



Algorithms for Knowledge and Information Extraction in Text with Wikipedia

Marco Ponza

Supervisor

Prof. Paolo Ferragina



Menu

0. Introduction

Knowledge and Information Extraction

1. Algorithms for Entity Relatedness
2. Algorithms for Entity and Fact Salience

Applications

3. Algorithms for Expert Finding
4. Future Research Directions

0

Introduction

Introduction

- ▷ Enhancing the humankind progress with new **intelligent** technologies



Tools that can afford general- or specific-domain **tasks** with **performance close or better than humans**

- ▷ Machines need of access, read and **understand** information stored in **data archives**



The **dominant** form on which information is produced every day by humans is still **Natural Language**



The New York Times

Introduction

Text Understanding

Easy for humans



Hard for machines



Introduction

Text Understanding

Leonardo is the scientist who painted Mona Lisa



Science



Italy



Leonardo da Vinci



Cartography



Art



Louvre



Mona Lisa (painting)



Renaissance



Florence

1

Map ambiguous words into the **real-world entities** they refer to as well as **contextualize** them together with **related** entities

Introduction

Text Understanding

Leonardo is the scientist who painted Mona Lisa



(“Leonardo”, “is”, “scientist”)

(“Leonardo”, “painted”, “Mona Lisa”)

2

Structure multiple **facts** (propositions) contained in the sentence

Triplets of **(subject, relation, object)**

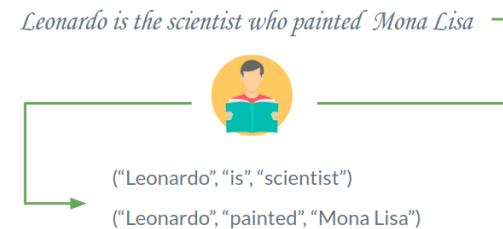
Introduction

Text Understanding



1

Map words into the **concepts** they refer to



2

Structure multiple **facts** (propositions) contained in the sentence
Triplets of (**subject, relation, object**)

How can we do that?

Humans can interpret words in a *larger context* hinging onto their **background** and **linguistic knowledge** (Gabrilovich, SIGIR'16)

Detect (1) unambiguous **entities** (2) **facts**, and (3) quantifying how much they are **related**

Introduction

Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) *A text is a vector of ambiguous keywords*

Limitations

- Curse of Dimensionality
- Synonymy and Polysemy problems of keywords
- No understanding of real-world entities
- Structure of the sentence is lost (no facts)

Introduction

Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) *A text is a vector of ambiguous keywords*
 - Limitations {
 - Curse of Dimensionality
 - Synonymy and Polysemy problems of keywords
 - No understanding of real-world entities
 - Structure of the sentence is lost (no facts)
- ▷ **LDA/LSI** (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations
 - A text is mapped into a latent space (vector of floating-points)*

Introduction

Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) A text is a vector of *ambiguous keywords*
- ▷ **Limitations** {
 - Curse of Dimensionality ✓
 - Synonymy and Polysemy problems of keywords ✓
 - No understanding of real-world entities ✗✓
 - Structure of the sentence is lost (no facts) ✗
- ▷ **LDA/LSI** (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations
A text is mapped into a *latent space* (vector of floating-points)

Introduction

Text Understanding



- ▷ Literature currently offers a number of solutions based on **BoW** (Harris, Word'54) A text is a vector of *ambiguous keywords*
- ▷ **Limitations** → **Need for a more structured and efficient algorithmic paradigm**
- ▷ LDA/LSI (Huffman, NIPS'10) and **Word Embeddings** (Mikolov, NIPS'13) overcome some limitations
- A text is mapped into a *latent space* (vector of floating-points)

Q. Curse of Dimensionality Q

Need for a more structured and efficient algorithmic paradigm

words ✓

s × ✓

x ×

Structure of the sentence is lost (no facts)

Introduction

- ▷ Need for a more structured algorithmic paradigm

Exploiting two
different resources



- World Knowledge
- Linguistic Knowledge

WIKIPEDIA
The Free Encyclopedia
Language Grammar

...thanks to recent advancements in the field of
Natural Language Processing:

- ▷ Entity Linking (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)

Introduction

- ▷ Need for a more structured algorithmic paradigm

Exploiting two
different resources



- World Knowledge
- Linguistic Knowledge

WIKIPEDIA
The Free Encyclopedia

Language Grammar

...thanks to recent advancements in the field of
Natural Language Processing:

- ▷ Entity Linking (Bunescu, EACL'06), (Scaiella, IEEE'12), (Piccinno, SIGIR'14)

Leonardo painted Mona Lisa

Introduction

- ▷ Need for a more structured algorithmic paradigm

Exploiting two
different resources



- World Knowledge
- Linguistic Knowledge

WIKIPEDIA
The Free Encyclopedia

Language Grammar

...thanks to recent advancements in the field of
Natural Language Processing:

