



DOMAIN-SPECIFIC KNOWLEDGE-GRAPH CONSTRUCTION AND APPLICATIONS (WITH ALMOST NO SUPERVISION)

Maya Ramanath

IIT-Delhi

What are Knowledge-bases?

Subject	Predicate	Object
Albert_Einstein	bornIn	Germany
Albert_Einstein	typeOf	Theoretical_Physicist
Albert_Einstein	wonAward	Nobel_prize_in_Physics
Titanic	directedBy	James_Cameron
Titanic	genre	Drama
Topological_sorting	typeOf	Graph_algorithm
Group_Testing	applicationOf	Coding_theory
Gigabit_ethernet	highSpeedFormOf	Ethernet
Iron_python	adaptationOf	Python
Barack_Obama	supports	Military_Action_in_Libya
Donald_Trump	opposes	Illegal_immigration
UK	supports	Brexit
Scotland	opposes	Brexit
Pseudomonas	foundIn	Sea
Pseudomonas	actsOn	Acetone
Hydrocarbon_degradation	occursIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

Facts about people, movies

Concepts in Computer Science and their relationship with each other

Opinions about important(?) political topics

Facts about entities in bio-chemical engineering

Turn the Web...

news, books, entertainment, scientific content, advertisements, photos, videos, sound, social networks, blogs, tweets, opinions, comments, spam, junk...

...into a Knowledge-base

well organized, crisp, machine-readable, browseable, searchable, self-maintaining encyclopedia human knowledge.

Outline

- Introduction
- Applications
- Construction
- Conclusion

Stake Holders

- Bob_Ney
- Conor_Cummins
- Avram_Grant
- Barack_Obama
- Léo_Apotheker
- Capital_Punishment
- Bill_Graves
- Phil_Wheatley
- Campbell_Smith
- Ken_Salazar
- Jack_Kevorkian
- Kenneth_Clarke
- Muammar_al-Gad
- Jon_Snow
- Jim_Davis_(cartoonist)
- Newt_Gingrich
- Philip_Noel-Baker
- Amr_Moussa
- Ted_Kennedy
- Kings_of_Dublin

Opinions

supportive

Topics

- Pickens US energy plan
- Cosmetic Surgery
- Alcohol
- US Renewable Electricity
- EU constitution reform
- Models, minimum wage
- Campaign Finance Reform
- Mini-nukes
- Asylum seekers, welcome
- Gay Marriages
- Factory Farming
- Online gambling, ban
- Drilling in the Arctic National Wildlife Refuge
- Democracy: Pace of Development
- Internet Censorship
- Flag Burning, Prohibition
- Hydrogen vehicles
- University Tuition Fees
- "Clean coal"
- Ransom Payments, criticism

Opinions

Check All	Uncheck All	Save	Delete	Opinion Holder	Opinion Polarity	Opinion Subject
<input type="checkbox"/>	view	Obama	oppose	current ban on offshore drilling		
<input type="checkbox"/>	view	Barack Obama II	support	Reduction in Dependence on Foreign Oil		
<input type="checkbox"/>	view	Barack Obama II	support	Arctic National Wildlife Refuge Drilling Amendment		
				Amendment on Certain Energy-related		

Opinion Holder	Opinion Polarity	Opinion Subject
Iron_python	adaptationOf	Python
Barack_Obama	supports	Military_Action_in_Libya
Donald_Trump	opposes	Illegal_immigration
UK	supports	Brexit
Scotland	opposes	Brexit
Pseudonomas	foundIn	Sea
Pseudonomas	actsOn	Acetone
Hydrocarbon_degradation	occursIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

George W. Bush & Abortion

(+ -)

(- +)

(+)

(-)

(++)

(--)

Joe Biden & Guantanamo Bay

(+)

(-)

(++)

(--)

Applications of KBs

1. Opinion Base


- “Explore the space of opinions”

2. TeKnowBase (Academic search)

- How can I search for research articles discussing algorithms for **autoencoders**?

