https://www.mpi-inf.mpg.de/www-2022-tutorial

Completeness, Recall, and Negation in Open-World Knowledge Bases



Simon Razniewski, Hiba Arnaout, Shrestha Ghosh, Fabian Suchanek







On the Limits of Machine Knowledge: Completeness, Recall and Negation in Web-scale Knowledge Bases

Simon Razniewski, Hiba Arnaout, Shrestha Ghosh, Fabian Suchanek

- 1. Introduction & Foundations (Simon) 15:45-16:00 CEST
- 2. Predictive recall assessment (Fabian) 16:00-16:20
- 3. Counts from text and KB (Shrestha) 16:20-16:40
- 4. Negation (Hiba) 16:40-17:00
- 5. Wrap-up (Simon) 17:00-17:10







Machine knowledge in action



your	numo	country
1901	Wilhelm Conrad Röntgen	Germany
1902	Hendrik Antoon Lorentz	Netherlands
1902	Pieter Zeeman	Netherlands
View 213	more rows	

https://www.research-in-germany.org > nobel-laureates -

German Nobel laureates - Research in Germany

J. Georg Bednorz: 1987 - Physics ... An unusual approach made Georg Bednorz a pioneer in the field of superconductivity – and **Physics Nobel Prize laureate** in ...



Machine knowledge in action



Machine knowledge is awesome

- Reusable, scrutable asset for knowledge-centric tasks
 - Semantic search & QA
 - Entity-centric text analytics
 - Distant supervision for ML
 - Data cleaning
- Impactful projects at major public and commercial players
 - Wikidata, Google KG, Microsoft Satori, ...
- Strongly rooted in semantic web community
 - Linked data, vocabularies, ontologies, indexing and querying,

But: Machine Knowledge is incomplete

Google	marie curie prize	S			×	٩		
	Q All 🖾 Image	s 🗉 News 🖓) Maps 🧷 Shopp	ing : More	Тс	pols		
	Awards/Mar	ie Curie						
							Nobel Prize	•
	Davy Medal	Matteucci Medal	Elliott Cresson Medal	Albert Medal	Actonian Prize	Willard Gibbs Award	• (2x)	•





Wikidata KB:

Semantic Web Journal has only published 84 articles ever

<u>https://scholia.toolforge.org/venue/Q15817015</u>

Only 7 papers ever deal with the topic "web science"

<u>https://scholia.toolforge.org/topic/Q579439</u>

But: Machine knowledge is one-sided



- In KB:
 - Nicola Tesla received title of IEEE fellow
 - Vietnam is a member of ASEAN
 - iPhone has 12MP camera
- Not in KB:
 - Nicola Tesla did not receive the Nobel Prize
 - Switzerland is not a member of the EU
 - *iPhone 12 has no headphone jack*

Why is this problematic? (1) Querying

- Decision making more and more data-driven
- Analytical queries paint wrong picture of reality
 - E.g., SW journal deemed too small
- Instance queries return wrong results
 - E.g., wrongly assuming certain topic is of no interest

Why is this problematic? (2) Data curation

- Effort priorization fundamental challenge in human-in-the-loop curation
 - Should we spend effort on obtaining data for SWJ or TWeb?
- Risk of effort duplication if not keeping track of completed areas
 - Spending effort on collecting data ... already present

Why is this problematic? (3) Summarization and decision making

Booking.com 🗳 Bathroom Safety & security Sellness 😵 Fitness Fire extinguishers Toilet paper Full body massage Towels CCTV outside property CCTV in common areas Additional charge Private bathroom Hand massage Additional charge Toilet Smoke alarms Head massage Additional charge Free toiletries 24-hour security Safety deposit box Hairdryer Shower Foot massage Additional charge (i) General Neck massage Additional charge E Bedroom Paid WiFi Back massage Additional charge Iinen Mini-market on site Spa/wellness packages Wardrobe or closet Vending machine (drinks) Steam room Alarm clock Designated smoking area Spa Facilities Air conditioning Room Amenities Light therapy ree room Soc No free WiFi! Facial treatments Flat-screen TV Ironing facilities Beauty Services Non-smoking rooms Satellite channels Sun loungers or beach chairs Radio Iron Pool/beach towels Air conditioning Telephone Hot tub/iacuzzi 🖉 тү Accessibility Massage Additional charge Pay-per-view channels Spa and wellness centre Visual aids: Tactile signs

Visual aids: Braille I ower bathroom sink

Higher level toilet

Toilet with grab rails

Wheelchair accessible

- On-site coffee house Chocolate or cookies
- Additional charge

Food & Drink

Ľ

×.

