Ontology Learning (with emphasis on refinement operator based learning)

Jens Lehmann, Johanna Völker

UNIVERSITÄT LEIPZIG

UNIVERSITÄT MANNHEIM

2010-09-02

Outline



1 Motivation and Definition

2 Overview of Ontology Learning Approaches

3 In Detail: Learning Definitions with Refinement Operators

Conclusions

Outline



2 Overview of Ontology Learning Approaches

3) In Detail: Learning Definitions with Refinement Operators

4 Conclusions

Definition: Ontology Learning

 "Ontology Learning is a subtask of information extraction. The goal of ontology learning is to (semi-)automatically extract relevant concepts and relations from a given corpus or other kinds of data sets to form an ontology." (Wikipedia, today)

Definition: Ontology Learning

- "Ontology Learning is a subtask of information extraction. The goal of ontology learning is to (semi-)automatically extract relevant concepts and relations from a given corpus or other kinds of data sets to form an ontology." (Wikipedia, today)
- "Ontology Learning is a mechanism for semi-automatically supporting the ontology engineer in engineering ontologies."
 A. D. Mädche. Ontology Learning for the Semantic Web. Dissertation. Universität Karlsruhe, 2001

Definition: Ontology Learning

- "Ontology Learning is a subtask of information extraction. The goal of ontology learning is to (semi-)automatically extract relevant concepts and relations from a given corpus or other kinds of data sets to form an ontology." (Wikipedia, today)
- "Ontology Learning is a mechanism for semi-automatically supporting the ontology engineer in engineering ontologies."
 A. D. Mädche. Ontology Learning for the Semantic Web. Dissertation. Universität Karlsruhe, 2001
- "Ontology Learning aims at the integration of a multitude of disciplines in order to facilitate the construction of ontologies, in particular ontology engineering and machine learning."
 A. D. Mädche, S. Staab. Ontology Learning. Handbook of Ontologies in Information Systems, 2004

Classification of Ontology Learning Data



sometimes heterogeneous sources of evidence (e.g., hyponymy [Snow et al. 2006], subsumption [Cimiano et al. 2005], [Manzano-Macho et al. 2008], [Buitelaar et al. 2008], disjointness [Völker et al. 2007])



Outline



2 Overview of Ontology Learning Approaches

3) In Detail: Learning Definitions with Refinement Operators

4 Conclusions

Ontology Learning Layer Cake [Cimiano 2006]



Patterns [Hearst 1992] for Class Subsumption

- NP such as {NP,}* {or|and} NP
 - "games such as baseball and cricket"
- NP {,NP}* {,} {and|or} other NP
 - "rabbits and other animals"
 - but: "rabbits and other pets"
- NP {,} including {NP,}* {or|and} NP
 - "fruits including apples and pears"
- NP {,} especially {NP,}* {or|and} NP
 - "Europeans, especially Italians"
 - but: "US presidents, especially democrats"



Patterns [Ogata and Collier 2004]

- NP is a NP
 - "A kangaroo is an animal living in Australia."
- a NP named called NP
 - "Japanese people like to play a game called Go."
- NP, NP
 - "Sencha, the most popular tea in Japan, ..."
- NP. The NP
 - "John loves his Ferrari. The car ..."
- Among NP, NP
 - "Among all musical instruments, violins are ..."
- NP except for other than NP
 - "Employees except for managers suffer from ..."



JAPE Rule

- GATE = General Architecture for Text Engineering
- written in Java
- mature, used worldwide
- JAPE = language for rapid prototyping and efficient implementation of shallow analysis methods
- can be used e.g. for domain specific patterns (financial blogs etc.)

JAPE Rule

```
rule: Hearst 1
(NounPhrase):superconcept
{SpaceToken.kind == space}
{Token.string=="such"}
{SpaceToken.kind == space}
{Token.string=="as"}
{SpaceToken.kind == space}
(NounPhrase):subconcept
):hearst1
-->
:hearst1.SubclassOfRelation = { rule = "Hearst1" },
:subconcept.Domain = { rule = "Hearst1" },
:superconcept.Range = { rule = "Hearst1" }
```

Lexical Context Similarity (e.g. [Cimiano and Völker 2005])

• "Columbus is the capital of the state of Ohio. Columbus has a population of about 700,000 inhabitants."

