

Contents

1	Linked Data	3
1.1	Resource Description Framework	3
1.2	Linked Open Data	5
2	Data Cubes	7
2.1	The Data Cube Vocabulary	8
3	Association Rules	9
4	AMIE Algorithm and Its Derivatives	10
4.1	AMIE	10
4.1.1	Language Bias	11
4.1.2	Measures of Significance	11
4.1.3	Confidence Measures	12
4.1.4	Algorithm	13
4.1.5	Refinement Operators	13
4.1.6	Count Projection Queries	13
4.1.7	In-Memory Database	13
4.2	AMIE+	14
4.2.1	Rule Refinement	14
4.2.2	Speeding Up Confidence Evaluation	15
4.3	RDFRules	15
4.3.1	Limitations of AMIE+	15
4.3.2	Faster Projection Counting	16
4.3.3	Processing of Numerical Attributes	16
4.3.4	Multiple Graphs	16
4.3.5	Improvements to Expressiveness of Rule Patterns	16
4.3.6	Top-k Approach	16
4.3.7	Support for the Lift Measure	17
4.3.8	Rule Clustering	17
4.3.9	Rule Pruning	18
5	RDFRules Reference Implementation	19
6	Leveraging a Combination of OLAP Cubes and Knowledge Graphs	20
7	Experiment	21
7.1	Czech Social Security Administration	21
7.2	Czech Statistical Office	24
7.3	Wikidata	25
7.4	Filtering the Observations	26

7.5	Slicing the Cubes	26
7.6	Linking	29
7.6.1	Sex Dimension Values	29
7.6.2	Reference Periods	29
7.6.3	Reference Areas	30
7.7	Discretization	31
7.8	Mining Tasks	32
7.8.1	Relation between the pension expenses and the policital alignment of the state's government	32
7.9	Discussion	32
	Conclusions	33
	List of References	34
	A Queries and Scripts	36

1. Linked Data

Linked Data is a set of best practices for publishing and connecting data on the Web structured in such a way that it is usable not only for human processing but also processable by machines. It builds upon the general architecture of the World Wide Web. Instead of creating links between particular documents from different sources in the case of the classic Web of documents, the Linked Data connects the representations of real-world objects or abstract concepts. Tim Berners-Lee expressed these best practises in four principles, known as the Linked Data principles.[2]

- Use URIs as names for things.
- Use HTTP URIs, so that people can look up those names.
- When someone looks up a URI, provide useful information, using the standards
- Include links to other URIs, so that they can discover more things.

URIs (Universal Resource Identifier) are used to identify real-world objects and abstract concepts. According to the second principle, information about the entity that the URI represents can be retrieved using HTTP protocol (so-called URI dereferencing). Based on the third principle, which advocates for a standard structure of data dereferenceable by URI, the Resource Description Framework (RDF) has been designed. Tim Berners-Lee also suggested 5-star deployment scheme for ranking published data according to the format in which it is published, comparing them based on their ability to be machine processed. It assigns one star to any data published and assigns the highest number of five stars to data, that is published as RDF, where the entities described are identified by an dereferenceable URI string and the data are connected to other data sources.

1.1 Resource Description Framework

RDF is a data model based on representing data as directed graphs. The basic building block of RDF structured data is a triple consisting of three parts called subject, predicate, and object. The subject is the URI representing the described resource. The object is either URI or literal value like string or number. The predicate specifies the type of relationship between the resources at the positions of subject and object. The predicate is always identified by URI. Predicate URIs come from vocabularies, intended to encompass various relations and concepts occurring in a certain domain.

Set of triples then establishes a RDF graph. URIs at the subject and object positions of the triples make nodes of the graph and each triple acts as an arc connecting the nodes. Type of the connecting is expressed by the predicate URI in the triple. Given the uniqueness of the URIs and their capability of being dereferenced and connected to URIs from various sources (the fourth Linked Data principle), one can imagine the linked data as one giant undivided

graph containing data from various topical domains, so-called Web of Data.

It is important to distinguish the model itself from its formats. RDF describes only an abstract structure of the data that has to be materialized into a certain format when the data is published on the Web. The first standard serialization format published together with RDF in [6] is RDF/XML. An example of two triples described in RDF/XML format is shown in the listing 1.1. The RDF/XML format suits well the use cases, where little human interaction with the data is expected because its syntax is difficult for a human to read and write compared to other formats. On the other hand, its XML background makes it a perfect format for data processed and generated solely by machines.

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:foaf="http://xmlns.com/foaf/0.1/">
  <rdf:Description rdf:about="http://example.com/john-johnson">
    <rdf:type rdf:resource="http://xmlns.com/foaf/0.1/Person"/>
    <foaf:name>John Johnson</foaf:name>
  </rdf:Description>
</rdf:RDF>
```

Listing 1.1: Example of RDF data described in RDF/XML format (Source: author)

One of the most used and most human-readable formats is Turtle (Tense RDF Triple Language).[1] It provides various shorthands, enabling to make the representation as brief as possible and thus suitable to be written by hand. The common part of URI strings can be prefixed, so only the decisive end of the URIs has to be stated. The symbol of the semicolon is used to divide pairs of predicate and object belonging to the same subject, so the subject does not have to be repeated. If the described triples share both subject and predicate, a comma can be used to divide the different objects of the triples. Usage of a prefix and the two symbols is shown in an example in the listing 7.5.

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
@prefix eg: <http://example.com/> .
eg:john-johnson rdf:type foaf:Person ;
                foaf:name "John Johnson" ;
                foaf:knows eg:john-jackson, eg:jack-johnson .
```

Listing 1.2: Example of RDF data described in Turtle format (Source: author)

Same as with the relational data model and SQL, RDF also needs a capable language for querying and manipulating the data. For this purposes, SPARQL was designed.[7] Example of a simple SELECT query written in SPARQL is shown in the listing 7.7. Similarly to SQL, the WHERE clause serves to limit the search place from which the result of the query is given. Content of the WHERE clause resembles the Turtle syntax. URIs, that can be bound to variable that occurs in the pattern stated in the WHERE clause, are contained in the output of the query if the variable is enumerated after the SELECT keyword. This

particular query would return all possible bindings for variable **name** ie. all names of persons, for whom the queried data states, that they know a certain John Jackson.

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX eg: <http://example.com/>

select ?name where {
    ?person rdf:type foaf:Person ;
           foaf:name ?name ;
           foaf:knows eg:john-jackson .
}
```

Listing 1.3: Example of a simple SPARQL query (Source: author)

1.2 Linked Open Data

The first activities with the goal of starting the publication of Linked Data on a global scale were conducted by the Semantic Web research community as part of the W3C Linking Open Data (LOD) project established in 2007.[5] The aim of the project was to identify datasets published under an open license and to publish them according to the Linked Data principles. All data sets that are published under an open license and are connected to other data sets are referred to as LOD cloud.

The content of the LOD spans across multiple domains. The website lod-cloud.net tracks the current state of published LOD data sets and divides the data sets into these categories, so-called subclouds: Cross-Domain, Geography, Government, Life Sciences, Linguistics, Media, Publications, Social Networking and User-Generated. A data set can fall into more than one category. Cross-Domain, general knowledge data sets play an important role of an intermediary through which unrelated data sets can be connected.

One of those data sets is Wikidata. It is a sister project of Wikipedia, founded and hosted by the Wikimedia Foundation. The data set is managed in an open and collaborative way. Everybody who is interested in expanding the knowledge base can create an account and start contributing. The website of the project provides an intuitive user interface for editing and creating data, so no technical skills beyond common usage of the Internet is needed. The data set currently contains over 93 million items edited by over 26 thousand active contributors. Every item of the dataset is allocated a unique identifier prefixed by the letter **Q**, so-called QID or Q number. The items are described by their statements corresponding to RDF triples. Predicates are in the context of Wikidata called properties and are prefixed by letter **P** similar to items. A sample of the triples contained in Wikidata in Turtle syntax shows listing 1.4.

An different approach is taken by the DBpedia project. Instead of relying on manual contribution of volunteers, DBpeditas data comes from an application of NLP extraction algorithms over plain text of Wikipedia's articles.

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix wd: <http://www.wikidata.org/entity/> .
@prefix wdt: <http://www.wikidata.org/prop/direct> .
@prefix schema: <http://schema.org/> .

wd:Q111 wdt:P361 wd:Q7879772
    rdf:label "Mars" ;
    schema:description "fourth planet from the Sun";
```

Listing 1.4: Wikidata content sample (Source: author)

2. Data Cubes

While relational databases with highly normalized data models fit well to situations where data is frequently modified, they can be quite cumbersome when being performed complex aggregating queries. Online Analytical Processing (OLAP) system fits better these purposes. In an OLAP system, numerical data is stored in a multidimensional data structure. The structure is comprised of hypothetical cells, which are identified by their assigned set of dimension values from each dimension of the structure. Each cell can contain zero to many numerical values, so-called measurements. The structure is referred to as OLAP Cube or Data Cube. The word *cube* implies exactly three dimensions, but its purpose is only to illustrate the multidimensionality of the structure.