Subject	Predicate	Object
Albert_Einstein	bornIn	Germany
Albert_Einstein	typeOf	Theoretical_Physicist
Albert_Einstein	wonAward	Nobel_prize_in_Physics
Titanic	directedBy	James_Cameron
Titanic	genre	Drama
Topological_sorting	typeOf	Graph_algorithm
Group_Testing	applicationOf	Coding_theory
Gigabit_ethernet	highSpeedFormOf	Ethernet
Iron_python	adaptationOf	Python
Barack_Obama	supports	Military_Action_in_Libya
Donald_Trump	opposes	Illegal_immigration
UK	supports	Brexit
Scotland	opposes	Brexit
Pseudonomas	foundIn	Sea
Pseudonomas	actsOn	Acetone
Hydrocarbon_degradation	occurIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

“Autoencoder”

 Semantic Scholar All Fields


About 34,000 results [Last Five Years](#) [Lit Reviews](#) [Has PDF](#) [Has Video](#) [More Filters](#)

Extracting and composing robust features with denoising **autoencoders**
[Pascal Vincent](#), [Hugo Larochelle](#), [Yoshua Bengio](#), [Pierre-Antoine Manzagol](#) • ICML • 2008
Previous work has shown that the difficulties in learning deep generative or discriminative models can be overcome by an initial unsupervised learning step that maps inputs to useful intermediate... [\(More\)](#)
🔥 299 📊 478 📺 1 [View on ACM](#) [View Slides](#) [Cite](#) [Save](#)

Dynamic Pooling and Unfolding Recursive **Autoencoders for Paraphrase Detection**
[Richard Socher](#), [Eric H. Huang](#), [Jeffrey Pennington](#), [Andrew Y. Ng](#), [Christopher D. Manning](#) • NIPS • 2011
Paraphrase detection is the task of examining two sentences and determining whether they have the same meaning. In order to obtain high accuracy on this task, thorough syntactic and semantic analysis... [\(More\)](#)
🔥 81 📊 74 [View PDF](#) [Cite](#) [Save](#)

Semi-Supervised Recursive **Autoencoders for Predicting Sentiment Distributions**
[Richard Socher](#), [Jeffrey Pennington](#), [Eric H. Huang](#), [Andrew Y. Ng](#), [Christopher D. Manning](#) • EMNLP • 2011
We introduce a novel machine learning framework based on recursive **autoencoders** for sentence-level prediction of sentiment label distributions. Our method learns vector space representations for... [\(More\)](#)
🔥 104 📊 118 [View PDF](#) [Cite](#) [Save](#)

Stacked Denoising **Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion**
[Pascal Vincent](#), [Hugo Larochelle](#), [Isabelle Lajoie](#), [Yoshua Bengio](#), [Pierre-Antoine Manzagol](#) • J. Mach. Learn. Res. • 2010
We explore an original strategy for building deep networks, based on stacking layers of denoising **autoencoders** which are trained locally to denoise corrupted versions of their inputs. The resulting... [\(More\)](#)
🔥 370 📊 551 [View PDF](#) [Cite](#) [Save](#)

 Semantic Scholar All Fields

About 28,400 results [Last Five Years](#) [Lit Reviews](#) [Has PDF](#) [Has Video](#) [More Filters](#)

Extracting and composing robust features with denoising **autoencoders**
[Pascal Vincent](#), [Hugo Larochelle](#), [Yoshua Bengio](#), [Pierre-Antoine Manzagol](#) • ICML • 2008
Previous work has shown that the difficulties in learning deep generative or discriminative models can be overcome by an initial unsupervised learning step that maps inputs to useful intermediate... [\(More\)](#)
🔥 299 📊 478 📺 1 [View on ACM](#) [View Slides](#) [Cite](#) [Save](#)

Semi-Supervised Recursive **Autoencoders for Predicting Sentiment Distributions**
[Richard Socher](#), [Jeffrey Pennington](#), [Eric H. Huang](#), [Andrew Y. Ng](#), [Christopher D. Manning](#) • EMNLP • 2011
We introduce a novel machine learning framework based on recursive **autoencoders** for sentence-level prediction of sentiment label distributions. Our method learns vector space representations for... [\(More\)](#)
🔥 104 📊 118 [View PDF](#) [Cite](#) [Save](#)

