

How and Why Computers Read the Web

Estevam R. Hruschka Jr.

Federal University of São Carlos and Carnegie Mellon University



Federal University of São Carlos

Carnegie Mellon University







Disclaimers

- Previous versions of this tutorial were presented at:
 - IBERAMIA2012 <u>http://iberamia2012.dsic.upv.es/tutorials/;</u>
 - WWW2013 <u>http://www2013.org/program/machine-reading-the-web/</u>
 - ECMLPKDD2015 <u>http://www.ecmlpkdd2015.org</u>
 - AAAI2016
 <u>http://www.aaai.org/Conferences/AAAI/2016/aaai16tutorials.php#SA</u>
 <u>4</u>
- Feel free to e-mail me (<u>estevam.hruschka@gmail.com</u>) with questions about this tutorial or any feedback/suggestions/criticisms. Your feedback can help improving the quality of these slides, thus, they are very welcome.







Disclaimers

- Due to time constraints, I do not intend to cover all the algorithms and publications related to DBPedia, YAGO, KnowItAll and NELL. What I do intend, instead, is to give an overview of all four projects and what is the main approach used in each project.
- DBPedia, YAGO, KnowItAll and NELL are not the only research efforts focusing on "Reading the Web". They were selected, to be presented in this tutorial, because they represent four different and very relevant approaches to this problem, but it does not mean they are the best (or the only relevant) ones.







- Machine Learning
- Machine Reading
- Reading the Web
 - YAGO
 - KnowltAll
 - NELL
 - DBPedia







Machine Learning

- Machine Reading
- Reading the Web
 - YAGO
 - KnowltAll
 - NELL
 - DBPedia







- Machine Learning
- Machine Reading
- Reading the Web
 - YAGO
 - KnowltAll
 - NELL
 - DBPedia







- Machine Learning
- Machine Reading
- Reading the Web
 - YAGO
 - KnowltAll
 - NELL
 - DBPedia



Federal University of São Carlos U-IST // fect 24 TC: an of two Machine Reading Bee L01Br Ca Ca Ch Auto-Text to Knowledge

Picture taken from [DARPA, 2012]

FAPESP

Carnegie Mellon University





Machine Reading



Auto-Text to Knowledge

Picture taken from [DARPA, 2012]



Carnegie Mellon University



Answer Wiki

No. Answers below note:

- Professional actors and directors need reference works to help them understand Shakespeare.
- It's easier to follow a Shakespeare play by watching it rather than reading it.
- Vocabulary, idiom and culture have changed since the 17th century. Seventeeth century British English is like a foreign language to many 21st century English speakers.
- Many people who say they understand Shakespeare are faking it because they don't want to appear dumb.
- Being able to read and understand Shakespeare is a particular kind of intelligence that you can develop.
- Reading the synopsis of a Shakespeare play can help you understand it better when you read the full play.









Reading Comprehension Understanding and Comprehension Psychology +1 🖍

Why don't I understand what I read?

I don't have dyslexia, and my English teacher says I am the best writer in my class (I am in the top stream). Yet still when I read something, I don't take it in.

I normally have to read a paragraph and then go back and read it again. One reading to hear it in my head, and the other to get the meaning. Why can't I simultaneously read and comprehend? Why must every sentence be read and then understood?









- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowltAll
 - NELL







- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







- Machine Learning
- Machine Reading
- Reading the Web
 - -DBPedia
 - YAGO
 - KnowltAll
 - NELL





ABOUT

DBpedia is a crowd-sourced community effort to extract structured information from Wikipedia and make this information available on the Web. DBpedia allows you to ask sophisticated queries against Wikipedia, and to link the different data sets on the Web to Wikipedia data. We hope that this work will make it easier for the huge amount of information in Wikipedia to be used in some new interesting ways. Furthermore, it might inspire new mechanisms for navigating, linking, and improving the encyclopedia itself. Further information

...





Carnegie Mellon University





- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - -YAGO
 - KnowltAll
 - NELL





The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web

Deutsch

Homepage

Departments:

Systems People

Research

About the Institute



Algorithms & Complexity Computer Vision and

Multimodal Computing

Applied Algorithmics Computer Graphics

Computational Biology &

Databases and Information

YAGO-NAGA

ANGIE HYENA JAVATOOLS

AIDA

AIDA

K2 NAGA

DATTV

The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web

Research Publications People

The YAGO-NAGA project started in 2006 with the goal of building a conveniently searchable, large-scale, highly accurate knowledge base of common facts in a machine-processible representation.

We have already harvested knowledge about millions of entities and facts about their relationships, from Wikipedia and WordNet with careful integration of these two sources. The resulting knowledge base, coined YAGO, has very high precision and is freely available. The facts are represented as RDF triples, and we have developed methods and prototype systems for querying, ranking, and exploring knowledge. Our search engine NAGA provides ranked answers to queries based on statistical models.

Several interlinked sub-projects are growing on the YAGO-NAGA basis. Our vision is a confluence of Semantic Web (Ontologies), Social Web (Web 2.0), and Statistical Web (Information Extraction) assets towards a comprehensive repository of human knowledge. Our methodologies combine concepts, models, and algorithms from several fields, including database systems, information retrieval, statistical learning, and logical reasoning.

max planck institut

informatik

AIDA is a method, implemented in an online tool, for disambiguating mentions of named entities that occur in natural-language text or Web tables.



Carnegie Mellon University





- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - -KnowItAll
 - NELL





W

KnowItAll



UNIVERSITY OF WASHINGTON | ABOUT US | CONTACT US | MY CSE | INTERNAL

Computer Science & Engineering





UNIVERSITY of WASHINGTON

News & Events	People	Education	Research	Current Students	Prospective Students	Faculty Candidates	Alumni	Industry Affiliates	Support CSE

KnowltAll



How can a computer accumulate a massive body of knowledge?
 What will Web search engines look like in ten years?

To address the questions above, the KnowItAll project has been developing a variety of domain-independent systems that extract information from the Web in an autonomous, scalable manner.



Carnegie Mellon University



UNIVERSITY OF WASHINGTON | ABOUT US | CONTACT US | MY CSE | INTERNAL

Computer Science & Engineering

UNIVERSITY of WASHINGTON

W



Open Information Extraction



How can a computer accumulate a massive body of knowledge? What will Web search engines look like in ten years?

To address these questions, the Open IE project has been developing a Web-scale information extraction system that reads arbitrary text from any domain on the Web, extracts meaningful information, and stores it in a unified knowledge base for efficient querying. In contrast to traditional information extraction, the Open Information Extraction paradigm attempts to overcome the knowledge acquisition bottleneck by extracting a large number of relations at once.

Demo: <u>TextRunner</u> extracted over 500,000,000 assertions from 100 million Web pages. Software: <u>ReVerb</u> Open Information Extraction Software and additional information. Data: <u>Horn-clause</u> inference rules learned by the Sherlock system. Demo: <u>Selectional Preferences from Web Text</u> compute admissible argument values for a relation.

Data: <u>10,000 Functional Relations learned from Web Text</u> predict the functionality of a phrase.





Q





- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll









Read the Web

Research Project at Carnegie Mellon University

Home Project

Project Overview

Resources & Data Publications

NELL

People

NELL: Never-Ending Language Learning

Can computers learn to read? We think so. "Read the Web" is a research project that attempts to create a computer system that learns over time to read the web. Since January 2010, our computer system called NELL (Never-Ending Language Learner) has been running continuously, attempting to perform two tasks each day:

 First, it attempts to "read," or extract facts from text found in hundreds of millions of web pages (e.g., playsInstrument(George Harrison, guitar)).



• Second, it attempts to improve its reading competence, so that tomorrow it can extract more facts from the web, more accurately.

So far, NELL has accumulated over 15 million candidate beliefs by reading the web, and it is considering these at different levels of confidence. NELL has high confidence in 1,888,984 of these beliefs — these are displayed on this website. It is not perfect, but NELL is learning. You can track NELL's progress below or <u>@cmunell on Twitter</u>, browse and download its <u>knowledge base</u>, read more about our <u>technical approach</u>, or join the <u>discussion group</u>.







Machine Learning

- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







Carnegie Mellon University

Machine Learning

• What is Machine Learning?

The field of Machine Learning seeks to answer the question

"How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?" [Mitchell, 2006]





Machine Learning

• What is Machine Learning?

a machine *learns* with respect to a particular:

- task T
- performance metric P
- type of experience E

if the system reliably improves its performance P at task T, following experience E. [Mitchell, 1997]





Machine Learning

- Examples of Machine Learning approaches for different tasks (T), performance metrics (P) an experiences (E)
- data mining
- autonomous discovery
- database updating
- programming by example
- Pattern recognition







Machine Learning

- Supervised Learning;
- Unsupervised Learning
- Semi-Supervised Learning













How supervised learning algorithms work [edit] Wikipedia definition

Federal University of São Carlo

Given a set of N training examples of the form $\{(x_1, y_1), ..., (x_N, y_N)\}$ such that x_i is the feature vector of the i-th example and y_i is its label (i.e., class), a learning algorithm seeks a function $g: X \to Y$, where X is the input space and Y is the output space. The function g is an element of some space of possible functions G, usually called the *hypothesis space*. It is sometimes convenient to represent g using a scoring function $f: X \times Y \to \mathbb{R}$ such that g is defined as returning the y value that gives the highest score: $g(x) = \arg \max_{y} f(x, y)$. Let F denote the space of scoring functions.

0.5	1	blue		
2.9	1.9	blue		
1.2	3.1	blue		
0.8	4.7	blue		
2.7	5.4	blue		
8.1	4.7	red		
8.3	6.6	red		
6.3	6.7	red		
8	9.1	red		
5.4	8.4	red		




























































U-ISI AT

0.5	1	blue
2.9	1.9	blue
1.2	3.1	blue
0.8	4.7	blue
2.7	5.4	blue
8.1	4.7	red
8.3	6.6	red
6.3	6.7	red
8	9.1	red
5.4	8.4	red



Federal University of São Carlos

U-ISI AT

0.5	1	
2.9	1.9	
1.2	3.1	
0.8	4.7	
2.7	5.4	
8.1	4.7	
8.3	6.6	
6.3	6.7	
8	9.1	
5.4	8.4	





(one simple anecdotal approach)



PESP





(one simple anecdotal approach)







(one simple anecdotal approach)







(one simple anecdotal approach)







uisca

(one simple anecdotal approach)













u isi an

(one simple anecdotal approach)



DESP







(one simple anecdotal approach)













(one simple anecdotal approach)







(one simple anecdotal approach)





Semi-supervised Learning

(one simple anecdotal approach)





Semi-supervised Learning

(one simple anecdotal approach)







(one simple anecdotal approach)







(one simple anecdotal approach)





Semi-supervised Learning







(one simple anecdotal approach)













(one simple anecdotal approach)







(one simple anecdotal approach)













(one simple anecdotal approach)







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL






- "The autonomous understanding of text" [Etzioni et al., 2007]
- "One of the most important methods by which human beings learn is by reading" [Clark et al., 2007], thus why not building machines capable of learning by reading?







- "The problem of deciding what was implied by a written text, of reading between the lines is the problem of inference." [Norvig, 2007]
- Typically, Machine Reading is different from Natural Language Processing alone





It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn? Worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at same te on the same and of Hedeby. The old man drows Blomkvist in by premising sond evidence against Wennerström. Blomkvist ag same pend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af same overing that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, bu Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at te on the same and of Hedeby. The old man drov a Blophovist in by promising solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco uncleOf inted with the men hires the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik. Af same)vering that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

ISC/a

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder. Blomkvist visits Henrik Vanger at te on the same and of Hedeby. same The old man drov a Blorikvist in by much sing solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist the extended Vanger family, most of whom resent beco uncleOf inted with the men hires his presence. He does, however, start a short lived affair with Cecilia, the niece of T enemvOf Af same vering that Salander has hacked into his co affair With persuade ASSIST him with research. They even¹ affairWith lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone sne meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same



ESWC Conference Series



Get all information about previous ESWC conferences

Brought to you by



The ESWC 2017 takes place from May 28th, 2016 to June 1st, 2017 in Portoroz, Slovenia.

About ESWC2017

The ESWC is a major venue for discussing the latest scientific results and technology innovations around semantic technologies. Building on its past success, ESWC is seeking to broaden its focus to span other relevant related research areas in which Web semantics plays an important role.

The goal of the Semantic Web is to create a Web of knowledge and services in which the semantics of content is made explicit and content is linked to both other content and services allowing novel applications to combine content from heterogeneous sites in unforeseen ways and support enhanced matching between users needs and content. This network of knowledge-based functionality will weave together a large network of human knowledge, and make this knowledge machine-processable to support intelligent behaviour by machines. Creating such an interlinked Web of knowledge which spans unstructured text, structured data (e.g. RDF) as well as multimedia content and services requires the collaboration of many disciplines, including but not limited to: Artificial Intelligence, Natural Language Processing, Databases and Information Systems, Information Retrieval, Machine Learning, Multimedia, Distributed Systems, Social Networks, Web Engineering, and Web Science. These complementarities are reflected in the outline of the technical program of the ESWC 2017; in addition to the standard research and in-use tracks, we will feature two special tracks putting particular emphasis on interdisciplinary research topics and areas that show the potential of exciting synergies





Series

Get all information about previous ESWC conferences

Brought to you by



About ESWC2017

The ESWC 2017 takes place from May 28th, 2016 to June 1st, 2017 in Portoroz, Slovenia.

The ESWC is a major venue for discussing the latest scientific results and technology innovations around semantic technologies. Building on its past success, ESWC is seeking to broaden its focus to span other relevant related research areas in which Web semantics plays an important role.

The goal of the Semantic Web is to create a Web of knowledge and services in which the semantics of content is made explicit and content is linked to both other content and services allowing novel applications to combine content from heterogeneous sites in unforeseen ways and support enhanced matching between users needs and content. This network of knowledge-based functionality will weave together a large network of human knowledge, and make this knowledge machine-processable to support intelligent behaviour by machines. Creating such an interlinked Web of knowledge which spans unstructured text, structured data (e.g. RDF) as well as multimedia content and services requires the collaboration of many disciplines, including but not limited to: Artificial Intelligence, Natural Language Processing, Databases and Information Systems, Information Retrieval, Machine Learning, Multimedia, Distributed Systems, Social Networks, Web Engineering, and Web Science. These complementarities are reflected in the outline of the technical program of the ESWC 2017; in addition to the standard research and in-use tracks, we will feature two special tracks putting particular emphasis on interdisciplinary research topics and areas that show the potential of exciting synergies































































- One important (*initial*) approach to machine reading is to extract facts from text and store them in a structured form.
- Facts can be seen as entities and their relations
- Ontology is one of the most common representation for the extracted facts







Machine Reading Some possible subtasks

- Named Entity Resolution/Recognition
- Relation Extraction
- Co-reference and Polysemy Resolution
- Relation Discovery
- Inference
- Knowledge Base
- Document/Sentence Understanding (Micro-Reading)







- Named Entity Resolution/Recognition
- Relation Extraction
- Co-reference and Polysemy Resolution
- Relation Discovery
- Inference
- Knowledge Base
- Document/Sentence Understanding (Micro-Reading)





Carnegie Mellon University

Machine Reading

• Named Entity Resolution/Recognition

Semi-structured data

The "Low-Hanging Fruit"

- Wikipedia infoboxes & categories
- HMTL lists & tables, etc.

– Free text

Federal University of São Car

- Hearst-patterns; clustering by verbal phrases
- Natural-language processing
- Advanced patterns & iterative bootstrapping



- Named Entity Recognition [Nadeau & Sekine, 2007]
 - term "Named Entity" coined for the Sixth Message Understanding Conference (MUC-6) (R. Grishman & Sundheim 1996).
 - important sub-tasks of IE called "Named Entity Recognition and Classification (NERC)".



Federal University of São Carlo



[Nadeau & Sekine, 2007]

- recognize information units like names, including person, organization and location names, and numeric expressions including time, date, money and percent expressions.
- In Machine Reading, many other entities: product, kitchen item, sport, etc.



Federal University of São Carlo





Named Entity Resolution

- Named Entity Resolution [Theobald & Weikum, 2012]
 - Which individual entities belong to which classes?
 - *instanceOf (Surajit Chaudhuri, computer scientists),*
 - instanceOf (BarbaraLiskov, computer scientists),
 - instanceOf (Barbara Liskov, female humans), ...





- Named Entities Recognition as a machine learning task.
 - Supervised Learning



Federal University of São Carlos

SC/ar





- Named Entities Recognition as a machine learning task.
 - Supervised Learning



Federal University of São Carlos

UFISC/ar





- Named Entities Recognition as a machine learning task.
 - Supervised Learning

Federal University of São Carlos

UFISC/ar







Carnegie Mellon University

Named Entity Recognition

- Named Entity Recognition as a Machine Learning task.
 - Supervised Learning

Federal University of São Carlos

IST /ar

- Possible features [Ratinov & Roth, 2009], [Khambhatla, 2004], [Zhou et. al. 2005]
 - Words "around" and including entities
 - POS (Part-Of-Speech)
 - Prefixes and suffixes
 - Capitalization
 - Number of words
 - Number of characters
 - First word, last word
 - gazetteer matches



- Named Entities Recognition as a machine learning task.
 - Supervised Learning

Federal University of São Carlos

ufisio/ar







Supervised Learning

Federal University of São Carlos

U-ISI AT







Supervised Learning

Federal University of São Carlos

U-ISI AT







Carnegie Mellon University

Named Entity Recognition [Bach & Badaskar, 2007]

• Supervised Learning using Kernels

Federal University of São Carl

- A Kernel defines similarity implicitly in a higher dimensional space
- Can be based on Strings, Word Sequences, Parse Trees, etc.
 - For strings similarity in number of common substrings (or subsequences)
 - Recommended reading on string kernels [Lodhi et. al., 2002]



Federal University of São Carlos

U-ISI AT

[Bach & Badaskar, 2007]

	Feature set Definition	Computational Complexity
Feature based Methods	Required to define a feature- set to be extracted after textual analysis. Good features arrived at by experimentation	Relatively lower
Kernel Methods	No need to define a feature- set. Similarity computed over a much larger feature space implicitly.	Relatively higher

IEA DESD



Semi-supervised
Approaches

Federal University of São Carlos

ufisio/ar

Bootstrap can
generate a large
number of patterns
and NE instances.







Semi-supervised
Approaches

Federal University of São Carlos

- ISC/a

Bootstrap can
generate a large
number of patterns
and NE instances.







Semi-supervised
Approaches

Federal University of São Carlos

U ISC/a

 Bootstrap can generate a large number of patterns and NE instances.











Federal University of São Carlos

U ISC/a






Federal University of São Carlos

ufisio/ar







Federal University of São Carlos

ufisio/ar







Federal University of São Carlos

ufisio/ar







Federal University of São Carlos

u isca



Unsupervised
 Approaches

Federal University of São Carlos

ufisio/ar

Bootstrap can
 generate a large
 number of patterns
 and NE instances.





• [Ratinov & Roth, 2009]

Federal University of São Carlos

U-ISI /an

	System	Resources Used	F_1
+	LBJ-NER	Wikipedia, Nonlocal Fea-	90.80
		tures, Word-class Model	
-	(Suzuki and	Semi-supervised on 1G-	89.92
	Isozaki, 2008)	word unlabeled data	
-	(Ando and	Semi-supervised on 27M-	89.31
	Zhang, 2005)	word unlabeled data	
-	(Kazama and	Wikipedia	88.02
	Torisawa, 2007a)		
-	(Krishnan and	Non-local Features	87.24
	Manning, 2006)		
-	(Kazama and	Non-local Features	87.17
	Torisawa, 2007b)		
+	(Finkel et al.,	Non-local Features	86.86
	2005)		





• [Pennington & Socher & Manning, 2014]

Federal University of São Carlos

ISC/ar

Model	Dev	Test	ACE	MUC7
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	88.7	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	93.2	88.3	82.9	82.2







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Polysemy Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







Relation Extraction

Semi-structured data

The "Low-Hanging Fruit"

- Wikipedia infoboxes & categories
- HMTL lists & tables, etc.

