

Profiling Linked Open Data with ProLOD



TECHNISCHE
UNIVERSITÄT
DARMSTADT

LiDDM: A Data Mining System for Linked Data



*Seminar aus maschinellem Lernen
Frederik Janssen, Dr. Heiko Paulheim*

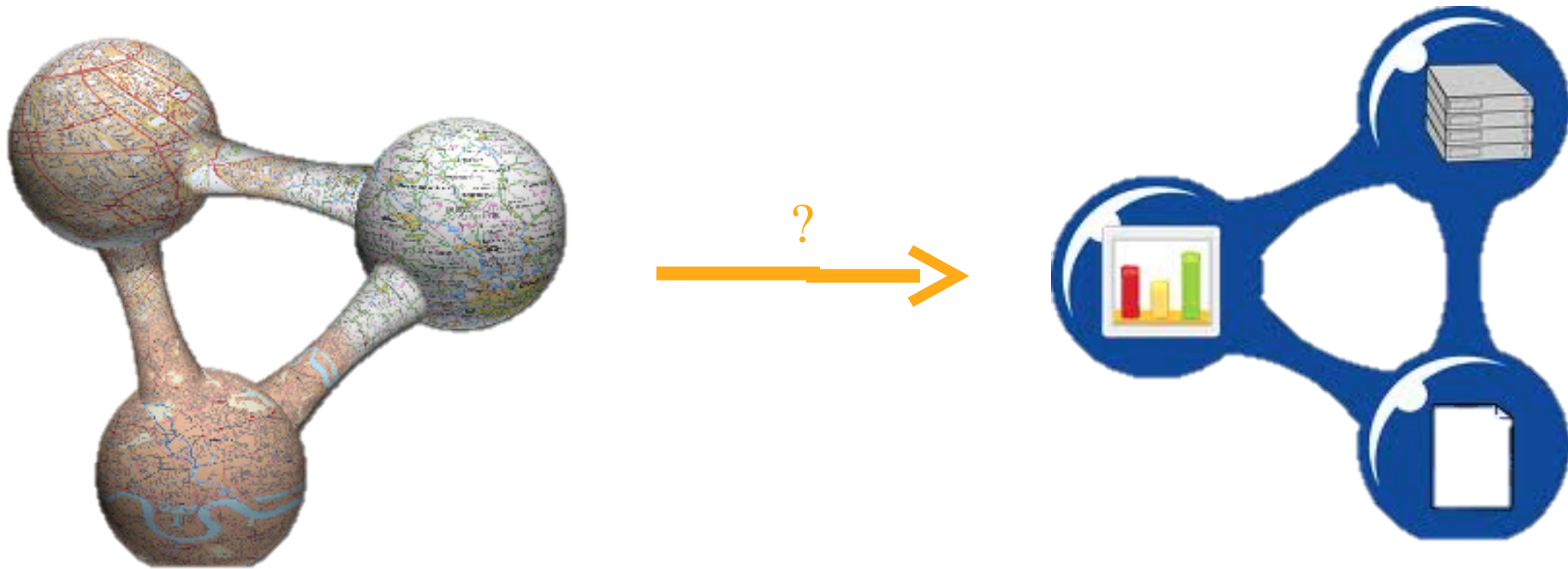
Profiling Linked Open Data with ProLOD

1. Linked Open Data
2. Profiling
3. Clustering and Labeling
4. Datatypes and Patterns

LiDDM: A Data Mining System for Linked Data

1. Einführung
2. Linked Data Data Mining Modell
3. Tool Umgebung
4. Ergebnisse

Linked Open Data



Traditionelle Profiling-methoden reichen nicht aus.