Stacked Denoising **Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion**
[Pascal Vincent](#), [Hugo Larochelle](#), [Isabelle Lajoie](#), [Yoshua Bengio](#), [Pierre-Antoine Manzagol](#) • J. Mach. Learn. Res. • 2010
We explore an original strategy for building deep networks, based on stacking layers of denoising **autoencoders** which are trained locally to denoise corrupted versions of their inputs. The resulting... [\(More\)](#)
🔥 370 📊 551 [View PDF](#) [Cite](#) [Save](#)

Dynamic Pooling and Unfolding Recursive **Autoencoders for Paraphrase Detection**
[Richard Socher](#), [Eric H. Huang](#), [Jeffrey Pennington](#), [Andrew Y. Ng](#), [Christopher D. Manning](#) • NIPS • 2011
Paraphrase detection is the task of examining two sentences and determining whether they have the same meaning in order to obtain high accuracy on this task, thorough syntactic and semantic analysis... [\(More\)](#)
🔥 81 📊 74 [View PDF](#) [Cite](#) [Save](#)

Query + aspect: A few results

Query	Application	Algorithm	Implementation
cryptography	Cryptographic Protection of Computer-Based Data Files	Cryptographic algorithm on multicore processor: A review	Design of a Reconfigurable Hardware for Efficient Implementation of Secret Key and Public Key Cryptography
autoencoder	Exploring autoencoders for unsupervised feature selection	Training Stacked Denoising Autoencoders for Representation Learning	Parallelizing the Sparse Autoencoder
neural_network	Application of Particle Swarm Optimization and RBF Neural Network in Fault Diagnosis of Analogue Circuits	Evolutionary Neural Network Based on New Ant Colony Algorithm	Design of FPGA based general purpose neural network
hashing	Sampling Based N-Hash Algorithm for Searching Frequent Itemset	An Alternative Analysis of the Open Hashing Algorithm	Low Power And Area Efficiency of SHA-1

Aspect-based Search [ECIR 2020]

- A LM-based method for aspect-aware search

$$MM(w) = \lambda P(w|a) + (1 - \lambda) P(w|q, a)$$

Prob. of term, given aspect

“Application”

Prob. of term, given both query and aspect

“Autoencoder”+

Topological_sorting	typeOf	Graph_algorithm
Group_Testing	applicationOf	Coding_theory
Gigabit_ethernet	highSpeedFormOf	Ethernet
Iron_python	adaptationOf	Python
Feature_selection	applicationOf	Autoencoders
Heap_sort	algorithmFor	Sorting

Approach	Algorithm			Application			Implementation		
	DCG@5	P@5	P@1	DCG@5	P@5	P@1	DCG@5	P@5	P@1
MM	6.27	0.70	0.75	2.64	0.45	0.47	2.33	0.44	0.40
QL+query	2.69	0.3	0.33	1.42	0.25	0.22	1.05	0.16	0.23
QL+query+aspect	5.03	0.56	0.59	2.38	0.41	0.35	1.92	0.30	0.43
QL+query+aspect+QE	5.12	0.58	0.61	2.5	0.43	0.41	2.29	0.37	0.49

$P|$

Applications of KBs

1. Opinion Base

- “Explore the space of opinions”