Distinct values of a dimension can be organized into a hierarchy, where a parent value is assigned to summarized measurements throughout its child values. An example of such a hierarchy could be a relationship of a product category and specific products belonging to this category. The depth of a dimension value in its hierarchy then determines the level of granularity the measurement values in a cell assigned to the dimension value are associated with. A cell with all dimension values at the lowest level in their hierarchies or in no hierarchy at all has the finest level of granularity. In a Data Cube consisting of only one cell, meaning each dimension of the cube has only one distinct value, the cell has the coarsest level of granularity.

Several operation can be performed on a Data Cube:

roll-up This operation aggregates data either by reduction of one or more dimensions or by climbing up a concept hierarchy for a dimension.

drill-down This operation transforms data to a more detailed level. It is the opposite of roll-up operation. Either a new dimension is added or the values are projected on a more granular level of a dimension.

slice and dice By slicing a cube only certain subset of the dimension values of one dimension is allowed in the resulting cube. Dicing means restricting dimension values across multiple dimensions.

pivot Pivoting means rotating the cube by its axis in order to change the view of the data.

If the values contained in the cube have an additive character (e.g. sales amount or a number of security incidents), the values can be rolled up or drilled down along any dimension. Not all facts are additive though (e.g. average temperature). The analytical process itself lies in performing the above-mentioned operations in order to find interesting insight into the data. By precomputing the aggregations of all possible subsets of dimensions from the cube on the finest level of granularity, the whole process can be accelerated.

2.1 The Data Cube Vocabulary

The principle of dimensions, measures and attributes are the basic building blocks of the standards and guidelines presented by the SDMX (Statistical Data and Metadata eXchange) initiative, that tries to standardise and modernise the exchange of statistical data. The World Wide Web Consortium's recommendation for representing multi-dimensional data in RDF is the Data Cube vocabulary. This vocabulary underlies the standards and guidelines of the SDMX initiative. It allows to publish the content of the cube together with information about its structure and its metadata. The structure of the vocabulary is shown in the picture 2.1.

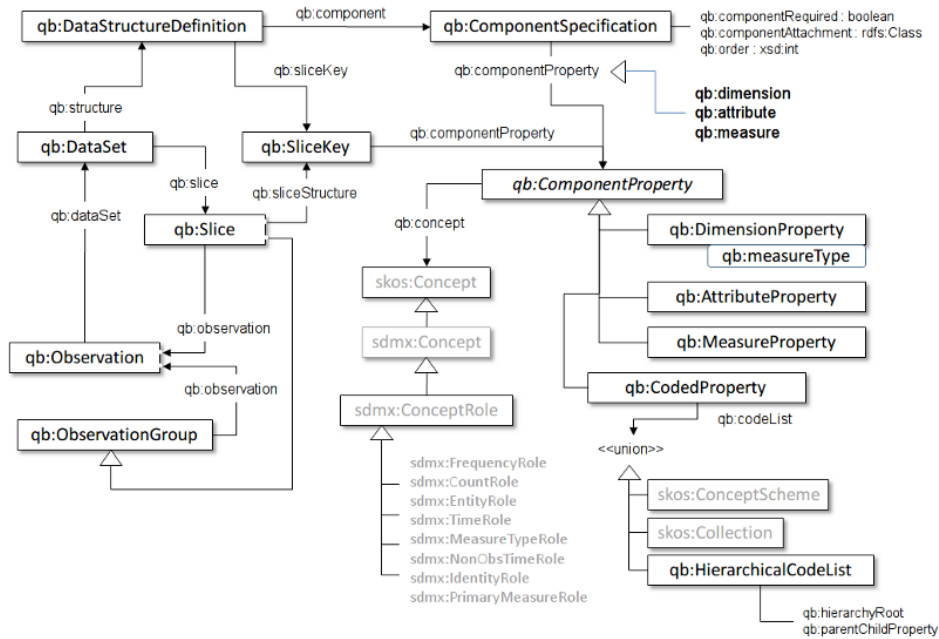


Figure 2.1: The Data Cube Vocabulary structure (Source: [3])

TODO SKOS

```
pk:PensionKindScheme a skos:ConceptScheme ;
    skos:prefLabel "Kinds of pensions"@en .
pkr:PK_VM a skos:Concept, pk:PensionKind ;
    skos:prefLabel "widower's pension"@en ;
    skos:inScheme pk:PensionKindScheme ;
    skos:notation "PK_VM" ;
    skos:altLabel "VM"@cs .
```

Listing 2.1: SKOS

3. Association Rules

4. AMIE Algorithm and Its Derivatives

Finding association rules in knowledge bases can serve several purposes. New facts that are not yet present in the dataset can be derived from the found regularities described by the rules. From such rules opposing facts present in the dataset can be deduced to be wrong. Mined rules can also help to understand the data better.

For mining rules from a graph database such as LOD datasets, Inductive Logic Programming (ILP) can be used. ILP work under the Closed World Assumption (CWA) meaning it supposed that both negative and positive statements are present in the data. However LOD operates under the Open World Assumption (OWA) ie. if a statement is not present in the data, it does not mean that this statement does not correspond to reality. Rules mined by ILP would not reflect this matter. Moreover ILP are not observed to be efficient over large datasets in the order of millions of statements, making it not a viable way to mine rules over real-world knowledge bases such as YAGO or Wikidata.

4.1 AMIE

An algorithm that is specifically designed to mine rules from data operating under OWA and consisting of binary predicates (just as Linked Data) is AMIE (**A**ssociation Rule **M**ining under **I**ncomplete **E**vidence).[4] AMIE mines rules in the form of Horn rule. Horn rule is an implication with conjunction of atoms on the left side, called body and a single atom on right side call head. We can imagine an atom as an RDF triple, where subject and object can be replaced by variables. Number of atoms in a rule indicates the length of the rule. In this work rules are represented with an infix notation. The AMIE literature use the Datalog notation commonly used in ILP domain. An example of a rule AMIE seeks to discover is shown below.

$$(?a \text{ worksIn } ?b) \wedge (?b \text{ hasHeadquartersIn } ?c) \Rightarrow (?a \text{ livesIn } ?c)$$

This rule states that any person lives in a place his or her company's headquarters. Length of this rule is 3 since it has two body atoms. The rule has 3 variables. When we substitute the variables by constants present in the examined data set, we get an *instantiation* of the rule. If all atoms of the instantiated rule appear in the data set, the head atom of the instantiation is one the *predictions* of the rule. Number of all instantiations of an atom that appear in the data set is called *size* of the atom.

4.1.1 Language Bias

In order to efficiently traverse the search space, AMIE subjects the rules to a particular language bias. Only the rules conforming the conditions stated below can be generated and further refined.

rules have to be connected A rule is connected when every atom in the rule shares every variable with another atom in the rule.

rules have to be closed A rule is closed when every variable appears at least twice in the rule.

rules cannot be reflexive reflexive rule contains at least one atom with identical subject and object variable or constant.

rule can be recursive Any predicate can occur more than once in a rule.

4.1.2 Measures of Significance

Support

For a chosen definition of a support measure for the AMIE algorithm it is crucial for the definition to have the property of monotonicity ie. by adding any new atom to the body of a rule, the support of the rule shall always decrease or remain the same. A naive way to count support of a rule would be to count all instantiations of the rule that appear in KB. Such definition would not comply the property of monotonicity, since an addition of a dangling atom to a rule would introduce a new variable multiplying the number of instantiations and thus the value of the support measure. By counting only all distinct pairs of subjects and objects in the head of all instantiations that appear in KB, the property of monotonicity is preserved:

$$supp(\vec{B} \Rightarrow r(x,y)) := \#(x,y) : \exists z_1 \dots z_n : \vec{B} \wedge r(x,y)$$

Head Coverage

Since the support is an absolute number, so the size of the examined data set has to be taken into account while defining this threshold. Plus If the defined support value is greater than a number of distinct triples containing a certain predicate, any rule containing this predicate in the head atom would be disregarded. Head Coverage is the relative expression of support. It is defined as support of a rule over the size of its head atom.

$$hc(\vec{B} \Rightarrow r(x,y)) := \frac{supp(\vec{B} \Rightarrow r(x,y))}{size(r)}$$

4.1.3 Confidence Measures

The above-mentioned measures describe a quantitative significance of the rule in relation to the examined data set. They quantify the true predictions of the rule but do not take into account the false predictions. Confidence is a way to measure the quality of a rule. Generally speaking, confidence is a ratio of true predictions of a rule to the sum of true predictions and the counterexamples. Number of true predictions can easily be expressed by the rule's support. Two different ways to count the counterexamples are discussed below.