Pets a

applic • /

Fruits Additional charge

Couples massage Additional charge

- Additional charge
- Fitness centre
- Sauna Additional charge
- Languages spoken
- English



Camera

- Pro 12MP camera system: Ultra Wide, Wide, and Telephoto cameras
- Ultra Wide: f/2.4 aperture and 120° field of view
- Wide: f/1.6 aperture
- Telephoto: f/2.2 aperture
- 2.5x optical zoom in, 2x optical zoom out; 5x optical zoom range
- Digital zoom up to 12x
- Night mode portraits enabled by LiDAR Scanner
- Portrait mode with advanced bokeh and Depth Control
- Portrait Lighting with six effects (Natural, Studio, Contour, Stage, Stage Mono, High-Key Mono)
- Dual optical image stabilization (Wide and Telephoto)
- Sensor-shift optical image stabilization
- Five-element lens (Ultra Wide): six-element lens (Telephoto); seven-element lens (Wide)
- Brighter True Tone flash with Slow Sync
- Panorama (up to 63MP)
- Sapphire crystal lens cover
- 100% Focus Pixels (Wide)
- Night mode (Ultra
- Deep Fusion (

No headphone jack

- Sensor-shift optical image stabilization for video (Wide)
- Optical image stabilization for video (Wide)
- 2.5x optical zoom in, 2x optical zoom out: 5x optical zoom range
- Digital zoom up to 7x
- Audio zoom
- Brighter True Tone flash
- OuickTake video
- Slo-mo video support for 1080p at 120 fps or 240 fps
- Time-lapse video with stabilization
- Night mode Time-lapse
- Extended dynamic range for video up to 60 fps
- Cinematic video stabilization (4K, 1080p, and 720p)
- Continuous autofocus video

Topic of this tutorial

How to know how much a KB knows?

How to = techniques How much knows = completeness/recall/coverage bookkeeping/estimation KB = General world knowledge repository

What this tutorial offers

- Logical foundations
 - Setting and formalisms for describing KB completeness (part 1)
- Predictive assessment
 - How (in-)completeness can be statistically predicted (Part 2)
- Count information
 - How count information enables (in-)completeness assessment (Part 3)
- Negation
 - How salient negations can be derived from incomplete KBs (Part 4)

Goals:

- 1. Systematize the topic and its facets
- 2. Lay out assumptions, strengths and limitations of approaches
- 3. Provide a practical toolsuite

15

What this tutorial is NOT about

- Knowledge base completion (KBC)
 - "How to make KBs more complete"
- Related: Understanding of completeness is needed to know when/when not to employ KBC
 - KBC naively is open-ended
 - → Understanding of completeness needed to "stop"
- But:
 - Heuristic, error-prone KBC not always desired
 - Completeness awareness != actionable completion
- Literature on knowledge graph completion, link prediction, missing value imputation, etc.
 - E.g., Rossi, Andrea, et al. <u>Knowledge graph embedding for link prediction: A comparative analysis</u> *TKDD 2021*

Beatles members: John Lennon 36% Paul McCartney 23% George Harrison 18% Bob Dylan 5% Ringo Starr 3% Elvis Presley 2% Yoko Ono 2%



On the Limits of Machine Knowledge: Completeness, Recall and Negation in Web-scale Knowledge Bases

Simon Razniewski, Hiba Arnaout, Shrestha Ghosh, Fabian Suchanek

- 1. Introduction & Foundations (Simon) 15:45-16:00 CEST
- 2. Predictive recall assessment (Fabian) 16:00-16:20
- 3. Counts from text and KB (Shrestha) 16:20-16:40
- 4. Negation (Hiba) 16:40-17:00
- 5. Wrap-up (Simon) 17:00-17:10