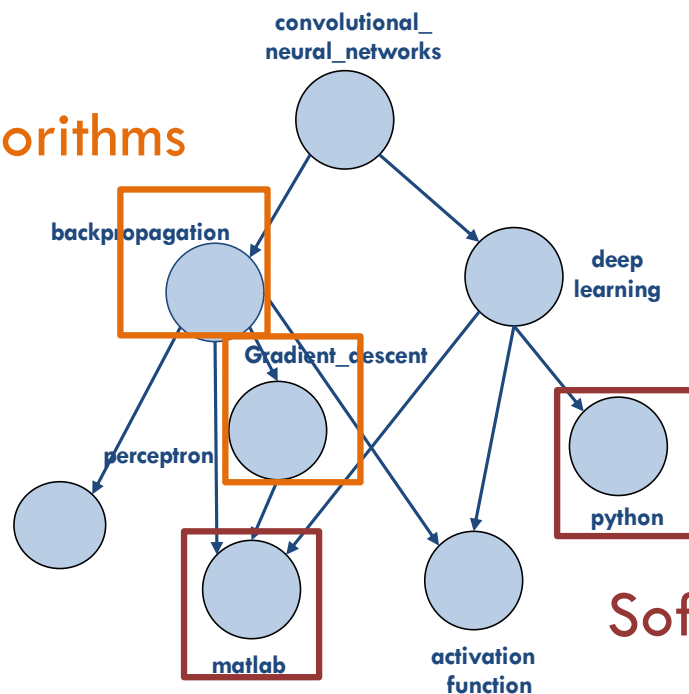
2. TeKnowBase (Academic search)

- How can I search for research articles discussing algorithms for **autoencoders**?
- What are the pre-requisites to learn about **convolutional neural networks**?

Subject	Predicate	Object
Albert_Einstein	bornIn	Germany
Albert_Einstein	typeOf	Theoretical_Physicist
Albert_Einstein	wonAward	Nobel_prize_in_Physics
Titanic	directedBy	James_Cameron
Titanic	genre	Drama
Topological_sorting	typeOf	Graph_algorithm
Group_Testing	applicationOf	Coding_theory
Gigabit_ethernet	highSpeedFormOf	Ethernet
Iron_python	adaptationOf	Python
Barack_Obama	supports	Military_Action_in_Libya
Donald_Trump	opposes	Illegal_immigration
UK	supports	Brexit
Scotland	opposes	Brexit
Pseudonomas	foundIn	Sea
Pseudonomas	actsOn	Acetone
Hydrocarbon_degradation	occurIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

Pre-requisites for “convolutional_neural_networks”

Algorithms



Software

Prerequisites for which aspect?

- **Software:** matlab, python
- **Activation functions:** sigmoid, relu, tanh, softmax
- **Applications:** face perception, object identification
- **Algorithms:** backpropagation, optimization, gradient descent

Subject	Predicate	Object
Topological_sorting	typeOf	Graph_algorithms
Group_Testing		
Gigabit_ethernet		
Iron_python		
Bubble_sort		
Sorting	typeOf	permutation
Sorting	typeOf	Computation
Feature_selection	applicationOf	Autoencoders
Bubble_sort	hasComplexity	$O(n^2)$
Backpropagation	algorithmFor	Neural_network
CNN	typeOf	Neural_network

Techniques	Precision
PreFace	0.76
QDMKB + RefD	0.636
RefD + TKB	0.68

Applications of KBs

1. Opinion Base

- “Explore the space of opinions”

2. TeKnowBase (Academic search)

- How can I search for research articles discussing algorithms for **autoencoders**?
- What are the pre-requisites to learn about **convolutional neural networks**?

3. Biochemical KB

- “Determine if a microbe is aerobic or anaerobic”

Subject	Predicate	Object
Albert_Einstein	bornIn	Germany
Albert_Einstein	typeOf	Theoretical_Physicist
Albert_Einstein	wonAward	Nobel_prize_in_Physics
Titanic	directedBy	James_Cameron
Titanic	genre	Drama
Topological_sorting	typeOf	Graph_algorithm
Group_Testing	applicationOf	Coding_theory
Gigabit_ethernet	highSpeedFormOf	Ethernet
Iron_python	adaptationOf	Python
Barack_Obama	supports	Military_Action_in_Libya
Donald_Trump	opposes	Illegal_immigration
UK	supports	Brexit
Scotland	opposes	Brexit
Pseudomonas	foundIn	Sea
Pseudomonas	actsOn	Acetone
Hydrocarbon_degradation	occursIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