Lexical Context Similarity (e.g. [Cimiano and Völker 2005])

- "Columbus is the capital of the state of Ohio. Columbus has a population of about 700,000 inhabitants."
- **Columbus** (capital (1), state (1), Ohio (1), population (1), inhabitant (1))
- City (country (2), state (1), inhabitant (2), mayor (1), attraction (1))
- Explorer (ship (1), sailor (2), discovery (1))



"most probably": City(Columbus)

Subcategorization Frames

- "Tina drives a Ford."
 - Person(Tina). Vehicle(Ford).
- "Her father drives a bus."
 - Father subclass-of Person
 - Bus subclass-of Vehicle
- subcat: drive(subj: person, obj: vehicle)
 - Person $\sqsubseteq \forall$ drive.Vehicle



✔ [Faure and Nédellec 1998], [Schutz and Buitelaar 2005], [Cimiano et al. 2006]

Text2Onto Perspective - NeOn Toolkit - F:\NeOnTool ile Edit Navigate Search Project Run Window H	ikit\workspace Helo	Text	20nto	
		ТСЛС		
T Workflow View 🕴 🚺 🖓 🗢 🗢 🏹	POM View S			
B 🔁 Algorithm	Concent Instance Similarity Su	Concent Instance (Similarity (Subclass)) Instance() Relation Disjoint		
😑 😂 Concept	Contract and and and and and	0.000	Canfidanaa	
Markflaur		individual	1.0	
WORKTIOW		content	1.0	
		communication	1.0	
Ontology Learning Methods		content	1.0	
oncology Leanning m	centous	content	1.0	
	knowledge base	content	1.0	
 VerticalRelationsConceptClassification 	🗹 designer	individual	1.0	
WordNetConceptClassification	discussion	communication	1.0	
e 😂 InstanceOf	personal	communication	1.0	
PatternInstanceClassification	✓ task	work	1.0	
😑 🗁 Relation	interoperability	quality	1.0	1
SubcatRelationExtraction	browsing	process	1.0	
Bisjont S	ubclassOf(So subclassOf(ftware_Agent, Software_Age	Computer_Bent, Technol	Program)(0.5 Logy)(0.5)
Corrous View 33	ubclassOf(So subclassOf(ftware_Agent, Software_Age	Computer_F ent, Technol	Program)(0.5 Logy)(0.5)
Disjont S	ubclassOf (So subclassOf (B) M report	ftware_Agent, Software_Age	Computer_F ent, Technol	Program)(0.5 Logy)(0.5)
Corpus View 12	ubclassOf(So subclassOf(subclassOf(software agent	ftware_Agent, Software_Age	Computer_F ent, Technol 0.5714285714285714 0.5 0.5	Program)(0.5 Logy)(0.5)
Disjont Corpus View 2 Corpus Corpus	ubclassOf (So subclassOf (software agent software agent software agent	ftware_Agent, Software_Age communication computer program technology method	Computer_F ent, Technol 0.5714285714 0.5 0.5 0.5	Program)(0.5 Logy)(0.5)
Disjont S Corpus View C Grócopus/corpus_sw1224967.bit Grócopus/corpus_sw1224967.bit Grócopus/corpus_sw12224967.bit	ubclassOf(So subclassOf(control of the second control of the sec	ftware_Agent, Software_Age communication communication communication communication	Computer_F ent, Technol 0.5 0.5 0.5 0.5 0.5 0.5	Program)(0.5 Logy)(0.5)
Depart S Corpus View 2 Gorpus/corpus_sw1228567.bt Gr(Corpus/corpus_w122850.bt Gr(Corpus/corpus_w128850.bt Gr(Corpus/corpu	ubclassof (So subclassof (So encort	ftware_Agent, Software_Age computerporam technology method communication impairage	Computer_F ent, Technol 05 05 05 05 05 05 05	Program)(0.5 Logy)(0.5)
Disjont S Corpus View C Grócopus/corpus_sw1229507.bt Grócopus/corpus_sw1229507.bt Grócopus/corpus_sw122250.bt Grócopus/corpus_sw122101.bt Grócopus/corpus_sw122101.bt Grócopus/corpus_sw122101.bt	ubclassOf(So subclassOf(Software agent Software a	ftware_Agent, Software_Age communication computer program technology method communication language information	Computer_H ent, Technol 05 05 05 05 05 05 05	Program)(0.5 Logy)(0.5)
Depart Depart Corpus View 22 Gorpus View 24 Gorpu	ubclassOf (So subclassOf (So subclassOf (9 software agent 9 software ag	ftware_Agent, Software_Age computer program technology metod communication terguage unguage unguage unguage unguage	Computer_B ent, Technol 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0 Logy)(0.5)
Disjont S Corpus View 2 Gropps Gropps GriCorpus/corpus_ph/12:950.ht GriCorpus/corpus_ph/32:1945.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1941.ht GriCorpus/corpus_ph/32:1942.ht GriCorpus/corpus_ph/32:1	ubclassof (So subclassof (So source agent Software agent Bargange Bargang Bargange Bargang Bargange Bargange Ba	ftware_Agent, Software_Age computer program webrology method communication targuage buoguage buoguage torowidge torowidge	Computer_F ent, Technol 05 05 05 05 05 05 05 05	Program)(0. Logy)(0.5)
Deport Deport Corpus View 32 GOrpus Viery 23	ubclassof(So subclassof(Softwareapert Softwareapert Usprage discussion Drowing Drowi	ftware_Agent, <u>Software_Age</u> <u>commutation</u> <u>commutation</u> <u>commutation</u> <u>isopage</u> information <u>isopage</u> information <u>isopage</u> information	Computer_F ent, Technol 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0.9
Disport Corpus View 22 Corpus View 22 Gr(corpus/corpus_psv1224507.bt) Gr(corpus/corpus_psv1224507.bt) Gr(corpus/corpus_psv1222520.bt) Gr(corpus/corpus_psv1221521.bt) Gr(corpus/corpus_psv1221521.bt) Gr(corpus/corpus_psv1221521.bt) Gr(corpus/corpus_psv1221521.bt) Gr(corpus/corpus_psv12645827.bt) Gr(corpus/corpus_psv12645827.bt) Gr(corpus/corpus_psv12645827.bt) Gr(corpus/corpus_psv12645827.bt) Gr(corpus/corpus_psv12645827.bt)	ubclassof (So subclassof (So subclassof (software agent software agent decusion browing browi	ftware_Agent, Software_Age computer program technology method communication larguage	Computer_E ent, Technol 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0.!
	ubclassof(So subclassof(So control Software agent Software agent U torruga U torruga	ftware_Agent, <u>Software_Agent</u> , <u>outputprogram</u> <u>outputprogram</u> <u>outputprogram</u> <u>outputprogram</u> <u>outputprogram</u> <u>outputprogram</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u> <u>inconside</u>	Computer_H ent, Technol 05/400/400/40 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0.9
Depart S Depart S Corpus Vew 22 Grops Gr(corpus/corpus_pw/32/350.bt Gr(corpus/corpus_pw/32/252.bt Gr(corp	ubclassof (So subclassof (So subclassof (software agent software agent bernage decusion browing bernage technology technology technology technology technology technology technology technology technology technology technology technology technology	ftware_Agent, Software_Age computer program technology method communication larguage larguage trowwidge trowwidge trowwidge trowwidge trowwidge trowwidge trowwidge trowwidge	Computer_H ent, Technol 05742871438714 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0.: Logy)(0.5)
	ubclassof(So subclassof(So software sport eschare sport eschare sport eschare sport eschare sport eschare sport eschare sport eschare sport eschare sport eschare eschare sport eschare esch	ftware_Agent, <u>Software_Agent</u> , <u>computer poyam</u> <u>computer poyam</u> <u>reformation</u> language insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge insolvidge	Computer_H ent, Technol 05 05 05 05 05 05 05 05 05 05 05 05 05	Program) (0.!
Disport Software Corpus View 12 Corpus View 12 GriCorpus/corpus_pw/12/34567.ht GriCorpus/corpus_pw/32/252.0.ht GriCorpus/corpus_pw/32/252.0.ht GriCorpus/corpus_pw/32/252.1.ht GriCorpus/corpus_pw/34/141.1.ht	ubclassof(So subclassof(So of software agent of software agent of software agent of browing the thorology of making of making	ftware_Agent, Software_Age computer program technology methodogy information torowidge tr	Computer_E ent, Technol 05 05 05 05 05 05 05 05 05 05 05 05 05	Program)(0.: Logy)(0.5)
	ubclassof(So subclassof(So concernent Source Sour	ftware_Agent, <u>Software_Agent</u> , <u>computer program</u> <u>computer program</u> <u>computer program</u> <u>inguage</u> inguage information inoviedge inovi	Computer_H ent, Technol 05/000/000/10 05 05 05 05 05 05 05 05 05 05 05 05 05	Program) (0.: Logy) (0.5)
Corpus View 12 Corpus Corpus	ubclassof(So subclassof(or recr or recr or software apert or software apert or software apert or software apert or software or recr or recr or recr or software or software or recr or recr or recr or recr or recr or recr or recr or software or recr or r	ftware_Agent, <u>Software_Age</u> <u>commission</u> <u>commission</u> <u>commission</u> <u>isopage</u> informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage informato isopage isop	Computer_H ent, Technol 03/1400/1400/14 05 05 05 05 05 05 05 05 05 05 05 05 05	Program) (0 Logy) (0.5)