- ▷ Entity Linking (Bunescu, EACL'06), (Scaeilla, IEEE'12), (Piccinno, SIGIR'14)

Leonardo painted Mona Lisa

Entities



Leonardo da Vinci

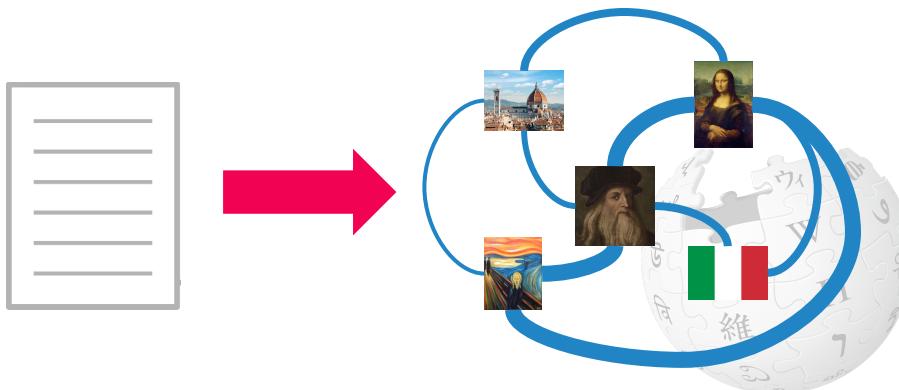


Mona Lisa (painting)

Introduction

- ▷ Need for a more structured algorithmic paradigm

...how? Applying **Graph Theory** to entity linking and open information extraction



Model a text as a
*Small Wikipedia
Graph!*

- Curse of Dimensionality →
- Synonymy and Polysemy problems of keywords →
- Understanding of real-world entities →
- Structured facts →

{ The graph is small
Wikipedia entities are unique and
they represent a real-world concept
OpenIE preserves subject-relation-
object structure

3

Algorithms for Expert Finding

WISER: A Semantic Approach for Expert Finding
in Academia based on Entity Linking
Paolo Cifariello, Paolo Ferragina, Marco Ponza

*Under review at Information
Systems (Elsevier Journal)*

Expert Finding

- ▷ Searching for experts with respect to an input topic
 - Extremely challenging task: Who is an expert?
The notion of expertise is hard to formalize as well as to be modeled (Balog, FTIR'12)



...so difficult that literature refers to expertise as “tacit knowledge”!

- Expertise is actually carried by people in their minds
- Machines have only one way to access to people expertise
 - Artifacts (e.g., papers, emails, ...) people write to share their expertise!

Experiments

Contributions

- ▷ New Expert Finding system



- Fully **unsupervised**
- Jointly combines classical retrieval techniques with the **Wikipedia KG** via **Entity Linking**



Indexing



Every authors' profile is modeled through a small Wikipedia graph...



Query Time



...used to design new profile-centric scoring strategies for the retrieval of experts!



Wikipedia Expertise Ranking Indexing





Wikipedia Expertise Ranking Indexing

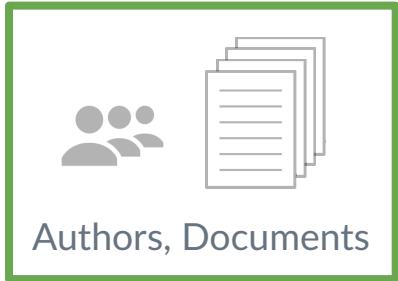


Authors, Documents





Wikipedia Expertise Ranking Indexing



 **elastic**
Documents
indexed
with Elasticsearch



 mongoDB

Indexing of pairs (Author, DocIDs)