Knowledge base - definition

Given set E (entities), L (literals), P (predicates)

- Predicates are positive or negated properties
 - bornIn, notWonAward, ...
- An assertion is a triple (s, p, o) $\in \mathbf{E} \times \mathbf{P} \times (\mathbf{E} \cup \mathbf{L})$
- A practically available KB K^a is a set of assertions
- The ``ideal'' (complete) KB is called \mathbf{K}^{i}
- Available KBs are incomplete: K^a ⊆ Kⁱ

Knowledge bases (KBs aka. KGs) subject-predicate-object triples about entities, attributes of and relations between entities predicate (subject, object) type (Marie Curie, physicist) subtypeOf (physicist, scientist) placeOfBirth (Marie Curie, Warsaw) residence (Marie Curie, Paris)

¬placeOfBirth (Marie Curie, France)

discovery (Polonium, 12345) discoveryDate (12345, 1898) discoveryPlace (12345, Paris) discoveryPerson (12345, Marie Curie)

atomicNumber (Polonium, 84) halfLife (Polonium, 2.9 y) + composite objects

taxonomic knowledge

factual knowledge

spatio-temporal & contextual knowledge

expert knowledge

KB incompleteness is inherent



Resulting challenges

Available KBs are incomplete K^a << Kⁱ

Available KBs hardly store negatives K^{a⁻} ≈ Ø

Formal semantics for incomplete KBs: Closed vs. open-world assumption

won			
name	award		
Brad Pitt	Oscar		
Marie Curie	Nobel Prize		
Berners-Lee	Turing Award		

	Closed-world assumption	Open-world assumption
won(BradPitt, Oscar)?	→ Yes	→ Yes
von(Pitt, Nobel Prize)?	→No	→ Maybe

- Databases traditionally employ closed-world assumption
- KBs (semantic web) necessarily operate under open-world assumption 21



The logicians way out – completeness metadata

 Need power to express both maybe and no

(Some paradigm which allows both open- and closed-world interpretation of data to co-exist)

• Approach: Completeness assertions [Motro 1989]



won(Pitt, Turing)? \rightarrow Maybe (OWA)

The power of completeness metadata

Know what the KB knows:

 \rightarrow Locally, K^a = Kⁱ

Absent assertions are really false:

 \rightarrow Locally, s $\neg \in K^a$ implies s $\neg \in K^i$

Completeness metadata: Formal view

Complete (won(name, award); award = 'Nobel')

Implies constraint on possible state of K^a and Kⁱ

wonⁱ(name, 'Nobel') \rightarrow won^a(name, 'Nobel')

(tuple-generating dependency)

Darari et al. <u>Completeness Statements about RDF Data</u> <u>Sources and Their Use for Query Answering</u> ISWC 2013 25

Cardinality assertions: Formal view

- "Nobel prize was awarded 603 times"
- \rightarrow |wonⁱ(name, 'Nobel') | = 603
- \rightarrow Allows counting objects in K^a
 - Equivalent count \rightarrow Completeness assertion
 - Otherwise, fractional coverage/recall information
 - "93% of awards covered"
- Grounded in number restrictions/role restrictions in Description Logics

B. Hollunder and F. Baader Qualifying Number Restrictions in Concept Languages KR 1991

Formal reasoning with completeness metadata



Where can completeness metadata come from?

- Data creators should pass them along as metadata
- Or editors should add them in curation steps



This is a complete list of compositions by Maurice Ravel,

:	28	Tout est lumière	soprano, mixed choir, and orchestra	1901	Prix de Rome competition
	29	Myrrha, cantata	soprano, tenor, baritone, and orchestra	1901	text: Fernand Beissier; • Prix de Rome competition
	31	Semiramis	cantata	1902	student competition;partially lost

• E.g., COOL-WD tool

Darari et al. <u>COOL-WD: A Completeness</u> <u>Tool for Wikidata</u> ISWC 2017 ²⁸

(i) cool-wd.inf.unibz.it/?p=Q22686				
COL-MAN In Analytics Query Search entity				
residence (P551)	White House	?		
country of citizenship (P27) United States of America				
child (P40)	Ivanka Trump			
	Donald Trump Jr.			
	Eric Trump	~		
	Tiffany Trump			
	Barron Trump			
field of work (P101)	politics			
	government			

But...