CWA and Standard Confidence

Standard confidence considers every fact that is not present in the examined dataset a false fact and thus a counterexample when predicted by a rule. Facts predicted by a rule is either present in the data set or it is not. Therefore the standard confidence is defined as the ratio of the number of true predictions of the rule to the number of all predictions of the rule.

$$conf(\vec{B} \Rightarrow r(x,y)) := \frac{supp(\vec{B} \Rightarrow r(x,y))}{\#(x,y) : \exists z_1 \dots z_n : \vec{B}}$$

This way of generating counterexamples fails to distinguish a false fact from an unknown fact. This conforms the CWA and it is traditionally used for association rule mining over transactional data where this assumption can be applied. For example if the data does not state, that I bought a bottle of milk last Wednesday, then I really did not buy it. AMIE, however, is intended to mine rules from data operating under OWA, so the usage of this measure is inappropriate.

PCA Confidence

For the PCA Confidence, Partial Completeness Assumption (PCA) is used for generating the counterexamples:

If $\langle s \ p \ o \rangle \in KBtrue$ then $\forall_{o'} : \langle s \ p \ o' \rangle \in (KBtrue \cup NEWtrue) \Rightarrow \langle s \ p \ o' \rangle \in KBtrue$.

Meaning that if we know any object for given predicate and subject, we know all triples of containing the predicate and subject together. This assumption is certainly true for predicates with high or complete functionality, such as birthdate or capital. A triple predicted by the measured rule is considered an counterexample only when triples with its combination of subject and predicate are present in the data set and none of those has the triple's object.

$$conf_{pca} := \frac{supp(\vec{B} \Rightarrow r(x,y))}{\#(x,y) : \exists z_1 \dots z_n, y' : \vec{B} \wedge r(x,y')}$$

4.1.4 Algorithm

Algorithm 1 AMIE algorithm

```
1: procedure AMIE( $x, y$ )
2:    $queue = [(?a\ r_1\ ?b), (?a\ r_2\ ?b) \dots (?a\ r_m\ ?b)]$ 
3:    $output = \langle \rangle$ 
4:   while  $\neg queue.isEmpty()$  do
5:      $rule = queue.dequeue()$ 
6:     if  $AcceptedForOutput(r, out, minConf)$  then
7:        $output.add(rule)$ 
8:     end if
9:     if  $length(rule) < maxLen$  then
10:       $R(rule) = Refine(rule)$ 
11:    end if
12:    for  $r_i \in R(rule)$  do
13:      if  $hc(r_i \geq minHC \ \& \ r_i \notin queue)$  then
14:         $queue.enqueue(r_i)$ 
15:      end if
16:    end for
17:  end while
18:  return  $output$ 
19: end procedure
```

4.1.5 Refinement Operators

O_D add dangling atom (with a fresh variable)

O_I add instantiated atom (with a constant)

O_C add closing atom (both arguments are shared variables)

4.1.6 Count Projection Queries

new relation r to the rule $B_1 \wedge \dots B_{n-1} \Rightarrow H$

find all relations that lead to a new rule that passes the min head coverage threshold.

4.1.7 In-Memory Database

query implementation

one fact index for each permutation of $\{S, P, O\}$

allowing to check existence of a triple in constant time

allowing to efficiently fetch instantiation of an atom

aggregated indexes S,P,O: store aggregated count of facts for each key of the fact indexes: P stores count of triples for each relation

Size Queries

$size(livesIn(x,y)) \rightarrow$ aggregated index **P**

1. look up: *livesIn* in **P**

$size(livesIn(x,USA)) \rightarrow$ fact index **POS**

1. look up: *livesIn* in **POS**
2. look up: USA in **OS** and count of subjects

Existence Queries

to determine whether there exists a binding for a conjunctive query

Select Queries

finding distinct instances of a variable, which is in a conjunction of atoms

Count Queries

compute count of bindings

for confidence of a rule

first fire SELECT then for each binding of x it instantiates the query and fires select query on variable y, adding up the count of instantiations

4.2 AMIE+

does not alter output in any way compared to AMIE.

1. refinements phase
2. confidence evaluation

4.2.1 Rule Refinement

Given a maximum rule length maxLen and a non-closed Horn rule of length maxLen-1, AMIE+ will refine it only if it is possible to close it before exceeding the length constraint.

For a not-yet-closed rule of length $\text{maxLen}-1$, AMIE+ will not apply the O_D , because this would result in a non-closed rule, which will be neither output nor refined.

If a rule contains more than two non-closed variables, AMIE+ will skip the application of O_C . O_C cannot close more than two variables.

Rules with more than one non-closed variable are not refined with instantiated atoms, because the addition of an instantiated atom can close at most one variable.

Rules with $\text{conf}_{\text{pca}} = 1$ are not further refined \rightarrow perfect rules

Simplyfing Projection Queries

addition of a dangling atom cannot reduce support when:

1. parent rule already contains atoms with the same relation as a dangling atom
2. these atoms have a variable in common with the dangling atom

Both rules have to have same support.

4.2.2 Speeding Up Confidence Evaluation

approximation of confidence, tends to overestimate (4% of errors)

turns days runtime to minutes

4.3 RDFRules

4.3.1 Limitations of AMIE+

little attention to data pre-processing and data post-processing

lack of various features which were found useful for mining rules from transactional data, such as support for additional interest measures (lift)

does not support multiple graphs

inability to process numerical data

absence of the top-k approach. In top-k approach, user is only returned the k rules with the highest values of a chosen measure.

coarse rule patterns. Without additional guidance by the user, the top-k approach often generates rules that reflect patterns in data that are obvious or uninteresting.

repetitive calculations during refinement

exhaustive calculations

4.3.2 Faster Projection Counting

reducing the number of calls to the binding functions

4.3.3 Processing of Numerical Attributes

4.3.4 Multiple Graphs

graph aware rules

$$(?a \langle wasBornIn \rangle ?b \langle YAGO \rangle) \Rightarrow (?a \text{ dbo : deatchPlace } ?b \langle DBpedia \rangle)$$

extension of fact indexes

PG, PSG, POG, PSOG

4.3.5 Improvements to Expressiveness of Rule Patterns

? pattern for any symbol

?_v pattern for any variable

?_c pattern for any constant

¬ negation

$$(?a \langle wasBornIn \rangle ?b) \wedge (?b ?_c ?_c \langle DBpedia \rangle) \Rightarrow (?a [\langle livesIn \rangle, \langle deadIn \rangle] ?b)$$

$$(?a \langle wasBornIn \rangle ?b) \wedge (?b \text{ dbo : cityOf } \langle USA \rangle \langle DBpedia \rangle) \Rightarrow (?a \langle deadIn \rangle ?b)$$

4.3.6 Top-k Approach

Top-k Confidence

increasing minConf may speed up the confidence calculation

$$conf(\vec{B} \Rightarrow H) = \frac{supp(\vec{B} \Rightarrow H)}{bsize(\vec{B} \Rightarrow H)}$$

if minConf is set, this inequality must apply:

$$bsize(\vec{B} \Rightarrow H) \leq \frac{supp(\vec{B} \Rightarrow H)}{minConf}$$

during calculation of bsize we can stop the calculation as soon as the value is greater than the ratio.