Learning from text and background knowledge via reasoning: "Washington is the capital of the US. (...) New York is the US capital of fashion."

Learning from text and background knowledge via reasoning: "Washington is the capital of the US. (...) New York is the US capital of fashion."

- extracted: hasCapital(US, New York); hasCapital(US, Washington)
- background knowledge: hasCapital is a functional property

Learning from text and background knowledge via reasoning: "Washington is the capital of the US. (...) New York is the US capital of fashion."

- extracted: hasCapital(US, New York); hasCapital(US, Washington)
- background knowledge: hasCapital is a functional property
- possible inferences:
 - $\bullet \ \, {\sf New York} = {\sf Washington}$
 - inconsistency (unique names assumption)
- logical contradictions can help to detect errors in automatically extracted information



Other Approaches

- Association rules and co-occurrence statistics
- WordNet: hyponymy \approx subsumption
 - hyponym(bank#1, institution#1)
 - Bank subclass-of Institution
- Noun phrase heuristics
 - "image processing software"
- Instance clustering (e.g. Columbus and Washington)
 - Hierarchical clustering of context vectors
- Knowledge Base Completion / Formal Concept Analysis (FCA)
 - asks knowledge engineer questions to complete a knowledge base
 - tool: OntoComp [Sertkaya et al.]