– Free text

- Hearst-patterns; clustering by verbal phrases
- Natural-language processing
- Advanced patterns & iterative bootstrapping







- Relation Extraction [Theobald & Weikum, 2012]
 - Which instances (pairs of individual entities) are there for given binary relations with specific type signatures?
 - hasAdvisor (JimGray, MikeHarrison)
 - hasAdvisor (HectorGarcia-Molina, Gio Wiederhold)
 - hasAdvisor (Susan Davidson, Hector Garcia-Molina)
 - graduatedAt (JimGray, Berkeley)
 - graduatedAt (HectorGarcia-Molina, Stanford)
 - hasWonPrize (JimGray, TuringAward)
 - bornOn (JohnLennon, 9Oct1940)
 - diedOn (JohnLennon, 8Dec1980)
 - marriedTo (JohnLennon, YokoOno)







[Bach & Badaskar, 2007]

- Extracting semantic relations between entities in text
- Relation extraction as a Machine Learning task.

- Supervised Learning









[Bach & Badaskar, 2007]

- Extracting semantic relations between entities in text
- Relation extraction as a Machine Learning task.

- Supervised Learning







[Bach & Badaskar, 2007]

- Extracting semantic relations between entities in text
- Relation extraction as a Machine Learning task.

- Supervised Learning







[Bach & Badaskar, 2007]

- Relation extraction as a Machine Learning task.
 - Supervised Learning
 - Possible features [Khambhatla, 2004], [Zhou et. al.
 2005]
 - Words between and including entities
 - Types of entities (person, location, etc)
 - Number of entities between the two entities, whether both entities belong to same chunk
 - # words separating the two entities
 - Path between the two entities in a parse tree







[Bach & Badaskar, 2007]

- Extracting semantic relations between entities in text
- Relation extraction as a classification task.

- Supervised Learning







[Bach & Badaskar, 2007]

- Extracting semantic relations between entities in text
- Relation extraction as a classification task.

- Supervised Learning







[Bach & Badaskar, 2007]

- Supervised Learning using Kernels
 - A Kernel defines similarity implicitly in a higher dimensional space
 - Can be based on Strings, Word Sequences, Parse Trees, etc.
 - For strings, similarity in number of common substrings (or subsequences)
 - Recommended reading on string kernels [Lodhi et. al., 2002]







[Bach & Badaskar, 2007]

	Feature set Definition	Computational Complexity
Feature based Methods	Required to define a feature- set to be extracted after textual analysis. Good features arrived at by experimentation	Relatively lower
Kernel Methods	No need to define a feature- set. Similarity computed over a much larger feature space implicitly.	Relatively higher

FAPESP



Semi-supervised
 Approaches

Federal University of São Carlos

u isca







Semi-supervised
 Approaches

Federal University of São Carlos

u isc⁄a







Carnegie Mellon University

Relation Extraction

Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar





Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar

 Bootstrap can generate a large number of patterns and relation instances.







Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar





Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar

Bootstrap can
 generate a large
 number of patterns
 and relation
 instances.











Federal University of São Carlos

u isca







Federal University of São Carlos

u isca



Unsupervised
 Approaches

Federal University of São Carlos

ufisio/ar

Bootstrap can
 generate a large
 number of patterns
 and relation
 instances.







- Supervised learning [Bunescu & Mooney, 2005]
- Distant and Partial Supervised [Angeli & Tibshirani & Wu & Manning, 2014]







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Polysemy Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)





Co-reference: expressions that refer to the same entity



Example (figure) taken from: http://nlp.stanford.edu/projects/coref.shtml



Federal University of São Ca

ISC/A



• Co-reference: expressions that refer to the same entity

Federal University of São C

"I voted for Nader because he was most aligned with my values," she said.

Example (figure) taken from: http://nlp.stanford.edu/projects/coref.shtml

within-document co-reference



Federal University of São Car

u isca

 Co-reference: expressions that refer to the same entity peca aligned w Example (figur niects/coref.shtml ocument c ference



Carnegie Mellon University

Co-Reference and Polysemy Resolution

• Co-reference: expressions that refer to the same entity



Example (figure) adapted from [Krishnamurthy & Mitchell, 2011]



Federal University of São Carl

u isc⁄a



• Co-reference: expressions that refer to the same entity



Example (figure) adapted from [Krishnamurthy & Mitchell, 2011]



Federal University of São Car

u isc⁄a





• Co-reference: expressions that refer to the same entity

Federal University of São Ca

U ISC/a





- Co-reference: expressions that refer to the same entity
- Which names denote which entities? [Theobald & Weikum, 2012]
 - means ("Lady Di", Diana Spencer),

Federal University of São

- means ("Diana Frances Mountbatten-Windsor", Diana Spencer), ...
- means ("Madonna", Madonna Louise Ciccone),
- means ("Madonna", Madonna(painting by Edward Munch)), ...

cross-document co-reference




 Polysemy: is the capacity for a sign (such as a word, phrase, or symbol) to have multiple meanings [Wikipedia]



Federal University of São Carlo

U-ISC/a



Federal University of São Ca

ISC/a

 Polysemy: is the capacity for a sign (such as a word, phrase, or symbol) to have multiple meanings [Wikipedia]





• Co-Reference and Polysemy

Federal University of São Carl

u isca



Example (figure) adapted from [Krishnamurthy & Mitchell, 2011]







- Co-reference and Polysemy:
 - Supervised Learning



Federal University of São Carlos





- Co-reference and Polysemy:
 - Supervised Learning

Federal University of São Carlos







- Co-reference and Polysemy:
 - Supervised Learning

Federal University of São Carlos









Co-Reference

Federal University of São Carlos

U-ISC/a

Resolution	Category	Feature	Source
Resolution.	Mention Types	Mention Type Pair	Annotation and tokens
	String Relations	Head Match	Tokens
– Supervised		Extent Match	Tokens
Loorning		Substring	Tokens
Learning		Modifiers Match	Tokens
		Alias	Tokens and lists
– Possible	Semantic	Gender Match	WordNet and lists
faaturaa		Number Match	WordNet and lists
reatures		Synonyms	WordNet
[Bongtson &		Antonyms	WordNet
[Dengtson a		Hypernyms	WordNet
Roth 2008]		Both Speak	Context
Noti, 2000]	Relative Location	Apposition	Positions and context
		Relative Pronoun	Positions and tokens
		Distances	Positions
	Learned	Anaphoricity	Learned
		Name Modifiers Predicted Match	Learned
	Aligned Modifiers	Aligned Modifiers Relation	WordNet and lists
	Memorization	Last Words	Tokens
	Predicted Entity Types	Entity Types Match	Annotation and tokens
		Entity Type Pair	WordNet and tokens





Co-Reference

Federal University of São Carlos

U-ISIO/

Posolution	Category	Feature	Source
Resolution.	Mention Types	Mention Type Pair	Annotation
	String Relations	Head Match	Toker
– Supervised		Extent Match	T
Loorning		Substring	columna
Leannig		Modifiers Match	
		Alias	ons and lists
— Possible	Semantic	Gender Match	WordNet and lists
faaturaa		Number Mat	WordNet and lists
reatures		Synony	WordNet
[Bangtson &		An.	WordNet
[Deligison &		2777	WordNet
Roth 2008]			Context
Notil, 2000]	Relative Loc	sition	Positions and context
		relative Pronoun	Positions and tokens
	- 101-	Distances	Positions
		Anaphoricity	Learned
		Name Modifiers Predicted Match	Learned
	.ed Modifiers	Aligned Modifiers Relation	WordNet and lists
	Memorization	Last Words	Tokens
	Predicted Entity Types	Entity Types Match	Annotation and tokens
		Entity Type Pair	WordNet and tokens



- Co-reference and Polysemy:
 - Supervised Learning

Federal University of São Carlos





• Supervised Learning using Kernels

Federal University of São Ca

ISCA

- A Kernel defines similarity implicitly in a higher dimensional space
- Can be based on Strings, Word Sequences, Parse Trees, etc.
 - For strings similarity in number of common substrings (or subsequences)
 - Recommended reading on string kernels [Lodhi et. al., 2002]





Semi-supervised
 Approaches

Federal University of São Carlos

u isca

Bootstrap can
 generate a large
 number of patterns
 and relation
 instances.







Semi-supervised
 Approaches

Federal University of São Carlo

u isca

 Bootstrap can generate a large number of patterns and relation instances.







Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar

 Bootstrap can generate a large number of patterns and relation instances.







Semi-supervised
 Approaches

Federal University of São Carlo

ufisio/ar

 Bootstrap can generate a large number of patterns and relation instances.





Federal University of São Carlo

ufisio/ar

 Semi-supervised X also know as Y Set of labeled Pattern Approaches **Examples** Bootstrap can generate a large Pair of Instances number of patterns Classifier and relation instances. **Apple Computer -**Set of labeled pairs of Apple **Instances Examples**



 Semi-supervised X also know as Y Set of labeled Pattern Approaches **Examples** Bootstrap can generate a large Pair of Pattern Instances number of patterns Classifier Classifier and relation instances. **Apple Computer -**Set of labeled pairs of Apple **Instances Examples**



Federal University of São Carlo

ufisio/ar







Federal University of São Carlo

U-ISI /an







Federal University of São Carlos

U-ISI /an



Semi-supervised
 Approaches

Federal University of São Carlos

ufisio/ar

 Bootstrap can generate a large number of patterns and relation instances.





Carnegie Mellon University

Co-Reference and Polysemy Resolution

Co-Reference Resolution: [Singh et al., 2011], [Krishnamurthy & Mitchell, 2011], [Dutta & Weikum, 2015]

• Polysemy Resolution:

Federal University of São Ca

[Krishnamurthy & Mitchell, 2011], [Galárraga et al., 2014]





- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







- Relation Discovery
 - Which new relations are there for given pair of entities?
 - hasAdvisor (JimGray, MikeHarrison)







- Relation Discovery
 - Which new relations are there for given pair of entities?
 - hasAdvisor (JimGray, MikeHarrison)
 - hasCoAuthor(HectorGarcia-Molina, Gio Wiederhold)







- Relation Discovery
 - Which new relations are there for given pair of entities?
 - hasAdvisor (JimGray, MikeHarrison)
 - hasCoAuthor(HectorGarcia-Molina, Gio Wiederhold)
 - graduatedAt (JimGray, Berkeley)







- Relation Discovery
 - Which new relations are there for given pair of entities?
 - hasAdvisor (JimGray, MikeHarrison)
 - hasCoAuthor(HectorGarcia-Molina, Gio Wiederhold)
 - graduatedAt (JimGray, Berkeley)
 - studiedAt (HectorGarcia-Molina, Stanford)
 - bornOn (JohnLennon, 9Oct1940)
 - releasedAlbum (JohnLennon, 10Dec1965)





Relation Discovery

Federal University of São Carlos

U-ISI AT







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







Inference

 Inference is the act or process of deriving logical conclusions from premises known or assumed to be true [Wikipedia]







Inference

• Manually crafted inference rules

- Automatically learned inference rules
- Data mining the Knowledge Base







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







Ontology Representation

Facts (RDF triples)

- 1: (Jim, hasAdvisor, Mike)
- 2: (Surajit, hasAdvisor, Jeff)
- 3: (Madonna, marriedTo, GuyRitchie)
- 4: (Nicolas, marriedTo, Carla)
- 5: (ManchesterU, wonCup, ChampionsLeague)

Reification:

"Facts about Facts":

- 6: (1, inYear, 1968)
- 7: (2, inYear, 2006)
- 8: (3, validFrom, 22-Dec-2000)
- 9: (3, validUntil, Nov-2008)
- 10: (4, validFrom, 2-Feb-2008)
- 11: (2, source, SigmodRecord)
- 12: (5, inYear, 1999)
- 13: (5, location, CampNou)
- 14: (5, source, Wikipedia)







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)





Document/Sentence Understanting (Micro-Read)

Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots.



Federal University of São Carlos

ufisio/ar







Document/Sentence Understanting (Micro-Read)

S1.) Alice caught the butterfly with the spots.





Document/Sentence Understanting (Micro-Read)

Nakashole and Mitchell, A

S1.) Alice caught the butterfly with the spots.S2.) Alice caught the butterfly with the net.

Federal University of São Carlos

u is dat







Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots. S2.) Alice caught the butterfly with the net. Relations:

Federal University of São Carlos

FISC/a

locatedIn(*city*, *country*): locatedIn(Saint Etienne, France), ...

animalHasCoatPattern(*animal, coatPattern*): animalHasCoatPattern(zebra, stripes), animalHasCoatPattern(leopard, spots), animalHasCoatPattern(butterfly, spots), ...

caught(instrument, animal):
caught(net, butterfly),
caught(fish, fishing rod), ...

Categories:

city(*city*): city(Saint Etienne), ...

animal(*animal*): animal(zebra), animal(butterfly), ...




Document/Sentence Understanting (Micro-Read)

S2.) Alice caught the butterfly with the net. **Relations:**

locatedIn(*city*, *country*): locatedIn(Saint Etienne, France), ... animalHasCoatPattern(*animal*, *coatPattern*): animalHasCoatPattern(zebra, stripes), animalHasCoatPattern(leopard, *Spots*), animalHasCoatPattern(butterfly, *spots*), ... caught(*instrument*, *animal*): caught(*inet*, butterfly), caught(fish, fishing rod), ...

Categories:

city(city):
city(Saint Etienne), ...

Federal University of São Carlos

uisca

animal(*animal*): animal(zebra), animal<mark>(bu</mark>tterfly), ...





Document/Sentence Understanting (Micro-Read)

Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots. S2.) Alice caught the butterfly with the net.

Federal University of São Carlos

u is dat







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - -DBPedia
 - YAGO
 - KnowltAll
 - NELL







DBPedia

http://wiki.dbpedia.org/



ABOUT

DBpedia is a crowd-sourced community effort to extract structured information from Wikipedia and make this information available on the Web. DBpedia allows you to ask sophisticated queries against Wikipedia, and to link the different data sets on the Web to Wikipedia data. We hope that this work will make it easier for the huge amount of information in Wikipedia to be used in some new interesting ways. Furthermore, it might inspire new mechanisms for navigating, linking, and improving the encyclopedia itself. Further information

...









DBPedia

Querying Wikipedia like a Database

Mapping Wikipedia semi-structured data into RDF triples







DBPedia

Querying Wikipedia like a Database

Mapping Wikipedia semi-structured data into RDF triples

Semi-structured data The "Low-Hanging Fruit"





It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn? Worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at same te on the same and of Hedeby. The old man drows Blomkvist in by premising sond evidence against Wennerström. Blomkvist ag same pend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af same overing that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, bu Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at te on the same and of Hedeby. The old man drov a Blophovist in by promising solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco uncleOf inted with the men hires the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik. Af same)vering that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

ISC/a

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder. Blomkvist visits Henrik Vanger at te on the same and of Hedeby. same The old man drov a Blorikvist in by much sing solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist the extended Vanger family, most of whom resent beco uncleOf inted with the men hires his presence. He does, however, start a short lived affair with Cecilia, the niece of T enemvOf Af same vering that Salander has hacked into his co affair With persuade ASSIST him with research. They even¹ affairWith lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone sne meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

- ISC/a







Carnegie Mellon University

DBPedia

- How to Read Wikipedia Semi-structured data? [Lehmann et al., 2014]
 - Parse Wikipedia Markup language
 - Overcome the lack of standard problem
 - Same properties might have different names
 - "Datebirth" and "Birth_date"
 - "Birthplace" and "Birth_place"
 - Instead of "Modeling the World", try to structure the available information

Federal University of São Carlos



DBPedia

About: Porto

An Entity of Type : Municipality, from Named Graph : http://dbpedia.org, within Data Space : dbpedia.org



Porto (Portuguese pronunciation: ['portu]), also known as Oporto in English, is the second-largest city in Portugal, after Lisbon, and one of the major urban areas in Southwestern Europe. The urban area of Porto, which extends beyond the administrative limits of the city, has a population of 1.4 million (2011) in an area of 389 km2 (150 sq mi), making it the second-largest urban area in Portugal. Porto Metropolitan Area, on the other hand, includes an estimated 1.8 million people.

Property	Value
dbo:abstract	Porto (Portuguese pronunciation: ['portu]), also known as Oporto in English, is the second-largest city in Portugal, after Lisbon, and one of the major urban areas in Southwestern Europe. The urban area of Porto, which extends beyond the administrative limits of the city, has a population of 1.4 million (2011) in an area of 389 km2 (150 sq mi), making it the second-largest urban area in Portugal. Porto Metropolitan Area, on the other hand, includes an estimated 1.8 million people. It is recognized as a Gamma-level global city by the Globalization and World Cities (GaWC) Study Group, the only Portuguese city besides Lisbon with such recognition.Located along the Douro river estuary in Northern Portugal, Porto is one of the oldest European centres, and its historical core was proclaimed a World Heritage Site by UNESCO in 1996. The western part of its urban area extends to the coastline of the Atlantic Ocean. Its settlement dates back many centuries, when it was an outpost of the Roman Empire. Its Latin name, Portus Cale, has been referred to as the origin of the name "Portugal", based on transliteration and oral evolution from Latin. In Portuguese the name of the city is spelled with a definite article ("o Porto"; English: the port). Consequently, its English name evolved from a misinterpretation of the oral pronunciation and referred to as Oporto in modern literature and by many speakers.One of Portugal's internationally famous exports, port wine, is named for Porto, since the metropolitan area, and in particular the adegas of Vila Nova de Gaia, were responsible for the production and export of the fortified wine. In 2014, Porto was elected The Best European Destination by the Best European Destinations Agency.
dbo:areaCode	(+351) 22
dbo:country	• dbr:Portugal
dbo:daylightSavingTimeZone	dbr:Western_European_Summer_Time
dbo:district	dbr:Porto_District
dbo:elevation	104.00000 (xsd:double)
dbo:government	 dbr:Concelho dbr:Câmara_municipal
dbo:governmentElevation	120.00000 (xsd:double)
dbo:governmentPosition	dbr:PortogovernmentPosition1
dbo:governmentRegion	dbr:Santo_Ildefonso_(Porto)
dbo:highestPosition	dbr:PortohighestPosition1
dbo:highestRegion	dbr:Bonfim_(Porto)
dbo:lowest	Sea level
dbo:lowestMountain	dbr:Atlantic_Ocean
dbo:lowestPlace	dbr:Atlantic_Ocean
dbo:maximumElevation	149.00000 (xsd:double)
dbo:mayor	• dbr:Rui_Moreira







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - -YAGO
 - KnowltAll
 - NELL





The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web

									Location	Press	Deutsch	
	P	max info	planck instit matik	tut							Q	
25ANNIVERSARY	HOME	INSTITUTE	NEWS	DEPARTMENTS	PUBLICATIONS	PEOPLE	SOFTWARE	SERVICES				
Departments Database	s and Inform	ation Systems	► Research ■	YAGO-NAGA								
Algorithms & Comple	xity		YAGO-NA	GA								
Computer Vision and Computing	Multimodal		Harvestin	g Searching ar	d Ranking Kno	wledge fro	m the Web					
Computational Biolog Algorithmics	y & Applied		The YAGO-NAG	GA project started in 2006	6 with the goal of building	g a conveniently	searchable, large-	AN				
Computer Graphics			scale, highly ac We have alread	curate knowledge base o ly harvested knowledge a	f common facts in a mac about millions of entities a	chine-processible and facts about th	representation. heir relationships,	la co		2		
Databases and Inform	ation Systen	ns	from Wikipedia base, coined Y/	and WordNet with carefu AGO, has very high preci-	l integration of these two sion and is freely availab	sources. The read	sulting knowledge represented as		1 E	han		
People			RDF triples, and	d we have developed me	thods and prototype syst	tems for querying	, ranking, and es based on					
Research			statistical mode	ls.	in the provided railited						~	
YAGO-NAGA			Several interline of Semantic We	ked sub-projects are grow b (Ontologies), Social W	ving on the YAGO-NAGA eb (Web 2.0), and Statis	A basis. Our vision tical Web (Inform	n is a confluence ation Extraction)					
AIDA			assets towards	a comprehensive reposit	ory of human knowledge	e. Our methodolog	gies combine					
AMIE			retrieval, statisti	ical learning, and logical r	reasoning.	anabase of otomo						