4.3.7 Support for the Lift Measure

$$lift(\vec{B} \Rightarrow H) = \frac{conf(\vec{B} \Rightarrow H)}{hconf(H)}$$

hconf

if $H = (?a \ p \ ?b)$ then

$$hconf(\vec{B} \Rightarrow H) = \frac{\#s : \exists \langle s, p, o \rangle \prec (?a \ p \ ?b)}{\#s : \exists \langle s, p, o \rangle \prec (?a \ ?r \ ?b)}$$

ratio between the number of distinct subjects bound with the predicate p and the number of distinct subjects in the whole KG.

if $H = (?a \ p \ C)$ then

$$hconf(\vec{B} \Rightarrow H) = \frac{\#s : \exists \langle s, p, o \rangle \prec (?a \ p \ C)}{\#s : \exists \langle s, p, o \rangle \prec (?a \ p \ ?b)}$$

if $H = (C \ p \ ?a)$ then

$$hconf(\vec{B} \Rightarrow H) = \frac{\#s : \exists \langle s, p, o \rangle \prec (C \ p \ ?a)}{\#s : \exists \langle s, p, o \rangle \prec (?a \ p \ ?b)}$$

4.3.8 Rule Clustering

$$sim(r_1, r_2) = \sum_{i=1}^m w_i * sim_i(R_{1,i}, R_{2,i})$$

$$\sum_{i=1}^m w_i = 1$$

Similarity Function of Atom Items

s can be substituted by o

$$sim(\langle s_1, p_1 \rangle, \langle s_2, p_2 \rangle) =$$

Similarity Function of Predicates

Similarity Function of Atoms

$$sim_a(A_1, A_2) = \frac{1}{3} [sim(\langle s_1, p_1 \rangle, \langle s_2, p_2 \rangle) + sim(\langle o_1, p_1 \rangle, \langle o_2, p_2 \rangle) + sim(p_1, p_2)]$$

Similarity Function of Rules

$$|r_1| \geq |r_2|$$

$$sim_r(r_1, r_2) = \frac{1}{r_1} \sum_{i=1}^{|r_1|} max(sim_a(A_i^{r_1}, A_1^{r_2}), \dots, sim_a(A_i^{r_1}, A_{|r_2|}^{r_2}))$$

4.3.9 Rule Pruning

data coverage pruning

ranking (order in which the rules enter the data coverage pruning algorithm)

Rule A is ranked higher than B if:

1. $conf(A) > conf(B)$
2. $conf(A) = conf(B)$ and $hc(A) > hc(B)$
3. rule A has a shorter body than the rule B

For each rule, the algorithm checks whether the rule correctly classifies at least one triple in the input KG

In AMIE+, it is often the case that a single triple is covered by multiple rules.

5. RDFRules Reference Implementation

6. Leveraging a Combination of OLAP Cubes and Knowledge Graphs

7. Experiment

This section describes an experiment of mining association rules from RDF data compiled of statistical data structured by the Data Cube Vocabulary and facts pulled from the Wikidata data set that was performed as an practical part of this work. The statistical data come from two sources. The first one is the Czech Social Security Administration and the second one is the Czech Statistical Office. Analysis was performed using the Scala API of the reference implementation of the RDRules algorithm. The following sections describe how the available data had to be preprocessed to give reasonable results in combination with KG data. The preprocessing was performed partly by the implementation's API itself, partly by performing SPARQL queries over the data. The described method can be taken as inspiration from when performing similar analysis ie. association rule mining task over the multidimensional data merged with loosely structured graph data.

7.1 Czech Social Security Administration

Czech Social Security Administration (CSSA) is a Czech public administration organisation responsible for collecting social security premiums and contributions to the state employment policy. Since 2015 the organization publishes its statistical yearbook datasets and other (vocabularies, code lists and datasets containing data concerning the internal operation of the organization) in the form LOD and became of the first Czech public institutions to do so. The yearbook statistical data sets are modelled using Data Cube Vocabulary. Their dimension values are represented by the SKOS vocabulary. The organization has published 73 datasets so far. All these datasets are downloadable as dumps¹ or accessible through a SPARQL endpoint². The CSSA's URI are dereferenceable.

The largest of the data cubes published is `cssa-d:duchodci-v-cr-krajich-okresech`³. From now on it will be denoted as CSSA1. It contains 368 118 observations structured spread over four dimensions: reference area⁴, reference period⁵, sex⁶ and pension kind⁷. Observations are assigned three measures: the average amount of pension⁸, the average age⁹ and the number of persons¹⁰. Each observation is assigned only one measure.

In this particular data set there are 102 distinct values of the dimension of reference area:

¹<http://data.cssz.cz/web/otevrena-data/katalog-otevrenych-dat>

²<http://data.cssz.cz/web/otevrena-data/sparql-query-editor/>

³<https://data.cssz.cz/resource/dataset/duchodci-v-cr-krajich-okresech>

⁴<https://data.cssz.cz/ontology/dimension/refArea>

⁵<https://data.cssz.cz/ontology/dimension/refPeriod>

⁶<https://data.cssz.cz/ontology/dimension/pohlavi>

⁷<https://data.cssz.cz/ontology/dimension/druh-duchodu>

⁸<https://data.cssz.cz/ontology/measure/prumerna-vyse-duchodu-v-kc>

⁹<https://data.cssz.cz/ontology/measure/prumerny-vek>

¹⁰<https://data.cssz.cz/ontology/measure/pocet-duchodcu>

```

@prefix qb: <http://purl.org/linked-data/cube#> .
@prefix cssa-om: <https://data.cssz.cz/ontology/measure/> .
@prefix cssa-d: <https://data.cssz.cz/resource/dataset/> .
@prefix cssa-od: <https://data.cssz.cz/ontology/dimension/> .

<https://data.cssz.cz/resource/observation/duchodci-v-cr-krajich-okresech/2017-12-31/
    prumerna-vyse-duchodu-v-kc/pk_srnvm/vc.35/m>
  a qb:Observation ;
  qb:dataSet cssa-d:duchodci-v-cr-krajich-okresech .
  qb:measureType cssa-om:prumerna-vyse-duchodu-v-kc ;
  cssa-od:druh-duchodu <https://data.cssz.cz/resource/pension-kind/PK_SRNVM_2010> ;
  cssa-od:pohlavi <https://data.cssz.cz/ontology/sdmx/code/sex-M> ;
  cssa-od:refArea <https://data.cssz.cz/resource/ruian/vusc/35>;
  cssa-od:refPeriod <https://data.cssz.cz/resource/reference.data.gov.uk/id/gregorian-day/2017-12-31>;
  cssa-om:prumerna-vyse-duchodu-v-kc 6622.0;

```

Listing 7.1: Example of an observation from CSSA1

14 regions (NUTS 3 administrative units, czech translation in singular nominative is *kraj*) including Prague, 77 districts (*okres*) also including Prague, 10 Prague districts (*správní obvod*) and a value representing the state in total. Each entity representing a reference area is assigned an unique numerical identifier which corresponds to this area's identifier in the official Registry of Territorial Identification, Addresses and Real Estate (RTIAR) runned by the State Administration of Land Surveying and Cadastre (SALSC). RTIAR codes are reference codes by law, so it is obligatory for CSSA to use them and have them correct.

```

@prefix skos: <http://www.w3.org/2004/02/skos/core#> .

<https://data.cssz.cz/resource/ruian/vusc/35>
  a <https://data.cssz.cz/ontology/ruian/Vusc> , skos:Concept ;
  <http://www.w3.org/2002/07/owl#sameAs> <https://linked.cuzk.cz/resource/ruian/vusc/35> ;
  skos:inScheme <https://data.cssz.cz/resource/ruian/ConceptScheme> ;
  skos:notation "VC.35" ;
  skos:prefLabel "Jihočeský kraj" .

```

Listing 7.2: Dereferenced proxy entity of the South Bohemian Region

There are official URIs of this registry but CSSA datasets do not used them directly. The entities for the dimension of reference area and other dimensions in the dataset CSSA1 and all other data sets of CSSA with the dimension of reference area work as so-called *proxy entities*. This means that instead of using the original code list item URIs directly as objects in the RDF triples, it uses their equivalents defined in the internal code lists. These equivalents are connected to the original URI by the `owl:sameAs` statement. These proxies can then contain data specific to CSSA e.g. labels. Another advantage of this is, that these URIs are dereferenceable to the CSSA domain and their versioning is under the control of CSSA and they can be easily redirect to a different equivalent code list. Previously the proxy entities of the reference area were directed to the unofficial transformation of the Opendata.cz initiative¹¹.