Lösung : Web-basiertes Program : **ProLOD**

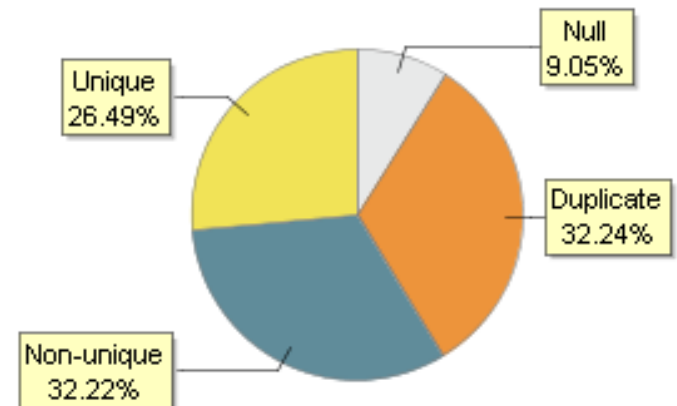


Erkennung einer Semantik von Spalten

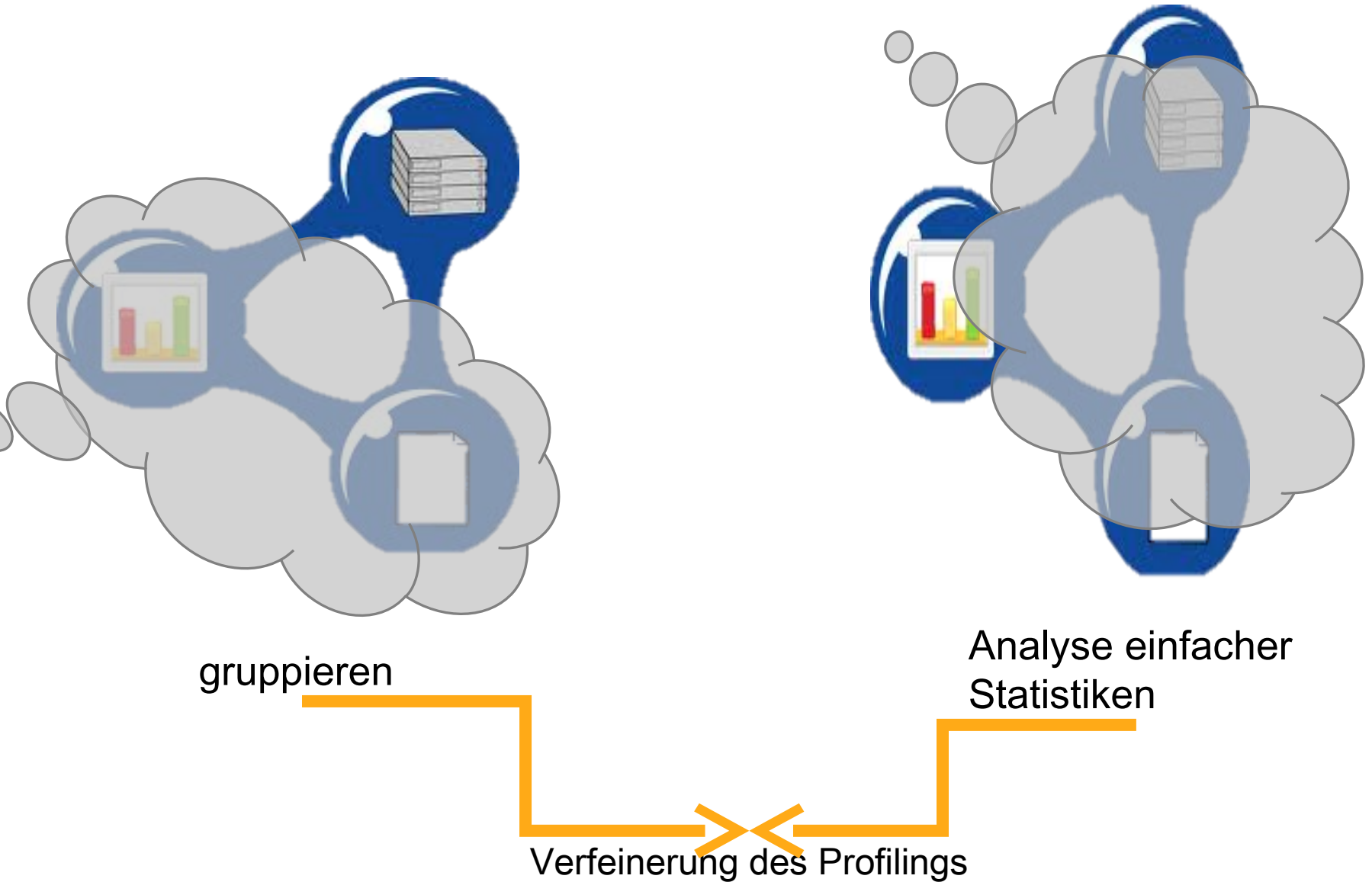
	Expression	Type	Domain	Non-null	Null	Unique	Distinct
<input type="checkbox"/>	src_name	string	specval pa...	139,317	0	130,464	132,575
<input type="checkbox"/>	src_gender	string	enum	15,780	123,537	0	4
<input type="checkbox"/>	src_birth_d...	string	day pattern	124,340	14,977	2,798	25,018
<input type="checkbox"/>	src_sin	string	pattern	134,830	4,487	134,766	134,798
<input type="checkbox"/>	src_card	string	long	55,904	83,413	15,800	35,630
<input type="checkbox"/>	src_address	string	enum patt...	9	139,308	5	7
<input type="checkbox"/>	src_primar...	integer		139,317	0	139,315	139,316

Zweck:

