Argument arg_i, arg_j **Output:** True if arg_i from r_i joins with arg_j from r_j , False otherwise

```

if  $r_i.arg_i = r_j.arg_j$  or  $subsumes(r_i.arg_i, r_j.arg_j)$  or  $subsumes(r_j.arg_j, r_i.arg_i)$  then
  | return true;
else
  | return false;

```

Algorithm 2: Checks valid join pairs for a given join patterns

List of Figures

List of Tables

Bibliography

- [1] Sergey Brin, Rajeev Rastogi, and Kyuseok Shim. Mining optimized gain rules for numeric attributes. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 135–144. ACM Press, 1999.
- [] Toon Calders, Christian W. Günther, Mykola Pechenizkiy, and Anne Rozinat. Using minimum description length for process mining. In *SAC*, pages 1451–1455, 2009.
- [2] Nada Lavrac and Saso Dzeroski. *Inductive Logic Programming: Techniques and Applications*. Ellis Horwood, New York, 1994. URL <http://www-ai.ijs.si/SasoDzeroski/ILPBook/>.
- [3] Szymon Jaroszewicz and Dan A. Simovici. Pruning redundant association rules using maximum entropy principle. In *In Advances in Knowledge Discovery and Data Mining, 6th Pacific-Asia Conference, PAKDD'02*, pages 135–147, 2002.
- [4] Thomas Neumann and Gerhard Weikum. The rdf-3x engine for scalable management of rdf data. *The VLDB Journal*, 19(1):91–113, February 2010. ISSN 1066-8888. doi: 10.1007/s00778-009-0165-y. URL <http://dx.doi.org/10.1007/s00778-009-0165-y>.