Determining the nature of microbes

- **Given:** corpus of research articles
- **Goal:** determine all microbes which are aerobic/anaerobic
- **Challenge:** Cannot extract “isAerobic” directly
- **Method:** Extract triples from the corpus and organise the information, *reason* about the “isAerobic” relation

Pseudonomas	foundIn	Sea
Pseudonomas	actsOn	Acetone
Hydrocarbon_degradation	occurIn	Sea
Hydrocarbon_degradation	carriedOn	Acetone

Interesting entity types

Microbes

Substrate

Enzyme

Process

Environment

Nutrient

Property

Outline

- Introduction
- Applications
- Construction
- Conclusion

General principles for domain-specific KB construction

Source identification

- Structured sources
 - Wikipedia, IMDB, domain-specific websites such as Webopedia...
- Unstructured sources
 - Newspaper articles, blogs, forums, reviews, tweets, comments...

Entities

- Extraction from structured sources
 - Design of appropriate heuristics and/or wrappers
- Extraction from unstructured sources
 - POS tagging + heuristics
- Types known/unknown
- Cleanup
 - Entity resolution to reconcile duplicate entity extractions

Relations

- From both structured and unstructured sources:
 - Known relations
 - Unknown relations
- Cleanup
 - Identifying and reconciling duplicate relationship names

Refinement

- “Complete” the KB by inferring new relations.
 - Simple string matching sometimes does the trick!
 - Graph embeddings are the method of choice (hype)
- Enhance the KB through feedback

Using the general principles

	Sources	Entity extraction techniques	Relation extraction techniques	Refinement
Opinion-base [WSDM 2012, CIKM 2012, CIKM workshops 2012]	Newspaper articles, Wikipedia, Debatepedia	<ol style="list-style-type: none"> Names of entities from YAGO, Debatepedia, <i>known types</i> Phrase extraction for topics, <i>unknown types</i> 	Known relations (“support”/“oppose”), unstructured sources	KB enhancement by feeding back new surface patterns for the relations
TeKnowbase [WWW workshops 2018, ISWC 2020, ECIR 2020, ECIR 2021]	Wikipedia, Webopedia, Techtarget, online textbooks	<ol style="list-style-type: none"> Wikipedia, Webopedia, TechTarget article names, <i>unknown types</i> Indexes of online textbooks, <i>unknown types</i> 	<ol style="list-style-type: none"> Known relations, structured sources KR, unstructured sources Unknown relations, SS UR, US 	<ol style="list-style-type: none"> Sub-string extraction <ul style="list-style-type: none"> “Graph algorithm” Embeddings for knowledge-base completion
Biochemical KB [Report available on request]	Research articles from the domain, BRENDA, MicrobesOnline, etc.	<ol style="list-style-type: none"> Expert-curated list Web-sources, <i>known types</i> POS-tagging + heuristics from the research articles 	<ol style="list-style-type: none"> Known relations, unstructured sources <ul style="list-style-type: none"> Training ML models after acquiring examples 	None
Event-base	Newspaper articles	<ol style="list-style-type: none"> Unknown entities, therefore, POS-tagging + heuristics 	<ol style="list-style-type: none"> Requires <i>nested triples</i>. Partially known relations, unstructured sources 	None

[End](#)

OpinioNet

Dictionary
creation

Relation
extraction

Entity and
facet
extraction

Subtopics

Hamid Karzai

Article Talk

From Wikipedia, the free encyclopedia

"Karzai" redirects here. For the surname, see *Karzai (surname)*.

Hamid Karzai (/ˈhæmɪd ˈkɑːrzaɪ/; Pashto/Persian: حامد کرزی, Pashto pronunciation: [ˈhɑmɪd kɑɾzɛˈiː]; born 24 December 1957) is an Afghan

Afghan peace process

Article Talk

From Wikipedia,

Peace process
1978 Saur Rev

Disarmament, demobilisation and reintegration (DDR), or disarmament, demobilisation, repatriation, reintegration and resettlement (DDRRR) are strategies used as a component of peace processes.