- Erkennung von einzigartigen Variablen
- Erstellung von Patterns



Profiling-LOD Methode



Proof-of-Concept

Ohne Ontologie

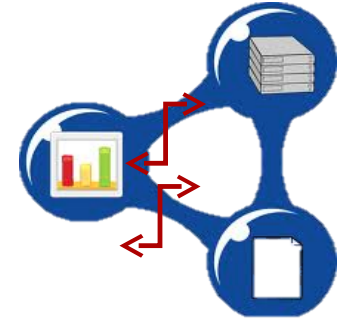
und



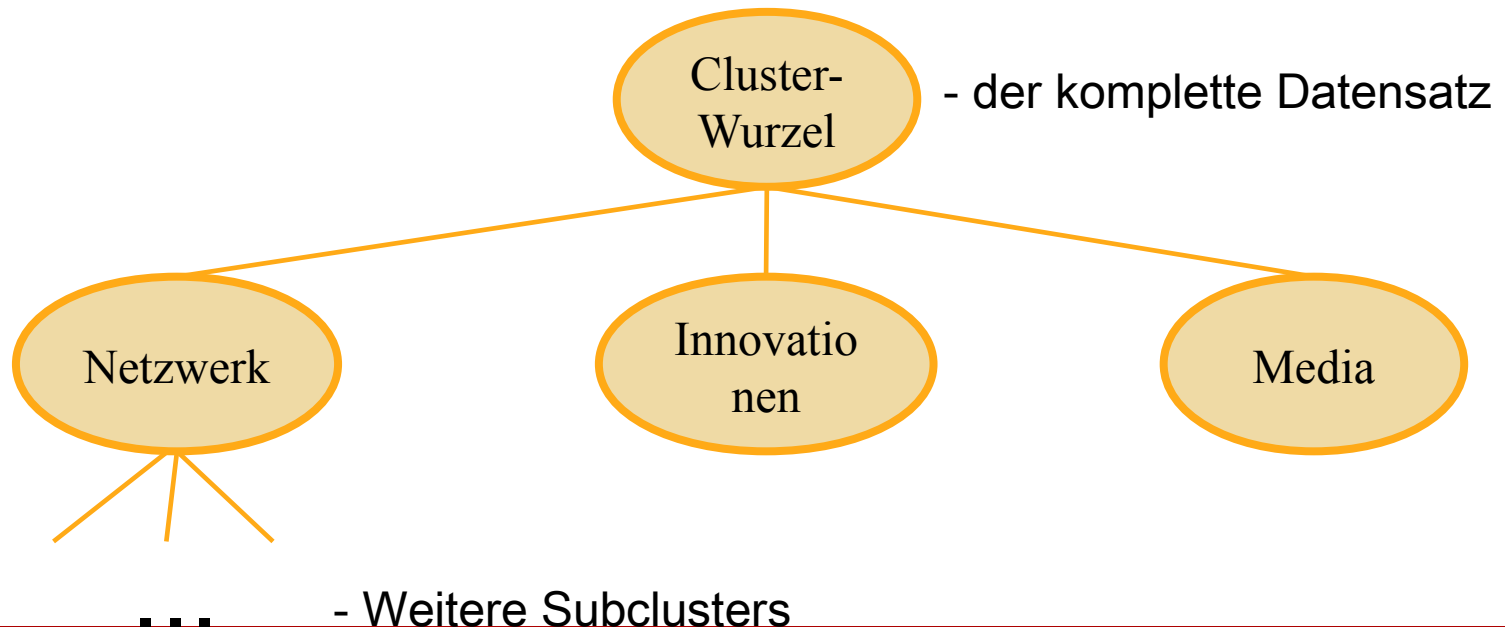
Clustering

Ein gutes Clustering:

semantikabhängig



Probleme: - die Clusteranzahl ist nicht bekannt.
- die Größe der Datensätze



Strategie bei einer großen Anzahl der Tripeln :

↓
Wurzelcluster
im Vorbereitungsschritt erstellen

↓
Die Grenzwerte der Unähnlichkeiten
gelten zur Bildung eines neuen Clusters

↓
K-Means Algorithmus

Labeling

TF-IDF um die m-wichtigsten Terme zu berechnen.

Beispiel für $m = 3$

Label	Sample subjects
<i>minister politician mayor</i>	Angela_Merkel Ted_Kennedy Presidency_of_George_Washington
<i>film directed starring</i>	Titanic_(1997_film) Metropolis_(film) Frankenstein_(1910_film)
<i>club football league</i>	FC_Bayern_Munich Liverpool_F.C. Los_Angeles_Galaxy

Für einen besseren Zusammenhang und
für das Verständnis der Tripelsätze findet man

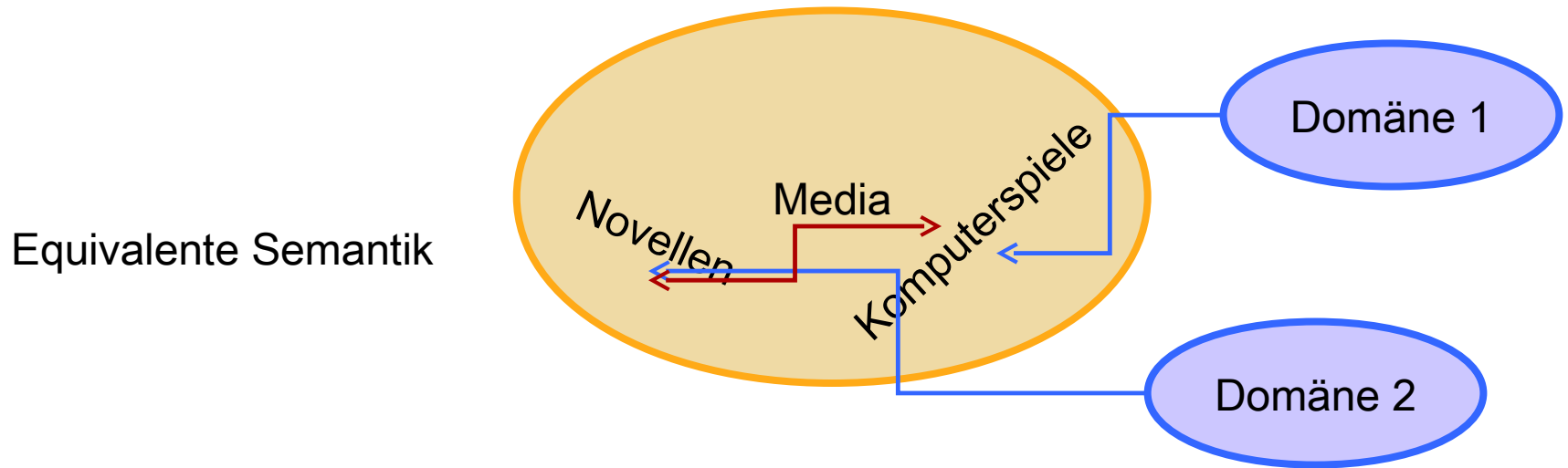
- die assoziativen Attribute
- die Assoziationsregeln im Cluster.

Als Ergebniss für ein Buch in “Media” kommen häufig author, isbn, genre vor.