Tools and Frameworks

Name	Institute	Authors
ASIUM	INRIA, Jouy-en-Josas	Faure and Nedellec 1999
TextToOnto	AIFB, University of Karlsruhe	Mädche and Volz 2001
HASTI	Amir Kabir University, Teheran	Shamsfard, Barforoush 2004
OntoLT	DFKI, Saarbrücken	Buitelaar et al. 2004
DOODLE	Shizuoka University	Morita et al. 2004
Text2Onto	AIFB, University of Karlsruhe	Cimiano and Völker 2005
OntoLearn	University of Rome	Velardi et al. 2005
OLE	Brno University of Technology	Novacek and Smrz 2005
OntoGen	Institute Jozef Stefan, Ljubljana	Fortuna et al., 2007
GALeOn	Technical University of Madrid	Manzano-Macho et al. 2008
DINO	DERI, Galway	Novacek et al. 2008
OntoLancs	Lancester University	Gacitua et al. 2008

Table: Lexical ontology learning: informal or semi-formal data (e.g. texts)

Tools and Frameworks

Name	Institute	Authors
YINGYANG	University of Bari	lannone 2006
DL-Learner	University of Leipzig	Lehmann 2006
RELExO	AIFB, University of Karlsruhe	Völker and Rudolph 2008
RoLExO	AIFB, University of Karlsruhe	Völker and Rudolph 2008
OntoComp	University of Dresden	Sertkaya 2008

Table: Logical Ontology Learning

Name	Institute	Authors
LeDA	AIFB, University of Karlsruhe	Völker et al. 2007
SOFIE	MPI, Saarbrücken	Suchanek et al. 2009

Table: Hybrid implementations

Problems and Challenges

- Homonymy and polysemy e.g. [Ovchinnikova et al. 2006]
 - "Peter is sitting on the **bank** in front of the **bank**."
 - "An interesting **book** is lying on the table."
- Semantics of adjectives
 - "red flower", "false friend"
- Empty heads e.g. [Völker et al. 2005], [Cimiano and Wenderoth 2005]
 - "Tuna is a **kind** of fish. The Southern Bluefin is one of the most endangered **types** of Tuna."
- Ellipsis and underspecification
 - "Mary started the book."
- Anaphora (e.g. pronouns) e.g. [Cimiano and Völker 2005]
 - "There is an apple on the table. It is red."

Problems and Challenges (ctd.)

- Metaphors and analogies e.g. [Gust et al. 2007]
 - "Live is a journey."
- Opinions, quotations and reported speech
 - "Tom thinks that dolphins are mammals."
- What should be represented as an individual? e.g. [Zirn et al. 2008]
 - "The kangaroo is an animal living in Australia."
- Class, relation (object property) or attribute (datatype property)?
 - "All elephants are grey."
 - "Easter monday is a national holiday."
- Knowledge is changing e.g. [Stojanovic 2004], [Zablith et al. 2009]
 - "Pluto is a planet."



Outline



Overview of Ontology Learning Approaches

③ In Detail: Learning Definitions with Refinement Operators

4 Conclusions

Learning OWL Class Expressions

- given:
 - background knowledge (particularly OWL/DL knowledge base)
 - positive and negative examples (particulary individuals in knowledge base)
- goal:
 - logical formula (particularly OWL Class Expression) covering positive examples and not covering negative examples



ILP and Semantic Web



- since early 90s Inductive Logic Programming
- only few approaches based on description logics
- Web Ontology Language (OWL) becomes W3C standard in 2004
- increasing number of RDF/OWL knowlegde bases, but ILP still mainly focuses on logic programs → research gap

ONTOLOGY LEARNING

Why ILP in the Semantic Web?

• Ontology Learning:

- $\bullet\,$ given class A in ${\cal K}$
- instances of A as positive examples
- non-instances as negative examples
- definitions can be learned if ABox data is available
- improvement of existing ML problem solutions
- direct usage of knowledge in the Semantic Web instead of conversion in e.g. horn clauses to apply ML methods



ontology network



ML problems

Refinement Operators - Definitions

- given a DL \mathcal{L} , consider the quasi-ordered space $\langle \mathcal{C}(\mathcal{L}), \sqsubseteq_{\mathcal{T}} \rangle$ over concepts of \mathcal{L}
- $\rho : C(\mathcal{L}) \to 2^{C(\mathcal{L})}$ is a downward \mathcal{L} refinement operator if for any $C \in C(\mathcal{L})$:

$$D \in \rho(C)$$
 implies $D \sqsubseteq_{\mathcal{T}} C$

- notation: Write $C \rightsquigarrow_{
 ho} D$ instead of $D \in
 ho(C)$
- example refinement chain in $\langle C(\mathcal{EL}), \sqsubseteq_{\mathcal{T}} \rangle$:

$$op \rightsquigarrow_{
ho}$$
 Male $\rightsquigarrow_{
ho}$ Male $\sqcap \exists extsf{hasChild}. op$



 start with most general concept (top)



