

# **Knowledge Harvesting from Web Sources**

### Part 2: Knowledge Base Applications: Search, Ranking, Disambiguation

### **Gerhard Weikum**

Max Planck Institute for Informatics http://www.mpi-inf.mpg.de/~weikum/



Source: DB & IR methods for knowledge discovery. Communications of the ACM 52(4), 2009

- Web 2011 contains more DB-style data than ever
- getting better at making structured content explicit: entities, classes (types), relationships
- but no (hope for) schema here!

# **Structure Now!**

IE-enriched Web pages with embedded entities and facts





# Knowledge bases with facts from Web and IE witnesses



# **Distributed Structure: Linking Open Data**



Source: Christian Bizer, Tom Heath, Tim Berners-Lee, Michael Hausenblas, WWW 2010 Workshop on Linked Data on the Web

# **Distributed Structure: Linking Open Data**



# **Entity Search**



... blending the styles of early J. S. Bach and late Igor Stravinsky that could win a ... by Josef Stalin and
 John Cage. Orchestrate it as Bayel might have. Perform the solo part. You will find a piano under



# Outline

### ✓ Motivation

### **Searching for Entities & Relations**

### ★ Informative Ranking

### ★ Entity-Name Disambiguation





# **RDF: Structure, Diversity, No Schema**

SPO triples (statements, facts):

(EnnioMorricone, bornIn, Rome) (Rome, locatedIn, Italy) (JavierNavarrete, birthPlace, Teruel) (Teruel, locatedIn, Spain) (EnnioMorricone, composed, l'Arena) (JavierNavarrete, composerOf, aTale)

bornIn (EnnioMorricone, Rome)

informatik

EnnioMorricone bornIn Rome

(uri1, hasName, EnnioMorricone) (uri1, bornIn, uri2) (uri2, hasName, Rome) (uri2, locatedIn, uri3)

 IocatedIn(Rome, Italy)

 Rome

 IocatedIn

 Italy

 IocatedIn

 City

 instanceOf

- **SPO triples**: Subject Property/Predicate Object/Value)
- pay-as-you-go: schema-agnostic or schema later
- RDF triples form fine-grained ER graph
- popular for Linked Data, comp. biology (UniProt, KEGG, etc.)
- open-source engines: Jena, Sesame, RDF-3X, etc.

# **Facts about Facts**

### <u>facts:</u>

- 1: (EnnioMorricone, composed, l'Arena)
- 2: (JavierNavarrete, composerOf, aTale)
- 3: (Berlin, capitalOf, Germany)
- 4: (Madonna, marriedTo, GuyRitchie)
- 5: (NicolasSarkozy, marriedTo, CarlaBruni)

### temporal facts:

- 6: (1, inYear, 1968)
- 7: (2, inYear, 2006)
- 8: (3, validFrom, 1990)
- 9: (4, validFrom, 22-Dec-2000)
- 10: (4, validUntil, Nov-2008)
- 11: (5, validFrom, 2-Feb-2008)

### provenance:

- 12: (1, witness, http://www.last.fm/music/Ennio+Morricone/)
- 13: (1, confidence, 0.9)
- 14: (4, witness, http://en.wikipedia.org/wiki/Guy\_Ritchie)
- 15: (4, witness, http://en.wikipedia.org/wiki/Madonna\_(entertainer))
- 16: (10, witness, http://www.intouchweekly.com/2007/12/post\_1.php)
- 17: (10, confidence, 0.1)
- temporal annotations, witnesses/sources, confidence, etc. can refer to reified facts via fact identifiers (approx. equiv. to RDF quadruples: Col × Sub × Prop × Obj)



# **SPARQL Query Language**

SPJ combinations of triple patterns (triples with S,P,O replaced by variable(s)) Select ?p, ?c Where { ?p instanceOf Composer . ?p bornIn ?t . ?t inCountry ?c . ?c locatedIn Europe . ?p hasWon ?a .?a Name AcademyAward . }

#### **Semantics:**

return all bindings to variables that match all triple patterns (subgraphs in RDF graph that are isomorphic to query graph)

+ filter predicates, duplicate handling, RDFS types, etc.