Stake-holders, Topics

Afghan President Says He Supports Talks With Taliban : The Two ...

www.npr.org/.../afghan-president-says-he-supports-talks-with-taliban

4 Jan 2012 – Experts say the plans for talks are a positive step toward a future peace in Afghanistan.

**Known relations,
unstructured source**

Hamid_Karzai

Afghan President

Karzai

Source	Overall	Support	Oppose	ALJ	BBC	CNN	GN	NY	WP
#Opinions	29648	16011	13637	970	3349	2650	6861	9877	5941
#Evaluated	2005	1121	884	364	349	308	301	364	319
Precision	0.724	0.70	0.75	0.72	0.66	0.78	0.788	0.74	0.69

Table 3.8: Opinions & precision results

War on

Afghanistan_War

Iraq_War

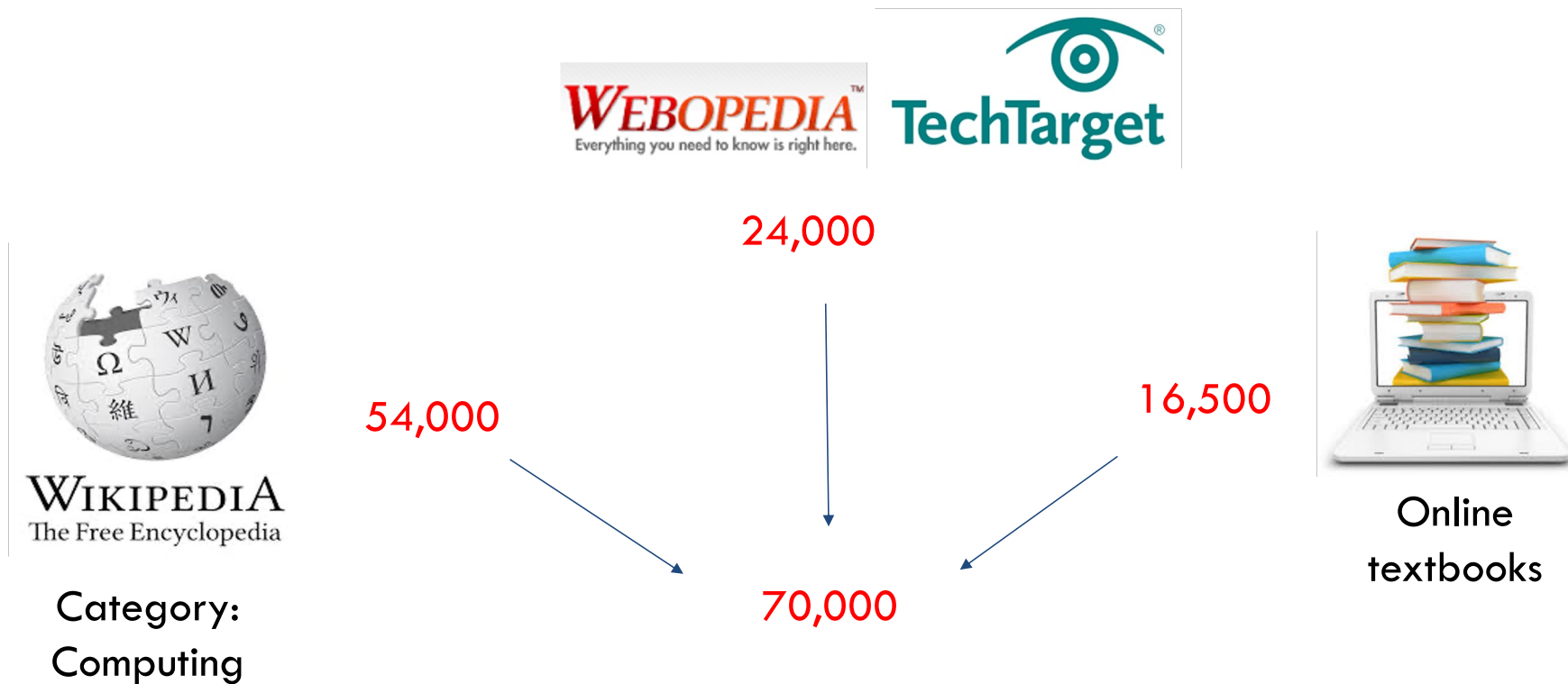
Afghan_Peace_Proc
ss

Drone_
Strikes

Suicide_Bombings

[End](#)