Rule	Confidence	Correlation Coefficient
<i>genre, isbn \Rightarrow author</i>	0.99	0.67
<i>isbn \Rightarrow author</i>	0.92	0.66
<i>isbn \Rightarrow author, genre</i>	0.83	0.66
<i>author, genre \Rightarrow isbn</i>	0.70	0.66
<i>author \Rightarrow isbn</i>	0.64	0.66
<i>author \Rightarrow genre, isbn</i>	0.58	0.67

Das Modell ist in der Lage, sowohl die positiven, als auch die negativen Assoziationsregeln zu erkennen, um die Abhängigkeiten zwischen Prädikaten in einem Cluster zu finden.

Zweite Möglichkeit : negative Abhängigkeiten.



author $\Rightarrow \neg$ developer und developer $\Rightarrow \neg$ author : Konfidenz 90%

Genauso mit name und title



X hat eine Verbindung zur Y durch Prädikat A

Je mehr Einheiten durch A und B verbunden sind, desto höher ist der Zusammenhang der inversen Prädikatenpaare.

z.B. für das Buch "Into the Wild" :

Into the Wild ^{debutWorks} Jon Krakauer

Jon Krakauer ^{author} Into the Wild

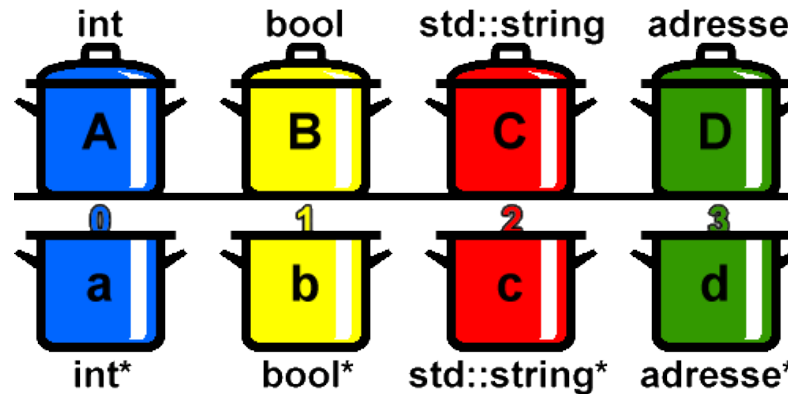
Tabelle der inversen Linkpaare :

<u>PredicateA</u>	<u>PredicateB</u>	Corr Coef	Frequency
before	after	0.239	28856
sisterStations	sisterStations	0.749	7494
precededBy	followedBy	0.830	7097
spouse	spouse	0.322	1964
before	before	-0.003	738
star	exoplanet	0.895	188

Wichtig für das semantische (Pre)-Clustering.