¹¹<https://linked.opendata.cz/dataset/cz-ruian>

Another dimension whose values work as proxy entities is the dimension of reference period. This dimension divides the observations into one year intervals. This applies to all other CSSA data cubes containing the reference period dimension. The data sets vary in the overall covered period. The last covered year of all data sets is the year 2019. The entities link to the `data.gov.uk` Time Intervals¹² OWL ontology. Usage of these entities is, however, not unified. In some data sets intervals are assigned an entity representing a year, in other they are assigned an entity representing the last day of the corresponding year. All of them are, however, representing a period of a whole year.

```
@prefix skos: <http://www.w3.org/2004/02/skos/core#> .

<https://data.cssz.cz/resource/reference.data.gov.uk/id/gregorian-year/2017>
  a skos:Concept ;
  <http://www.w3.org/2002/07/owl#sameAs> <http://reference.data.gov.uk/id/gregorian-year/2017> ;
  skos:inScheme <https://data.cssz.cz/ontology/years/YearsScheme> ;
  skos:notation "2017" ;
  skos:prefLabel "2017" .
```

Listing 7.3: Dereferenced proxy entity of the year 2017

The CSSA1 dataset uses two distinct schemes for the pension kinds, because in 2010, the official categorization of pensions was changed in the Czech legislation. Only the observations assigned to year 2008 are divided according to the old pension scheme. All other year's observations correspond to the new pension scheme. So one can not simply multiply the numbers of distinct values for each dimension and the number of measures to get the total number of observations for this particular cube. The URIs of both pension kind schemes are suffixed either by `_2008` (31 of them) for the old scheme or `_2010` (37 of them) for the new scheme. Not all entities correspond to a particular kind of pension. Some of them represent an aggregation over related pension kinds or simply an aggregation over all of the pension kinds.

```
@prefix skos: <http://www.w3.org/2004/02/skos/core#> .

<https://data.cssz.cz/resource/pension-kind/PK_SRNVN_2010>
  a <https://data.cssz.cz/ontology/pension-kinds/PensionKind> , skos:Concept ;
  skos:altLabel "SRNVN"@cs ;
  skos:exactMatch <https://data.cssz.cz/resource/pension-kind/PK_SRNVN> ;
  skos:inScheme <https://data.cssz.cz/ontology/pension-kinds/PensionKindScheme_2010> ;
  skos:notation "PK_SRNVN" ;
  skos:prefLabel "Starobní důchod SRN vyplácený v souběhu s vdoveckým důchodem"@cs .
```

Listing 7.4: Dereferenced pension kind

The dimension of sex consists of three distinct values: dimension of male pension, dimension of female pensions and its total. The values are proxy entities linking to the SDMX representations of sexes¹³

¹²<http://old.datahub.io/dataset/data-gov-uk-time-intervals>

¹³<http://purl.org/linked-data/sdmx/2009/code>

7.2 Czech Statistical Office

The Czech Statistical Office (CZSO) is the main public organization responsible for collecting and analyzing statistical data in the Czech Republic. This organization is for example responsible for the state's census. Data about demography, economics, education, health care etc. are made available on the organization's website¹⁴ in a form of interactive spreadsheet builder. Thanks to the Opendata.cz initiative this data sets are made available as LOD modelled by the Dacube Vocabulary in the initiative's catalogue. The data is also hosted as a SPARQL endpoint¹⁵. 8 of these data sets have a dimension of reference period. Each data set's dump file can be downloaded from the catalogue¹⁶.

```
@prefix qb: <http://purl.org/linked-data/cube#> .
@prefix czso: <http://data.czso.cz/ontology/> .

<http://data.czso.cz/resource/observation/job-applicants-and-unemployment-rate/CZ0513/2009-12-31/T> a
    qb:Observation ;
    czso:refArea <http://ruian.linked.opendata.cz/resource/okresy/3505> ;
    czso:refPeriod <http://reference.data.gov.uk/id/gregorian-day/2009-12-31> ;
    czso:sex sdmx-code:sex-T ;
    czso:neumisteniUchazeciOZamestnani 9692.0 ;
    czso:dosazitelniNeumisteniUchazeciOZamestnani 9528.0 ;
    czso:podilNezamestnanych 7.92 ;
    czso:pocetVolnychMist 569.0 ;
    qb:dataSet <http://data.czso.cz/resource/dataset/job-applicants-and-unemployment-rate> .
```

Listing 7.5: Example of an observation from the CZSO data sets

Observations in CZSO data cubes are assigned multiple measures. Their URIs are not dereferenceable. Their dimension value URIs do not work as proxy entities. The dimension of reference area uses entities of an above-mentioned initiative's unofficial RTIAR transformation¹⁷ with its own SPARQL endpoint¹⁸. The measured values relate to regions and district. They do not contain observations related to the whole state. The proxy entities of the reference area dimension values for the CSSA data sets previously linked to this code list.

For time intervals representation the CZSO data cubes also use the the `data.gov.uk` Time Intervals OWL ontology. They only do so directly unlike the CSSA data cubes. The data cubes vary their time span. The earliest recorded values are for the year 2005. The latest values are for the year 2013. There are two data sets that contain values for both the earliest and latest year mentioned meaning they cover a period of 9 years: **czso-deaths-by-selected-causes-of-death**¹⁹ and **czso-job-applicants-and-unemployment-rate**²⁰

¹⁴<https://vdb.czso.cz/vdbvo2/faces/en/index.jsf?page=uziv-dotaz>

¹⁵<http://linked.opendata.cz/sparql>

¹⁶The download link URLs are, however, broken and return HTTP status code 404. To get the dump file, word *dumps* has to be substituted with word *soubor* (czech word for *file*). For example, the dump file of the data set **czso-job-applicants** is available at <https://linked.opendata.cz/soubor/czso-job-applicants.trig>

¹⁷<https://linked.opendata.cz/dataset/cz-ruian>

¹⁸<https://ruian.linked.opendata.cz/sparql>

¹⁹<https://linked.opendata.cz/dataset/czso-deaths-by-selected-causes-of-death>

²⁰<https://linked.opendata.cz/dataset/czso-job-applicants-and-unemployment-rate>

Just as with the CSSA data cubes, some of the CZSO data cubes contain the dimension of sex consisting of three distinct values: dimension of male pension, dimension of female pensions and its total. The values used are the SDMX representations of sexes themselves.

7.3 Wikidata

Wikidata data set contains data about political representation of countries, their administrative areas and municipalities. For the purposes of this experiment, such data concerning the Czech Republic was extracted from the data set. In the Czech Republic, regions and municipalities²¹ are being assigned a government that emerges from elections. In Wikidata data set, there exist records of who was or still is head of this local government, including the head of the state government. Records of these *head of government* roles are given a time period of validity of this role by stating the date of this role's start and optionally the end of this role when it is not a current area government head anymore. For the persons who hold or held the office, the affiliation to a political party is stated in the data also with the start and end date of this affiliation. The entities of the political parties are assigned their political alignment (left, center, far-right etc.). In the sample of Wikidata data set's content below it is stated that since 2008 till 2016 the head of the government of the South Moravian Region was Michal Hašek who since 1998 is a member of the Czech Social Democratic Party, which has the centre-left political alignment.

```
@prefix wd: <http://www.wikidata.org/entity/> .
@prefix p: <http://www.wikidata.org/prop/> .

wd:Q192697 rdfs:label "South Moravian Region"@en ;
  p:P6 [ p:P6 wd:Q6835752 ; p:P580 "2008-11-21T00:00:00Z" ; p:P582 "2016-11-16T00:00:00Z" ] .

wd:Q6835752 rdfs:label "Michal Hašek"@en ;
  p:P102 [ p:P102 wd:Q341148 ; p:P580 "1998-01-01T00:00:00Z" ] .

wd:Q341148 rdfs:label "Czech Social Democratic Party"@en ;
  p:P1387 [ p:P1387 wd:Q737014 ] .

wd:Q737014 rdfs:label "centre-left"@en .
```

Listing 7.6: Description of XXX (Source: author)

When adequately preprocessed, this data can be utilized to find relations of statistical data described in the data cubes of CSSA and CZSO and the political cycle in the country. For example a rule can be found that states, that if in any year, the head of the Czech Republic's government was a member of a left-leaning political party, the pension expenses for one-off allowance to pensions were above average. A query that extracts this data from the Wikidata's SPARQL endpoint is listed in A.2. This data, however, cannot be used for mining such rules yet. Measures in the data cubes are recorded on year to year bases. For the governmental roles and political affiliations it is only known the start date and end date.

²¹Not districts though.

It is necessary to transform these triples into set of triples stating that a governmental role or the political affiliation *applies* for a certain year. The edge years are a bit tricky because the role or the membership was not valid for the whole year it started or ended. To facilitate the query and to generate more triples I decided to generate the triples for the edge years as well. The SPARQL query that constructs the *appliesToRefPeriod* triples from the extracted data is listed in A.3. The triples stating the start dates and end dates are no longer needed and do not have to be loaded into the RDFSRules mining task.

7.4 Filtering the Observations

The values for the city of Prague are duplicated. The city is assigned an entity both as a region and as a district. The administrative area of the Czech Republic's capital is given a special status by the Act No. 131/2000 Coll., on the Capital City of Prague and does not in fact fall into neither of those categories. Nonetheless, Prague is assigned an identifier both as a district (3100) and as a region (19) in the RTIAR registry. When it comes to total population of the area (1 324 227 as of 2020), it is comparable to other czech regions. The least populated is the Karlovy Vary Region with around 300 000 inhabitants and the most populated: the Central Bohemian Region has around 1 300 000 inhabitants. Its surface area is on the other hand comparable with the districts. As the statistics about pensioners are certainly more correlated with the population rather than the surface area, the observation allocated to the dimension value of Prage as being district would be filtered out to maintain the commensurability along values measured for districts (see 7.5).