Wikipedia Expertise Ranking Indexing

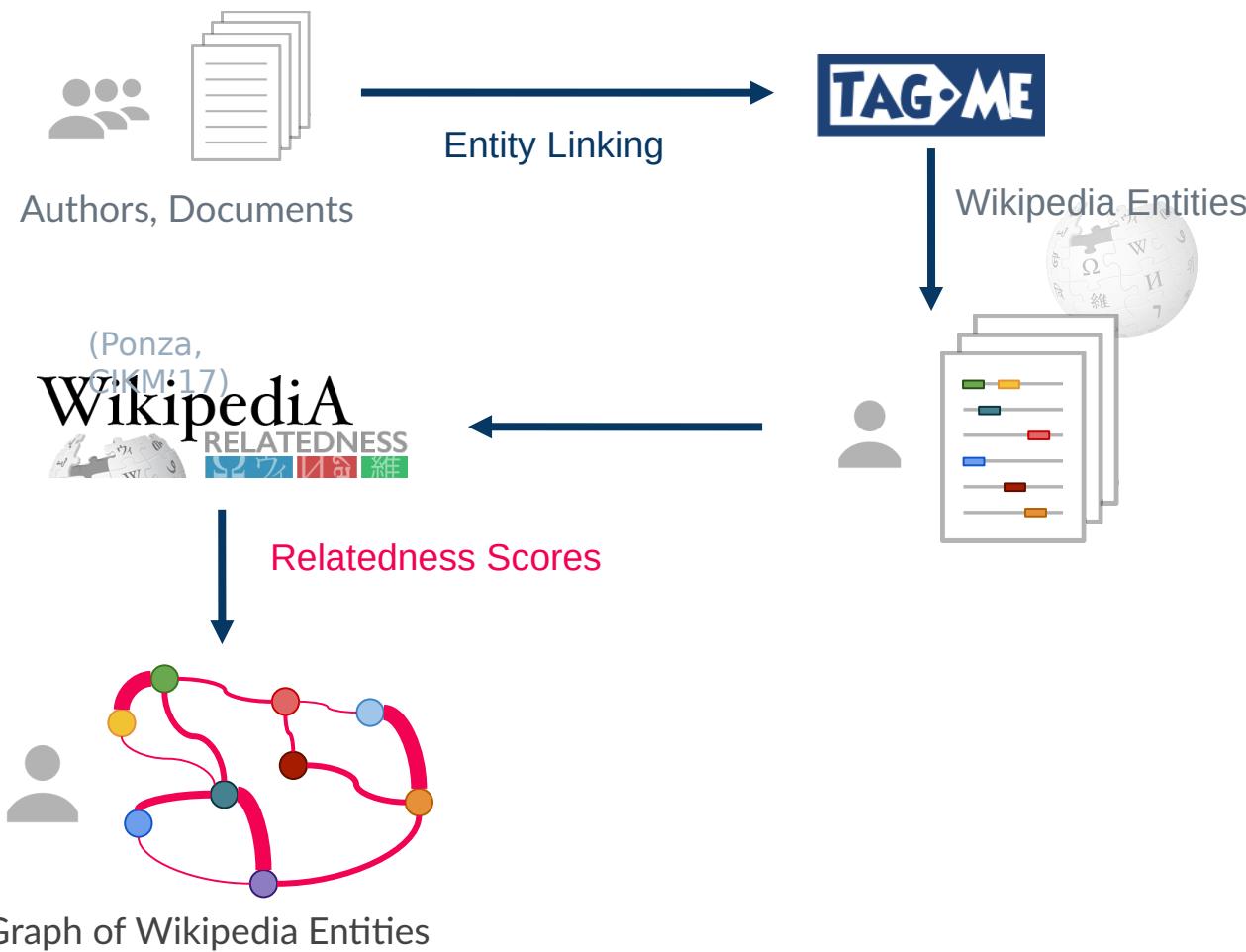


Authors, Documents



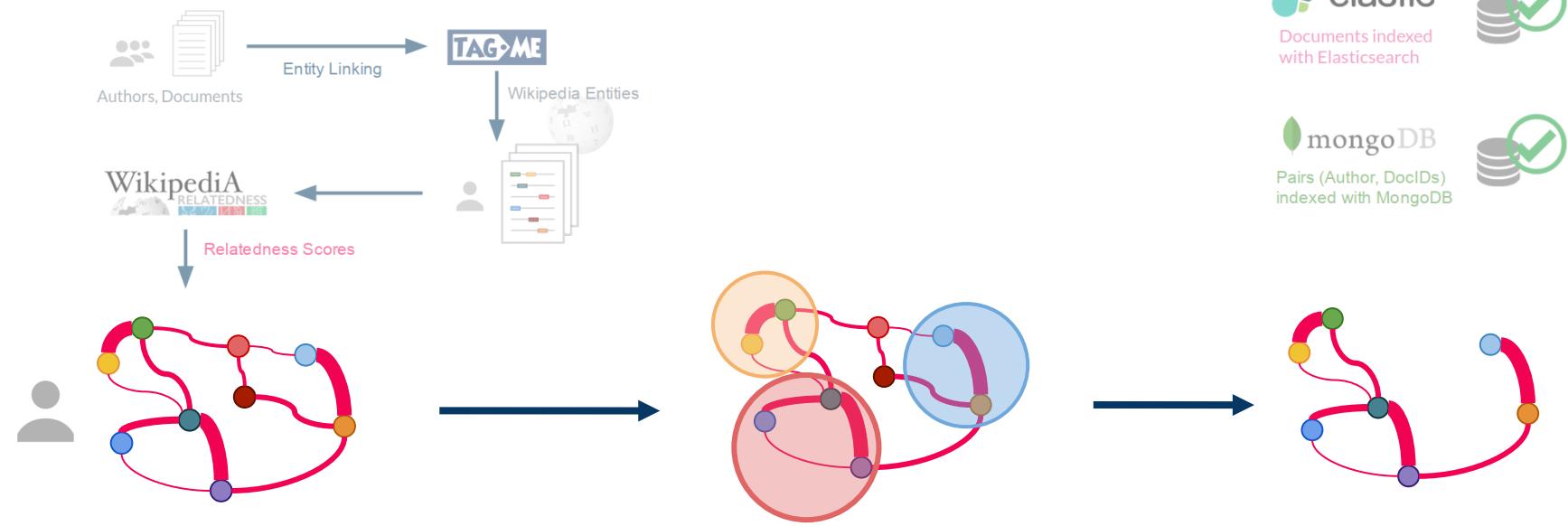


Wikipedia Expertise Ranking Indexing





Wikipedia Expertise Ranking Indexing



Graph of Wikipedia Entities

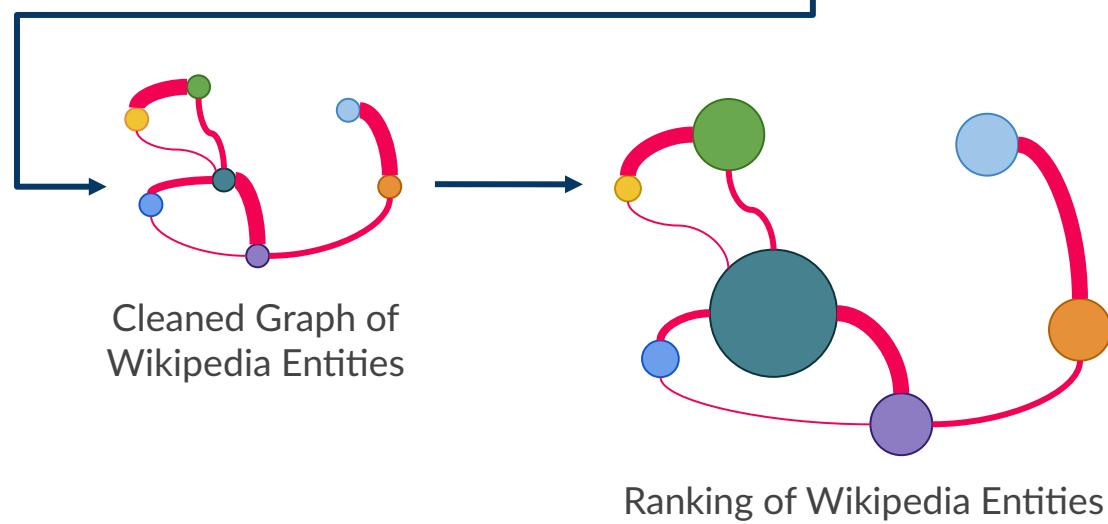