Constructing TeKnowbase: Dictionary of entities



TeKnowbase [WWW Workshops 2019]

*Dataset available from <https://github.com/prajnaupadhyay/TeKnowbase>

Structured sources, known relations

typeOf	subtopicOf		applicationOf		algorithmFor		techniqueIn																														
<div>List of data structures</div> <div>Trees [edit]</div> <div>Main article: <i>Tree (data structure)</i></div> <div>Binary trees [edit]</div> <div><ul style="list-style-type: none">AA treeAVL treeBinary search treeBinary treeCartesian treeLeft-child right-sibling binary treeOrder statistic treePagodaRandomized binary search treeRed-black tree</div>	<div>List of set theory topics</div> <div>Articles on individual set theory topics</div> <div><ul style="list-style-type: none">Algebra of setsAxiom of choice<ul style="list-style-type: none">Axiom of countable choiceAxiom of dependent choiceZorn's lemmaBoolean-valued modelBurali-Forti paradoxCantor's back-and-forth method</div>		<div>Coding theory</div> <div>From Wikipedia, the free encyclopedia</div> <div>Contents [hide]</div> <div><ol style="list-style-type: none">History of coding theorySource coding<ul style="list-style-type: none">2.1 Definition2.2 Properties2.3 Principle2.4 ExampleChannel coding<ul style="list-style-type: none">3.1 Linear codes<ul style="list-style-type: none">3.1.1 Linear block codes3.1.2 Convolutional codesCryptographical codingLine codingOther applications of coding theory<ul style="list-style-type: none">6.1 Group testing6.2 Analog coding</div>		<div>Derivative-free optimization</div> <div>Algorithms [edit]</div> <div><ul style="list-style-type: none">DONEBayesian optimizationCoordinate descent and adaptive coordinate dCuckoo searchEvolution strategies, Natural evolution strategiSNES)Genetic algorithmsLIPO algorithmMCS algorithmNelder-Mead method</div>		<div>Document classification</div> <div>Techniques [edit]</div> <div>Automatic document classification techniques include:</div> <div><ul style="list-style-type: none">Expectation maximization (EM)Naive Bayes classifiertf-idfInstantaneously trained neural networksLatent semantic indexingSupport vector machines (SVM)Artificial neural networkK-nearest neighbour algorithmsDecision trees such as ID3 or C4.5</div>																														
<table><tr><td>Binary_sea rch_tree</td><td>typeOf</td><td>Binary_ tree</td></tr><tr><td>Binary_tree</td><td>typeOf</td><td>Tree</td></tr><tr><td>Tree</td><td>typeOf</td><td>Data_st ructure</td></tr></table>	Binary_sea rch_tree	typeOf	Binary_ tree	Binary_tree	typeOf	Tree	Tree	typeOf	Data_st ructure	<table><tr><td>Boolean- valued_mo del</td><td>subtopic Of</td><td>Set_theory</td></tr><tr><td>Algebra_o f_sets</td><td>subtopic Of</td><td>Set_theory</td></tr></table>	Boolean- valued_mo del	subtopic Of	Set_theory	Algebra_o f_sets	subtopic Of	Set_theory	<table><tr><td>Group_test ing</td><td>appl Of</td><td>Coding_th eory</td></tr><tr><td>Analog_co ding</td><td>appl Of</td><td>Coding_th eory</td></tr></table>	Group_test ing	appl Of	Coding_th eory	Analog_co ding	appl Of	Coding_th eory	<table><tr><td>Cuckoo_se arch</td><td>algo For</td><td>Derivative -free_opt</td></tr><tr><td>MCS_algo rithm</td><td>algo For</td><td>Derivative -free_opt</td></tr></table>	Cuckoo_se arch	algo For	Derivative -free_opt	MCS_algo rithm	algo For	Derivative -free_opt	<table><tr><td>Naïve- Bayes_classifi er</td><td>techIn</td><td>Document_cl assification</td></tr><tr><td>K- nearest_neig hbor_algorith m</td><td>techIn</td><td>Document_cl assification</td></tr></table>	Naïve- Bayes_classifi er	techIn	Document_cl assification	K- nearest_neig hbor_algorith m	techIn	Document_cl assification
Binary_sea rch_tree	typeOf	Binary_ tree																																			
Binary_tree	typeOf	Tree																																			
Tree	typeOf	Data_st ructure																																			
Boolean- valued_mo del	subtopic Of	Set_theory																																			
Algebra_o f_sets	subtopic Of	Set_theory																																			
Group_test ing	appl Of	Coding_th eory																																			
Analog_co ding	appl Of	Coding_th eory																																			
Cuckoo_se arch	algo For	Derivative -free_opt																																			
MCS_algo rithm	algo For	Derivative -free_opt																																			
Naïve- Bayes_classifi er	techIn	Document_cl assification																																			
K- nearest_neig hbor_algorith m	techIn	Document_cl assification																																			