- start with most general concept (top)
- operator specialises concept



- start with most general concept (top)
- operator specialises concept
- heuristic assigns score using pos/neg examples



- start with most general concept (top)
- operator specialises concept
- heuristic assigns score using pos/neg examples
- continue until termination criterion is met
Learning with Refinement Operators



- start with most general concept (top)
- operator specialises concept
- heuristic assigns score using pos/neg examples
- continue until termination criterion is met

learning algorithm

An ${\mathcal L}$ downward refinement operator ρ is called

• finite iff $\rho(C)$ is finite for any concept $C \in \mathcal{C}(\mathcal{L})$



- finite iff $\rho(C)$ is finite for any concept $C \in \mathcal{C}(\mathcal{L})$
- redundant iff there exist two different ρ refinement chains from a concept C to a concept D.



- finite iff $\rho(C)$ is finite for any concept $C \in \mathcal{C}(\mathcal{L})$
- redundant iff there exist two different ρ refinement chains from a concept C to a concept D.
- proper iff for $C, D \in \mathcal{C}(\mathcal{L})$, $C \rightsquigarrow_{\rho} D$ implies $C \not\equiv_{\mathcal{T}} D$



- finite iff $\rho(C)$ is finite for any concept $C \in \mathcal{C}(\mathcal{L})$
- redundant iff there exist two different ρ refinement chains from a concept C to a concept D.
- proper iff for $C, D \in \mathcal{C}(\mathcal{L})$, $C \rightsquigarrow_{\rho} D$ implies $C \not\equiv_{\mathcal{T}} D$
- complete iff for $C, D \in C(\mathcal{L})$ with $D \sqsubset_{\mathcal{T}} C$ there is a concept E with $E \equiv_{\mathcal{T}} D$ and a refinement chain $C \rightsquigarrow_{\rho} \cdots \rightsquigarrow_{\rho} E$
- weakly complete iff for any concept C with C □_T ⊤ we can reach a concept E with E ≡_T C from ⊤ by ρ.



- finite iff ho(C) is finite for any concept $C \in \mathcal{C}(\mathcal{L})$
- redundant iff there exist two different ρ refinement chains from a concept C to a concept D.
- proper iff for $C, D \in \mathcal{C}(\mathcal{L})$, $C \rightsquigarrow_{\rho} D$ implies $C \not\equiv_{\mathcal{T}} D$
- complete iff for $C, D \in C(\mathcal{L})$ with $D \sqsubset_{\mathcal{T}} C$ there is a concept E with $E \equiv_{\mathcal{T}} D$ and a refinement chain $C \rightsquigarrow_{\rho} \cdots \rightsquigarrow_{\rho} E$
- weakly complete iff for any concept C with C □_T ⊤ we can reach a concept E with E ≡_T C from ⊤ by ρ.



- Properties indicate how suitable a refinement operator is for solving the learning problem:
 - Incomplete operators may miss solutions
 - Redundant operators may lead to duplicate concepts in the search tree
 - Improper operators may produce equivalent concepts (which cover the same examples)
 - For infinite operators it may not be possible to compute all refinements of a given concept
- We researched properties of refinement operators in Description Logics
- Key question: Which properties can be combined?

Refinement Operator Property Theorem

Theorem

 $\begin{aligned} \text{Maximal sets of properties of } \mathcal{L} \text{ refinement operators which can be} \\ \text{combined for } \mathcal{L} \in \{\mathcal{ALC}, \mathcal{ALCN}, \mathcal{SHOIN}, \mathcal{SROIQ}\}: \end{aligned}$

- { weakly complete, complete, finite}
- { weakly complete, complete, proper }
- § {weakly complete, non-redundant, finite}
- { weakly complete, non-redundant, proper}
- {non-redundant, finite, proper}

"Foundations of Refinement Operators for Description Logics",
 J. Lehmann, P. Hitzler, ILP conference, 2008
 "Concept Learning in Description Logics Using Refinement Operators",

J. Lehmann, P. Hitzler, Machine Learning journal, 2010

LEHMANN, VÖLKER (LEIPZIG+MANNHEIM)