People

Weikum, Gerhard Hoffart, Johannes Suchanek, Fabian Berberich, Klaus Kuzey, Erdal Yahya, Mohamed Elbassuoni, Shady Kasneci, Gjergji Lewis-Kelham, Edwin de Melo, Gerard Nakashole, Ndapandula Neumann 1, Thomas Preda, Nicoleta Sozio, Mauro Taneva, Bilyana Wang, Yafang Yosef, Mohamed Amir

ANGIE

EVIN

HYENA

Javatools

IBEX



The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web









- Yet Another Great Ontology YAGO
- Main Goal: building a conveniently searchable, large-scale, highly accurate knowledge base of common facts in a machine-processable representation







Carnegie Mellon University

YAGO

						Location	Press	Deutsch
	a planck institut ormatik							Q
25ANNIVERSARY HOME INSTITUTE	NEWS DEPARTMENTS	PUBLICATIONS	PEOPLE	SOFTWARE	SERVICES			
Departments Databases and Information Systems	s ▶ Research ▶ YAGO-NAGA ▶ YA	GO						
Overview	YAGO: A High-Quality Kno	wledge Base						
Demo	. .							
Downloads	Overview							
Statistics	YAGO is a huge semantic knowledge ba Currently, YAGO has knowledge of more	ase, derived from <u>Wikiped</u> e than 10 million entities (lia WordNet and like persons, org	GeoNames. anizations, cities,				
Publications	etc.) and contains more than 120 million	facts about these entities	5.				\neg	
Linking	YAGO is special in several ways: 1 The accuracy of YAGO has been m	anually evaluated proving	n a confirmed ac	curacy of 95%	Уċ	انال	\cup	
Archive	Every relation is annotated with its c 2 YAGO combines the clean taxonom	onfidence value.	ness of the Wiki	ipedia category	select	t knowle	edge	
Acknowledgements	system, assigning the entities to mo	re than 350,000 classes.		poula outogoly				
FAQ	YAGO is an ontology that is anchore facts and entities.	ed in time and space. YAC	GO attaches a te	mporal dimension ar	nd a spacial dim	ension to ma	any of its	
	 In addition to a taxonomy, YAGO has 	s thematic domains such	as "music" or "s	cience" from WordNe	et Domains.			
	YAGO extracts and combines entities	es and facts from 10 Wikip	pedias in differen	it languages.				
	YAGO is developed jointly with the DBM	leb group at <u>Télécom Par</u>	risTech Universit	<u>¥</u> .				

FAPESP





- Turn Web into Knowledge Base [Weikum et al., 2009]
 - Building a comprehensive Knowledge Base of human knowledge
 - knowledge from Wikipedia and WordNet
 - the ontology check itself for precision







- The knowledge base is automatically constructed from Wikipedia
- Each article in Wikipedia becomes an entity in the kb (e.g., since Leonard Cohen has an article in Wikipedia, LeonardCohen becomes an entity in YAGO).







- Main page Contents Featured content Current events Random article Donate to Wikipedia
- Interaction Help About Wikipedia Community portal Recent changes Contact Wikipedia
- Toolbox
- Print/export

 Languages Afrikaans للعربية Asturianu Български Bosanski Català Česky Dansk Deutsch Eesti Eλληνικά Español

Leonard Cohen

From Wikipedia, the free encyclopedia

Leonard Norman Cohen, CC GOQ (born 21 September 1934) is a Canadian singer-songwriter, musician, poet, and novelist. His work often explores religion, isolation, sexuality, and interpersonal relationships.^[1] Cohen has been inducted into the American Rock and Roll Hall of Fame and both the Canadian Music Hall of Fame and the Canadian Songwriters Hall of Fame. He is also a Companion of the Order of Canada, the nation's highest civilian honour. In 2011 Cohen received Prince of Asturias Award for literature.

In a speech at Cohen's induction into the Rock and Roll Hall of Fame on March 10, 2008, Lou Reed described Cohen as belonging to the "highest and most influential echelon of songwriters."^[2]

The critic Bruce Eder wrote an assessment of Cohen's overall career in popular music, writing, "[Cohen is] one of the most fascinating and enigmatic. . .singer/songwriters of the late '60s. . . [and] has retained an audience across four decades of music-making. . . Second only to Bob Dylan (and perhaps Paul Simon) [in terms of influence], he commands the attention of critics and younger musicians more firmly than any other musical figure from the 1960s who is still working at the outset of the 21st century."^[3]

The Academy of American Poets has commented more broadly on Cohen's overall career in the arts, including his work as a poet, novelist, and songwriter, stating that "[Cohen's] successful blending of poetry, fiction, and music is made most clear in *Stranger Music: Selected Poems and Songs*, published in 1993, which gathered more than two hundred of Cohen's poems . . .several novel excerpts, and almost sixty song lyrics. . .While it may seem to some that Leonard Cohen departed from the literary in pursuit of the musical, his fans continue to embrace him as a Renaissance man who straddles the elusive artistic borderlines." ^[4]

FAPESP

Contents [hide]
Early life
Poetry and novels
Recording career
3.1 1960s and 1970s
3.2 1980s
3.3 1990s
3.4 2000s
3.4.1 Post-monastery records: Ten New Songs, Dear Heather and Anjani's Blue Alert
3.4.2 Lawsuits and financial troubles
3.4.3 Book of Longing
3.5 2008-2010 World Tour
3.5.1 2008 tour



Back	ground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
Instruments	Vocals, guitar, piano, keyboards, synthesizer



Free Text

Leonard Cohen

From Wikipedia, the free encyclopedia

Featured content Order of Canada, the nation's highest civilian honour. In 2011 Cohen received Prince of Asturias Award for literature. Current events

Random article Donate to Wikipedia

WIKIPEDIA The Free Encyclopedia

Main page

Contents

- Interaction Help About Wikipedia Community portal Recent changes Contact Wikipedia
- Toolbox
- Print/export

 Languages Afrikaans العربية Asturianu Български Bosanski Català Česky Dansk Deutsch Eesti Ελληνικά Español

Leonard Norman Cohen, CC GOQ (born 21 September 1934) is a Canadian singer-sor often explores religion, isolation, sexuality, and interpersonal relationships.^[1] Cohen ha Hall of Fame and both the Canadian Music Hall of Fame and the Canadian Songwrite

poet, and novelist. His work d into the American Rock and Roll e is also a Companion of the

In a speech at Cohen's induction into the Rock and Roll Hall of Fame on March 10, 2008, Lou Reed described Cohen as belonging to the "highest and most influential echelon of songwriters."[2]

The critic Bruce Eder wrote an assessment of Cohen's overall career in popular music, writing, "[Cohen is] one of the most fascinating and enigmatic. . . singer/songwriters of the late '60s. . . [and] has retained an audience across four decades of music-making. . . Second only to Bob Dylan (and perhaps Paul Simon) [in terms of influence], he commands the attention of critics and younger musicians more firmly than any other musical figure from the 1960s who is still working at the outset of the 21st century."[3]

The Academy of American Poets has commented more broadly on Cohen's overall career in the arts, including his work as a poet, novelist, and songwriter, stating that "[Cohen's] successful blending of poetry, fiction, and music is made most clear in Stranger Music: Selected Poems and Songs, published in 1993, which gathered more than two hundred of Cohen's poems . . . several novel excerpts, and almost sixty song lyrics. . . While it may seem to some that Leonard Cohen departed from the literary in pursuit of the musical, his fans continue to embrace him as a Renaissance man who straddles the elusive artistic borderlines."^[4]

Contents [hide] 1 Early life 2 Poetry and novels 3 Recording career 3.1 1960s and 1970s 3.2 1980s 3.3 1990s 3.4 2000s 3.4.1 Post-monastery records: Ten New Songs, Dear Heather and Anjani's Blue Alert 3.4.2 Lawsuits and financial troubles 3.4.3 Book of Longing 3.5 2008-2010 World Tour 3.5.1 2008 tour



Back	ground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
nstruments	Vocals, guitar, piano, keyboards, synthesizer

FAPESP



1934) is a Can

l relation

nard C

dles the elusiv

Main page Contents

Featured content Current events Random article Donate to Wikipedia

- Interaction Help About Wikipedia Community portal Recent changes Contact Wikipedia
- Toolbox
- Print/export

 Languages Afrikaans العربية Asturianu Български Bosanski Català Česky Dansk Deutsch Eesti Ελληνικά

Español

WIKIPEDIA The Free Encyclopedia

Leonard Norman Co often explores religion, is Hall of Fame and both the C Order of Canada, the nation's h

Leonard Cohen

From Wikipedia, the free en

In a speech at Cohen's induction into "highest and most influential echelon of

The critic Bruce Eder wrote an assessment of enigmatic. . .singer/songwriters of the late '60 Bob Dylan (and perhaps Paul Simon) [in t any other musical figure from the 1960

The Academy of American Poets, novelist, and songwriter, stating Selected Poems and Song almost sixty song lyrics, continue to embrace

1 Early life 2 Poetry and novels 3 Recording career 3.1 1960s and 1970s 3.2 1980s 3.3 1990s 3.4 2000s 3.4.1 Post-monastery records: Ten New Songs, Dear Heather and Anjani's Blue Alert 3.4.2 Lawsuits and financial troubles 3.4.3 Book of Longing 3.5 2008-2010 World Tour 3.5.1 2008 tour

et, and novelist. His work o the American Rock and Roll e is also a Companion of the ard for literature.

Free

Text

Reed described Cohen as belonging to the

writing, "[Cohen is] one of the most fascinating and cross four decades of music-making. . . Second only to on of critics and younger musicians more firmly than rv."^[3]

> e arts, including his work as a poet, made most clear in Stranger Music: ems . . .several novel excerpts, and ursuit of the musical, his fans

Leonard Cohen



Back	ground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
Instruments	Vocals, guitar, piano, keyboards, synthesizer





pet, and novelist. His work o the American Rock and Roll e is also a Companion of the ard for literature.

leed described Cohen as belonging to the

cross four decades of music-making. . . Second only to on of critics and younger musicians more firmly than

> e arts, including his work as a poet, made most clear in Stranger Music: ursuit of the musical, his fans

Leonard Cohen



Back	ground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
nstruments	Vocals, guitar, piano, keyboards, synthesizer

WIKIPEDIA The Free Encyclopedia

- Main page Contents Featured content Current events Random article Donate to Wikipedia
- Interaction Help About Wikipedia Community portal Recent changes Contact Wikipedia
- Toolbox
- Print/export

 Languages Afrikaans العربية Asturianu Български Bosanski Català Česky Dansk Deutsch Eesti Ελληνικά Español

Leonard Cohen From Wikipedia, the free eng

Leonard Norman Co often explores religion, is Hall of Fame and both the C Order of Canada, the nation's h

In a speech at Cohen's induction into "highest and most influential echelon of

The critic Bruce Eder wrote an assessment enigmatic. . .singer/songwriters of the late '60s Bob Dylan (and perhaps Paul Simon) [in te any other musical figure from the 1960

The Academy of American Poets novelist, and songwriter, stating Selected Poems and Songs almost sixty song lyrics. continue to embrace

1 Early life 2 Poetry and novels 3 Recording career 3.1 1960s and 1970s 3.2 1980s 3.3 1990s 3.4 2000s 3.4.1 Post-monastery records: Ten New Songs, Dear Heather and Anjani's Blue Alert 3.4.2 Lawsuits and financial troubles 3.4.3 Book of Longing 3.5 2008-2010 World Tour 3.5.1 2008 tour

1934) is a Can al relation

hard C

dles the elusiv

Free

Text

writing, "[Cohen is] one of the most fascinating and VV."[3]

oems . . .several novel excerpts, and

Carnegie Mellon University

InfoBox







Wikipedia InfoBox

Bac	kground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
Instruments	Vocals, guitar, piano, keyboards, synthesizer

FAPESP





YAGO Semi-structured data

Wikipedia InfoBox

The "Low-Hanging Fruit"

Bac	kground information
Birth name	Leonard Norman Cohen
Born	21 September 1934 (age 78) Montreal, Quebec, Canada
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop
Occupations	Musician, singer-songwriter, poet, novelist
Instruments	Vocals, guitar, piano, keyboards, synthesizer

FAPESP





YAGO Semi-structured data

Wikipedia InfoBox

The "Low-Hanging Fruit"

Background information				
Birth name	Leonard Norman Cohen			
Born	21 September 1934 (age 78) Montreal, Quebec, Canada			
Genres	Folk, folk rock, rock, pop rock, spoken word, synthpop			
Occupations	Musician, singer-songwriter, poet, novelist			
Instruments	Vocals, guitar, piano, keyboards, synthesizer			

FAPESP







 Certain categories are exploited to deliver type information (e.g., the article about Leonard Cohen is in the category Canadian male poets, so he becomes a Canadian poet).

Categories: Leonard Cohen | 1934 births | Anglophone Quebec people | Ashkenazi Jews | Canadian buskers | Canadian folk singers | Canadian male singers | Canadian Music Hall of Fame inductees | Canadian male novelists | Canadian people of Lithuanian-Jewish descent | Canadian people of Polish-Jewish descent | Canadian singer-songwriters | Canadian soft rock musicians | Canadian Zen Buddhists | Companions of the Order of Canada | Buddhists of Jewish descent | Genie and Canadian Screen Award winning people | Governor General's Award winning poets | Grammy Award winners | Grand Officers of the National Order of Quebec | Jewish Canadian musicians | Jewish songwriters | Jewish poets | Jewish singers | Jewish Canadian writers | Juno Award winners | Governor General's Performing Arts Award winners | Kohanim | Living people | McGill University alumni | Musicians from Montreal | Mystics | Rinzai Buddhists | Rock and Roll Hall of Fame inductees | People from Westmount, Quebec | Writers from Montreal | Fellows of the American Academy of Arts and Sciences | Converts to Buddhism | Canadian Buddhists | 20th-century Canadian novelists | 20th-century Canadian singers | 21st-century Canadian singers | 20th-century Canadian poets | Canadian male poets | Columbia Records artists | Genie Award winners for Best Achievement in Music - Original Song | Canadian screenwriters | Male screenwriters













- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes
- thematic categories (such as Canadian poetry) are different from conceptual Categories







- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes
- thematic categories (such as Canadian poetry) are different from conceptual Categories






- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes
- thematic categories (such as Canadian poetry) are different from conceptual Categories





- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes
- thematic categories (such as Canadian poetry) are different from conceptual Categories







Carnegie Mellon University

- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes
- thematic categories (such as Canadian poetry) are different from conceptual Categories





- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes







- For each category of a page [Hoffart et al., 2012]
 - Using shallow parsing, determine the head word of the category name. In the example of Canadian poets, the head word is poets.
 - If the head word is in plural, then proposes the category as a class and the article entity as an instance
 - Link the class to the WordNet taxonomy (most frequent sense of the head word in WordNet)
- only countable nouns can appear in plural form
- only countable nouns can be ontological classes







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







- About 100 manually defined relations
 - wasBornOnDate
 - locatedIn
 - hasPopulation
- Categories and infoboxes are exploited to deliver facts (instances of relations).







- Manually defined patterns that map categories and infobox attributes to fact templates
 - infobox attribute born=Montreal, thus wasBornIn(LeonardCohen, Montreal)
- Pattern-based extractions resulted in 2 million extracted entities and 20 million facts





- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







- Based on declarative rules (stored in text files)
- The rules take the form of subject-predicateobject triples, so that they are basically additional facts
- There are different types of rules







Carnegie Mellon University

- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.





Carnegie Mellon University

- Factual rules: definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.





- Factual rules: definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the lorn
 knc
 <
- Re cle
- **Representation** the

article titles, and even other regular elements in the source such as headings, links, or references.







- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.







- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.







- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.







- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wil arti suc

Representation





- **Factual rules:** definition of all relations, their domains and ranges, and the definition of the classes that make up the YAGO hierarchy of literal types.
- Implication rules: express that if certain facts appear in the knowledge base, then another fact shall be added. Horn clause rules.
- **Replacement rules:** for interpreting micro-formats, cleaning up HTML tags, and normalizing numbers.
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.



- Factual rules: definition of all relations, their domains and range Trans
 The TAGE
 The TAGE
- Imp
 kno
 clau
 Bep
 Imp
 in the property of the prop
- Extraction rules: apply primarily to patterns found in the Wikipedia infoboxes, but also to Wikipedia categories, article titles, and even other regular elements in the source such as headings, links, or references.

Federal University of São Carlos

I IST A





- Automatically verifies consistency
 - Check uniqueness of functional arguments
 - spouse(x,y) \land diff(y,z) $\Rightarrow \neg$ spouse(x,z)
 - Check domains and ranges of relations
 - spouse(x,y) \Rightarrow female(x)
 - spouse(x,y) \Rightarrow male(y)
 - spouse(x,y) \Rightarrow (f(x) \land m(y)) \lor (m(x) \land f(y))







- Automatically verifies consistency
 - Check uniqueness of functional arguments
 - spouse(x,y) \land diff(y,z) $\Rightarrow \neg$ spouse(x,z)
 - Check domains and ranges of relations
 - spouse(x,y) \Rightarrow female(x)
 - spouse(x,y) ⇒ male(y)
 - spouse(x,y) \Rightarrow (f(x) \land m(y)) \lor (m(x) \land f(y))

Inference







- Automatically verifies consistency
 - Hard Constraint
- hasAdvisor(x,y) ∧ graduatedInYear(x,t) ∧ graduatedInYear(y,s) ⇒ s < t

Soft Constraint

 firstPaper(x,p) ∧ firstPaper(y,q) ∧ author(p,x) ∧ author(p,y)) ∧ inYear(p) > inYear(q) + 5years ⇒ hasAdvisor(x,y) [0.6]







- Automatically verifies consistency
 - Hard Constraint
- hasAdvisor(x,y) \land graduatedInYear(x,t) \land graduatedInYear(y,s) \Rightarrow s < t

Soft Constraint

 firstPaper(x,p) ∧ firstPaper(y,q) ∧ author(p,x) ∧ author(p,y)) ∧ inYear(p) > inYear(q) + 5years ⇒ hasAdvisor(x,y) [0.6]

Inference



Carnegie Mellon University





- Ontology Representation
 - Entities and Relations of public interest
 - Format: TSV, RDF, XML, N3, Web Interface
 - Learns
 - Instances and patterns from Wikipedia;
 - Taxonomy from WordNet;
 - Geotags information from Geonames.







Carnegie Mellon University

- Named Entity Resolution/Extraction [Theobald & Weikum, 2012]
 - Based on rules and patterns extracted from Wikipedia
 - Disambiguation is a relevant issue
 - Semi-structured data
 - The "Low-Hanging Fruit"
 - Wikipedia infoboxes & categories
 - HMTL lists & tables, etc.





- Named Entity Resolution/Extraction [Theobald & Weikum, 2012]
 - Based on rules and patterns extracted from Wikipedia
 - Natural Language Processing
 - HMTL lists & tables, etc.