I also chose to filter out all observations regarding year 2008. For 2008 the pension kinds are structured according to a different scheme than any other year and it is hard to assume compatibility for the URIs that only differ in the year's suffix. 2008 scheme contains penkind kinds that 2010 does not end vice versa. It could be possible to just cut the cube so that the year's 2008 become one cube and the other years the other one, but a cube concerning only one reference period has not got much value. Both filters can be performed in a single SPARQL query. In the CSSA1 data set, discarding the Prague as a District entity removes 3 609 observations. The year 2008 contained 28 458 observations. 336 330 observations are contained in the query's result making up 91,4% of the unfiltered data set.

7.5 Slicing the Cubes

It is in part aimed to mine rules, in which the head atom's predicate is one of the cube's measures. In order to ensure, that such rules can achieve a reasonable support, the numerical values at the position of object in the measure triples have to be discretized and replaced by intervals. Irrespective to a chosen discretization approach, it is inadmissible to discretize values belonging to different disproportionate contexts. For example, we cannot create intervals for the number of pensioners from values measured for both regions and districts together. A

```

PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX qb: <http://purl.org/linked-data/cube#>
PREFIX cssa-d: <https://data.cssz.cz/resource/dataset/>
PREFIX cssa-od: <https://data.cssz.cz/ontology/dimension/>
PREFIX cssa-rd: <https://data.cssz.cz/resource/ruian/okresy/>
PREFIX cssa-op: <https://data.cssz.cz/ontology/pension-kinds/>

CONSTRUCT {
    ?observation ?p ?o
}
WHERE {
    GRAPH cssa-d:duchodci-v-cr-krajich-okresech {
        ?observation qb:dataSet cssa-d:duchodci-v-cr-krajich-okresech ;
        cssa-od:druh-duchodu ?druh ;
        ?p ?o .
        NOT EXISTS {
            ?observation cssa-od:refArea cssa-rd:3100 .
        }
    }
    GRAPH cssa-d:pomocne-ciselniky {
        ?druh skos:inScheme cssa-op:PensionKindScheme_2010 .
    }
}

```

Listing 7.7: SPARQL query to filter the CSSA1 data set (Source: author)

district is a lower administrative unit. It belongs to a lower level in the concept hierarchy and it is assumed that its numbers of pensioners are of a different order of magnitude than those for regions or for the whole state. Same applies for values of dimensions sex and pension kind of the described data set. The values of the reference period represent even time intervals so the commensurability is assumed.

One way to solve is to slice a preprocessed cube having disproportionate dimension values into a set of smaller subcubes, in which the dimension values belong to the same level of a concept hierarchy. Measured values can be then discretized into intervals in each subcube separately. Number of subcubes that the main cube has to be divided into depends on the number dimensions and the number of levels in each dimension's hierarchies. Also when the commensurability cannot be expected among dimension values on the same level of their hierarchy, it is a good idea to *make a cut* for each dimension value. For example, the dataset `cssa-d:vydaje-na-duchody-v-cr`²² capturing costs on pensions in the Czech Republic by year and kind of pension contains 10 distinct values of the dimension of pension kind (not considering the scheme used only for year 2008). The cube would have to be divided into 10 subcubes for each value of the pension kind dimension. The CSSA1 data set would have to be divided into 222 subcubes. It has 37 pension kinds. Dimension of sex has two hierarchy levels: each sex and total. The dimension of reference area is considered to have 3 hierarchy levels: State's total, regions and Prague districts combined with the regional districts since they are comparable in number of inhabitants.

²²<https://data.cssz.cz/resource/dataset/vydaje-na-duchody-v-cr>

$$37 \times 2 \times 3 = 222 \text{ subcubes}$$

The construction of a subcube from a *master* cube can be performed by a SPARQL CONSTRUCT query. An example of such query is shown in the listing 7.8. This query filters the triples of the CZSO data cube **czso-job-applicants-and-unemployment-rate** to create a smaller cube of statistics about job applicants and unemployment rate for districts by sex. Notice how the reference area values corresponding to districts are distinguished. After the reference area values are linked (see 7.6) to their CSSA counterparts, the ontology provided with the CSSA data sets can be reused.

```
PREFIX qb: <http://purl.org/linked-data/cube#>
PREFIX sdmx-c: <http://purl.org/linked-data/sdmx/2009/code#>
PREFIX czso: <http://data.czso.cz/ontology/>
PREFIX czso-rd: <http://data.czso.cz/resource/dataset/>
PREFIX cssa-rd: <https://data.cssz.cz/resource/dataset/>
PREFIX cssa-or: <https://data.cssz.cz/ontology/ruian/>
PREFIX owl: <http://www.w3.org/2002/07/owl#>

CONSTRUCT { ?observation ?p ?o }
WHERE {
  GRAPH czso-rd:job-applicants-and-unemployment-rate {
    ?observation qb:dataSet czso-rd:job-applicants-and-unemployment-rate ;
    ?p ?o ;
    czso:refArea ?refAreaCZSO .
    NOT EXISTS {
      ?observation czso:sex sdmx-c:sex-T .
    }
  }
  ?refAreaCSSZ owl:sameAs ?refAreaCZSO .
  GRAPH cssa-rd:pomocne-ciselniky { ?refAreaCSSZ a cssa-or:Okres }
}
```

Listing 7.8: SPARQL query to create a subcube (Source: author)

For every subcube a similar query has to be created and performed over the master cube. For a cube that has to be divided into a small number of subcubes it is plausible to write (and save it for the documentation and repeatability purposes) and perform these queries manually. But there are cubes for which this would involve an hours long work. At the same time, it is an trivial activity that can easily be automated. For this preprocessing step for the CSSA1 dataset, a Scala script was written that creates 222 distinct SPARQL queries that construct 222 subcubes, saves each query to a text file and also creates a shell script that triggers all queries and saves a result of each query to a distinct file in the turtle format. The script is listed in A.1. This, however, still requires writing such script for each preprocessed data cube and solve the problem of a time consuming resolution of the queries.

7.6 Linking

In order to find rules describing relations across multiple sources (meaning data cubes of CZSO, data cubes of CSSA and Wikidata triples) the entities either have to be assigned the same URIs or to be connected by the `owl:sameAs` statements. The shared dimensions of the CSSA and CZSO data cubes are the reference area, reference period and sex. The dimension values URIs used for these dimensions differ not only institution from institution but also data cube from data cube from the same institution (In the CSSA data cubes, reference period is represented by an entity of a year and of the last day of the year as well). Linking of equivalent dimension values of the three dimensions is done by creating `owl:sameAs` statements.

7.6.1 Sex Dimension Values

It was already mentioned that the CSSA data cubes use proxy entities linking to the SDMX representations of sexes, whereas the CZSO data cubes use these representations directly. So the linking statements are already provided with the CSSA code list file. To extract these very triples, the query listed below can be used. These triples can be then loaded into a mining task involving mining from data cubes contain the dimension of sex instead of loading the whole code list file.

```
PREFIX owl: <http://www.w3.org/2002/07/owl#>

CONSTRUCT { ?cssaSex owl:sameAs ?sdmxSex }
WHERE {
  GRAPH <https://data.cssz.cz/resource/dataset/pomocne-ciselniky> {
    ?cssaSex a <https://data.cssz.cz/ontology/sdmx/code/Sex> ; owl:sameAs ?sdmxSex
  }
}
```

Listing 7.9: Linking the sex dimension values

7.6.2 Reference Periods

In CSSA data cubes, values from two concept schemes are used for representing the year intervals: a years scheme and a days scheme. The entities in the schemes are proxy entities linking to the `data.gov.uk` Time Intervals ontology. CZSO data cubes use the ontology's day scheme concepts directly. That means that every year is represented by three distinct URIs in the data cubes so two `owl:sameAs` statements are required for each year. A query was written that generates theses statements for every year entity in the CSSA code list:

```

PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>

CONSTRUCT {
    ?cssaYear owl:sameAs ?cssaDay .
    ?cssaYear owl:sameAs ?dataGovDay .
    ?cssaDay owl:sameAs ?dataGovDay .
}
WHERE {
    GRAPH <https://data.cssz.cz/resource/dataset/pomocne-ciselniky> {
        ?cssaYear skos:inScheme <https://data.cssz.cz/ontology/years/YearsScheme>
        BIND (REPLACE(str(?cssaYear), ".*(\\d{4})", "$1") as ?cssaYearValue)
        ?cssaDay skos:inScheme <https://data.cssz.cz/ontology/days/DaysScheme>
        FILTER (REGEX(str(?cssaDay), ".*day.*12-31"))
        BIND (REPLACE(str(?cssaDay), ".*(\\d{4})-12-31", "$1") as ?cssaDayValue)
        ?cssaDay owl:sameAs ?dataGovDay
        FILTER (?cssaYearValue = ?cssaDayValue)
    }
}