ONTOLOGY LEARNING

Refinement Operator Property Theorem

- no ideal refinement in OWL and many description logics
- indicates that learning in DLs is hard
- algorithms need to counteract disadvantages
- goal: develop operators close to theoretical limits

$$\rho(C) = \begin{cases} \{\bot\} \cup \rho_{\top}(C) & \text{if } C = \top \\ \rho_{\top}(C) & \text{otherwise} \end{cases}$$

$$\begin{cases} \emptyset & \text{if } C = \bot \\ \{C_1 \sqcup \cdots \sqcup C_n \mid C_i \in M_B \ (1 \le i \le n)\} & \text{if } C = \top \\ \{A' \mid A' \in sh_{\downarrow}(A)\} & \text{if } C = \Lambda \ (A \in N_C) \\ \cup \{A \sqcap D \mid D \in \rho_B(\top)\} & \text{if } C = \neg A \ (A \in N_C) \\ \cup \{\neg A \sqcap D \mid D \in \rho_B(\top)\} & \text{if } C = \neg A \ (A \in N_C) \\ \cup \{\neg A \sqcap D \mid D \in \rho_B(\top)\} & \text{if } C = \exists r.D \\ \cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\} & \text{if } C = \exists r.D \\ \cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\} & \text{if } C = \forall r.D \\ \cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\} & \text{if } C = \forall r.D \\ \cup \{\forall r.L \mid A = ar(r), E \in \rho_A(D)\} & \text{if } C = \forall r.D \\ \cup \{\forall r.L \mid B = A \in N_C \ and \ sh_{\downarrow}(A) = \emptyset\} & \text{if } C = C_1 \sqcap \cdots \sqcap C_n \\ D \in \rho_B(C_i), 1 \le i \le n\} & (n \ge 2) \\ \{C_1 \sqcup \cdots \sqcup C_{i-1} \sqcup D \sqcup C_{i+1} \sqcup \cdots \sqcup C_n \mid \text{if } C = C_1 \sqcup \cdots \sqcup C_n \\ D \in \rho_B(C_i), 1 \le i \le n\} & (n \ge 2) \\ \cup \{(C_1 \sqcup \cdots \sqcup C_n) \sqcap D \mid B \in P_B(\top)\} & \text{if } C = C_1 \sqcup \cdots \sqcup C_n \\ D \in \rho_B(\top)\} & \text{if } C = C_1 \sqcup \cdots \sqcup C_n \end{cases}$$

Base Operator (excerpt)



Base Operator (excerpt)

 $\{\exists r.E \mid A = ar(r), E \in \rho_A(D)\} \quad \text{if } C = \exists r.D$ $\cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\}$ $\cup \{\exists s.D \mid s \in sh_{\downarrow}(r)\}$

Examples:

 $\exists \texttt{takesPartIn.SocialGathering} \rightsquigarrow$

∃takesPartIn.Meeting

$$\{\exists r.E \mid A = ar(r), E \in \rho_A(D)\} \quad \text{if } C = \exists r.D$$
$$\cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\}$$
$$\cup \{\exists s.D \mid s \in sh_{\downarrow}(r)\}$$

Examples:

$$\{\exists r.E \mid A = ar(r), E \in \rho_A(D)\} \quad \text{if } C = \exists r.D$$
$$\cup \{\exists r.D \sqcap E \mid E \in \rho_B(\top)\}$$
$$\cup \{\exists s.D \mid s \in sh_{\downarrow}(r)\}$$

Examples:

- ρ_{\downarrow} is complete
- ρ_{\downarrow} is infinite, e.g. there are infinitely many refinement steps of the form:

$$\top \rightsquigarrow_{\rho_{\downarrow}} C_1 \sqcup C_2 \sqcup C_3 \sqcup \ldots$$

- ρ_{\downarrow} not proper, but can be extended to a proper operator ρ_{\downarrow}^{cl} (refinements more expensive to compute)
- ρ_{\downarrow} is redundant: $\forall r_1.A_1 \sqcup \forall r_2.A_1 \rightsquigarrow_{\rho_{\downarrow}} \forall r_1.(A_1 \sqcap A_2) \sqcup \forall r_2.A_1$ $\downarrow^{\downarrow}_{\gtrsim}$ $\forall r_1.A_1 \sqcup \forall r_2.(A_1 \sqcap A_2) \rightsquigarrow_{\rho_{\downarrow}} \forall r_1.(A_1 \sqcap A_2) \sqcup \forall r_2.(A_1 \sqcap A_2)$

 "A Refinement Operator Based Learning Algorithm for the ALC Description Logic", J. Lehmann, P. Hitzler, ILP conference, 2008
 "Concept Learning in Description Logics Using Refinement Operators", J. Lehmann, P. Hitzler, Machine Learning journal, 2010

LEHMANN, VÖLKER (LEIPZIG+MANNHEIM)

ONTOLOGY LEARNING

OCEL

- uses ρ for top down search
- OCEL is complete it always find a solution if one exists
- highly configurable, e.g. flexible target language, termination criteria and heuristics
- implements redundancy elimination technique with polynomial complexity wrt. search tree size based on ordered negation normal form
- can handle infinite refinement operators by stepwise length-limited horizontal expansion





- length of child concepts limited by horizontal expansion (he)
- ρ (infinite) is applicable



- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)



- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)



- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)



- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)



- length of child concepts limited by horizontal expansion (he)
- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)



- length of child concepts limited by horizontal expansion (he)
- ρ (infinite) is applicable
- he influences heuristic (Bias towards short concepts - Occam's Razor, higher diversity)

Scalability: Reasoning

```
\mathcal{K} = \{ \texttt{Male} \sqsubseteq \texttt{Person}, \\ \texttt{OnlyMaleChildren}(a), \\ \texttt{Person}(a), \texttt{Male}(a_1), \texttt{Male}(a_2), \\ \texttt{hasChild}(a, a_1), \texttt{hasChild}(a, a_2) \}
```

- given \mathcal{K} , we want to learn a description of OnlyMaleChildren
- C = Person □ ∀hasChild.Male appears to be a good solution, but a is not an instance of C under OWA
- idea: dematerialise \mathcal{K} using standard (OWA) DL reasoner, but perform instance checks using CWA
- closer to intuition and provides order of magnitude performance improvements
- optimised for thousands of instance checks on a static knowledge base

Scalability: Stochastic Coverage Computation

Heuristics often require expensive instance checks or retrieval, e.g.:

$$\frac{1}{2} \cdot \left(\frac{|R(A) \cap R(C)|}{|R(A)|} + \sqrt{\frac{|R(A) \cap R(C)|}{|R(C)|}} \right)$$

Scalability: Stochastic Coverage Computation

Heuristics often require expensive instance checks or retrieval, e.g.:

$$\frac{1}{2} \cdot \left(\frac{a}{|R(A)|} + \sqrt{\frac{a}{b}} \right)$$

- replace $|R(A) \cap R(C)|$ und |R(C)| by variables *a* and *b* we want to estimate
- Wald-Method for computing the 95% confidence interval
- first estimate *a*, then the whole expressions
- method can be applied to various heuristics

Scalability: Stochastic Coverage Computation

Heuristics often require expensive instance checks or retrieval, e.g.:

$$\frac{1}{2} \cdot \left(\frac{a}{|R(A)|} + \sqrt{\frac{a}{b}} \right)$$

- replace $|R(A) \cap R(C)|$ und |R(C)| by variables *a* and *b* we want to estimate
- Wald-Method for computing the 95% confidence interval
- first estimate a, then the whole expressions
- method can be applied to various heuristics
- in tests on real ontologies up to 99% less instance checks and algorithm up to 30 times faster
- low influence on learning results empirically shown in 380 learning problems on 7 real ontologies (differs by ca. $0, 2\% \pm 0, 4\%$)

Scalability: Fragment Extraction



"Learning of OWL Class Descriptions on Very Large Knowledge Bases", Hellmann, Lehmann, Auer, Int. Journal Semantic Web Inf. Syst, 2009

- lack of evaluation standards in OWL/DL learning
- procedure: convert existing benchmarks to OWL (time consuming, requires domain knowledge)
- measure predictive accuracy in ten fold cross validation
- part 1: evaluation against other OWL/DL learning systems
- part 2: evaluation against other ML systems (carcinogenesis problem)
- part 3: evaluation of ontology enginering

Evaluation: Accuracy



cross validation accuracy in % (Durchschnitt über 6 Benchmarks)

- Collection of 6 Benchmarks
- OCEL often stat. significantly better than other algorithms for most benchmarks

Evaluation: Readability



• YinYang generates significantly longer solutions

Evaluation: Runtime



Carcinogenesis

- goal: predict whether chemical compounds cause cancer
- Why?
 - more than 1000 new substances each year
 - substances can often only be tested via long and expensive experiments on rats and mice
- background knowledge:
 - database of US National Toxicology Program (NTP)
 - converted from Prolog to OWL



"Obtaining accurate structural alerts for the causes of chemical cancers is a problem of great scientific and humanitarian value." (A. Srinivasan, R.D. King, S.H. Muggleton, M.J.E. Sternberg 1997)

Carcinogenesis



- very challenging problem: low accuracy, high standard deviation
- OCEL stat. sign. better than most other approaches

Ontology Learning Evaluation

- 5 PhD studens
- 5 real ontologies in different domains
- 998 decision of each test person for 92 classes
- in 35% of the cases accepted suggestions for ontology enhancements
- problem: ontology quality, modelling errors (unsatisfiable classes, disjunction and conjunction confused etc.)
DL-Learner Project

- DL-Learner Open-Source-Projekt: http://dl-learner.org, http://sf.net/projects/dl-learner
- extensible platform for different learning problems and algorithms
- Interfaces: command line, GUI, Web-Service
- supports common OWL formats
- allows different reasoners (via OWL API, DIG, OWLLink)
- sourceforge.net (Open Source Portal): 4000 Downloads
- mloss.org (ML & Open Source Software): 1600 Downloads



LEHMANN, VÖLKER (LEIPZIG+MANNHEIM)