Patterns

Variablenverteilung in den Datentypen



Im Bereich der Datentypen verteilt

Patterns

Normalisierte Patterns



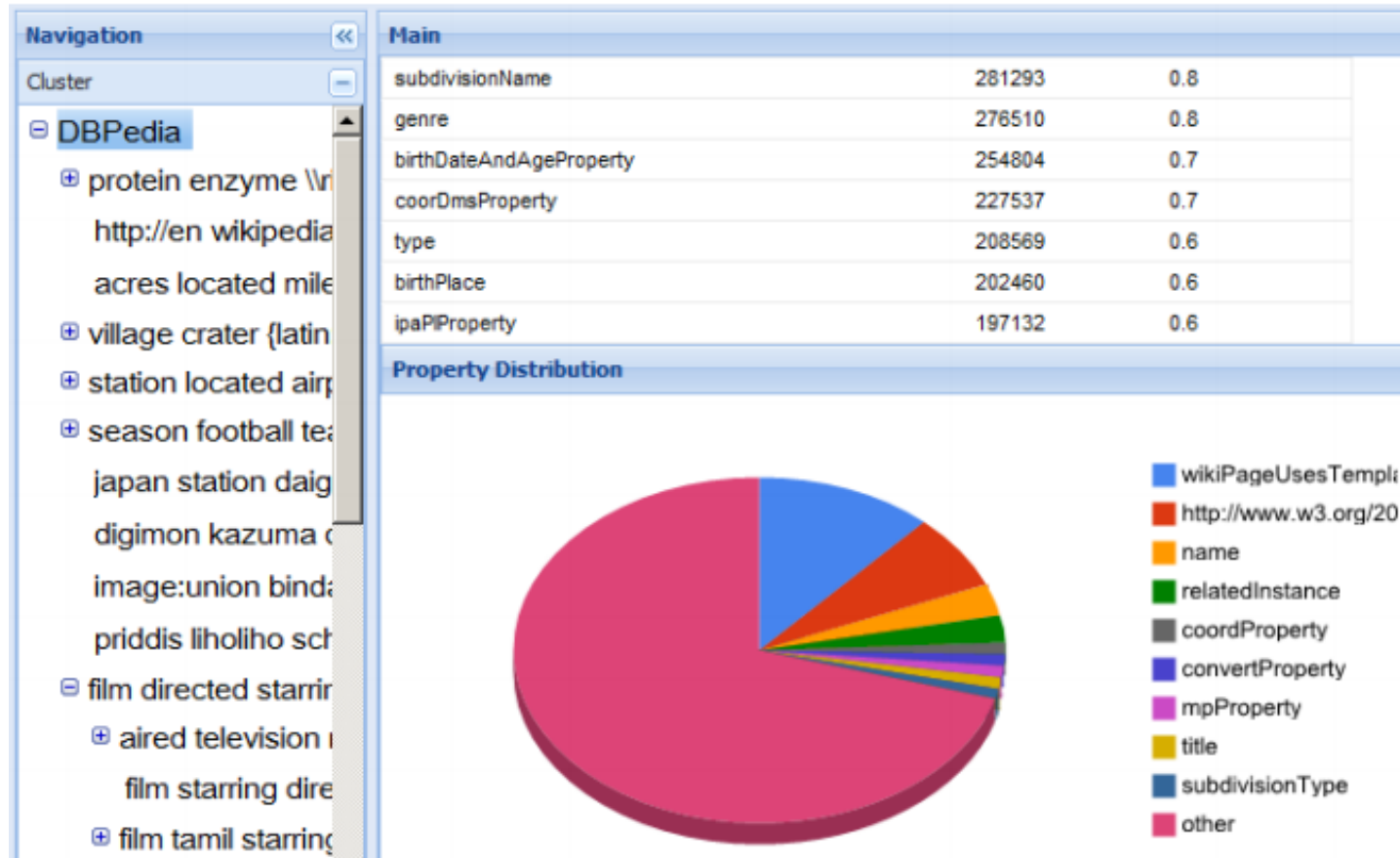
Preprocessing

1. Tripel umgewandelt in der Datenbank gespeichert
2. Clustering & Labeling
3. Jedem Tripel ein Pattern bzw. Datentyp zugewiesen



Aufwand : 1 Tag mit DBPedia

Realtime Profiling



A Data Mining System for Linked Data

extrem wachsende verlinkte Daten:

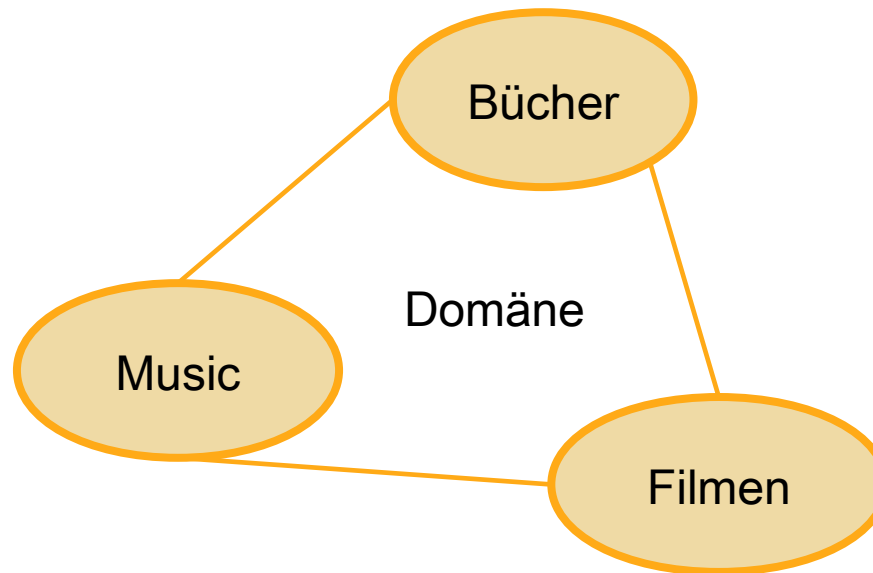


- Sammeln von verschiedenen Quellen
- Integrieren
- Für Statistiken verfeinern

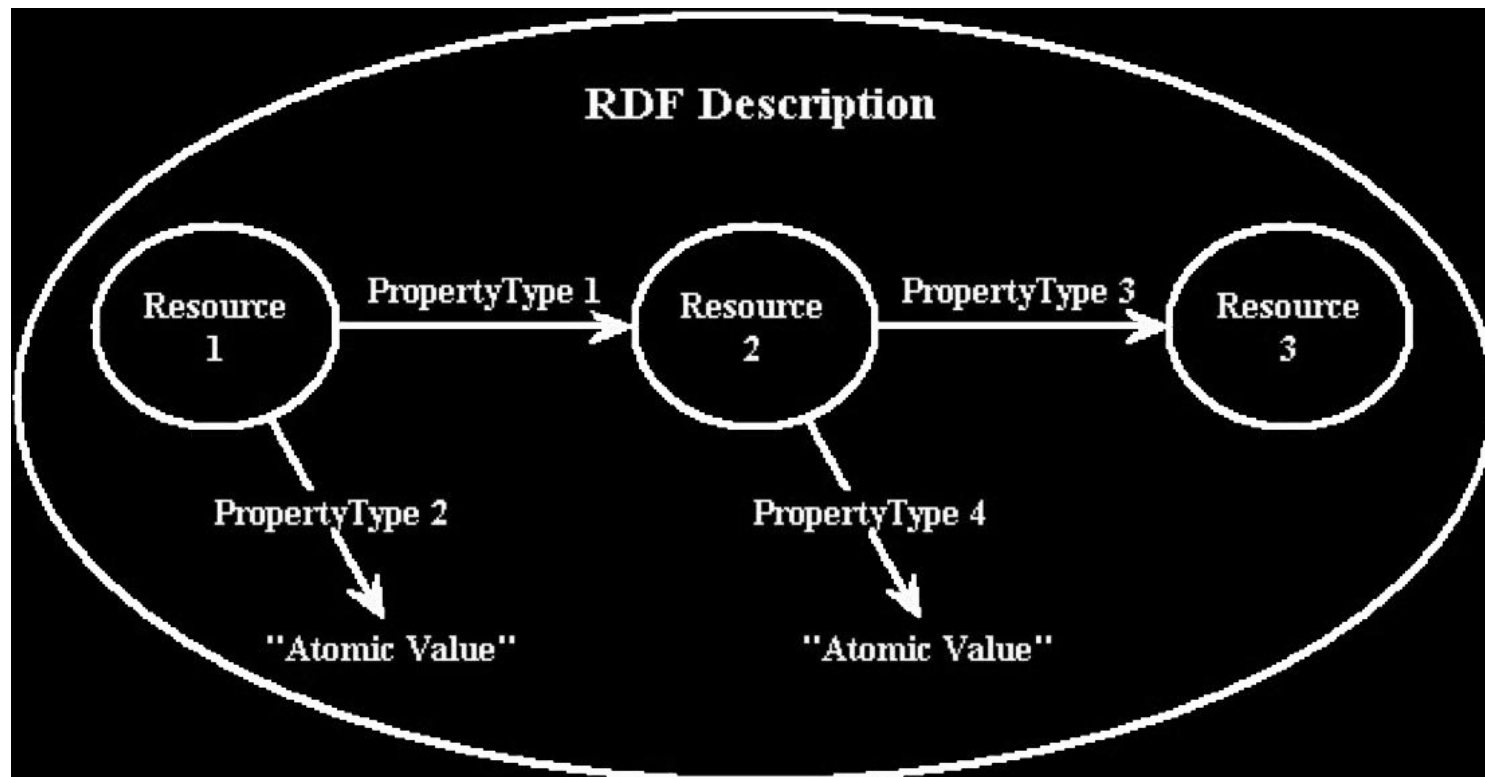
Deswegen hat man dafür ein Modell entwickelt