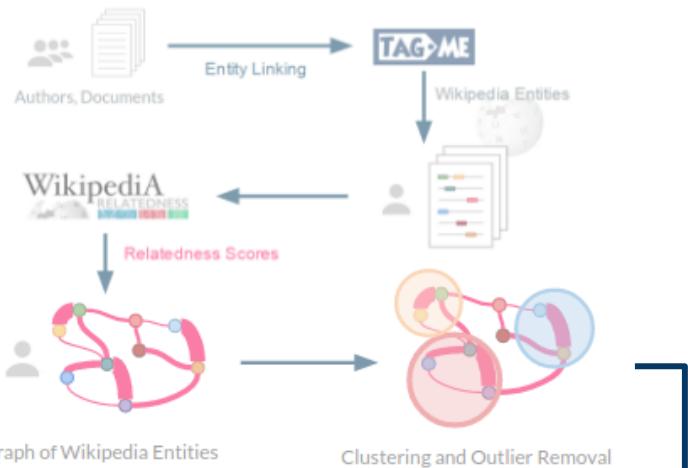
Clustering and Outlier Removal
of Wikipedia Entities

Cleaned Graph of
Wikipedia Entities

- ▷ **HDBScan Algorithm (McInnes, IEEE'17)**
- ▷ **Conservative Approach:**
 $\geq 20\%$ of nodes marked as outliers implies no cleaning