More relation extractions

Structured Source, Unknown relation

Template:Databases

From Wikipedia, the free encyclopedia

V • T • E	
Types	Object-oriented (comparison) • Relational (comparison) • Database
Concepts	Database • ACID • Armstrong's axioms • CAP theorem • CRU
Objects	Relation (table • column • row) • View • Transaction • Transa
Components	Concurrency control • Data dictionary • JDBC • XQJ • ODBC •

View	Object	Databases
Data_dictionary	Component	Databases

Unstructured Source, Known relation

synonymOf
 “is abbreviation for”
 “bfs (breadth-first_search)”
 “is short for”

BFS	synonymOf	Breadth-first_search
RegEx	synonymOf	Regular_expression

Unstructured Source, Unknown relation

Step	Action	No. of triples	Problem
1	Run OpenIE	400,000	Non-technical entities in the result
2	Remove entities not in dictionary	300,000	Long phrases containing entity
3	Retain triples with 50% match	100,000	Generic entities
4	Remove triples where argument start with “The”	3520	Accuracy still around 65%

gigabit_ethernet	is_a_high_speed_for	ethernet
ironpython	is_an_adaption_of	python
utc	uses	gregorian_calendar

Refinement using KG embeddings

Xbasic	Typeof	Programming_language
Binomial_heap	Typeof	Tree
Palmdos	Typeof	Operating_system
Delayed_column_generation	Typeof	Convex_programming

TeKnowbase statistics

No. of unique entities	70,285
No. of unique relations	2,574
Taxonomic relations (typeOf)	27,078
Total no. of triples	146,657
No. of overlapping entities with DBPedia	17,987
No. of overlapping entities with Freebase	34,785
No. of triples extracted from Wikipedia	99,357
No. of triples extracted from Unstructured sources	3,506

#	Relation (rows 1–5)	# Evaluated triples	Accuracy
1.	typeOf	515	99.0% \pm 0.8%
2.	terminologyOf	676	98.9% \pm 0.7%
3.	synonymOf	70	100% \pm 0.0%
4.	subTopicOf	42	91.3% \pm 8.2%
5.	conceptIn	334	95.4% \pm 2.1%
6.	<i>Unstructured sources</i>	435	63.2% \pm 3.7%
7.	<i>Inferencing with NTN</i>	428	64.2% \pm 4.5%

Conclusion

- Applications

- Opinion-base (*heat maps, correlated opinions, flip-floppers*)
- TeKnowbase (*academic search, pre-requisite construction*)
- Biochemical KB (*reasoning – incomplete*)

- Construction

- General principles (*most important: identify good sources, construct a dictionary of interesting entities*)

Thanks: Rawia Awadallah, Srikanta Bedathur, Tanmoy Chakraborty,
Prajna Upadhyay, Manas Joshi, Shubham Singla, Ashutosh Bindal,
Manjeet Kumar, Gerhard Weikum