komplex und dynamisch.

- Fundament für Linked Data bei Web
- verknüpft die Daten aus verschiedenen Domänen



- bietet verschiedene Datasätze in RDF-Format:



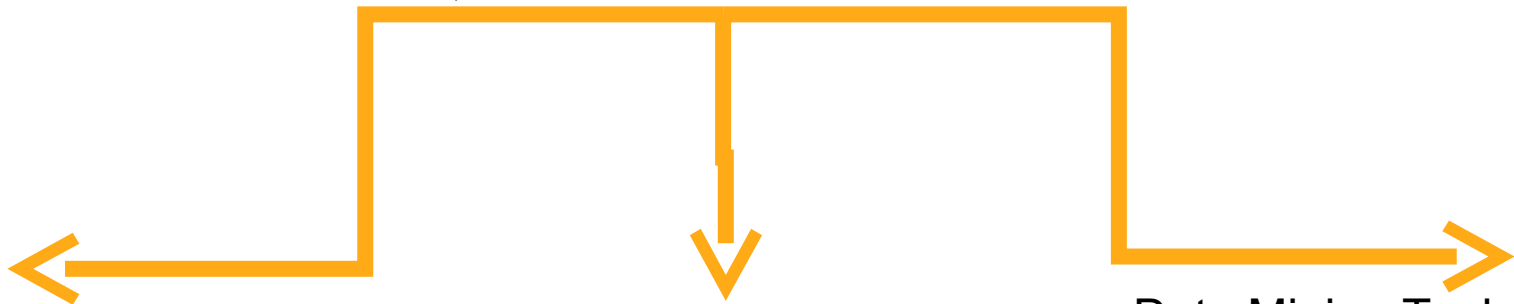
Wordl FactBook

Data.gov

DBPedia

- Informationen aus der reellen Welt
- Vorhersage für künftige Statistiken

Die wertvollen, versteckten Informationen extrahieren.