ONTOLOGY LEARNING

2010-09-02 60 / 63

- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks

- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin



- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin
 - OntoWiki Plugin

OntoWiki -	DL Learner - Learnt Class Expressions				
User Extras Help Debug Search for Resources	Add Class Engression				
Knowledge Bases	Suggested Equivalence Classes for customer requirement				
SoftWiki Ontology for Requirements Engineering	accuracy toge at suggested class expressions				
Seebis EOAE Profile	O 100% togge details is created by some customer				
00000100010000	O 100% toggle details requirement and is created by some customer				
Classes	O 90% togge datable requirement and is created by some (customer or government)				
string image system requirement > abstract source	Boys Boys Boys Covered Instances (DO%): Use A Lift System Resources As People Life Obtains: Life Obtains Life Obtains Life Obtains Create Anderin COLDuragin Lise A Lift System Resources As People Life Obtains Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College Create Resources Anderin College College				
abstract requirement scenario vrequirement quality requirement	Additional Instances: Build A Secure Login System Build Login System Build Network Login System Create Database storface Create Network Interface Use Of toons Cadiculations Technical Database (Usability)				
performance requirement	O 80% toggle details requirement and is created by only customer				
functional requirement	O 75% toggle datable requirement and is created by some (author or customer)				
design requirement	O 73% toggle details requirer List instances				
derived requirement	Create Instance				
customer requirement	The suggestions were generated by Collete Resource the				
allocated requirement	DL-Learner plugin, see the OntoWiki				
goal	Learn Equivalent Class Expression				
► reference point	Learn Duper Class expression				
► abstract comment					

- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin
 - OntoWiki Plugin
 - ORE



- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin
 - OntoWiki Plugin
 - ORE
- Recommendation/Navigatior
 - moosique.net

All I metal Sear	ch	Search Playlist	Recommendations	More Info Help	
Currently playing:					
Rouler Pinder — Formidable				< → →	
Plavlist					
1 haymse					
You can delete entries from the playlist by clicking the small x on the right and change their order by clicking on the small up- and down-arrows.					
1. O Rouler Pinder – Formidable				000	
2. O Rouler Pinder – Laska			000		
3. O Rouler Hinder – Les pates (intro)			000		
4. O Rouler Pinder — Les pates			000		
5. © Rouler Pinder – Bratou			000		
7. O Rouler Pinder – Pacifiste			000		
8. O Rouler Pinder — J'ai baisé (intro)			000		
9. 🕞 Rouler Pinder — J'ai baisé				000	
10. 💿 Rouler Pinder — Chilé				000	
 11. Rouler Pinder — Formidable Tror 	npette			000	
Delete all					
Recently Listened to					

- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin
 - OntoWiki Plugin
 - ORE
- Recommendation/Navigation
 - moosique.net
 - DBpedia Navigator



Navigation Suggestions

- "classical" ML problems
 - carcinogenesis
 - other biomedical tasks
- Ontology Learning
 - Protégé Plugin
 - OntoWiki Plugin
 - ORE
- Recommendation/Navigation
 - moosique.net
 - DBpedia Navigator
- other/external:
 - ISS (Gerken et al.)
 - Learning in Probabilistic DLs (Ochoa Luna et al.)
 - TIGER Corpus Navigator (Hellmann et al.)

(III		
Lehmann, Völker	(Leipzig+Mannheim)	ONTOLOGY LEARNING

TIGER Corpus Navigator see here for data license reset						
Search	Learned Concept (Sentence and has Token some VVPP and has Token some (ADV and next Token some VAFIN)) Accuracy:1.					
Fulltext Search search Lemma Search werden search						
Search Results	Matching Comment These are finite auxiliary verbs. ("(du) bist", "(wir) werden".					
show	Learning Input					
Classified Instances (displaying result 1-47 of 5299)						
hide ** Ich olaube kaum . daß mit seinem . nala . etwas	Nur 1,4 Milliarden Mark selen gestrichen worden . x					
undplomatischen Still im Weißen Haus dem Land ein + - Gefallen getan wäre . ** Es ist wirklich schwerz u sagen , weiche Positionen er einnimmt , da er sich noch nicht konkret geäußert + -	Zuletzt war am 25. Dezember das Wohnmobil des Göttinger Oberstadtdirektors Hermann Schierwater in Brand gesteckt worden . Negative Samples					
nat *, beklagen volkswirte . Sie haben den längst überfälligen Bruch mit der	Er wird zum direkten Partner , der umweglos alle und x					

Outline



Overview of Ontology Learning Approaches

3) In Detail: Learning Definitions with Refinement Operators

4 Conclusions

Conclusions

- Ontology Learning is a diverse research area involving several research disciplines (NLP, Machine Learning, Ontology Engineering)
- approaches vary in used data sources and the expressiveness of the created ontologies
- refinement operator based learning as one method for learning definitions (with applications outside of learning ontologies)
- new Wiki (under construction): http://ontology-learning.net

• new ontology learning book in 2011