Carnegie Mellon University





- Nam Natural Language eobald & Weiku
 - Bai Processing m Wikipedia
 - Disambiguation is a relevant issue

Machine Learning

Carnegie Mellon University



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn? Worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at same te on the same and of Hedeby. The old man drows Blomkvist in by premising sond evidence against Wennerström. Blomkvist ag same pend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af same overing that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, bu Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at te on the same and of Hedeby. The old man drov a Blophovist in by promising solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco uncleOf inted with the men hires the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik. Af same)vering that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

ISC/a

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder. Blomkvist visits Henrik Vanger at te on the same and of Hedeby. same The old man drov a Blorikvist in by much sing solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist the extended Vanger family, most of whom resent beco uncleOf inted with the men hires his presence. He does, however, start a short lived affair with Cecilia, the niece of T enemvOf Af same vering that Salander has hacked into his co affair With persuade ASSIST him with research. They even¹ affairWith lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone sne meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

- ISC/a







- YAGO2: Exploring and Querying World Knowledge in Time, Space, Context, and Many Languages
 - New relations specifically designed to cover time, space and context
 - Wikipedia translated pages as sources for other languages







Carnegie Mellon University

- YAGO3 [Mahdisoltani & Biega & Suchanek, 2015]
 - an extension of the YAGO knowledge base;
 - built from the Wikipedias in multiple languages.
 - fuses the multilingual information with the English WordNet
 - categories, infoboxes, and Wikidata, to learn the meaning of infobox attributes across languages
 - 10 different languages
 - precision of 95%-100% in the attribute mapping
 - enlarges YAGO by 1m new entities and 7m new facts.





- More on YAGO:
 - Very nice tutorials:
 - "<u>Knowledge Bases for Web Content Analytics</u>" at WWW 2015, Florence, May 2015.
 - <u>"Semantic Knowledge Bases from Web Sources</u>" at IJCAI 2011, Barcelona, July 2011 <u>"Harvesting Knowledge from</u> <u>Web Data and Text</u>" at CIKM 2010, Toronto, October 2010 <u>"From Information to Knowledge: Harvesting Entities and</u> <u>Relationships from Web Sources</u>" at PODS 2010, Indianapolis, June 2010
 - Project Website:
 - http://www.mpi-inf.mpg.de/yago-naga/




• More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

Homepage	The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web						
About the Institute	Present Bulling Banda						
Departments:	Research Publications People						
Algorithms & Complexity							
Computer Vision and	The YAGO-NAGA project started in 2006 with the goal of building a conveniently searchable, large-scale, highly accurate						
Multimodal Computing	knowledge base of common facts in a machine-processible representation.						
Computational Biology &							
Applied Algorithmics	We have already harvested knowledge about millions of entities and facts about their relationships, from Wikipedia and WordNet						
Computer Graphics	with careful integration of these two sources. The resulting knowledge base, coined YAGO, has very high precision and is freely						
Databases and Information	available. The facts are represented as RDF triples, and we have developed methods and prototype systems for querying,						
Systems	ranking, and exploring knowledge. Our search engine NAGA provides ranked answers to queries based on statistical models.						
People							
Research	Several interlinked sub-projects are growing on the YAGO-NAGA basis. Our vision is a confluence of Semantic Web						
YAGO-NAGA	(Ontologies), Social Web (Web 2.0), and Statistical Web (Information Extraction) assets towards a comprehensive repository of						
AIDA	human knowledge. Our methodologies combine concepts, models, and algorithms from several fields, including database systems, information retrieval, statistical learning, and logical reasoning.						
AMIE >							
ANGIE							
DEANNA							
EVIN	AIDA						
HYENA	AIDA is a method implemented in an online tool for disambiguating mentions of named entities that occur in natural-language text or Web tables						
JAVATOOLS	AIDA is a method, implemented in an online tool, for disambiguating mentions of named entities that occur in natural-language text of web tables.						
K2							
NAGA							
PATTY	AMIE (Association Rule Mining under Incomplete Evidence in Ontological Knowledge Bases) is a joint project with the Ontologies group.						
PRAVDA							
PROSPERA	ANGIE						
RDF3X							
SOFIE	ANGIE is an active knowledge system for interactive exploration.						
UWN							
YAGO	DEANNA						

FAPESP





• More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

Homepage	The YAGO-NAGA Project: Harvesting, Searching, and Ranking Knowledge from the Web						
About the Institute	Becervela Bublications Becela						
Departments:	Research Publications People						
Algorithms & Complexity							
Computer Vision and	The YAGO-NAGA project started in 2006 with the goal of building a conveniently searchable, large-scale, highly accurate						
Multimodal Computing	knowledge base of common facts in a machine-processible representation.						
Computational Biology &							
Applied Algorithmics	We have already harvested knowledge about millions of entities and facts about their relationships, from Wikipedia and WordNet						
Computer Graphics	with careful integration of these two sources. The resulting knowledge base, coined YAGO, has very high precision and is freely						
Databases and Information	available. The facts are represented as RDF triples, and we have developed methods and prototype systems for querying,						
Systems	ranking, and exploring knowledge. Our search engine NAGA provides ranked answers to queries based on statistical models.						
People							
Research	Several Interlinked sub-projects are growing on the YAGO-NAGA basis. Our vision is a confluence of Semantic Web						
YAGO-NAGA	(Ontologies), Social Web (Web 2.0), and Statistical Web (Information Extraction) assets towards a comprehensive repository of						
AIDA	numan knowledge. Our methodologies combine concepts, models, and algorithms from several fields, including database						
AMIE 🦻	systems, information retrieval, statistical learning, and logical reasoning.						
ANGIE							
DEANNA							
EVIN	AIDA						
HYENA	AIDA is a method, implemented in an online tool, for disambiguating mentions of named entities that occur in natural-language text or Web tables.						
JAVATOOLS	,						
K2	AMTE						
NAGA							
PATTY	AMIE (Association Rule Mining under Incomplete Evidence in Ontological Knowledge Bases) is a joint project with the Ontologies group.						
PRAVDA							
PROSPERA	ANGIE						
RDF3X							
SOFIE	ANGIE is an active knowledge system for interactive exploration.						
UWN							
YAGO	DEANNA						

FAPESP





• More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

							Location	Press	Deutsch
	x planck instit ormatik	ut							C
	ormatik								
25ANNIVERSARY HOME INSTITUT	E NEWS	DEPARTMENTS	PUBLICATIONS	PEOPLE	SOFTWARE	SERVICES			
Departments Databases and Information System	is ▶ Research ▶	YAGO-NAGA > AM	IE						
Algorithms & Complexity	AMIE								
Computer Vision and Multimodal Computing	Association	Rule Mining unde	r Incomplete Evide	ence in Ontol	ogical Knowled	ge Bases			
Computational Biology & Applied Algorithmics	This project is	developed jointly with	the <u>DBWeb team</u> of <u>Tél</u>	écom ParisTech					
Computer Graphics	AMIE is a system the data. For ex	m that extracts supported ample the rule:	d and confident logical ru	les from a knowle	edge base (KB). Log	ical rules encod	e frequent co	orrelations	3 in
Databases and Information Systems	?x <haschild> ?c ?y <haschild> ?c => ?x <ismarriedto> ?y</ismarriedto></haschild></haschild>								
People	states that peop prediction, irreg	le having children in con ularities detection, autom	nmon are frequently marr natic schema generation,	ied. Logical rules ontologies recon	have potential in a l ciliation, etc. AMIE c	broad range of a an mine these p	applications so patterns in m	such as da edium-siz	ata .ed
Research	KBs, several orders of magnitude faster than state-of-the-art approaches to mine logical rules from KBs. The first application of AMIE uses logical rules to address the problem of incompleteness in KBs (particularly web-extracted KBs)								
YAGO-NAGA									
AIDA	Decele								
AMIE	People								
ANGIE	Galárraga, I	_uis	Teflioudi, Christina	Hose	, Katia	Such	nanek, Fabia	n	
DEANNA	0,							_	
EVIN									
HIGGINS	Results								
HYENA									
IBEX	Runtime i	nformation							
Javatools K2	AMIE can extra AMIE, that inclu	ct closed horn rules fror des a set of runtime enh	n medium-sized ontologi ancements. AMIE and A	ies in a few minu MIE+ can sort ar	ites. We report the in the indication of the second s	runtimes for AM ort, head cover	IE+, the late age, standar	est versior d confider	n of nce





More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

AMIE

Association Rule Mining under Incomplete Evidence in Ontological Knowledge Bases

This project is developed jointly with the DBWeb team of Télécom ParisTech.

AMIE is a system that extracts supported and confident logical rules from a knowledge base (KB). Logical rules encode frequent correlations in the data. For example the rule:

?x <hasChild> ?c ?y <hasChild> ?c => ?x <isMarriedTo> ?y

states that people having children in common are frequently married. Logical rules have potential in a broad range of applications such as data prediction, irregularities detection, automatic schema generation, ontologies reconciliation, etc. AMIE can mine these patterns in medium-sized KBs, several orders of magnitude faster than state-of-the-art approaches to mine logical rules from KBs. The first application of AMIE uses logical rules to address the problem of incompleteness in KBs (particularly web-extracted KBs)

?X <hasChild> ?C ?Y <hasChild> ?C => ?X <isMarriedTo> ?Y







More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

AMIE

Associatio

This project

AMIE is a sys

the data. For

?x <hasChild

states that pe

Machine

ent correlations in

ons such as data

Learning in medium-sized prediction, irre KBs, several orders of magnitude faster than state-of-the-art approaches to mine logical rules from KBs. The first application of AMIE uses logical rules to address the problem of incompleteness in KBs (particularly web-extracted KBs)

?X <hasChild> ?C ?Y <hasChild> ?C => ?X <isMarriedTo> ?Y







• More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

AMIE

Associatio

This project

AMIE is a sys

the data. For

?x <hasChild

states that pe

Machine

ent correlations in

Learning ons such as data prediction, irre in medium-sized KBs, several orders of magnitude faster than state-of-the-art approaches to mine logical rules from KBs. The first application of AMIE uses logical rules to address the problem of incompleteness in KBs (particularly web-extracted KBs)

?X <hasChild> ?C ?Y <hasChild> ?C => ?X <isMarriedTo> ?Y

Inference







• More on YAGO (http://www.mpi-inf.mpg.de/yago-naga/)

* YAGO3 is out, and combines information from Wikipedias in 10 different languages: https://suchanek.name/work/publications/cidr2015.pdf

* Automatically find rules such as "if you are married then you live in the same city as your spouse" in YAGO: http://www.mpi-inf.mpg.de/departments/databases-and-informationsystems/research/yago-naga/amie/

* Plans to add commercial products to YAGO: http://www.mpi-inf.mpg.de/departments/databases-and-informationsystems/research/yago-naga/ibex/







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - -KnowItAll
 - NELL





W

KnowItAll



UNIVERSITY OF WASHINGTON | ABOUT US | CONTACT US | MY CSE | INTERNAL

Computer Science & Engineering





UNIVERSITY of WASHINGTON

News & Events	People	Education	Research	Current Students	Prospective Students	Faculty Candidates	Alumni	Industry Affiliates	Support CSE

KnowltAll



How can a computer accumulate a massive body of knowledge?
 What will Web search engines look like in ten years?

To address the questions above, the KnowItAll project has been developing a variety of domain-independent systems that extract information from the Web in an autonomous, scalable manner.





UNIVERSITY OF WASHINGTON | ABOUT US | CONTACT US | MY CSE | INTERNAL

Computer Science & Engineering

UNIVERSITY of WASHINGTON

W



Open Information Extraction



How can a computer accumulate a massive body of knowledge? What will Web search engines look like in ten years?

To address these questions, the Open IE project has been developing a Web-scale information extraction system that reads arbitrary text from any domain on the Web, extracts meaningful information, and stores it in a unified knowledge base for efficient querying. In contrast to traditional information extraction, the Open Information Extraction paradigm attempts to overcome the knowledge acquisition bottleneck by extracting a large number of relations at once.

Demo: <u>TextRunner</u> extracted over 500,000,000 assertions from 100 million Web pages. Software: <u>ReVerb</u> Open Information Extraction Software and additional information. Data: <u>Horn-clause</u> inference rules learned by the Sherlock system. Demo: <u>Selectional Preferences from Web Text</u> compute admissible argument values for a relation.

Data: <u>10,000 Functional Relations learned from Web Text</u> predict the functionality of a phrase.





Q

KnowItAII: Open Information Extraction



UFE CAT

Turing Center

Investigating problems at the crossroads of natural language processing, data mining, Web search, and the Semantic Web.



	HOME	PEOPLE	PUBLICATIONS	PRESS	EVENTS	JOBS	CONTACT	
				Research Sc Niranian Bala	subramanian			
				Prof. Mausam				
				Michael Schm	iitz			
				Dr. Stephen S	Soderland			
People				Affiliated Fa	culty:			
copie				Prof. Emily Be	ender			
Director:				Prof. Jeff Bilm	<u>ies</u>			
Prof. Oren Etzioni				Prof. Pedro D	omingos irchhoff			
Advisory Board:				Prof. Mari Ost	rendorf			
Prof. Emily Bender				Prof. Dan We	ld			
Prof. David Notkin					-			
Dr. Jonathan Pool				Students:				
Assistant Director:				Tony Fader	ensen			
Patrick Allen				Jeff Huang				
				Tom Lin				
Events Committee:				Alan Ritter				
Christie Evans				Stefan Schoe	nmackers			
Dr. Jonathan Pool				Fei Wu				
				Alumni:				
				Michele Bank	o (Microsoft, 2009)			
				Mike Cafarella	(University of Mich	igan, 2009)	00)	
				Michael Skipp	per (Google 2008)	n University, 20	08)	
				Dr. Ana-Maria	Popescu (Yahoo. 2	007)		
				Kobi Reiter (C	Google, 2007)	,		
				Prof. Alex Yat	es (Temple Universi	ty, 2007)		







KnowItAll

- Motivation: New Paradigm for Search [Etzioni, 2008]
 - The future of Web Search
 - Read the Web instead of retrieving Web pages to perform Web Search







Carnegie Mellon University

KnowItAll

- Information Extraction (IE) + tractable inference
 - IE(sentence) = who did what?
 - speaker(Sheila Mcllraith, IJCAI-16)
 - Inference = uncover implicit information
 - Will Pittsburgh Steelers be champions again?
- Open Information Extraction [Banko et al., 2007]





- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)





[Banko et al., 2007]

Open IE systems avoid specific nouns and verbs

Federal University of São Carl

- Extractors are unlexicalized—formulated only in terms of:
 - syntactic tokens (e.g., part-of-speech tags)
 - closed-word classes (e.g., of, in, such as).
- Open IE extractors focus on generic ways in which relationships are expressed in English
 - naturally generalizing across domains.



Open Information Extraction [Banko et al., 2007]

 Open IE extractors focus on generic ways in which relationships are expressed in English

naturally generalizing across domains.





Federal University of São Carlos



- Open IE systems are traditionally based on three steps [Etzioni et al., 2011]:
 - **1. Label:** Sentences are automatically labeled with extractions using heuristics or distant supervision.



Federal University of São Carlo



• Open IE systems are traditionally based on three steps [Etzioni et al., 2011]:

Federal University of São Carlo

IST /ar

- **1. Label:** Sentences are automatically labeled with extractions using heuristics or distant supervision.
- 2. Learn: A relation phrase extractor is learned using a sequence-labeling graphical model (e.g., CRF).





- Open IE systems are traditionally based on three steps [Etzioni et al., 2011]:
 - 1. Label: Sentences are automatically labeled with
 - ext - 2. L Inference sequence-labeling graphical model (e.g., CRF).
 - 3. Extract: given a sentence as input, identifies a candidate pair of NP arguments (Arg1, Arg2) from the sentence, and then uses the learned extractor to label each word between the two arguments as part of the relation phrase or not.



Federal University of São Car



- TextRunner [Banko et al., 2007] was the first OIE system;
- OIE became the main focus of the KnowItAll project;
- Two main problems:

Federal University of São Carl

- incoherent extractions;
- uninformative relations



incoherent extractions

Federal University of São Carlos

U-ISI /an

Sentence	Incoherent Relation	
The guide contains dead links	contains omits	
and omits sites.		
The Mark 14 was central to the	was central torpedo	
torpedo scandal of the fleet.	_	
They recalled that Nungesser	recalled began	
began his career as a precinct		
leader.		





uninformative relations

Federal University of São Carlos

UFISC/ar

is	is an album by, is the author of, is a city in
has	has a population of, has a Ph.D. in, has a cameo in
made	made a deal with, made a promise to
took	took place in, took control over, took advantage of
gave	gave birth to, gave a talk at, gave new meaning to
got	got tickets to see, got a deal on, got funding from







TextRunner was based on

Federal University of São Carlos

ufisio/ar

		Simplified
Relative		Lexico-Syntactic
Frequency	Category	Pattern
37.8	Verb	E ₁ Verb E ₂
		X established Y
22.8	Noun+Prep	E ₁ NP Prep E ₂
		X settlement with Y
16.0	Verb+Prep	E ₁ Verb Prep E ₂
		X moved to Y
9.4	Infinitive	E_1 to Verb E_2
		X plans to acquire Y
5.2	Modifier	E ₁ Verb E ₂ Noun
		X is Y winner
1.8	Coordinate _n	E_1 (and , - :) E_2 NP
		X-Y deal
1.0	$Coordinate_v$	E_1 (and ,) E_2 Verb
		X, Y merge
0.8	Appositive	$E_1 \text{ NP (: ,)? } E_2$
		X hometown : Y
	-	





Carnegie Mellon University

OIE: the second generation

• New syntactic constraint based on POS tag patterns

```
V | VP | VW^*P

V = \text{verb particle? adv?}

W = (\text{noun} | \text{adj} | \text{adv} | \text{pron} | \text{det})

P = (\text{prep} | \text{particle} | \text{inf. marker})
```

- simple verb phrase (e.g., invented)
- verb phrase followed immediately by a preposition or particle (e.g., located in)
- verb phrase followed by a simple noun phrase and ending in a preposition or particle (e.g., has atomic weight of)
- multiple possible matches, then the longest possible match is chosen.





- New lexical constraint to separate valid relation phrases from over-specified relation phrases
- The lexical constraint is based on the intuition that a valid relation phrase should take many distinct arguments in a large corpus.





- New OIE System: ReVerb [Fader et al., 2011]
 - Input: a POS-tagged and NP-chunked sentence
 - Output: a set of (x,r,y) extraction triples
 - Based on two extraction algorithm:
 - 1. Relation Extraction: based on the new constraints
 - 2. Argument Extraction: For each relation phrase r identified in Step 1, find the nearest noun phrase x to the left and the nearest noun phrase y to the right of r in s.



Federal University of São Carlo





• New OIE System: ReVerb [Fader et al., 2011]



Figure 2: REVERB identifies correct relation phrases with substantially higher precision and recall than state-of-the-art open extractors, including WOE^{parse} that uses patterns learned over dependency parse paths.







Carnegie Mellon University

OIE: the second generation

Open Information Extraction				
Argument 1: porto	Relation:	Argument 2:	All 1 Q Search	

286 answers from 886 sentences (cached)

Do you mean...