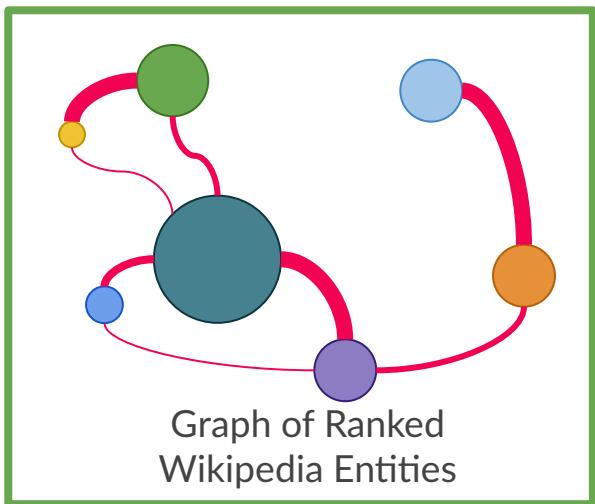
Wikipedia Expertise Ranking Indexing



- ▷ PageRank Algorithm
- ▷ Teleport vector instantiated by taking into account the **frequency** of an **entity** annotated in the **documents** of an author



Wikipedia Expertise Ranking Indexing



 **elastic**
Documents indexed with Elasticsearch



 mongoDB
Pairs of (Author, DocIDs)





Wikipedia Expertise Ranking Indexing

Indexing Completed!



Pairs of (Author, DocIDs)
+
Pairs of (Author, Graph of Ranked
Wikipedia Entities)



Wikipedia Expertise Ranking

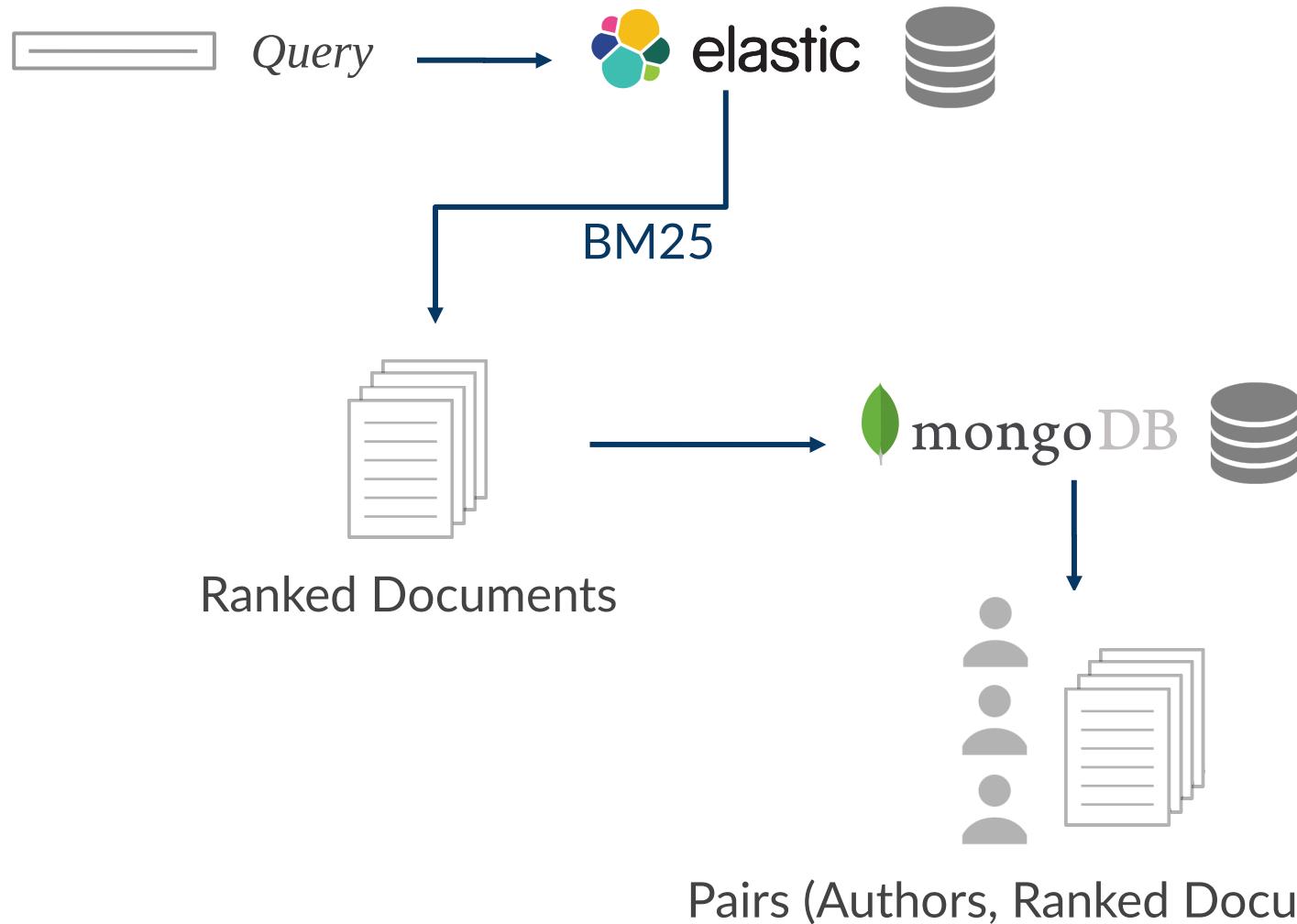
Query Time: Two Strategies

- ▷ Jointly combine **two** different **Authors' Scoring Strategies**
 - Document-Centric
 1. Retrieve relevant documents
 2. Score each author wrt documents' rank (BM25)
 - Profile-Centric
 1. Retrieve relevant authors (wrt query Wikipedia entities)
 2. Score each author wrt entities relevance



Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy





Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



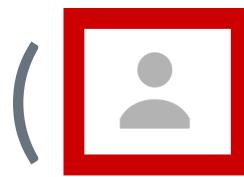


Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



Reciprocal Rank
(Macdonald,
CIKM'08)



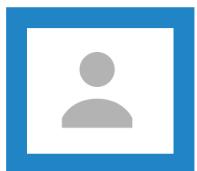
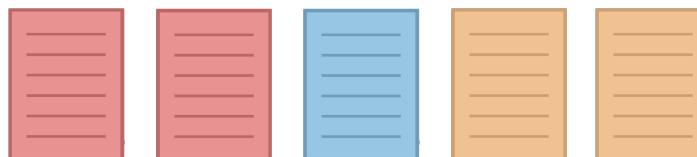
$$\frac{1}{1} + \frac{1}{2}$$





Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



Reciprocal Rank
(Macdonald,
CIKM'08)

$$\left(\boxed{\text{User}} \right) \Rightarrow \frac{1}{1} + \frac{1}{2}$$



1.5



Wikipedia Expertise Ranking

Query Time: Document-Centric Strategy



Reciprocal Rank

$$\left(\boxed{\text{User}} \right) \Rightarrow 1.5$$

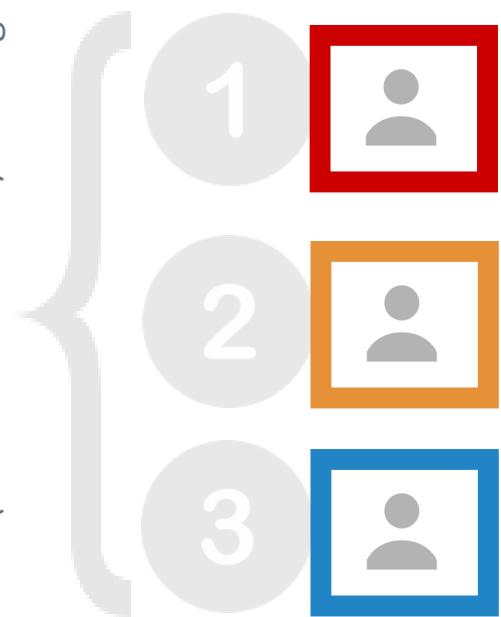
Reciprocal Rank

$$\left(\boxed{\text{User}} \right) \Rightarrow 0.3$$

Reciprocal Rank

$$\left(\boxed{\text{User}} \right) \Rightarrow 0.4$$

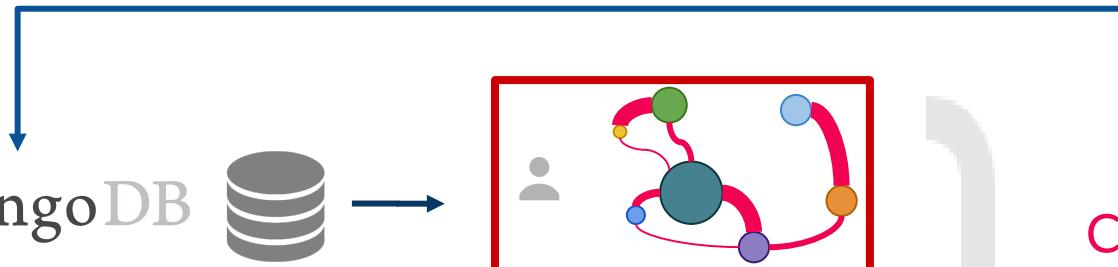
Final (Document-Centric) Ranking





Wikipedia Expertise Ranking

Query Time: Profile-Centric Strategy



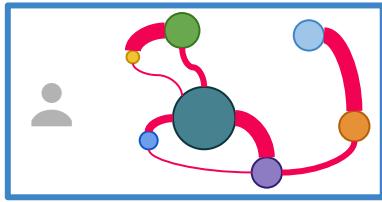
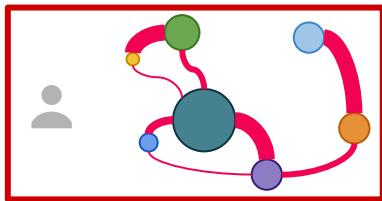
Candidate experts and their Wikipedia-based **profiles** matching the query's entities