Requires human effort

- Soliciting metadata more demanding than data
- Automatically created KBs do not even have editors

Remainder of this tutorial:

How to automatically acquire information about what a KB knows

Takeaway Part 1: Foundations

- KBs are pragmatic collections of knowledge
 - Issue 1: Inherently incomplete
 - Issue 2: Hardly store negative knowledge
- Open-world assumption (OWA) as formal interpretation leads to counterintuitive results
- Metadata about completeness or counts as way out

Next: How to use predictive models to derive completeness metadata

On the Limits of Machine Knowledge: Completeness, Recall and Negation in Web-scale Knowledge Bases

Simon Razniewski, Hiba Arnaout, Shrestha Ghosh, Fabian Suchanek

- 1. Introduction & Foundations (Simon) 15:45-16:00 CEST
- 2. Predictive recall assessment (Fabian) 16:00-16:20
- 3. Counts from text and KB (Shrestha) 16:20-16:40
- 4. Negation (Hiba) 16:40-17:00
- 5. Wrap-up (Simon) 17:00-17:10







Wrap-up: Take-aways



- 1. KBs are incomplete and limited on the negative side
- 2. Predictive techniques work from a surprising set of paradigms
- 3. Count information a prime way to gain insights into completeness/coverage
- 4. Salient negations can be heuristically materialized
- 5. Relative completeness tangible alternative

Wrap-up: Recipes

Ab-initio KB construction

- 1. Intertwine data and metadata collection
- 2. Human insertion: Provide tools
- 3. Automated extraction: Learn from extraction context

• KB curation

- 1. Exploit KB-internal or textual cardinality assertions
- 2. Inspect statistical properties on density or distribution
- 3. Compute overlaps on pseudo-random samples

Open research questions

- 1. How are entity, property and fact completeness related?
- 2. How to distinguish salient negations from data modelling issues?
- 3. How to estimate coverage of knowledge in pre-trained language models?
- 4. How to identify most valuable areas for recall improvement?

Wrap-up: Wrap-up

- KBs major drivers of knowledge-intensive applications
- Severe limitations concerning completeness and coverage-awareness
- This tutorial: Overview of problem, techniques and tools to obtain awareness of completeness

Takeaway Part 1: Foundations

- KBs are pragmatic collections of knowledge
 - Issue 1: Inherently incomplete
 - Issue 2: Hardly store negative knowledge
- Open-world assumption (OWA) as formal interpretation leads to counterintuitive results
- Metadata about completeness or counts as way out

https://www.mpi-inf.mpg.de/www-2022-tutorial

Takeaway: Predictive recall assessment

Using statistical techniques, we can predict more or less

- the recall of facts
 - are we missing objects for a subject?
 - do all subjects have an attribute in the real world?
 - does a text enumerate all objects for a subject?
- the recall of entities
- is the distribution of entities representative?
- how many entities are in the real world?

Takeaway: negation Takeaway: Counts from text and KB 64 Count information comes in two variants 1. Current KBs lack negative knowledge Counting predicates - store integer counts 0 Enumerating predicates - store entities Rising interest in the explicit addition of negation to OW KB. ٠ Count information in text 2 Negations highly relevant in many applications including: occurs as cardinals, ordinals, non-numeric noun phrases Commercial decision making (e.g., hotel booking) 0 General-domain question answering systems (e.g., is Switzerland a member occurs with compositional cues of the EU?) Count information in KBs 3. Methodologies include: is expressed in two variants 0 Statistical inference o occurs semantically related count predicates Text extraction Pretrained LMs. Count information 4. can enrich KB 0 highlight inconsistencies