Porto (78)



Porto (Portuguese pronunciation: [portu]), also known as Oporto in English, is the second largest city in Portugal, after Lisbon, and one of the major urban areas in Southern Europe. Its administrative limits (an area of 41.66 km²/16 sq.mi) include a population of 237,584 (2011) inhabitants distributed within 15 civil parishes. The urbanized area of Porto, which extends beyond the administrative limits of the city, has a population of 1.3 million (2011) in an area of 389 km (150 sq mi),...read more

F.C. Porto (14)

Futebol Clube do Porto MH IH MH OM (Portuguese pronunciation: [futeboł 'klube du 'portu]) (Euronext: FCP), commonly known as FC Porto, Porto, or FCP, is a Portuguese multi-sports club from the city of Porto, in the northern region of the country. Although they successfully compete in a number of different sports, FC Porto is mostly known for its association football team. Founded in Porto in September 28, 1893, it is one of the "Três Grandes" (Big Three, in English), football clubs in...read more



Table extracted from [Etzioni et al., 2011]

Federal University of São Carlos

SC/ar

Category	Potterns	Frequency Arg1	Frequency Arg?
Division		Frequency Argi	Frequency Arg2
Basic NP	NN, JJ NN, etc	65%	60%
		Chicago was founded in 1833.	Calcium prevents osteoporosis.
Prepositional	NP PP ⁺	19%	18%
Attachments		The forest in Brazil	Lake Michigan is one of the five
		is threatened by ranching.	Great Lakes of North America.
List	NP (,NP)* ,?	15%	15%
	and/or NP	Google and Apple are	A galaxy consists of
		headquartered in Silicon Valley.	stars and stellar remnants.
Independent	(that WP WDT)?	0%	8%
Clause	NP VP NP	Google will acquire YouTube,	Scientists estimate that 80%
		announced the New York Times.	of oil remains a threat.
Relative	NP (that WP WDT)	<1%	6%
Clause	VP NP?	Chicago, which is located in Illinois,	Most galaxies appear to be
		has three million residents.	dwarf galaxies, which are small.

Table 4: Taxonomy of arguments for binary relationships. In each sentence, the argument is bolded and the relational phrase is italicized. Multiple patterns can appear in a single argument so percentages do not need to add to 100. In the interest of space, we omit argument structures that appear in less than 5% of extractions. Upper case abbreviations represent noun phrase chunk abbreviations and part-of-speech abbreviations.







• New OIE System: ArgLearner [Etzioni et al., 2011]



Figure 4: ARGLEARNER's system architecture.





• New OIE System:

Federal University of São Carlo

• ReVerb + ArgLearner = R2A2 [Etzioni et al., 2011]



REVERB.





• New OIE System:

Federal University of São Carlos

IS CA

• ReVerb + ArgLearner = R2A2 [Etzioni et al., 2011]

Free text Hearst-patterns; clustering by verbal phrases Natural-language processing Advanced patterns & iterative bootstrapping ("Dual Iterative Pattern Relation Extraction")

Fi Reverd.





It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

This slide was adapted from [Hady et al., 2011]



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn? Worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

This slide was adapted from [Hady et al., 2011]

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at same te on the same and of Hedeby. The old man drows Blomkvist in by premising sond evidence against Wennerström. Blomkvist ag same pend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af same overing that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, bu Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

This slide was adapted from [Hady et al., 2011]
It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at te on the same and of Hedeby. The old man drov a Blophovist in by promising solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco uncleOf inted with the men hires the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik. Af same)vering that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

This slide was adapted from [Hady et al., 2011]

Federal University of São Carlo

ISC/a

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder. Blomkvist visits Henrik Vanger at te on the same and of Hedeby. same The old man drov a Blorikvist in by much sing solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist the extended Vanger family, most of whom resent beco uncleOf inted with the men hires his presence. He does, however, start a short lived affair with Cecilia, the niece of T enemvOf Af same vering that Salander has hacked into his co affair With persuade ASSIST him with research. They even¹ affairWith lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone sne meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

This slide was adapted from [Hady et al., 2011]



It's about the disap scion of one of the determined to know **Blomkvist** visits He The old man dray of Blomkvist ag same assignment: the dis home to several ger beco uncleOf inte his presence. He do Af same)vering t him with research. to Lisbeth who trea discover that Harri

ance forty years ago of Harriet V hiest families in Sweden, and a truth about what he believ s her Vanger at and te un he same akvist in by memory evidence er family hi nd 2 riting th Harriet some wns aran rt owners in Vai ions of e extended Vange h the me t lived affair with owever, sta c affairW alander lovers, bat Blomk even¹ affairWith

oung cle, der. edeby. nst Wennerström. as a cover for the real vears earlier. Hedeby is Enterprises<mark>.</mark> Bl<mark>omkvist</mark> ily, most of whom resent ia, the niece 💕 🎹 enemyOf bersuade **assis**i same has trouble getting close Itimately the two is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

This slide was adapted from [Hady et al., 2011]







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)







More on KnowItAll

• Open IE is now included in the Stanford CoreNLP package: http://nlp.stanford.edu/software/openie.shtml (ACL 2015 paper is linked to from this page).







More on KnowltAll

• Open IE is now included in the Stanford CoreNLP package: http://nlp.stanford.edu/software/openie.shtml (ACL 2015 paper is linked to from this page).



The Stanford Natural Language Processing Group

people publications research blog software teaching local

Software > Stanford OpenIE

Stanford Open Information Extraction

About | Download | Usage | Support | Questions | Release history |

About

Open information extraction (open IE) refers to the extraction of structured relation triples from plain text, such that the *schema* for these relations does not need to be specified in advance. For example, *Barack Obama was born in Hawaii* would create a triple (Barack Obama; was born in; Hawaii), corresponding to the open domain relation "was born in". This software is a Java implementation of an open IE system as described in the paper:

Gabor Angeli, Melvin Johnson Premkumar, and Christopher D. Manning. Leveraging Linguistic Structure For Open Domain Information Extraction. In *Proceedings of the Association of Computational Linguistics (ACL)*, 2015.

The system first splits each sentence into a set of entailed clauses. Each clause is then maximally shortened, producing a set of entailed shorter sentence fragments. These fragments are then segmented into OpenIE triples, and output by the system. An illustration of the process is given for an example sentence below:











• Open IE 4.0: http://knowitall.github.io/openie/

W Getting Started		ç o			
Open IE		Download ZIP File	Download TAR Ball	View On GitHub	
Quality information extraction at web scale.					
View the Project on GitHub knowitall/openie					
Description					
Open IE is an open information extractor that find relationships in sent	ence. Here is an exa	ample sentenc	e and extraction	IS:	
U.S. president Barack Obama gave his inaugural address on January 20,	2013.				
(Barack Obama; is president of; U.S.) (Barack Obama; gave; [his inaugural address, on Janua	nry 20, 2013])				
In this example, the second part of the tuple is the relationship betweer	the first argument	t and the seco	nd arguments.		
	ESD		Carn	egie Mellon Un	





More on KnowItAll

 Open IE in Chinese: https://www.semanticscholar.org/paper/19f83e125aa638419



Figure 1: The parsed tree of a Chinese sentence.







More on KnowltAll

• Open IE in Chinese:

https://www.semanticscholar.org/paper/19f83e125aa638419 7b7a1fb138f0de531ab73eb

Chinese	Open IE	Precision	Recall	F1
Exact	ReVerb	0.5820	0.0987	0.1688
Match	CORE	0.5579	0.3291	0.4140
Relation	ReVerb	0.8361	0.1425	0.2435
Only	CORE	0.8463	0.5000	0.6286

Table 1: Performance evaluation on Chinese Open IE.







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll
 - NELL







Outline

- Machine Learning
- Machine Reading
- Reading the Web
 - DBPedia
 - YAGO
 - KnowItAll









Never-Ending Language Learner



Joint work with Carnegie Mellon Read The Web Project Group (http://rtw.ml.cmu.edu/rtw/) and MaLL (Machine Learning Lab) from Federal University of São Carlos (http://www.dc.ufscar.br/MaLL/MaLL.html)







Humans learn many things, for years, and become better learners over time

Why not machines?







Never-Ending Learning

We'll never really understand learning until we build machines that

- learn many different things,
- over years,
- and become better <u>learners</u> over time.







Never-Ending Learning

- We'll never produce natural language <u>understanding</u> systems until we have systems that react to arbitrary sentences by saying one of:
- I understand, and already knew that
- I understand, and didn't know, but accept it
- I understand, and disagree because ...





NELL: Never-Ending Language Learner

•Inputs:

Federal University of São Carl

- initial ontology
- handful of examples of each predicate in ontology
- the web
- occasional interaction with human trainers

•The task:

- run 24x7, forever
- -• each day:
- extract more facts from the web to populate the initial ontology
- 2. learn to read (perform #1) better than yesterday





• Goal:

Federal University of São Carlo

ISC/ar

- • run 24x7, forever
- • each day:
- 1. extract more facts from the web to populate given ontology
- 2. learn to read better than yesterday
- Today...
- Running 24 x 7, since January, 2010
- Input:
- ontology defining ~800 categories and relations
- • 10-20 seed examples of each
- 1 billion web pages (ClueWeb Jamie Callan)
- Result:
- continuously growing KB with +90.000,000 extracted beliefs (different levels of confidence)





E

Se



http://rtw.ml.cmu.edu

Read the Web

Research Project at Carnegie Mellon University

Home	Project Overview	Resources & Data	Publications	People	

NELL: Never-Ending Language Learning

Can computers learn to read? We think so. "Read the Web" is a research project that attempts to create a computer system that learns over time to read the web. Since January 2010, our computer system called NELL (Never-Ending Language Learner) has been running continuously, attempting to perform two tasks each day:

 First, it attempts to "read," or extract facts from text found in hundreds of millions of web pages (e.g., playsInstrument(George_Harrison, guitar)).



Browse the Knowledge Base!

Carnegie Mellon University

• Second, it attempts to improve its reading competence, so that tomorrow it can extract more facts from the web, more accurately.

So far, NELL has accumulated over 15 million candidate beliefs by reading the web, and it is considering these at different levels of confidence. NELL has high confidence in 1,471,011 of these beliefs — these are displayed on this website. It is not perfect, but NELL is learning. You can track NELL's progress below or <u>@cmunell on Twitter</u>, browse and download its <u>knowledge base</u>, read more about our <u>technical approach</u>, or join the <u>discussion group</u>.

NELL: Never-Ending Language Learner

Recently-Learned Facts Lewitter

date learned confidence iteration instance 100.0 🐊 🖏 chase bank is a credit union 10-apr-2016 989 91.4 2 8 unnamed subpeak of long valley peak is a mountain 27-apr-2016 991 98.9 28 27 currency funding is an economic sector 989 10-apr-2016 90.5 🖄 🖏 10-apr-2016 snow mountains grassland melomys is a mammal 989 91.3 🍰 🖏 10-apr-2016 the paramus post is a blog 989 100.0 🍃 🖏 15-may-2016 <u>aashto</u> is <u>headquartered</u> in the country <u>u</u> s 994 93.8 🖾 🖏 division is an organization also known as us department 02-may-2016 992 100.0 🖄 🖏 dome is an attraction made of bronze 15-may-2016 994 100.0 🖄 🖏 15-may-2016 jeff bezos is the CEO of amazon 994 99.6 🍰 ኛ tomato is an agricultural product that is usually cooked with salt 27-apr-2016 991



Refresh









Learning to Read and Learning by Reading

1. Classify noun phrases (NP's) by category







Paris Pittsburgh Seattle Cupertino







Paris Pittsburgh Seattle Cupertino



mayor of arg1 live in arg1









mayor of arg1 live in arg1





















































Key Idea 1: Coupled semi-supervised training of many functions



Federal University of São Carlos

ISC/ar



hard (underconstrained) semi-supervised learning problem much easier (more constrained) semi-supervised learning problem







Coupled Training Type 1: Co-training, Multiview, Coregularization [Blum & Mitchell; 98]

[Blum & Mitchell; 98] [Dasgupta et al; 01] [Ganchev et al., 08] [Sridharan & Kakade, 08] [Wang & Zhou, ICML10]





Coupled Training Type 1: Co-training, Multiview, Coregularization [Blum & Mitchell; 98]

[Dasgupta et al; 01] [Ganchev et al., 08] [Sridharan & Kakade, 08] [Wang & Zhou, ICML10]



Constraint: $f_1(x_1) = f_2(x_2)$





Coupled Training Type 1: Co-training, Multiview, Coregularization [Blum & Mitchell; 98]

[Dasgupta et al; 01] [Ganchev et al., 08] [Sridharan & Kakade, 08] [Wang & Zhou, ICML10]

If f₁, f₂ PAC learnable, X₁, X₂ conditionally indep Then PAC learnable from <u>unlabeled</u> data and weak initial learner

Constraint: $f_1(x_1) = f_2(x_2)$

and disagreement between f₁, f₂ bounds error of each







Type 1 Coupling Constraints in NELL





Federal University of São Carlos

UFIST AT



Coupled Training Type 2:

Structured Outputs, Multitask, Posterior Regularization, Multilabel

Learn functions with the same input, different outputs, where we know some constraint

[Daume, 2008] [Bakhir et al., eds. 2007] [Roth et al., 2008] [Taskar et al., 2009] [Carlson et al., 2009]




Coupled Training Type 2:

Structured Outputs, Multitask, Posterior Regularization, Multilabel

Learn functions with the same input, different outputs, where we know some constraint

[Daume, 2008] [Bakhir et al., eds. 2007] [Roth et al., 2008] [Taskar et al., 2009] [Carlson et al., 2009]



Constraint: $\Phi(f_1(x), f_2(x))$





Coupled Training Type 2:

Structured Outputs, Multitask, Posterior Regularization, Multilabel

Learn functions with the same input, different outputs, where we know some constraint

[Daume, 2008] [Bakhir et al., eds. 2007] [Roth et al., 2008] [Taskar et al., 2009] [Carlson et al., 2009]



Effectiveness ~ probability that $\Phi(Y_1, Y_2)$ will be violated by incorrect f_i and f_k

Constraint: $\Phi(f_1(x), f_2(x))$







Type 2 Coupling Constraints in NELL







Multi-view, Multi-Task Coupling





Federal University of São Carlos

u istolati





Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation







Learning Relations between NP's









Learning Relations between NP's









Type 3 Coupling: Argument Types Constraint: f3(x1,x2) → (f1(x1) AND f2(x2))



playsSport(NP1,NP2) → athlete(NP1), sport(NP2)





Pure EM Approach to Coupled Training



- E: jointly estimate latent labels for each function of each unlabeled example
- M: retrain all functions, based on these probabilistic labels

Scaling problem:

- E step: 20M NP's, 1014 NP pairs to label
- **M** step: 50M text contexts to consider for each function \rightarrow 10¹⁰ parameters to retrain
- even more URL-HTML contexts..





NELL's Approximation to EM

E' step:

- Consider only a growing subset of the latent variable assignments
- category variables: up to 250 NP's per category per iteration
- relation variables: add only if confident and args of correct type
- this set of explicit latent assignments *IS* the knowledge base

M' step:

- Each view-based learner retrains itself from the updated KB
- "context" methods create growing subsets of contexts







NELL Architecture







Never-Ending Language Learning

arg1_was_playing_arg2_arg2_megastar_arg1_arg2_icons_arg1 arg2 player named arg1 arg2 prodigy arg1 arg1 is the tiger woods of arg2 arg2 career of arg1 arg2_greats_as_arg1 arg1_plays_arg2 arg2_player_is_arg1 arg2 legends arg1 arg1 announced his retirement from arg2 arg2 operations chief arg1 arg2 player like arg1 arg2_and_golfing_personalities_including_arg1_arg2_players_like_arg1 arg2 greats like arg1 arg2 players are steffi graf and arg1 arg2 great arg1 arg2 champ arg1 arg2 greats such as arg1 arg2_professionals_such_as_arg1_arg2_hit_by_arg1 arg2_greats_arg1 arg2_icon_arg1_arg2_stars_like_arg1_arg2_pros_like_arg1 arg1 retires from arg2 arg2 phenom arg1 arg2 lesson from arg1 arg2 architects robert trent jones and arg1 arg2 sensation arg1 arg2 pros arg1 arg2 stars venus and arg1 arg2 hall of famer arg1 arg2 superstar arg1 arg2 legend arg1 arg2 legends such as arg1 arg2_players_is_arg1_arg2_pro_arg1_arg2_player_was_arg1 arg2_god_arg1_arg2_idol_arg1_arg1_was_born_to_play_arg2 arg2_star_arg1_arg2_hero_arg1 arg2_players_are_arg1 arg1_retired_from_professional_arg2_arg2_legends_as_arg1 arg2 autographed by arg1 arg2 champion arg1



Predicate	Feature	Weight
mountain	LAST=peak	1.791
mountain	LAST=mountain	1.093
mountain	FIRST=mountain	-0.875
musicArtist	LAST=band	1.853
musicArtist	POS=DT_NNS	1.412
musicArtist	POS=DT_JJ_NN	-0.807
newspaper	LAST=sun	1.330
newspaper	LAST=university	-0.318
newspaper	POS=NN_NNS	-0.798
university	LAST=college	2.076
university	PREFIX=uc	1.999
university	LAST=state	1.992
university	LAST=university	1.745
university	FIRST=college	-1.381
visualArtMovement	SUFFIX=ism	1.282
visualArtMovement	PREFIX=journ	-0.234
visualArtMovement	PREFIX=budd	-0.253

Predicate	Web URL	Extraction Template
academicField athlete bird bookAuthor	http://scholendow.ais.msu.edu/student/ScholSearch.Asp http://www.quotes-search.com/d_occupation.aspx?o=+athlete http://www.michaelforsberg.com/stock.html http://lifebehindthecurve.com/	<pre> [X] - - <option>[X]</option> [X] by [Y] –</pre>





Never-Ending Language Learning

arg1_was_playing_arg2 arg2_megastar_arg1 arg2_icons arg1 arg2 player named arg1 arg2 prodigy arg1 arg1 is the tiger woods of arg2 arg2 career of arg1 arg2_greats_as_arg1 arg1_plays_arg2 arg2_player_is_arg1 arg2 legends arg1 arg1 announced his retirement from arg2 arg2 operations chief arg1 arg2 player like arg1 arg2_and_golfing_personalities_including_arg1_arg2_players_like_arg1 arg2 greats like arg1 arg2 players are steffi graf and arg1 arg2 great arg1 arg2 champ arg1 arg2 greats such as arg1 arg2_professionals_such_as_arg1_arg2_hit_by_arg1 arg2_greats_arg1 arg2_icon_arg1_arg2_stars_like_arg1_arg2_pros_like_arg1 arg1 retires from arg2 arg2 phenom arg1 arg2 lesson from arg1 arg2 architects robert trent jones and arg1 arg2 sensation arg1 arg2 pros arg1 arg2 stars venus and arg1 arg2 hall of famer arg1 arg2 superstar arg1 arg2 legend arg1 arg2 legends such as arg1 arg2_players_is_arg1_arg2_pro_arg1_arg2_player_was_arg1 arg2_god_arg1_arg2_idol_arg1_arg1_was_born_to_play_arg2 arg2_star_arg1_arg2_hero_arg1 arg2_players_are_arg1 arg1_retired_from_professional_arg2_arg2_legends_as_arg1 arg2 autographed by arg1 arg2 champion arg1



Predicate	Feature	Weight
mountain	LAST=peak	1.791
mountain	LAST=mountain	1.093
mountain	FIRST=mountain	-0.875
musicArtist	LAST=band	1.853
musicArtist	POS=DT_NNS	1.412
musicArtist	POS=DT_JJ_NN	-0.807
newspaper	LAST=sun	1.330
newspaper	LAST=university	-0.318
newspaper	POS=NN_NNS	-0.798
university	LAST=college	2.076
university	PREFIX=uc	1.999
university	LAST=state	1.992
university	LAST=university	1.745
university	FIRST=college	-1.381
visualArtMovement	SUFFIX=ism	1.282
visualArtMovement	PREFIX=journ	-0.234
visualArtMovement	PREFIX=budd	-0.253

Predicate	Web URL	Extraction Template
academicField athlete bird bookAuthor	http://scholendow.ais.msu.edu/student/ScholSearch.Asp http://www.quotes-search.com/d_occupation.aspx?o=+athlete http://www.michaelforsberg.com/stock.html http://lifebehindthecurve.com/	<pre> [X] - - <option>[X]</option> [X] by [Y] –</pre>





Never-Ending Language Learning

arg1_was_playing_arg2 arg2_megastar_arg1 arg2_icons arg1 arg2 player named arg1 arg2 prodigy arg1 arg1 is the tiger woods of arg2 arg2 career of arg1 arg2_greats_as_arg1 arg1_plays_arg2 arg2_player_is_arg1 arg2 legends arg1 arg1 announced his retirement from arg2 arg2 operations chief arg1 arg2 player like arg1 arg2_and_golfing_personalities_including_arg1_arg2_players_like_arg1 arg2 greats like arg1 arg2 players are steffi graf and arg1 arg2 great arg1 arg2 champ arg1 arg2 greats such as arg1 arg2_professionals_such_as_arg1_arg2_hit_by_arg1 arg2_greats_arg1 arg2_icon_arg1_arg2_stars_like_arg1_arg2_pros_like_arg1 arg1 retires from arg2 arg2 phenom arg1 arg2 lesson from arg1 arg2 architects robert trent jones and arg1 arg2 sensation arg1 arg2 pros arg1 arg2 stars venus and arg1 arg2 hall of famer arg1 arg2 superstar arg1 arg2 legend arg1 arg2 legends such as arg1 arg2_players_is_arg1_arg2_pro_arg1_arg2_player_was_arg1 arg2_god_arg1_arg2_idol_arg1_arg1_was_born_to_play_arg2 arg2_star_arg1_arg2_hero_arg1 arg2_players_are_arg1 arg1_retired_from_professional_arg2_arg2_legends_as_arg1 arg2 autographed by arg1 arg2 champion arg1



Predicate	Feature	Weight
mountain	LAST=peak	1.791
mountain	LAST=mountain	1.093
mountain	FIRST=mountain	-0.875
musicArtist	LAST=band	1.853
musicArtist	POS=DT_NNS	1.412
musicArtist	POS=DT_JJ_NN	-0.807
newspaper	LAST=sun	1.330
newspaper	LAST=university	-0.318
newspaper	POS=NN_NNS	-0.798
university	LAST=college	2.076
university	PREFIX=uc	1.999
university	LAST=state	1.992
university	LAST=university	1.745
university	FIRST=college	-1.381
visualArtMovement	SUFFIX=ism	1.282
visualArtMovement	PREFIX=journ	-0.234
visualArtMovement	PREFIX=budd	-0.253

Predicate	Web URL	Extraction Template
academicField	http://scholendow.ais.msu.edu/student/ScholSearch.Asp	[X] -
athlete	http://www.quotes-search.com/d_occupation.aspx?o=+athlete	-
bird	http://www.michaelforsberg.com/stock.html	<option>[X]</option>
bookAuthor	http://lifebehindthecurve.com/	[X] by [Y] –







If coupled learning is the key idea, how can we get new coupling constraints?







Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances







- first order, probabilistic horn clause constraints
- 0.93 athletePlaysSport(?x,?y) :- athletePlaysForTeam(?x,?z), teamPlaysSport(?z,?y)
- connects previously uncoupled relation predicates
- – infers new beliefs for KB



Federal University of São Carlo





- 0.95 athletePlaysSport(?x,basketball) :- athleteInLeague(?x,NBA)
- 0.93 athletePlaysSport(?x,?y) :- athletePlaysForTeam(?x,?z)
- teamPlaysSport(?z,?y)
- 0.91 teamPlaysInLeague(?x,NHL) :- teamWonTrophy(?x,Stanley_Cup)
- 0.90 athleteInLeague(?x,?y):- athletePlaysForTeam(?x,?z), teamPlaysInLeague(?z,?y)
- 0.88 cityInState(?x,?y) :- cityCapitalOfState(?x,?y), cityInCountry(?y,USA)
- 0.62* newspaperInCity(?x,New_York) :- companyEconomicSector(?x,media), generalizations(?x,blog)



Federal University of São Carlo

IST /ar





Learned Probabilistic Horn Clause Rules









Learned Probabilistic Horn Clause Rules

0.93 playsSport(?x,?y) ← playsForTeam(?x,?z), teamPlaysSport(?z,?y)









NELL Architecture







Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts





Distinguish Text Tokens from Entities

[Jayant Krishnamurthy]



Coreference Resolution:

- Co-train classifier to predict coreference as f(string similarity, extracted beliefs)
- Small amount of supervision: ~10 labeled coreference decisions
- Cluster tokens using f as similarity measure



Federal University of São Carlos

u istolati



Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology







Key Idea 3: Automatically Extending the Ontology















































• Find open triangles in the Graph















• open triangles









• open triangles









• open triangles








• open triangles









• open triangles









• open triangles









• open triangles









• open triangles









• Compute the number of common neighbors $\aleph(u, w)$









- Compute the number of common neighbors $\aleph(u, w)$
- Compute the cumulative number of instances for the categories nodes Σ $\aleph(u, w)$









- Compute the number of common neighbors $\aleph(u, w)$
- Compute the cumulative number of instances for the categories nodes Σ $\aleph(u, w)$
- $N_{\Lambda_c(u_c,w_c)}$ is the number of open triangles for categories u and w.







$$\aleph_c(u_c,w_c) = \sum lpha(u,w) - N_{\Lambda_c(u_c,w_c)}$$







$$\aleph_c(u_c,w_c) = \sum leph(u,w) - N_{\Lambda_c(u_c,w_c)}$$

- If $\aleph_c(u_c, w_c) > \xi$ then create the new relation
- $\xi = 10$ (empirically)







$$\aleph_c(u_c, w_c) = \sum \aleph(u, w) - N_{\Lambda_c(u_c, w_c)}$$

- If $\aleph_c(u_c, w_c) > \xi$ then create the new relation
- $\xi = 10$ (empirically)
- Name the new relation based on text







OntExt

Contexts/ Contexts	may cause	can cause	can lead to	to treat	for treatment of	medication
may cause	0.176	0.074	0.030	0.015	0.011	0.000
can cause	0.051	0.150	0.039	0.018	0.013	0.010
can lead to	0.034	0.064	0.189	0.019	0.021	0.018
to treat	0.006	0.011	0.007	0.109	0.043	0.015
for treatment of	0.005	0.008	0.008	0.045	0.086	0.023
medication	0.000	0.011	0.009	0.030	0.036	0.111



FAPESP



NELL: sample of self-added relations

- athleteWonAward
- animalEatsFood

Federal University of São Carlos

uferen

- languageTaughtInCity
- clothingMadeFromPlant
- beverageServedWithFood
- fishServedWithFood
- athleteBeatAthlete
- athleteInjuredBodyPart
- arthropodFeedsOnInsect
- animalEatsVegetable
- plantRepresentsEmotion
- foodDecreasesRiskOfDisease

- clothingGoesWithClothing
- bacteriaCausesPhysCondition
- buildingMadeOfMaterial
- emotionAssociatedWithDisease
- foodCanCauseDisease
- agriculturalProductAttractsInsect
- arteryArisesFromArtery
- countryHasSportsFans
- bakedGoodServedWithBeverage
- beverageContainsProtein
- animalCanDevelopDisease
- beverageMadeFromBeverage



Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology
- 6. Learn to infer relation instances via targeted random walks







[Lao, Mitchell, Cohen, EMNLP 2011]

Pittsburgh

<u>Feature = Typed Path</u> CityInState, CityInstate⁻¹, CityLocatedInCountry

Feature Value 0.32

Logistic Regresssion Weight







[Lao, Mitchell, Cohen, EMNLP 2011]



<u>Feature = Typed Path</u> CityInState, CityInstate⁻¹, CityLocatedInCountry

Feature Value 0.32









[Lao, Mitchell, Cohen, EMNLP 2011]



<u>Feature = Typed Path</u> CityInState, CityInstate⁻¹, CityLocatedInCountry

Feature Value 0.32

Logistic Regresssion Weight





<u>Feature = Typed Path</u> CityInState, CityInstate⁻¹, CityLocatedInCountry

Feature Value 0.32

Logistic Regresssion Weight













Carnegie Mellon University

Logistic



CityInState, CityInstate⁻¹, CityLocatedInCountry 0.8 AtLocation⁻¹, AtLocation, CityLocatedInCountry

Feature Value 0.32 0.20

Weight

















•••









Random walk inference: learned rules

•CityLocatedInCountry(*city, country*):

- •8.04 cityliesonriver, cityliesonriver⁻¹, citylocated incountry
- •5.42 hasofficeincity⁻¹, hasofficeincity, citylocated incountry
- •4.98 cityalsoknownas, cityalsoknownas, citylocatedincountry
- •2.85 citycapitalofcountry, citylocated incountry⁻¹, citylocated incountry
- •2.29 agentactsinlocation⁻¹, agentactsinlocation, citylocated incountry
- •1.22 statehascapital⁻¹, statelocatedincountry
- •0.66 citycapitalofcountry
 - 7 of the 2985 learned rules for CityLocatedInCountry





Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology
- 6. Learn to infer relation instances via targeted random walks
- 7. Vision: connect NELL and NEIL





New Direction: Integrate Vision with Text

•The problem:

uisca

- Many things not learnable from text
- •New direction:
- integrate NELL with NEIL (Never Ending Image Learner) [Gupta, Chen, 2013]
- NELL gives noun phrases it understands to NEIL
- NEIL collects images associated with these, and analyzes
- NELL, NEIL cotraining







NEIL / NELL Polysemy: Bass









NEIL / NELL Polysemy: Bass









NEIL / NELL Polysemy: Bass



FAPESP





Building the Knowledge Graph by Reading

- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology
- 6. Learn to infer relation instances via targeted random walks
- 7. Vision: connect NELL and NEIL
- 8. Mutilingual NELL (English, Portuguese, Spanish, ...)



Recently_Learned Facts twitter

Recently learned beliefs (from English text)

instance	iteration	date learned	confidence	
<u>actimmune</u> is a <u>product</u>	890	11-dec-2014	100.0	D 5
<u>dogwood_drive</u> is a <u>street</u>	892	30-dec-2014	100.0	_2 T
<u>the_news_progress</u> is a <u>newspaper</u>	892	30-dec-2014	100.0	_2 ¢
university_of_washington is a train station	892	30-dec-2014	100.0	_2 ¢
<u>iranian_real</u> is a <u>currency</u>	892	30-dec-2014	91.5	_2 ¢
lotronex is a drug worked on by glaxosmithkline	892	30-dec-2014	93.8	_2 ¢
<u>peter_finch</u> <u>starred in</u> the movie <u>network</u>	892	30-dec-2014	100.0	_2 ¢
<u>bmw</u> is a specific automobile maker dealer <u>in tampa_bay</u>	893	02-jan-2015	100.0	_2 ¢
<u>jeremy</u> is a person who <u>died at the age of</u> 5	895	22-jan-2015	98.4	D 5
johannes_brahms is a person <u>born on</u> the date <u>n1833</u>	895	22-jan-2015	100.0	D 🖓
Recently learned beliefs (from Portuguese text)				

instance	iteration	date learned	confidence
william_maxwell is a South American person	105	17-oct-2016	99.4 🏖 ኛ
emily_procter is an actor	105	17-oct-2016	97.7 🍃 ኛ
<u>five_petalled_flowers</u> is a <u>plant</u>	105	17-oct-2016	99.3 🕼 ኛ
detergent_compatible_protein_assay is a chemical	107	19-oct-2016	99.5 🗳 ኛ
southern_us_ is a geopolitical entity that is an organization	106	18-oct-2016	98.3 🏖 ኛ
edgar_degas is part of the visual art movement realismo	110	21-oct-2016	100.0 🍃 ኛ
suzano is a sports team that plays in the league superliga_de_volei_volei	110	21-oct-2016	100.0 🏖 ኛ
adobe has acquired macromedia	110	21-oct-2016	99.9 🗳 ኛ
ryan_giggs plays for the team manchester_united	110	21-oct-2016	100.0 🍃 ኛ
ponta_grossa is a city located in the state or province paran®	110	21-oct-2016	99.2 🗳 ኛ







How to Read the Web in Many Languages?









NELL: Never-Ending Language Learner









NELL: Never-Ending Language Learner









NELL: Never-Ending Language Learner








NELL: Never-Ending Language Learner









NELL: Never-Ending Language Learner



Multilingual NELL







NELL: Never-Ending Language Learner







Federal University of São Carlos

ISC/ar



Picture Adapted from http://www.xlike.org



Federal University of São Carlos

ISC/ar



Picture Adapted from http://www.xlike.org



Federal University of São Carlos

ISCA



Picture Adapted from http://www.xlike.org





Picture Adapted from http://www.xlike.org





Picture Adapted from http://www.xlike.org











Federal University of São Carlos

u is dat

Multilingual Reading The Web





Federal University of São Carlos

u is dat

Multilingual Reading The Web





Federal University of São Carlos

u is dat









Federal University of São Carlos

uferer.







Federal University of São Carlos

uferer.

Federal University of São Carlos





Federal University of São Carlos





Federal University of São Carlos





Federal University of São Carlos









Federal University of São Carlos

U-ISI AT





Federal University of São Carlos

U-ISI AT





Federal University of São Carlos

UFIST AT









Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

Input:





Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

Carnegie Mellon University



Input:



Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016



Input:

Positive Examples: About 4000 manually Defined equivalence relations



Hernandez-Gonzalez, Hruschka Jr.



Federal University of São Carlos Merging Concepts Learned in Different Merging Concepts Learned in Different Hernandez-Gonzalez, Hruschka Jr. Languages Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016 Mitchell, 2016



Figure 1: Graphical representation of the generation of negative examples. Given a positive pair (USA, equivalentTo, EUA), the relationship (USA, bor-derWith, Mexico) is followed to reach an entity with the same category (Mex-ico, Country), generating the negative example (Mexico, notEquivalentTo, EUA).



Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016



Input:

1.4M automatically translated literal strings pairs



Hernandez-Gonzalez, Hruschka Jr.





Merging Concepts Learned in Different

Languages

Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

ML Task:

binary classifier (yes/no) for a given "equivalentTo" Relation candidate

The Classifier:

lf

sameCategory(concept1,concept2), and differentLang(concept1,concept2), and higherThanTreshold(score)







Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

3 Approaches:

Personalized Page Rank – PPR Path Ranking Algorithm – PRA Dictionary-based – Baseline





Merging Concepts Learned in Different Languages Hernandez-Gonzalez,

Hruschka Ir and Mitchell, 2016

PPR configuration:

- 2,000 random walks
- walk length: 5
- probability of restart: 0.01

PRA configuration:

- breadth-first search
- depth: 2



Merging Concepts Learned in Different

Languages

Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

	English						Portuguese					
Category	In graph			In positive pairs			In graph			In positive pairs		
animal	12436	591	746	36	23	27	101	63	97	36	14	35
$\operatorname{country}$	6031	443	460	106	93	93	153	94	136	106	87	103
city	18893	4437	5311	460	237	263	5767	483	1404	460	138	282
movie	7008	712	831	42	38	38	368	64	132	42	25	40
person	6693	2898	3050	403	395	399	621	611	614	403	304	376
writer	18911	1707	2143	61	39	40	114	26	63	61	23	37
actor	28361	794	1421	512	139	167	1870	129	793	512	36	208
\mathbf{sport}	5022	205	381	109	66	75	153	34	125	109	29	94
global	1909339	66239	96086	4126	2112	2331	30401	5119	12930	4126	1827	2565

Table 1: For both language subgraphs, the total number of nodes in a graph and the number of nodes included in positive pairs are shown. The three columns in each vertical subdivision show counts for the original graph, the pruned graph and the SVO+pruned graph, respectively from left to right. Each line shows the counts for the considered categories (in the last line, all the categories).





Figure 2: PR curves comparing both proposals with the dictionary as a classifier. Each figure displays the results of the three approaches with the examples of a specific category of the ontology using the pruned graph.





Hernandez-Gonzalez, Hruschka Jr. and Mitchell, 2016

Originally Sparse Graphs

	E	nglish	Portuguese			
Category	Pruned graph	SVO+pruned graph	Pruned graph	SVO+pruned graph		
animal	13.32 ± 47.12	392.02 ± 1859.43	1.57 ± 1.17	27.23 ± 58.58		
country	22.65 ± 68.60	289.97 ± 970.55	5.04 ± 15.99	180.79 ± 820.71		
city	5.00 ± 22.31	134.72 ± 1176.05	1.71 ± 4.33	32.03 ± 135.13		
movie	1.60 ± 1.69	88.95 ± 891.37	1.33 ± 0.79	34.77 ± 177.64		
person	2.99 ± 6.76	244.10 ± 1394.07	1.18 ± 0.66	16.62 ± 101.67		
writer	2.72 ± 4.75	43.50 ± 419.41	1.15 ± 0.46	6.40 ± 12.04		
actor	2.19 ± 2.12	77.21 ± 819.10	1.67 ± 1.41	6.96 ± 11.40		
\mathbf{sport}	19.94 ± 132.71	256.49 ± 1430.06	3.88 ± 6.31	26.98 ± 83.30		
global	4.48 ± 28.95	275.49 ± 1715.53	1.81 ± 3.93	51.26 ± 199.03		

Table 2: Mean out-degree value and associate standard deviation of the nodes in the pruned and SVO+pruned graphs of both languages. Each line shows the values for the considered categories (in the last line, all the categories).


Figure 3: PR curves comparing both proposals with the dictionary as a classifier. Each figure displays the results of the three approaches with the examples of a specific category of the ontology using the graph pruned after populating it with new relationships inferred from a SVO corpus.





and Mitchell, 2016

- denser SVO+pruned graph enhances the performance of the PRA-based proposal
- PPR is less regular and changes considerably among categories





- Dictionary is very competitive in categories we can find translations for their instances. In category sport, the performance was very good, but not in category Movie
- PRA often considers the "canBeTranslatedAs" path as strong feature, thus it can be seen as internally using the dictionary approach. But, can help when the Dictionary brings bad results (see country and person).







- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology
- 6. Learn to infer relation instances via targeted random walks
- 7. Vision: connect NELL and NEIL
- 8. Mutilingual NELL (Portuguese)
- 9. CrossLingual NELL
- 10. Learn to microread single sentences and single documents
- 11. Self reflection, self-directed learning
- 12. Goal-driven reading: predict, then read to corroborate/correct
- 13. Make NELL learn by conversation (e.g, Twitter)
- 14. Add a robot body, or mobile phone body, to NELL





- 1. Classify noun phrases (NP's) by category
- 2. Classify NP pairs by relation
- 3. Discover rules to predict new relation instances
- 4. Learn which NP's (co)refer to which latent concepts
- 5. Discover new relations to extend ontology
- 6. Learn to infer relation instances via targeted random walks
- 7. Vision: connect NELL and NEIL
- 8. Mutilingual NELL (Portuguese)
- 9. CrossLingual NELL
- 10. Learn to microread single sentences and single documents
- 11. Self reflection, self-directed learning
- 12. Goal-driven reading: predict, then read to corroborate/correct
- 13. Make NELL learn by conversation (e.g, Twitter)
- 14. Add a robot body, or mobile phone body, to NELL



NELL is here



Knowledge on Demand

NELL Knowledge Base Browser

Federal University of São Carlos

CMU Read the Web Project

categories

relations

- everypromotedthing
- abstractthing
 - creativework
 - book
 - poem
 - lyrics
 - musicalbum
 - musicsong
 - televisionshow
 - movie
 - visualartform
 - species
 - animal
 - vertebrate
 - bird fish
 - reptile
 - mammal amphibian
 - invertebrate
 - arthropod
 - insect
 - crustacean
 - arachnid
 - mollusk

To browse the knowledge base:

- Click on a category (or relation) from the list in the left-hand panel. This will bring up a list of facts that NELL has read that are relevant to that category (or relation).
- By default, facts are sorted by NELL's confidence that they are true. You may also sort alphabetically, by iteration, or by the date at which that fact was first read on the Web. To do so, simply click on the corresponding column heading.
- You may also search entities in the KnowledgeBase using the search box in the upper-right.
- Click on an entity (noun phrase) to bring up a detailed view of all the facts that are known about it.
- The "facts" that are shown in light grey (like this) are things that NELL has found some weak evidence for somewhere the Web, but doesn't quite believe to be true.
- For each fact in the detailed view, we also present some "source" information, describing which subsystems (e.g., CPL, SEAL, CMC, RL) were used in contributing to NELL's understanding of this fact. This includes the system iteration, confidence, and date at the time it was read, plus some details (e.g., web page links or text patterns).