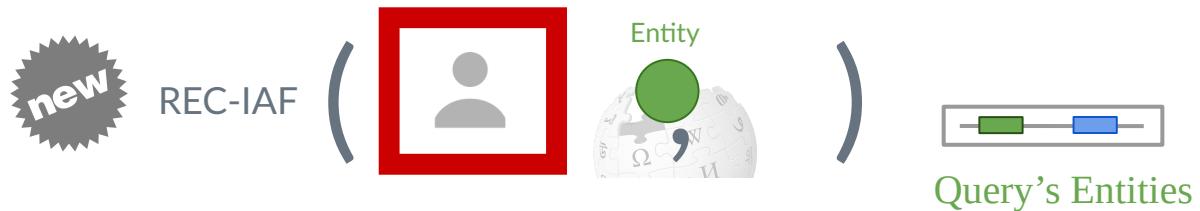
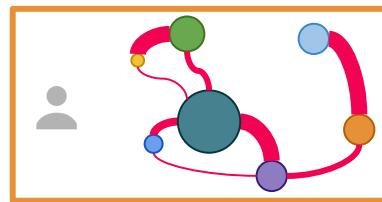


Wikipedia Expertise Ranking

Query Time: Profile-Centric Strategy



⋮

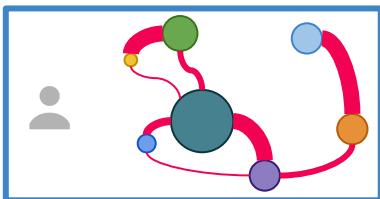
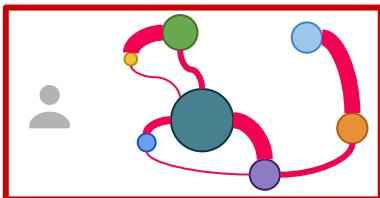


- ▷ Score each author's entity (matched in the input query)
- ▷ Combination of multiple scores of the entity:
 - Document Frequency
 - Confidence (provided by TagMe)
 - Inverse document frequency
 - PageRank in the author's profile

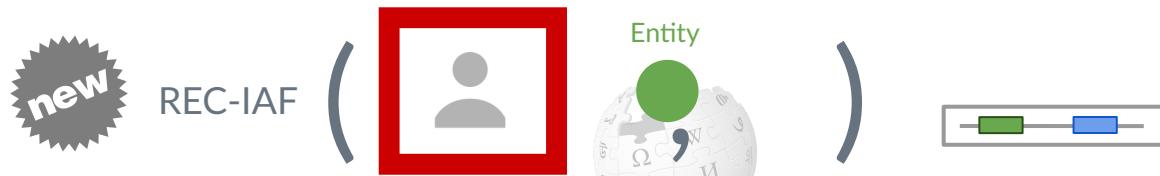
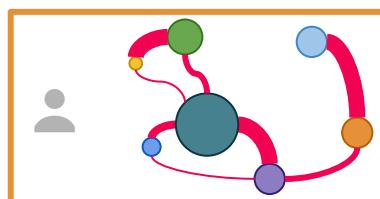


Wikipedia Expertise Ranking

Query Time: Profile-Centric Strategy



⋮



REC-IAF is computed for each entity!

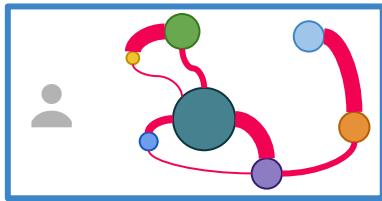
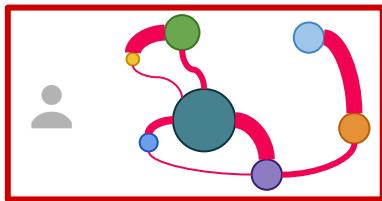
REC-IAF (User Profile, Entity)
REC-IAF (User Profile, Entity)

The *average* is the ranking of the authors!

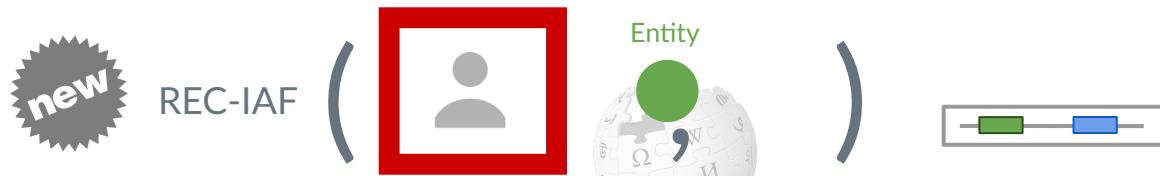
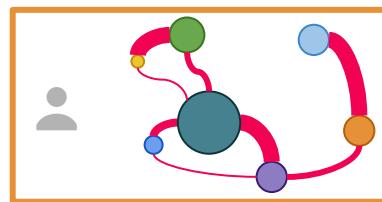


Wikipedia Expertise Ranking

Query Time: Profile-Centric Strategy



⋮



REC-IAF is computed for each entity!