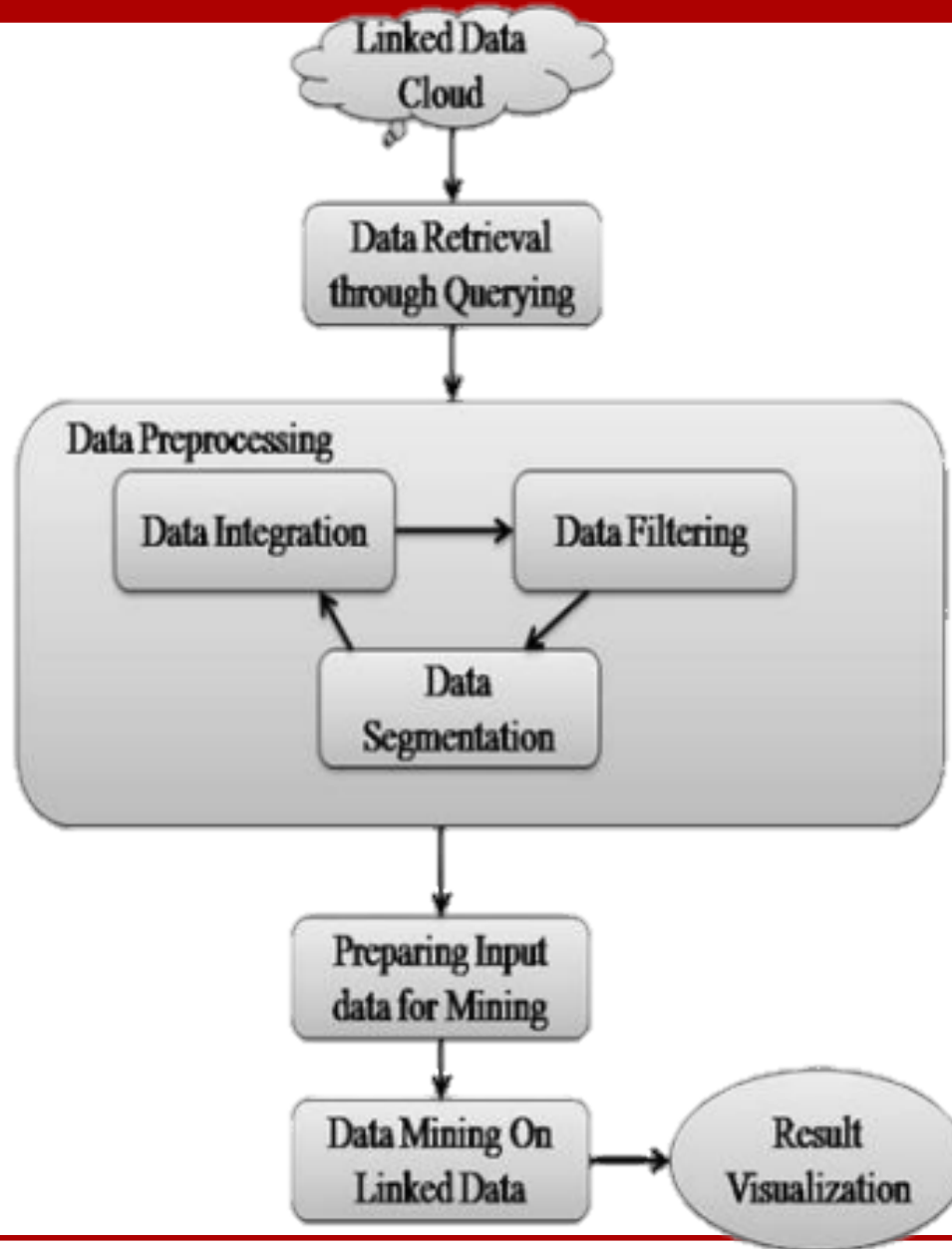


Daten aus Linked Data Cloud holen

KDD durchführen

Data Mining Techniken
(assoziationen, Clustering)

visualisieren

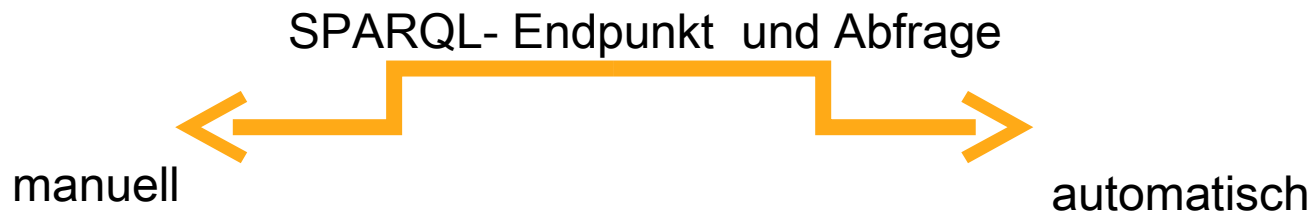


Tool Umgebung

Linked Data Data Mining Tool - LiDDMT



Schritt 1.



Abfrage Builder liefert die möglichen Prädikaten und eine spezifizierte Abfrage.

FOR COMPLEX QUERIES

CrossProduct
 CrossProduct
 Append at the End

Query 1 colNo

Query 2 colNo

ShowResult

Add Another Query

continue

Skip if you have just one query

RESULT-- QUERY 1

70.2	Aruba	100
69	Antigua and Barbuda	
78.5	United Arab Emirate	
53	Afghanistan	318
67.9	Algeria	333
67.7	Azerbaijan	812
66.6	Albania	360
69.3	Armenia	297
71.2	Andorra	718
53.5	Angola	122
63.5	American Samoa	
64.4	Argentina	403
67.4	Australia	204
67.5	Austria	819

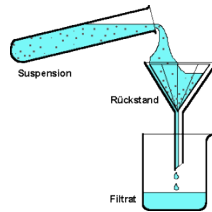
RESULT-- QUERY 2

Final Fantasy: The Spirits Within 2001-07-1		
LÃ@on	1994-11-18	France
Manufacturing Consent: Noam Chomsky a		
Manufacturing Consent: Noam Chomsky a		
Manufacturing Consent: Noam Chomsky a		
Manufacturing Consent: Noam Chomsky a		
Mulholland Drive 2001		
O Brother, Where Art Thou? 2000-05-1		
O Brother, Where Art Thou? 2000-05-1		
Original Sin	2001	NotFOUN
Star Trek: Generations 1994-11-1		
Star Trek: First Contact 1996-11-2		
Stargate	1994-10-28	France
The Rock	1996-06-07	NotFOUN

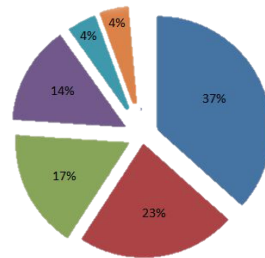
RESULT--CURRENT STATE AFTER MERGING BOTH THE QUERIES

Australia	20434176	Manufacturing Consent: Noam Chomsky and the Media	
Australia	20434176	Babe	1995-08-04 Australia
Australia	20434176	Shine	1996-01-21 Australia
Australia	20434176	Moulin Rouge!	2001-05-09 Australia
Australia	20434176	Holy Smoke!	1999-09-04 Australia
Australia	20434176	Muriel's Wedding	1995-03-10 Austr
Australia	20434176	The Adventures of Priscilla, Queen of the Desert 1994-0	
Australia	20434176	Babe: Pig in the City	1998-11-25 Austr
Australia	20434176	Mighty Morphin' Power Rangers: The Movie 1995-0	
Australia	20434176	Chopper	2000-08-03 Australia

Schritt 3. Fillterung



Schritt 4.