For more technical details on the NELL system and how it reads the Web, see our AAAI 2010 paper.

NEW: Knowledge on demand:

Try our new Ask NELL service to see what NELL can read and infer on the fly.

Carnegie Mellon University

log in | preferences | help/instructions | feedback

Search



log in | preferences | help/instructions | feedback

NELL: Never-Ending Language Learner •NELL is grown enough for new steps

Knowledge on Demand – Ask NELL

NELL Knowledge Base Browser

Federal University of São Carlos

CMU Read the Web Project

ISC/a

categories relat	Ask NELL:
everypromotedthing	You can now ask NELL what it believes about any noun phrase (e.g., rocking chair, chocolate). Try it!
 abstractining creativework 	What categories does belong to? Answer
 book poem 	
 lyrics 	
 musicalbum musicsong televisionshow 	What is NELL Doing?
 movie visualartform 	NELL is looking up your input noun phrase in its knowledge base, and also attempting to infer additional beliefs about it on the fly (by reasoning from other beliefs, and reading more). Therefore, it might take a minute or two.
species animal vertebrate	Underlying API
 bird fish reptile mammal amphibian invertebrate arthropod 	The demos above are based on a public machine-friendly web-based API that returns a JSON object in response to an HTTP GET request. This underlying API is somewhat more complicated to use, and we offer both <u>detailed documentation</u> and a <u>test UI</u> for developers.

- crustacean
- arachnid
- mollusk

Search





NELL Knowledge Base Browser Search CMU Read the Web Project log in | preferences | help/instructions | feedback categories relations **Microreading on Demand Demo** everypromotedthing Enter text to microread in the box below: abstractthing event convention musicfestival protestevent meetingeventtitle conference mlconference weatherphenomenon sportsevent sportsgame race olympics grandprix crimeorcharge earthquakeevent election bombingevent Microread militaryeventtype Output Format: JSON \$ militaryconflict







NELL Knowledge Base Browser

CMU Read the Web Project

Search log in | preferences | help/instructions | feedback

categories relations everypromotedthing abstractthing event convention musicfestival protestevent meetingeventtitle conference • mlconference weatherphenomenon sportsevent sportsgame race olympics grandprix crimeorcharge earthquakeevent election bombingevent Microread militaryeventtype Output Format: JSON \$ militaryconflict

Microreading on Demand Demo

Enter text to microread in the box below:

The statement by US Secretary of State Colin Powell about the Venezuelan Government triggered yesterday reactions from members of the ruling party . He stressed : `` Venezuela 's domestic policy is the exclusive concern of Venezuelans because of our historical tradition and because it is a constitutional mandate . " He added : `` There is no reason for any country to request the agreement of others to determine its conduct and adopt its decisions . " However , he ruled out that the statement by the US Foreign Secretary will lead to a change in the policies of President Chavez . Deputy Cilia Flores , chief of bloc of the Fifth Republic Movement -LSB- MVR -RSB- , used the same tone to express her opinion . She added : Venezuela designs its own policies in an autonomous manner and , in any event , it is up to Venezuelans to decide them . We respect the sovereignty of all countries and this is why we demand respect for ours . " According to her , the remarks by the US Government spokesman were uncalled-for . She stated : `` We are quite respectful of the decisions made by Washington , even if we do not agree with them . We demand equal treatment for our foreign policy . " The opposition seeks intervention or blockade Journalists Ibeyise Pacheco, Marta Colomina, Patricia Poleo, and Marianela Salazar disclosed this video in recent days . He took one further step and invited Powell to visit Venezuela to see what is really happening . He accused the four journalists of plotting with powerful economic groups that presumably contributed the funds to purchase the video from an intelligence agent. Oddly, he defended the





Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots.









Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots.







Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots.S2.) Alice caught the butterfly with the net.







Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots. S2.) Alice caught the butterfly with the net. Relations:

locatedln(*city*, *country*): locatedln(Saint Etienne, France), ...

animalHasCoatPattern(*animal, coatPattern*): animalHasCoatPattern(zebra, stripes), animalHasCoatPattern(leopard, spots), animalHasCoatPattern(butterfly, spots), ...

caught(instrument, animal):
caught(net, butterfly),
caught(fish, fishing rod), ...

Categories:

city(*city*): city(Saint Etienne), ...

animal(*animal*): animal(zebra), animal(butterfly), ...







Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots. S2.) Alice caught the butterfly with the net. Relations:

IocatedIn(*city*, *country*): IocatedIn(Saint Etienne, France), ... animalHasCoatPattern(*animal*, *coatPattern*): animalHasCoatPattern(zebra, stripes),

animalHasCoatPattern(leopard, *Spots*),

animalHasCoatPattern(butterfly, **Spots**), ...

caught(instrument, animal):
caught(net, butterfly),
caught(fish, fishing rod), ...

Categories:

city(city):
city(Saint Etienne), ...

animal(*animal*): animal(zebra), animal(butterfly), ...





Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots. S2.) Alice caught the butterfly with the net. Relations:

IocatedIn(*city*, *country*): IocatedIn(Saint Etienne, France), ... animalHasCoatPattern(*animal*, *coatPattern*): animalHasCoatPattern(zebra, stripes),

animalHasCoatPattern(leopard, *Spots*),

animalHasCoatPattern(butterfly, *Spots*), ...

caught(instrument, animal):

caught(*net*, butterfly), caught(fish, fishing rod), ...

Categories:

city(*city*): city(Saint Etienne), ...

animal(*animal*): animal(zebra), animal(butterfly), ...





Nakashole and Mitchell, 2015

S2.) Alice caught the butterfly with the net. Relations:

locatedIn(*city*, *country*): locatedIn(Saint Etienne, France), ... animalHasCoatPattern(*animal*, *coatPattern*): animalHasCoatPattern(zebra, stripes), animalHasCoatPattern(leopard, *SpotS*), animalHasCoatPattern(butterfly, *spotS*), ... caught(*instrument*, *animal*): caught(*inet*, butterfly), caught(fish, fishing rod), ...

Categories:

city(city):
city(Saint Etienne), ...

animal(*animal*): animal(zebra), animal<mark>(bu</mark>tterfly)<mark>, ...</mark>







Nakashole and Mitchell, 2015

S1.) Alice caught the butterfly with the spots.S2.) Alice caught the butterfly with the net.







NELL Knowledge Base Browser

CMU Read the Web Project

Search log in | preferences | help/instructions | feedback

categories relations everypromotedthing abstractthing event convention musicfestival protestevent meetingeventtitle conference • mlconference weatherphenomenon sportsevent sportsgame race olympics grandprix crimeorcharge earthquakeevent election bombingevent Microread militaryeventtype Output Format: JSON \$ militaryconflict

Microreading on Demand Demo

Enter text to microread in the box below:

The statement by US Secretary of State Colin Powell about the Venezuelan Government triggered yesterday reactions from members of the ruling party . He stressed : `` Venezuela 's domestic policy is the exclusive concern of Venezuelans because of our historical tradition and because it is a constitutional mandate . " He added : `` There is no reason for any country to request the agreement of others to determine its conduct and adopt its decisions . " However , he ruled out that the statement by the US Foreign Secretary will lead to a change in the policies of President Chavez . Deputy Cilia Flores , chief of bloc of the Fifth Republic Movement -LSB- MVR -RSB- , used the same tone to express her opinion . She added : Venezuela designs its own policies in an autonomous manner and , in any event , it is up to Venezuelans to decide them . We respect the sovereignty of all countries and this is why we demand respect for ours . " According to her , the remarks by the US Government spokesman were uncalled-for . She stated : `` We are quite respectful of the decisions made by Washington , even if we do not agree with them . We demand equal treatment for our foreign policy . " The opposition seeks intervention or blockade Journalists Ibeyise Pacheco, Marta Colomina, Patricia Poleo, and Marianela Salazar disclosed this video in recent days . He took one further step and invited Powell to visit Venezuela to see what is really happening . He accused the four journalists of plotting with powerful economic groups that presumably contributed the funds to purchase the video from an intelligence agent. Oddly, he defended the





(PP (IN in) (S (VP (VBG handling) (NP (NNS relations)) (PP (IN with) (NP (DT the) (NNP United) (NNPS States))))))))))) (. .))) Type: con-parse Score: 0.0000 Value: (ROOT (S (NP (PRP It)) (VP (VBZ has) (VP (VBN been) (ADJP (JJ vague) (PP (IN about) (NP (NP (NP (NP venezuela) (POS 's)) (NNS intentions)) (PP (VBG regarding) (NP (NP (DT the) (NN protection)) (PP (IN of) (NP (NNS prices)))) (CC and) (NP (NP (DT the) (NN war)) (PP (IN on) (NP (NN terrorism))))))))))) (. .)))

Type: dep-parse Value: det(statement-2, The-1) nsubj(triggered-14, statement-2) case(Secretary-5, by-3) compound(Secretary-5, US-4) nmod(statement-2, Secretary-5) case(Powell-9, of-6) compound(Powell-9, State-7) compound(Powell-9, Colin-8) nmod(Secretary-5, Powell-9) case(Government-13, about-10) det(Government-13, the-11) amod(Government-13, Venezuelan-12) nmod(Powell-9, Government-13) root(ROOT-0, triggered-14) compound(reactions-16, yesterday-15) dobj(triggered-14, reactions-16) case(members-18, from-17) nmod(triggered-14, members-18) case(party-22, of-19) det(party-22, the-20) amod(party-22, ruling-21) nmod(members-18, party-22) punct(triggered-14, .-23)







NELL Architecture









NELL Architecture









Carnegie Mellon University

Conversing Learning





Federal University of São Carlos

Conversing Learning •Basic Steps:

- Decide which task is going to be asked
- Determine who are the oracles the ML system is going to consult
- Propose a method of conversation with oracles, often humans
- Determine how to feedback the ML system with the community inputs





•Basic Steps:

- Decide which task is going to be asked
- Determine who are the oracles the ML system is going to consult
- Propose a method of conversation with oracles, often humans
- Determine how to feedback the ML system with the community inputs





Decide which task is going to be asked

- Learned facts
- Learned Inference Rules
- Metadata (mainly for automatically extending the ontology)







Conversing Learning •Basic Steps:

- Decide which task is going to be asked
- Determine who are the oracles the ML system is going to consult
- Propose a method of conversation with oracles, often humans
- Determine how to feedback the ML system with the community inputs





•who are the oracles the ML system is going to consult

- •Yahoo! Answers
 - very popular on the Web
 - a lot of metadata to harvest
- •Twitter
 - millions of users worldwide
 - a system that was not designed to work as a QA environment
- •Both web communities have API to connect to their database











Federal University of São Carlos

Conversing Learning •Basic Steps:

- Decide which task is going to be asked
- Determine who are the oracles the ML system is going to consult
- Propose a method of conversation with oracles, often humans
- Determine how to feedback the ML system with the community inputs





- Propose a method of conversation with oracles, often humans
- Macro Question-Answering
- •For each posted question:
 - Ask for yes/no simple answers
 - Try to understand every answer
 - Discard answers too difficult to understand
 - Conclude based only on fully understood answers





Federal University of São Carlos

Conversing Learning •Basic Steps:

- Decide which task is going to be asked
- Determine who are the oracles the ML system is going to consult
- Propose a method of conversation with oracles, often humans
- Determine how to feedback the ML system with the community inputs





•how to feedback the ML system with the community inputs?

- •Suggested actions to NELL:
 - Synonym/co-reference resolution
 - Automatically update the Knowledge Base
 - Feature Selection/Weighting







- •Some Initial Results with First Order Rules:
- Take top 10% of rules from Rule Learner
- 60 rules were converted into questions and asked with both the regular and the Yes/No question approach
- The 120 questions received a total of 350 answers.







•Some Initial Results with First Order Rules:

- Rule extracted from NELL in PROLOG format
 stateLocatedInCountry(x,y):statehascapital(x,z), citylocatedincoutry(z,y)
- converted into question:

•Is this statement always true? If state X has capital Z and city Z is located in country Y then state X is located in country Y.







Conversing Learning - FOR

- •Question: (Yes or No?) If athlete Z is member of team X and athlete Z plays in league Y, then team X plays in league Y.
- Twitter answers sample:
- No. (Z in X) \land (Z in Y) \rightarrow (X in Y)
- Yahoo! Answers sample:

• NO, Not in EVERY case. Athlete Z could be a member of football team X and he could also play in his pub's Friday nights dart team. The Dart team could play in league Y (and Z therefore by definition plays in league Y). This does not mean that the football team plays in the darts league!














Question Type	TP	TN	FP	FN	Prec.	Recall	Acc.	F1
Regular Individual	41	31	7	39	0.85	0.51	0.61	0.64
Regular Combined	27	2	17	7	0.61	0.79	0.54	0.69
Yes/No Individual	62	53	14	45	0.81	0.57	0.66	0.67
Yes/No Combined	27	7	11	11	0.71	0.71	0.60	0.71
Best Answers	19	28	3	29	0.86	0.39	0.59	0.54

Carnegie Mellon University





•Some Initial Results with Facts Validation:







•Some Initial Results with Facts Validation:







•Some Initial Results with Facts Validation:

Relations	outliers	answers
AthletePlaysInLeague & Players	9	72
TeamPlaysSport & TeamPlaysInLeague	20	144
TeamPlaysSport & TeamWonTrophy	53	386







•Some Initial Results with Facts Validation:

Relations	outliers	answers
AthletePlaysInLeague & Players	9	72
TeamPlaysSport & TeamPlaysInLeague	20	144
TeamPlaysSport & TeamWonTrophy	53	386

Manchester United







•Some Initial Results with Metadata:

- Question: Could you please give me some examples of clothing?
- Answer 01: Snowshoes, rain ponchos, galoshes, sunhats, visors, scarves, mittens, and wellies are all examples of weather specific clothing!
- Answer 02: pants
- Answer 03: Training shoes can be worn by anyone for any purpose, but the term means to train in sports





Conversing Learning •Some Initial Results with Metadata:

- Users replied with 552 seeds for 129 categories
 Total of 5900 promotions with seeds created by NELL's developers
- Total of 5300 promotions with seeds extracted from answers of Twitter users (similar precision)







- •Some Initial Results with Metadata:
- For Relation Discovery Components
 - Symmetry: Is it always true that if a person P1 is neighbor of a person P2, then P2 is neighbor of P1?
 - Anti-symmetry: Is it always true that if a person P1 is the coach of a person P2, then P2 is not coach of P1?







- •Some Initial Results with Metadata:
- Feature Weighting/Selection for CMC
 - Logistic Regression features are based on noun phrase morphology
 - (true or false) hotel names tend to be compound noun phrases having "hotel" as last the word.
 - (true or false) a word having "burgh" as sufix (ex.
 Pittsburgh) tend to be a city name.







•On going and future work

- Asking to the right community and to the right person
- Asking the right thing to maximize the results with minimum questions (multi-view Active Learning)
- Better Question-Answering methods
- Asking in different languages and explore time zones.







- •Knowledge Base Validation in NELL
- Human Supervision: RTW group members;
- **Conversing Learning:** NELL can autonomously talk to people in web communities and ask for help
- Web Querying: NELL can query the Web on specific facts to verify correctness, or to predict the validity of a new fact;
- Hiring Labelers: NELL can autonomously hire people (using web services such as Mechanical Turk) to label data and help the system to validate acquired knowledge.



Federal University of São Carlo





NELL Architecture





Carnegie Mellon University





Motivation

Samadi, Veloso and Blum, 2013





Samadi, Veloso and Blum, 2013

healthyFood(shrimp)?

Federal University of São Carlos

HSC/ar



Given the input time, which CBIs should be extracted?



Carnegie Mellon University

510





NELL: Never-Ending Language Learner

Samadi, Veloso and Blum, 2013

•OpenEval in the last iteration:

-academicfield 0.8976357986206526 Environmental Anthropology.

-Several excellent textbooks and readers in environmental anthropology have now appeared, establishing a basic survey of the field.





Federal University of São Carlos

NELL: Never-Ending Language Learner

Samadi, Veloso and Blum, 2013

•OpenEval in the last iteration:

-academicfield 0.912473775634353 Anesthesiology.

-The Department of Anesthesiology is committed to excellence in clinical service, education, research and faculty development.





Federal University of São Carlos

NELL: Never-Ending Language Learner

Samadi, Veloso and Blum, 2013

•OpenEval in the last iteration:

-worksfor 0.9845774661303888 (charles osgood, cbs).

-Charles Osgood, often referred to as CBS News' poet-in-residence, has been anchor of "CBS News Sunday Morning" since 1994.





It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn, worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."



It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

Blomkvist visits Henrik Vanger at his estate on the tiny island of Hedeby. The old man draws Blomkvist in by promising solid evidence against Wennerström. Blomkvist agrees to spend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

After discovering that Salander has hacked into his computer, he persuades her to assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in turn? Worries that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at same te on the same and of Hedeby. The old man drows Blomkvist in by premising sond evidence against Wennerström. Blomkvist ag same pend a year writing the Vanger family history as a cover for the real assignment: the disappearance of Vanger's niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist becomes acquainted with the members of the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik.

Af same overing that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, bu Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep background investigations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder.

Federal University of São Carlo

- ISC/a

Blomkvist visits Henrik Vanger at te on the same and of Hedeby. The old man drov a Blophovist in by promising solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist beco uncleOf inted with the men hires the extended Vanger family, most of whom resent his presence. He does, however, start a short lived affair with Cecilia, the niece of Henrik. Af same)vering that Salander has hacked into his computer, he persuade same assist him with research. They eventually become lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone she meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries, is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same

Federal University of São Carlo

ISC/a

It's about the disappearance forty years ago of Harriet Vanger, a young scion of one of the wealthiest families in Sweden, and about her uncle, determined to know the truth about what he believes was her murder. Blomkvist visits Henrik Vanger at te on the same and of Hedeby. same The old man drov a Blorikvist in by much sing solid evidence against Wennerström. Blomkvist ag same year writing the Vanger family history as a cover for the real assignment: the disappearance of V owns niece Harriet some 40 years earlier. Hedeby is home to several generations of Vangers, all part owners in Vanger Enterprises. Blomkvist the extended Vanger family, most of whom resent beco uncleOf inted with the men hires his presence. He does, however, start a short lived affair with Cecilia, the niece of T enemvOf Af same vering that Salander has hacked into his co affair With persuade ASSIST him with research. They even¹ affairWith lovers, but Blomkvist has trouble getting close to Lisbeth who treats virtually everyone sne meets with hostility. Ultimately the two discover that Harriet's brother Martin, CEO of Vanger Industries is secretly a serial killer. A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill." tu: same



It's about the disap scion of one of the determined to know **Blomkvist** visits He The old man dray of Blomkvist ag same assignment: the dis home to several ger beco uncleOf inte his presence. He do Af same)vering t him with research. to Lisbeth who trea discover that Harri

ance forty years ago of Harriet V hiest families in Sweden, and a truth about what he believ s her Vanger at and te un he same akvist in by memory evidence er family hi nd 2 riting th Harriet some wns aran rt owners in Vai ions of e extended Vange h the me t lived affair with owever, sta c affairW alander lovers, bat Blomk even¹ affairWith

oung cle, der. edeby. nst Wennerström. as a cover for the real vears earlier. Hedeby is Enterprises<mark>.</mark> Bl<mark>omkvist</mark> ily, most of whom resent ia, the niece 💕 🎹 enemyOf bersuade **assis**i same has trouble getting close Itimately the two is secretly a serial killer.