$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$
$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$

$\{ \}$ avg = 0.23

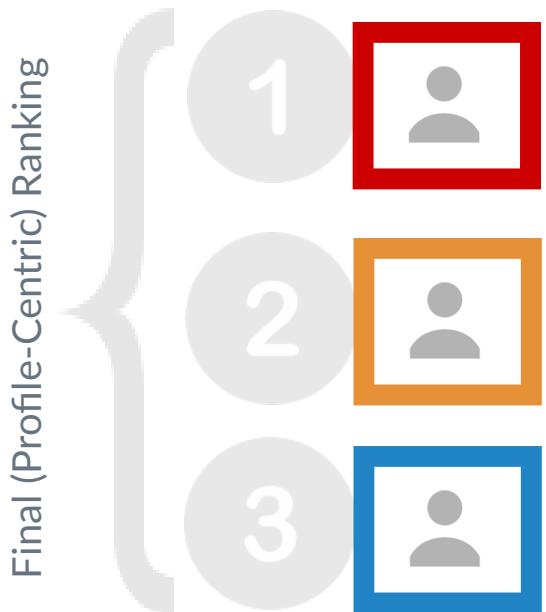
$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$
$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$

$\{ \}$ avg = 0.12

$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$
$$\text{REC-IAF} \left(\begin{array}{c} \text{User} \\ \text{Entity} \end{array} \right)$$

$\{ \}$ avg = 0.43

Final (Profile-Centric) Ranking



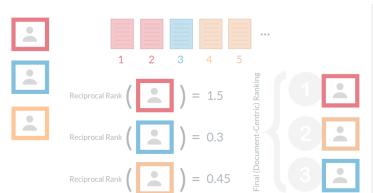


Wikipedia Expertise Ranking

Query Time: Data Fusion

- ▷ We have two different rankings

- Document-Centric Ranking



- Profile-Centric Ranking



Final ranking of experts is given by the **Reciprocal Rank** (Macdonald, CIKM'08) between the **Product** of these two ranking scores

$$\text{Final Score} \left(\boxed{\text{User}} \right) = \frac{1}{\text{Doc-Cent Rank}(\boxed{\text{User}})} \cdot \frac{1}{\text{Prof-Cent Rank}(\boxed{\text{User}})}$$

Experiments

Benchmark

- ▷ TU Dataset (Berendsen, DBWIR'13)
 - ~31K documents (*largest available*)
 - ~1K researchers
 - ~1K test queries
 - Human-assessed Ground-Truth
- ▷ Other systems
 - JM Model (Balog, SIGIR'06)
 - Based on Frequency statistics between (Author, Keywords)
 - Log-Linear (Van Gysel, WWW'16)
 - Based Deep Learning (each author's profile is represented with an embedding vector)
 - Ensamble
 - Product Reciprocal Rank between JM Model and Log-Linear

Experiments

Results

Method	MAP	MRR	P@5	P@10	NDCG@100
JM Model	0.253	0.302	0.108	0.081	0.394
Log-Linear	0.287	0.363	0.134	0.092	0.425
Ensamble	0.331	0.402	0.156	0.105	0.477
 WISER	0.385	0.459	0.163	0.105	0.513

+5.4%

+5.7%

+0.7%

+3.6%

from the University of Pisa



URL: <https://wiser.mkapp.it>



Search for expertise...



Search by Expertise

Search by Name

Search by Department

- ▷ ~1.5K Authors
- ▷ ~65K Documents (papers' abstracts)
- ▷ ~35K Research Topics
- ▷ More than 1K queries and ~2K profiles view in few months
- ▷ Currently used by UniPi's Technological Transfer Office



Select the range of years to analyze.



Ra...	Entity	Count	Doc. count	Years	Wiser score
1	Recurrent neural network	26	15	1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
2	Artificial neural network	44	24	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
3	Tree (data structure)	24	17	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013	
4	Machine learning	16	14	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
5	Recursive neural network	16	16	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013	
6	Quantitative structure-activi...	30	16	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013	
7	Echo state network	14	13	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
8	Mathematical model	39	28	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
9	Prediction	19	17	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
10	Generative model	9	7	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
11	Group representation	22	14	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
12	Empiricism	15	15	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
13	Data model	8	8	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
14	Data set	20	14	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	
15	Reservoir computing	8	8	2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018	

Previous

Page 1 of 24

15 rows

Next



Main Topics Main Areas **Stream Graph** Tag Cloud Publications Survey

Select range of years to examine. Select how many topics for each range

6

4