```

Listing 7.10: Linking the reference periods

7.6.3 Reference Areas

The proxy entities of the reference area in CSSA data set used to link to the same entities that are used by the CZSO data sets. This linking is no longer present in the CSSA's code list but can be retrieved from the Opendata.cz's SPARQL endpoint. The query listed below returns 114 linking triples for all districts (including the Prague district entity), all regions (including Prague), Prague districts and the entity of the whole state. The linking with the Wikidata's entities had to be performed manually.

```

PREFIX owl: <http://www.w3.org/2002/07/owl#>

CONSTRUCT {
    ?cssaArea owl:sameAs ?odArea
}
WHERE {
    ?cssaArea a ?class ; owl:sameAs ?odArea .
    FILTER (?class IN (
        <https://data.cssz.cz/ontology/ruian/Okres>,
        <https://data.cssz.cz/ontology/ruian/Vusc>,
        <https://data.cssz.cz/ontology/ruian/SpravniObvod>,
        <https://data.cssz.cz/ontology/ruian/Stat>
    ))
}

```

Listing 7.11: Linking the reference periods

7.7 Discretization

The RDRules implementation provides discretization functionality. The discretization tasks are, however, federated to the implemented discretization algorithms of the EasyMiner-Discretization library²³. The equifrequent and the equisize discretization were chosen for the purposes of this experiment. The count of intervals that is set to be created from the set of measured values while performing the equifrequent discretization determines how many a which rules are generated by the RDRules algorithm. If the the values are discretized into a small number of intervals, more rules with should be generated but the measure values become coarse and information is lost. If creating too many intervals, more specific rules should be found but they happen to have lower support. When performing the equisize discretization, coarser rules are found for the intervals created for a higher support.

To avoid guessing, which number of intervals and which minimal support suits best the preprocessed data, multiple discretizations were performed with different parameters for both discretization algorithms. The minimal support thresholds were calculated as absolute numbers of various percentages of observations in the discretized cubes. As it was already mentioned, the preprocessed cube has to be cut into subcubes with commensurable observations, measures in theses subcube have to be discretized separately and only after that the triples can be merged and performed the mining tasks on.

That means that the number of overall measurements multiplies by the number of distinct discretizations. A situation has to be avoided, when the instantiations of variable representing observations are involving observations not only from various subcubes. The approach to solve this problem no matter how many measures are assigned to each observation will be shown on a sample of data below.

```
@prefix qb: <http://purl.org/linked-data/cube#> .

<o1> qb:dataSet <original-dataset> ;
    <dimension1> <dimension1value1> ; <dimension2> <dimension2value1> ;
    <measure1> 25000 ;
    <measure2> 3 .

<o2> qb:dataSet <original-dataset> ;
    <dimension1> <dimension1value2> ; <dimension2> <dimension2value2> ;
    <measure1> 10000 ;
    <measure2> 10 .
```

Listing 7.12: XYZ

Each application of a discretization algorithm will create a new measurement triple for each measure and observation with an object of the assigned interval based on the discretization algorithm and the parameter. The objects in triples assigning the observations to a data set will be changed to point to the particular subcube. In the example below two discretizations for each measure were performed on the two observations. In the example the same pair of

²³<https://github.com/KIZI/EasyMiner-Discretization>

discretizations were performed on the two distinct measures, but that does not have to be so. Assigning multiple triples of the same measure is fine as far as the measure is differently discretized.

```
@prefix qb: <http://purl.org/linked-data/cube#> .

<o1> qb:dataSet <subcube1> ;
    <dimension1> <dimension1value1> ; <dimension2> <dimension2value1> ;
    <measure1> <subcube1_ef3_measure1_3>, <subcube1_es10_measure1_2> ;
    <measure2> <subcube1_ef3_measure2_2>, <subcube1_es10_measure2_1> .

<o2> qb:dataSet <subcube2> ;
    <dimension1> <dimension1value2> ; <dimension2> <dimension2value2> ;
    <measure1> <subcube2_ef3_measure1_3> , <subcube2_es10_measure1_2> ;
    <measure2> <subcube2_ef3_measure2_1> , <subcube2_es10_measure2_1> .
```

Listing 7.13: XYZ

The discretized subcubes can be then merged into a single data set and be performed mining tasks on. In each rule it has to be ensured, that the set of observations is limited to a certain subcube. For each cube in a rule the body of a rule has to contain an atom of a pattern (`?o qb:dataSet AnyConstant`).

An alternative to creating subcubes based on the concept hierarchy in the master cube's dimensions is to use an unsupervised clustering algorithm (eg. k-means) to divide the observations into subcubes based on the proximity of their measures. After that the workflow is the same, the measures are discretized in the generated subcubes separately and then the subcubes are merged. But this brings two problems when used for the RDRules algorithm.

1. There is no clear interpretation of the generated subcubes, because their observation can belong to different levels of the concept hierarchy in a dimension. Unlike with the previous approach where a generated subcube could be described as for example *Population in districts by age category*.
2. For each distinct measure in the master cube a set of subcubes would have to be generated. So if the observations are assigned multiple measures (as the CZSO observations are) the number of observations would multiply with the number of measures and the observation URIs would have to be distinguished as one observation cannot be assigned to multiple data sets by the `qb:dataSet` triple.

7.8 Mining Tasks

7.8.1 Relation between the pension expenses and the political alignment of the state's government

7.9 Discussion

Conclusions

List of References

- [1] David Beckett et al. *RDF 1.1 Turtle*. 2014-02. URL: <http://www.w3.org/TR/turtle/>.
 - [2] Tim Berners-Lee. *Linked Data - Design Issues*. 2006. URL: <http://www.w3.org/DesignIssues/LinkedData.html>.
 - [3] Richard Cyganiak and Dave Reynolds. *The RDF Data Cube Vocabulary*. W3C Recommendation. <https://www.w3.org/TR/2014/REC-vocab-data-cube-20140116/>. W3C, 2014-01.
 - [4] Luis Antonio Galárraga et al. ‘AMIE: Association rule mining under incomplete evidence in ontological knowledge bases’. In: 2013-05, pp. 413–422. DOI: 10.1145/2488388.2488425.
 - [5] Tom Heath and Christian Bizer. *Linked Data: Evolving the Web into a Global Data Space*. Morgan & Claypool, 2011. ISBN: 9781608454303. URL: <http://linkeddatabook.com/>.
 - [6] Frank Manola and Eric Miller. *RDF Primer*. 2004-02. URL: <http://www.w3c.org/TR/rdf-primer/>.
 - [7] Eric Prud’hommeaux and Andy Seaborne. *SPARQL Query Language for RDF*. 2008-01. URL: <http://www.w3.org/TR/rdf-sparql-query/>.
-