A 24-year-old computer hacker sporting an assortment of tattoos and body piercings supports herself by doing deep backgrou headOf gations for Dragan Armansky, who, in tu same ies that Lisbeth Salander is "the perfect victim for anyone who wished her ill."







- Named Entity Resolution/Extraction
- Relation Extraction
- Co-reference and Synonym Resolution
- Relation Discovery
- Inference
- Knowledge Base Representation
- Document/Sentence Understanding (Micro-Reading)





More on NELL

http://rtw.ml.cmu.edu/rtw/publications

Read the Web

Federal University of São Carlos

Research Project at Carnegie Mellon University







Industrial Approaches

- Google Knowledge Vault:
- "The Knowledge Vault is substantially bigger than any previously published structured knowledge repository, and features a probabilistic inference system that computes calibrated probabilities of fact correctness".
- 271M facts with probability of being true at (or above) 0.9.
- It uses 16 different information extraction systems to extract (subject, predicate, object) knowledge triples from webpages.





Industrial Approaches

- IBM Watson:
- From Wikipedia: "Watson is a question answering computer system capable of answering questions posed in natural language, (...) The computer system was specifically developed to answer questions on the quiz show Jeopardy! (...) Watson had access to 200 million pages of structured and unstructured content consuming four terabytes of disk storage including the full text of Wikipedia, but was not connected to the Internet during the game. (...) Watson consistently outperformed its human opponents on the game's signaling device, but had trouble in a few categories, notably those having short clues containing only a few words."





estevam.hruschka@gmail.com

Thank you very much!

and thanks to all people from NELL, KnowItAII, YAGO and DBPedia projects for very nice discussions and suggestions. And thanks to FAPESP (Fundação de Amparo a Pesquisa do Estado de Sao Paulo), and thanks to Federal University of São Carlos and thanks to Carnegie Mellon University.



Carnegie Mellon University





• [Angeli et al., 2014] Angeli, Gabor, Julie Tibshirani, Jean Y. Wu, and Christopher D. Manning. "Combining distant and partial supervision for relation extraction." In*Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014.

- [Appel & Hruschka Jr., 2011] Appel, A.P., Hruschka Jr., E.R.: Prophet a link-predictor to learn new rules on Nell. In: Proceedings of the 2011 IEEE 11th International Conference on Data Mining Workshops. pp. 917–924. ICDMW '11, IEEE Computer Society, Washington, DC, USA (2011)
- Auer, S.; Bizer, C.; Kobilarov, G.; Lehmann, J.; Cyganiak, R.; and Ives, Z. 2007. *Dbpedia: A nucleus for a web of open data*. Springer.
- [Bach & Badaskar, 2007] Bach, Nguyen, and Sameer Badaskar. "A survey on relation extraction."*Language Technologies Institute, Carnegie Mellon University* (2007).
- [Banko, et al., 2007] Michele Banko, Michael J. Cafarella, Stephen Soderland, Matthew Broadhead, Oren Etzioni: Open Information Extraction from the Web. IJCAI 2007: 2670-2676
- Banko, M., and Etzioni, O. 2007. Strategies for lifelong knowledge extraction from the web. In *Proc. of K-CAP*.
- Biega, J.; Kuzey, E.; and Suchanek, F. M. 2013. In- side yago2s: A transparent information extraction architecture. In *Proceedings of the 22nd international conference on World Wide Web companion*, 325–328. International World Wide Web Conferences Steering Committee.
- Bizer, C.; Lehmann, J.; Kobilarov, G.; Auer, S.; Becker, C.; Cyganiak, R.; and Hellmann, S. 2009. Dbpedia-a crystallization point for the web of data. *Web Semantics: science, services and agents on the world wide web* 7(3):154– 165.
- [Bengtson & Roth, 2008] Bengtson, E., & Roth, D. (2008, October). Understanding the value of features for coreference resolution. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing* (pp. 294-303). Association for Computational Linguistics.





References and Bibliography

- [Carlson et al., 2010b] Carlson, A., Betteridge, J., Kisiel, B., Settles, B., Jr., E.R.H., Mitchell, T.M.: Toward an architecture for never-ending language learning. In: Proceedings of the Twenty-Fourth Conference on Artificial Intelligence (AAAI 2010).
- Carlson, A.; Betteridge, J.; Wang, R. C.; Hruschka Jr., E. R.; and Mitchell, T. M. 2010b. Coupled semisupervised learn- ing for information extraction. In *Proc. of WSDM*.
- [Clark et al., 2007] Peter Clark, Phil Harrison, John Thompson, Rick Wojcik, Tom Jenkins, David Israel, Reading to Learn: An Investigation into Language Understanding. The 2007 AAAI Spring Symposium. Published by The AAAI Press, Menlo Park, California, 2007.
- [DARPA, 2012] DARPA Machine Reading Program, http://www.darpa.mil/Our_Work/I2O/Programs/Machine_Reading.aspx.

- Dong, X.; Gabrilovich, E.; Heitz, G.; Horn, W.; Lao, N.; Murphy, K.; Strohmann, T.; Sun, S.; and Zhang, W. 2014. Knowledge vault: A web-scale approach to proba- bilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, 601–610. ACM.
- Dong, X. L.; Gabrilovich, E.; Murphy, K.; Dang, V.; Horn, W.; Lugaresi, C.; Sun, S.; and Zhang, W. 2015. Knowledgebased trust: Estimating the trustworthiness of web sources. *Proceedings of the VLDB Endowment* 8(9):938–949.
- [Dutta & Weikum, 2015]Dutta, S., & Weikum, G. (2015). Cross-document co-reference resolution using samplebased clustering with knowledge enrichment. *Transactions of the Association for Computational Linguistics*, *3*, 15-28.
- [Etzioni et al., 2007] Oren Etzioni, Michele Banko, and Michael J. Cafarella, Machine Reading. The 2007 AAAI Spring Symposium. Published by The AAAI Press, Menlo Park, California, 2007.
- [Etzioni et al., 2011] Oren Etzioni, Anthony Fader, Janara Christensen, Stephen Soderland, and Mausam "Open Information Extraction: the Second Generation". *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI 2011).*







- [Etzioni, 2008] Oren Etzioni. 2008. Machine reading at web scale. In *Proceedings of the international conference* on Web search and web data mining (WSDM '08). ACM, New York, NY, USA, 2-2.
- [Fader et al., 2011] Anthony Fader, Stephen Soderland, and Oren Etzioni. "Identifying Relations for Open Information Extraction". *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing (EMNLP 2011)*
- [Fern, 2008] Xiaoli Z. Fern, CS 434: Machine Learning and Data Mining, School of Electrical Engineering and Computer Science, Oregon State University, Fall 2008.
- Ferrucci, D. A.; Levas, A.; Bagchi, S.; Gondek, D.; and Mueller, E. T. 2013. Watson: beyond jeopardy! *Artif. Intell.* 199:93–105.
- Ferrucci, D. 2010. Build watson: an overview of deepqa for the jeopardy! challenge. In *Proceedings of the 19th interna- tional conference on Parallel architectures and compilation techniques*, 1–2. ACM.
- Ferrucci, D. A. 2012. Introduction to this is watson. *IBM Journal of Research and Development* 56(3.4):1–1.
- [Galárraga et al., 2014] Luis Galárraga, Geremy Heitz, Kevin Murphy, and Fabian M. Suchanek. 2014. Canonicalizing Open Knowledge Bases. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management* (CIKM '14). ACM, New York, NY, USA, 1679-1688.
- Gardner, M.; Talukdar, P. P.; Kisiel, B.; and Mitchell, T. 2013. Improving learning and inference in a large knowledge-base using latent syntactic cues. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP 2013)*.
- Gardner, M.; Talukdar, P.; Krishnamurthy, J.; and Mitchell, T. 2014. Incorporating vector space similarity in random walk inference over knowledge bases. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Lan-guage Processing (EMNLP)*. Doha, Qatar.: Association for Computational Linguistics.





Carnegie Mellon University



- Gerber, D.; Esteves, D.; Lehmann, J.; Bu hmann, L.; Usbeck, R.; Ngonga Ngomo, A.-C.; and Speck, R. 2015. Defacto

 temporal and multilingual deep fact validation. Web Seman-tics: Science, Services and Agents on the World
 Wide Web.
- [Grishman & Sundheim, 1996] Grishman, Ralph; Sundheim, B. 1996. Message Understanding Conference 6: A Brief History. In *Proc. International Conference on Computational Linguistics*.
- [Hady et al., 2011] Hady W. Lauw, Ralf Schenkel, Fabian Suchanek, Martin Theobald, and Gerhard Weikum[,] "Semantic Knowledge Bases from Web Sources" at IJCAI 2011, Barcelona, July 2011
- Hoffart, J.; Suchanek, F. M.; Berberich, K.; Lewis-Kelham, E.; de Melo, G.; and Weikum, G. 2011. Yago2: exploring and querying world knowledge in time, space, context, and many languages. In *Proceedings of the 20th international conference companion on World wide web*, WWW '11, 229–232. New York, NY, USA: ACM.
- [Hoffart et al., 2012] Johannes Hoffart, Fabian Suchanek, Klaus Berberich, Gerhard Weikum. YAGO2: A Spatially and Temporally Enhanced Knowledge Base from Wikipedia. Special issue of the Artificial Intelligence Journal, 2012
- [Hoffart et al., 2015] Johannes Hoffart, Nicoleta Preda, Fabian M. Suchanek, and Gerhard Weikum. Knowledge Bases for Web Content Analytics. WWW 2015 Tutorial.
- [Hruschka et al., 2012] E. R. Hruschka Jr. and M. C. Duarte and M. C. Nicoletti. Coupling as Strategy for Reducing Concept-Drift in Never-ending Learning Environments. Fundamenta Informaticae, IOS Press, 2012.
- [Krishnamurthy & Mitchell, 2011] Krishnamurthy, J., Mitchell, T.M.: Which noun phrases denote which concepts. In: Proceedings of the Forty Ninth Annual Meeting of the Association for Computational Linguistics (2011)



References and Bibliography

- [Lao et al., 2011] Lao, N., Mitchell, T., Cohen, W.W.: Random walk inference and learning in a large scale knowledge base. In: Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing. pp. 529–539. Associa- tion for Computational Linguistics, Edinburgh, Scotland, UK. (July 2011), http://www.aclweb.org/anthology/D11-1049
- Lehmann, J., Isele, R., Jakob, M., Jentzsch, A., Kontokostas, D., Mendes, P. N., ... & Bizer, C. (2014).
 DBpedia-a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web Journal*, *5*, 1-29.
- Lehmann, J.; Schu ppel, J.; and Auer, S. 2007. Discovering unknown connections-the dbpedia relationship finder. *CSSW* 113:99–110.
- Lin, T.; Mausam; and Etzioni, O. 2012a. Entity link- ing at web scale. In *Proceedings of the Joint Workshop* on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction (AKBC-WEKEX), 84–88. Montré al, Canada: Association for Computational Linguistics.
- Lin, T.; Mausam; and Etzioni, O. 2012b. No noun phrase left behind: Detecting and typing unlinkable entities. In *Pro- ceedings of the 2012 Joint Conference on Empirical Meth- ods in Natural Language Processing and Computational Natural Language Learning*, EMNLP-CoNLL '12, 893–903. Stroudsburg, PA, USA: Association for Computational Linguistics.
- Lin, T. 2013. *Leveraging Knowledge Bases in Web Text Processing*. Ph.D. Dissertation, University of Washington.
- Lukovnikov, D.; Kontokostas, D.; Stadler, C.; Hellmann, S.; and Lehmann, J. 2014. Dbpedia viewer-an integrative inter-face for dbpedia leveraging the dbpedia service eco system. In *Proceedings of the 7th Workshop on Linked Data on the Web*.





Carnegie Mellon University



- [Lodhi et al., 2002] Lodhi, H., Saunders, C., Shawe-Taylor, J., & Cristianini, N. (2002). Text classification using string kernels. Journal of Machine Learning Research (pp. 419–444).
- [Mahdisoltani et al., 2015] Mahdisoltani, F., Biega, J., & Suchanek, F. M. YAGO3: A Knowledge Base from Multilingual Wikipedias. 7th Biennial Conference on Innovative Data Systems Research, CIDR 2015
- [Mitchell, 1997] Tom M. Mitchell, Machine Learning. McGraw-Hill, 1997.

- [Mitchell, 2006] Tom M. Mitchell, The Discipline of Machine Learning, my perspective on this research field, July 2006 (http://www.cs.cmu.edu/~tom/pubs/MachineLearning.pdf).
- Mitchell, T.; Cohen, W.; Hruschka, E.; Talukdar, P.; Bet- teridge, J.; Carlson, A.; Dalvi, B.; Gardner, M.; Kisiel, B.; Krishnamurthy, J.; Lao, N.; Mazaitis, K.; Mohamed, T.; Nakashole, N.; Platanios, E.; Ritter, A.; Samadi, M.; Settles, B.; Wang, R.; Wijaya, D.; Gupta, A.; Chen, X.; Saparov, A.; Greaves, M.; and Welling, J. 2015. Never-ending learn- ing. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI-15)*.
- [Mohamed et al., 2011] Mohamed, T.P., Hruschka, Jr., E.R., Mitchell, T.M.: Discovering relations between noun categories. In: Proceedings of the Conference on Empirical Methods in Nat- ural Language Processing. pp. 1447–1455. EMNLP '11, Association for Computa- tional Linguistics, Stroudsburg, PA, USA (2011)
- Morsey, M.; Lehmann, J.; Auer, S.; and Ngomo, A.-C. N. 2011. Dbpedia sparql benchmark–performance assessment with real queries on real data. In *The Semantic Web–ISWC 2011*. Springer. 454–469.
- Morsey, M.; Lehmann, J.; Auer, S.; Stadler, C.; and Hell- mann, S. 2012. Dbpedia and the live extraction of structured data from wikipedia. *Program* 46(2):157–181.
- [Nadeau & Sekine, 2007] Nadeau, David, and Satoshi Sekine. "A survey of named entity recognition and classification." *Lingvisticae Investigationes* 30.1 (2007): 3-26.





References and Bibliography

- [Navarro et al., 2013] Navarro, L. F. and Appel, A. P. and Hruschka Jr., E. R., GraphDB Storing Large Graphs on Secondary Memory. In New Trends in Databases and Information. Advances in Intelligent Systems and Computing, Springer, 177-186, 2013.
- Nakashole, N.; Weikum, G.; and Suchanek, F. M. 2012. Discovering and Exploring Relations on the Web. In *Demo at VLDB*.
- Nickel, M.; Murphy, K.; Tresp, V.; and Gabrilovich, E. 2015. A review of relational machine learning for knowl- edge graphs: From multi-relational link prediction to au-
- tomated knowledge graph construction. *arXiv:1503.00759*.
- [Norvig, 2007] Peter Norvig, Inference in Text Understanding. The 2007 AAAI Spring Symposium. Published by The AAAI Press, Menlo Park, California, 2007.
- [Pedro & Hruschka Jr., 2012a] Saulo D.S. Pedro and Estevam R. Hruschka Jr., Conversing Learning: active learning and active social interaction for human supervision in never-ending learning systems. Xiii Ibero-american Conference On Artificial Intelligence, IBERAMIA 2012, 2012.
- [Pedro & Hruschka Jr., 2012b] Pedro, S.D.S., Hruschka Jr., E.R.: Collective intelligence as a source for machine learning self-supervision. In: Proc. of the 4th International Workshop on Web Intelligence and Communities. WIC12, NY, USA, ACM (2012) 5:1–5:9
- [Pedro et al., 2013] Saulo D.S. Pedro, Ana Paula Appel, and Estevam R. Hruschka, Jr. Autonomously reviewing and validating the knowledge base of a never-ending learning system. In *Proceedings of the 22nd international conference on World Wide Web companion* (WWW '13 Companion), 1195-120, 2013.







References and Bibliography

- [Pennington et al., 2014] Pennington, Jeffrey, Richard Socher, and Christopher D. Manning. "Glove: Global vectors for word representation." *Proceedings of the Empiricial Methods in Natural Language Processing (EMNLP 2014)* 12 (2014): 1532-1543.
- [Raghunathan et al., 2010] Raghunathan, K., Lee, H., Rangarajan, S., Chambers, N., Surdeanu, M., Jurafsky, D., & Manning, C. (2010, October). A multi-pass sieve for coreference resolution. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing* (pp. 492-501). Association for Computational Linguistics.
- [Ratinov & Ruth, 2009] Ratinov, Lev, and Dan Roth. "Design challenges and misconceptions in named entity recognition." *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*. Association for Computational Linguistics, 2009.
- [Recasens et al., 2013] Recasens, M., Can, M., & Jurafsky, D. (2013). Same Referent, Different Words: Unsupervised Mining of Opaque Coreferent Mentions. In *Proceedings of NAACL-HLT* (pp. 897-906).
- Samadi, M.; Veloso, M. M.; and, M. 2013. Openeval: Web information query evaluation. In AAAI. Citeseer.
- [Settles, 2011] Settles, B.: Closing the loop: Fast, interactive semi-supervised annotation with queries on features and instances. In: Proc. of the EMNLP'11, Edinburgh, ACL (2011) 1467–1478 5.






References and Bibliography

- [Singh et al., 2011] Sameer Singh, Amarnag Subramanya, Fernando Pereira, and Andrew McCallum. 2011. Large-scale cross-document coreference using distributed inference and hierarchical models. In*Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies Volume 1* (HLT '11), Vol. 1. Association for Computational Linguistics, Stroudsburg, PA, USA, 793-803.
- Suchanek, F., and Weikum, G. 2013. Knowledge harvesting in the big-data era. In Proceedings of the 2013 international conference on Management of data, 933–938. ACM.
- Suchanek, F. M.; Hoffart, J.; Kuzey, E.; and Lewis-Kelham, E. 2013. Yago2s: Modular highquality information extrac- tion with an application to flight planning. In *BTW*, 515–518.
- Suchanek, F. M.; Kasneci, G.; and Weikum, G. 2007. Yago: A Core of Semantic Knowledge. In *16th international World Wide Web conference (WWW 2007)*. New York, NY, USA: ACM Press.
- Talukdar, P. P.; Wijaya, D.; and Mitchell, T. 2012. Cou- pled temporal scoping of relational facts. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining (WSDM)*.





Carnegie Mellon University

References and Bibliography

• [Theobald & Weikum, 2012] Martin Theobald and Gerhard Weikum. From Information to Knowledge: Harvesting Entities and Relationships from Web Sources. Tutorial at PODS 2012

Federal University of São Carlos

IS Char

- Tseng, Y.-H.; Lee, L.-H.; Lin, S.-Y.; Liao, B.-S.; Liu, M.-J.; Chen, H.-H.; Etzioni, O.; and Fader, A. 2014. Chinese open relation extraction for knowledge acquisition. In *EACL*.
- [Verma & Hruschka Jr., 2012] S. Verma and E. R. Hruschka Jr. Coupled Bayesian Sets Algorithm for Semi-supervised Learning and Information Extraction. In Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD), 2012.
- [Wang & Cohen, 2007] Richard C. Wang and William W. Cohen: Language-Independent Set Expansion of Named Entities using the Web. In *Proceedings of IEEE International Conference on Data Mining* (ICDM 2007), Omaha, NE, USA. 2007.
- [Weikum et al., 2009] G. Weikum, G., Kasneci, M. Ramanath, F. Suchanek. DB & IR methods for knowledge discovery. Communications of the ACM 52(4), 2009.
- Wijaya, D.; Talukdar, P. P.; and Mitchell, T. 2013. Pidgin: Ontology alignment using web text as interlingua. In *Pro- ceedings of the Conference on Information and Knowledge Management (CIKM 2013)*. San Francisco, USA: Associa- tion for Computing Machinery.