Schritt 5. Die Daten im geeigneten Format für Mining speichern : ARFF – Attribute-Relation File Format

LIDDMT

STATISTICAL ANALYSIS

☐ Use the current file for analysis

☒ Use some other file for analysis

Look In: Project_Total_K.V.N.PAVAN

- LiDDM Data Mining System for Linked Data
- Presentations during Internship
- Project Source Folder
- Research Paper Files
- Some ARFF FILES
- case study 1 factbook percapita.txt
- country data for films.txt
- factbook test for movie.txt
- LiDDM Data Mining System for Linked Data.r
- LIDDMT Demo Movie.avi

File Name:

Files of Type: All Files

Open Cancel

Visualize

OPTIONS

Associate

Classify

Associate

Cluster

Apriori

car false

Class Index -1

Delta 0.05

Lower Bound Min Support 0.1

Metric Type Confidence

Min Metric 0.9

Number of Rules 10

Output Item Sets false

Remove all missing columns false

Ergebnisse

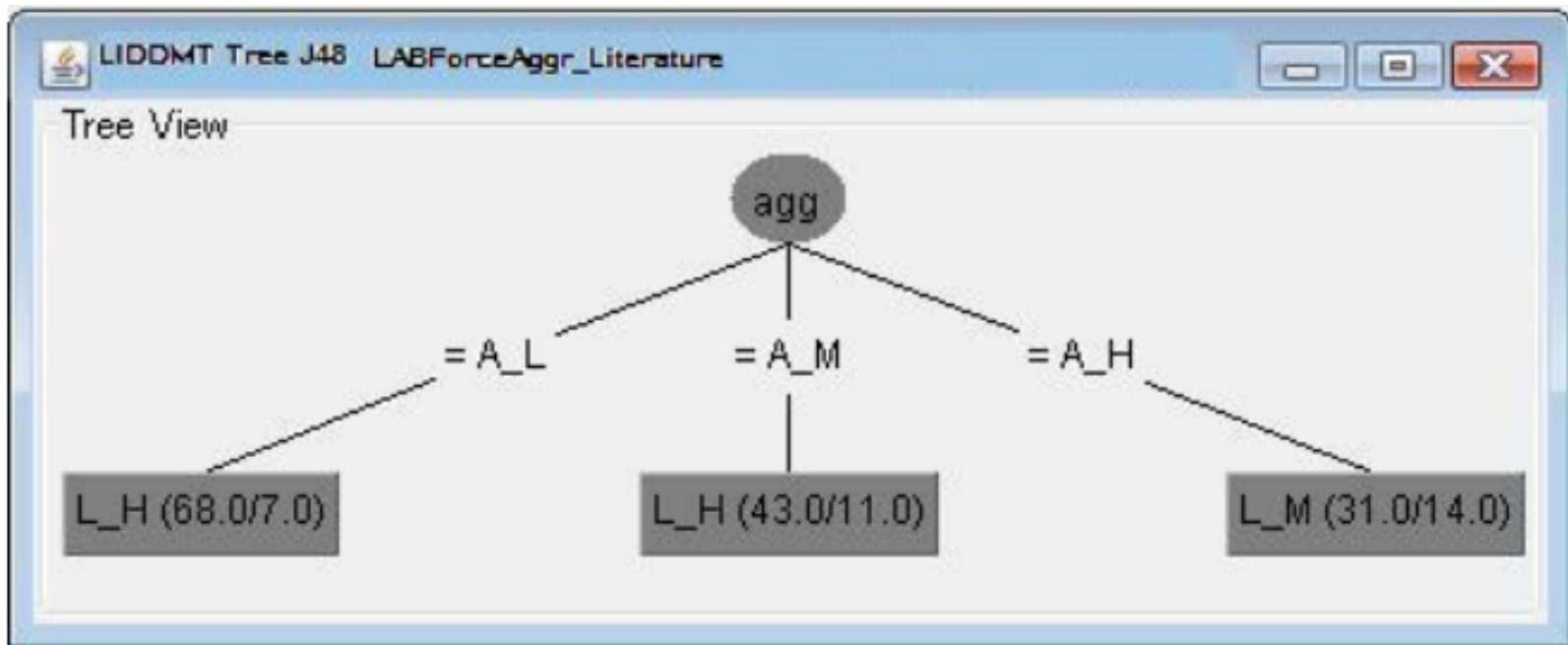
World FactBook für Bildungsniveau der Arbeitskräfte in der Landwirtschaft, in der Wirtschaft und im Services-Bereich.

Attribute sind in drei Klassen segmentiert;

A_L: agriculture low
A_M: agriculture medium
A_H: agriculture high.

Bildungsniveau in drei Klassen segmentiert:

low	0 - 50
medium	50 - 85
high	85 - 100





Danke für Ihre Aufmerksamkeit!





Autoren:

“Profiling Linked Open Data with ProLOD” : Christoph Böhm, Felix Naumann, Ziawasch Abedjan, Dandy Fenz, Toni Grütze, Daniel Hefenbrock, Matthias Pohl, David Sonnabend

“LiDDM: A Data Mining System for Linked Data” : Venkata Narasimha Pavan Kappara, Ryutaro Ichise , O.P. Vyas

Referenzen:

<http://www.highscore.de/cpp/einfuehrung/zeiger.html>

<http://ex-ample.blogspot.com/2011/07/tthe-weka-environment-weka-is.html>

<http://jena.sourceforge.net/>