Appendices

A. Queries and Scripts

```
import scala.io.Source
import java.io.PrintWriter
import scala.sys.process._

val byRegion = """
?observation cssz-dimension:refArea ?refAreaCSSA .
GRAPH <https://data.cssz.cz/resource/dataset/pomocne-ciselniky> {
  ?refAreaCSSA a <https://data.cssz.cz/ontology/ruian/Vusc> .
}
"""

val byDistrict = """
?observation cssz-dimension:refArea ?refAreaCSSA .
GRAPH <https://data.cssz.cz/resource/dataset/pomocne-ciselniky> {
  ?refAreaCSSA a ?class .
  FILTER (?class IN (
    <https://data.cssz.cz/ontology/ruian/Okres>,
    <https://data.cssz.cz/ontology/ruian/SpravniObvod>
  ))
}
"""

val stateTotal = """
?observation cssz-dimension:refArea <https://data.cssz.cz/resource/ruian/staty/1> .
"""

val bySex = """
NOT EXISTS {
  ?observation cssz-dimension:pohlavi <https://data.cssz.cz/ontology/sdmx/code/sex-T
    >
}
"""

val bothSexes = """
?observation cssz-dimension:pohlavi <https://data.cssz.cz/ontology/sdmx/code/sex-T> .
"""

val template = """
PREFIX qb: <http://purl.org/linked-data/cube#>
PREFIX cssz-dimension: <https://data.cssz.cz/ontology/dimension/>

CONSTRUCT {
  ?observation ?p ?o
}
WHERE {

  ?observation qb:dataSet <https://data.cssz.cz/resource/dataset/duchodci-v-cr-
    krajich-okresech> .
  ?observation ?p ?o .
  %s
  ?observation cssz-dimension:druh-duchodu <%s> .
  %s
}
"""

val commandTemplate = "../../../jena/bin/arq --data \"../../../data/pensions-filtered.ttl\""
```

```

--data \"../..../data/pomocne-ciselniky.trig\" --query \"queries/%s.rq\" > \"results
/%s.ttl\"
val source = Source.fromFile(\"pensionkinds.txt\")
val lines = source.getLines().toArray
val pw = new PrintWriter(s\"script.sh\")
lines.foreach(line => {

    val kindName = line.replaceAll(\"https://data.cssz.cz/resource/pension-kind/\", \"\").
        replaceAll(\"_2010\", \"\")
    val districtTotalQuery = template.format(bothSexes, line, byDistrict)
    val districtTotalFileName = s\"pensions-by-district-total-$kindName\"
    val districtTotalPw = new PrintWriter(s\"queries/$districtTotalFileName.rq\")
    districtTotalPw.print(districtTotalQuery)
    districtTotalPw.close()
    pw.println(commandTemplate.format(districtTotalFileName, districtTotalFileName))

    val districtBySexQuery = template.format(bySex, line, byDistrict)
    val districtBySexFileName = s\"pensions-by-district-by-sex-$kindName\"
    val districtBySexPw = new PrintWriter(s\"queries/$districtBySexFileName.rq\")
    districtBySexPw.print(districtBySexQuery)
    districtBySexPw.close()
    pw.println(commandTemplate.format(districtBySexFileName, districtBySexFileName))

    val regionTotalQuery = template.format(bothSexes, line, byRegion)
    val regionTotalFileName = s\"pensions-by-region-total-$kindName\"
    val regionTotalPw = new PrintWriter(s\"queries/$regionTotalFileName.rq\")
    regionTotalPw.print(regionTotalQuery)
    regionTotalPw.close()
    pw.println(commandTemplate.format(regionTotalFileName, regionTotalFileName))

    val regionBySexQuery = template.format(bySex, line, byRegion)
    val regionBySexFileName = s\"pensions-by-region-by-sex-$kindName\"
    val regionBySexPw = new PrintWriter(s\"queries/$regionBySexFileName.rq\")
    regionBySexPw.print(regionBySexQuery)
    regionBySexPw.close()
    pw.println(commandTemplate.format(regionBySexFileName, regionBySexFileName))

    val totalTotalQuery = template.format(bothSexes, line, stateTotal)
    val totalTotalFileName = s\"pensions-total-total-$kindName\"
    val totalTotalPw = new PrintWriter(s\"queries/$totalTotalFileName.rq\")
    totalTotalPw.print(totalTotalQuery)
    totalTotalPw.close()
    pw.println(commandTemplate.format(totalTotalFileName, totalTotalFileName))

    val totalBySexQuery = template.format(bySex, line, stateTotal)
    val totalBySexFileName = s\"pensions-total-by-sex-$kindName\"
    val totalBySexPw = new PrintWriter(s\"queries/$totalBySexFileName.rq\")
    totalBySexPw.print(totalBySexQuery)
    totalBySexPw.close()
    pw.println(commandTemplate.format(totalBySexFileName, totalBySexFileName))
})
pw.close()

```

Listing A.1: Scala script for creating SPARQL queries for preprocessing the CSSA1 data set

```

PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX wd: <http://www.wikidata.org/entity/>
PREFIX wdt: <http://www.wikidata.org/prop/direct/>
PREFIX p: <http://www.wikidata.org/prop/>

CONSTRUCT {

```

```

?area p:P31 wd:Q5153359.
?area p:P6 ?headRole ;
    p:P131 ?superiorArea .
?headRole p:P6 ?person ;
    p:P580 ?headRoleStartDate ;
    p:P582 ?headRoleEndDate .
?person p:P102 ?partyRole .
?partyRole p:P580 ?partyRoleStartDate ;
    p:P582 ?partyRoleEndDate ;
    p:P102 ?party .
?party p:P1387 ?alignmentLabel .
}
WHERE {
    {
        ?area wdt:P31 wd:Q38911 .
    } UNION {
        ?area wdt:P31|p:P31 wd:Q5153359 .
        OPTIONAL {
            ?area wdt:P131|p:P131 ?superiorArea
        }
    } UNION {
        wd:Q3342946 wdt:P1269 ?area
    }
    OPTIONAL {
        ?area p:P6|wdt:P6 ?headRole .
        ?headRole ps:P6 ?person ;
            pqv:P580|pq:P580 ?headRoleStartDate .
        OPTIONAL {
            ?headRole pqv:P582|pq:P582 ?headRoleEndDate .
        }
        OPTIONAL {
            ?person p:P102|wdt:102 ?partyRole .
            OPTIONAL {
                ?partyRole pqv:P580|pq:P580 ?partyRoleStartDate .
            }
            OPTIONAL {
                ?partyRole pqv:P582|pq:P582 ?partyRoleEndDate .
            }
            OPTIONAL {
                ?partyRole ps:P102 ?party .
                ?party wdt:P1387|p:P1387 ?alignment .
                ?alignment rdfs:label ?alignmentLabel
                FILTER (lang(?alignmentLabel) = "en")
            }
        }
    }
}
}

```

Listing A.2: Extracting the Wikidata’s political data about the Czech Republic

```

PREFIX cssa-ody: <https://data.cssz.cz/ontology/days/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX p: <http://www.wikidata.org/prop/>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX dt: <http://kizi.vse.cz/novp19/diploma-thesis/>

CONSTRUCT {
    ?headRole dt:appliesToRefPeriod ?refPeriod .
    ?partyRole dt:appliesToRefPeriod ?refPeriod .
}
WHERE {

```

```

GRAPH <https://data.cssz.cz/resource/dataset/pomocne-ciselniky> {
  ?refPeriod skos:inScheme cssa-ody:DaysScheme .
  FILTER REGEX(str(?refPeriod), ".*year.*")
  FILTER REGEX(str(?refPeriod), "^.*\d{4}-12-31")
  BIND(REPLACE(str(?refPeriod), ".*(\d{4}).*", "$1") as ?refPeriodBind)
}

?area p:P6 ?headRole .
?headRole p:P580 ?headRoleStartDate .
FILTER ( datatype(?headRoleStartDate) = xsd:string)
BIND(REPLACE(str(?headRoleStartDate), "^(\\d{4}).*", "$1") as ?
      headRoleStartDateBind)
OPTIONAL {
  ?headRole p:P582 ?headRoleEndDate .
  FILTER ( datatype(?headRoleEndDate) = xsd:string)
  BIND(REPLACE(str(?headRoleEndDate), "^(\\d{4}).*", "$1") as ?
        headRoleEndDateBind)
}
FILTER (
  (bound(?headRoleEndDateBind) && ?refPeriodBind >= ?headRoleStartDateBind && ?
    refPeriodBind <= ?headRoleEndDateBind)
  ||
  (!bound(?headRoleEndDateBind) && ?refPeriodBind >= ?headRoleStartDateBind && ?
    headRoleStartDateBind > "2000")
)
OPTIONAL {
  ?headRole p:P6 ?person .
  ?person p:P102 ?partyRole .
  OPTIONAL {
    ?partyRole p:P580 ?partyRoleStartDate .
    FILTER ( datatype(?partyRoleStartDate) = xsd:string)
    BIND(REPLACE(str(?partyRoleStartDate), "^(\\d{4}).*", "$1") as ?
          partyRoleStartDateBind)
  }
  OPTIONAL {
    ?partyRole p:P582 ?partyRoleEndDate .
    FILTER ( datatype(?partyRoleEndDate) = xsd:string)
    BIND(REPLACE(str(?partyRoleEndDate), "^(\\d{4}).*", "$1") as ?
          partyRoleEndDateBind)
  }
  FILTER (
    (bound(?partyRoleEndDateBind) && ?refPeriodBind >= ?partyRoleStartDateBind
      && ?refPeriodBind <= ?partyRoleEndDateBind)
    ||
    (!bound(?partyRoleEndDateBind) && ?refPeriodBind >= ?
      partyRoleStartDateBind && ?partyRoleStartDateBind > "2000")
    ||
    (!bound(?partyRoleStartDateBind) && !bound(?partyRoleEndDateBind))
  )
}
}

```

Listing A.3: Creating the **appliesToRefPeriod** triples