Select Distinct ?c Where { ?p instanceOf Composer . ?p bornIn ?t . ?t inCountry ?c . ?c locatedIn Europe . ?p hasWon ?a .?a Name ?n . ?p bornOn ?b . Filter (?b > 1945) . Filter(regex(?n, "Academy") . }

# **Querying the Structured Web**

Structure but no schema: SPARQL well suited

flexible subgraph matching

wildcards for properties (relaxed joins): Select ?p, ?c Where { ?p instanceOf Composer . ?p ?r1 ?t . ?t ?r2 ?c . ?c isa Country . ?c locatedIn Europe . }

Extension: transitive paths [K. Anyanwu et al.: WWW'07] Select ?p, ?c Where { ?p instanceOf Composer . ?p ??r ?c . ?c isa Country . ?c locatedIn Europe . PathFilter(cost(??r) < 5) . PathFilter (containsAny(??r,?t ) . ?t isa City . }

Extension: regular expressions [G. Kasneci et al.: ICDE'08] Select ?p, ?c Where { ?p instanceOf Composer . ?p (bornIn | livesIn | citizenOf) locatedIn\* Europe . }



# **Querying Facts & Text**

Problem: not everything is triplified

- Consider witnesses/sources
   (provenance meta-facts)
- Allow text predicates with each triple pattern (à la XQ-FT)

Semantics: triples match struct. pred. witnesses match text pred.

Location



European composers who have won the Oscar, whose music appeared in dramatic western scenes, and who also wrote classical pieces ?

Select ?p Where { ?p instanceOf Composer . ?p bornIn ?t . ?t inCountry ?c . ?c locatedIn Europe . ?p hasWon ?a .?a Name AcademyAward . ?p contributedTo ?movie [western, gunfight, duel, sunset] . ?p composed ?music [classical, orchestra, cantata, opera] . }



# **Querying Facts & Text**

Problem: not everything is triplified

- Consider witnesses/sources (provenance meta-facts)
- Allow text predicates with each triple pattern (à la XQ-FT)

### **Grouping of** keywords or phrases **boosts expressiveness** 1858-04-23

#### French politicians married to Italian singers?

Select ?p1, ?p2 Where { ?p2 instanceOf singer [Italy]. ?p1 marriedTo ?p2. }

Select ?p1, ?p2 Where { ?p1 instanceOf politician [France]. ?p1 instanceOf ?c1 [France, politics]. ?p2 instanceOf ?c2 [Italy, singer]. ?p1 marriedTo ?p2. }

#### **CS** researchers whose advisors worked on the Manhattan project?

Select ?r, ?a Where { ?r Inst Ofore [compute [cscriptde] science]. ?a Wpo2 KecolO[Max hklitatah attajept]oject]. ?r RpsAdvisor ?a. }

## **Relatedness Queries**

Schema-agnostic keyword search (on RDF, ER graph, relational DB) becomes a special case



Relationship between Angela Merkel, Jim Gray, Dalai Lama?

Select ??p1, ??p2, ??p3 Where { AngelaMerkel ??p1 MaxPlanck. MaxPlanck ??p2 DalaiLama . DalaiLama ??p3 AngelaMerkel . } ?e3 ?r3 ?c3 ["Dalai Lama"] . ?e1 ??p1 ?e2 . ?e1 ??p1 ?e2 . ?e3 ??p3 ?e1 . }



#### Querying Temporal Facts [Y. Wang et al.: EDBT'10] [O. Udrea et al.: TOCL'10]

Problem: not all facts hold forever (e.g. CEOs, spouses, ...)

- Consider temporal scopes of reified facts
- Extend Sparql with temporal predicates
  - 1: (BayernMunich, hasWon, ChampionsLeague)
  - 2: (BorussiaDortmund, hasWon, ChampionsLeague)
  - 3: (1, validOn, 23May2001) 4: (1, validOn, 15May1974)
  - 5: (2, validOn, 28May1997)

6: (OttmarHitzfeld, manages, BayernMunich)

7: (6, validSince, 1Jul1998) 8:(6, validUntil, 30Jun2004)

When did a German soccer club win the Champions League? Select ?c, ?t Where { ?c isa soccerClub . ?c inCountry Germany . ?id1: ?c hasWon ChampionsLeague . ?id1 validOn ?t . }

Managers of German clubs who won the Champions League? Select ?m Where { isa soccerClub . ?c inCountry Germany . ?id1: ?c hasWon ChampionsLeague . ?id1 validOn ?t . ?id2: ?m manages ?c . ?id2 validSince ?s . ?id2 validUntil ?u . [?s,?u] overlaps [?t,?t] . }

#### max planck institut

# **Querying with Vague Temporal Scope**

#### 8 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 20

- Problem: user's temporal interest is often imprecise
  - Consider temporal phrases as text conditions
  - Allow approximate matching and rank results wisely

### **German Champion League winners in the nineties?**

Select ?c Where { ?c isa soccerClub . ?c inCountry Germany . ?c hasWon ChampionsLeague [nineties] . }

### **Soccer final winners in summer 2001?**

Select ?c Where { ?c isa soccerClub . ?id: ?c matchAgainst ?o [final] . ?id winner ?c ["summer 2001"] . }



# **Take-Home Message (Querying)**

Don't re-invent the wheel:

**SPARQL** is there already, entity search is special case

• Extensions should be conceptually simple

meta-fact & text predicates naturally embedded

• Ease of use (progammability) is crucial:

doing well for API, UI is harder



# Outline

### ✓ Motivation

### ✓ Searching for Entities & Relations

★ Informative Ranking

### ★ Entity-Name Disambiguation





# **Ranking Criteria**

### **Confidence:**

Prefer results that are likely correct

- accuracy of info extraction
- trust in sources (authenticity, authority)

### Informativeness:

Prefer results with salient facts Statistical estimation from:

- frequency in answer
- Frequency on Web
- frequency in query log

**Diversity:** Prefer variety of facts

### **Conciseness:**

Prefer results that are tightly connected

- size of answer graph
- cost of Steiner tree

bornIn (Jim Gray, San Francisco) from "Jim Gray was born in San Francisco" (en.wikipedia.org)

livesIn (Michael Jackson, Tibet) from "Fans believe Jacko hides in Tibet" (www.michaeljacksonsightings.com)

q: Einstein isa ? Einstein isa scientist Einstein isa vegetarian

q: ?x isa vegetarian Einstein isa vegetarian Whocares isa vegetarian

E won ... E discovered ... E played ... E won ... E won ... E won ... E won ...

Einstein won NobelPrize Bohr won NobelPrize

Einstein isa vegetarian Cruise isa vegetarian Cruise born 1962 Bohr died 1962

# **Ranking Approaches**

### **Confidence:**

Prefer results that are likely correct

- accuracy of info extraction
- trust in sources (authenticity, authority)

empirical accuracy of IE PR/HITS-style estimate of trust combine into: max { accuracy (f,s) \* trust(s) | s ∈ witnesses(f) }

### **Informativeness:**

Prefer results with salient facts Statistical LM with estimations from:

- frequency in answer
- > frequency in corpus (e.g. Web)
- frequency in query log

#### **Diversity:**

Prefer variety of facts

### **Conciseness:**

Prefer results that are tightly connected
➢ size of answer graph
➢ cost of Steiner tree

**PR/HITS**-style entity/fact ranking [V. Hristidis et al., S.Chakrabarti, ...]

or

IR models: tf\*idf ... [K.Chang et al., ...] Statistical Language Models

**Statistical Language Models** 

#### graph algorithms (BANKS, STAR, ...) [J.X. Yu et al., S.Chakrabarti et al., B. Kimelfeld et al., G.Kasneci et al., ...]

## Statistical Language Models (LM's)

[Maron/Kuhns 1960, Ponte/Croft 1998, Hiemstra 1998, Lafferty/Zhai 2001]

C



"God rolls dice in places where you can't see them" (Stephen Hawking)

• each doc d<sub>i</sub> has LM: generative prob. distr. with params  $\theta_i$ 

• query q viewed as sample from LM( $\theta_1$ ), LM( $\theta_2$ ), ...

 $LM(\theta_2)$ 

- estimate likelihood P[ q | LM(θ<sub>i</sub>) ] that q is sample of LM of doc d<sub>i</sub> (q is "generated by" d<sub>i</sub>)
- rank by descending likelihoods (best "explanation" of q)

**d**₁

 $\mathbf{d}_2$ 

# LM: Doc as Model, Query as Sample



# document d: sample of M used for parameter estimation

max planck institut



## **Some LM Basics**

independ. assumpt.  $s(d,q) = P[q | d] = \prod_{i} P[q_{i} | d]$  $\sim \sum_{i} \log \frac{tf(i,d)}{\sum_{i} tf(k,d)}$ simple MLE: overfitting  $s(d,q) = \lambda P[q \mid d] + (1 - \lambda) P[q]$ mixture model for smoothing  $\sim \sum_{i} \log \left\{ \lambda \frac{tf(i,d)}{\sum_{k} tf(k,d)} + (1-\lambda) \frac{df(i)}{\sum_{k} df(k)} \right\}$ **P[q]** est. from log or corpus  $\sim \sum_{i} \log \left( 1 + \frac{tf(i,d)}{\sum_{k} tf(k,d)} \cdot \frac{1-\lambda}{\lambda} \frac{\sum_{k} df(k)}{df(i)} \right) rank by ascending "improbability"$  $\sim KL(q \mid d) = \sum_{i} P[i \mid q] \log \frac{P[i \mid q]}{P[i \mid d]}$ **KL divergence** (Kullback-Leibler div.) aka. relative entropy

Precompute per-keyword scores

max planck institut informatik

- Store in postings of inverted index
- Score aggregration for (top-k) multi-keyword query

efficient implementation

**Entity Search with LM Ranking** [Z. Nie et al.: WWW'07, H. Fang et al.: ECIR'07, P. Serdyukov et al.: ECIR'08, ...]

#### query: keywords → answer: entities

 $s(e,q) = \lambda P[q | e] + (1 - \lambda) P[q] \sim \prod \frac{P[q_i | e_i]}{P[q_i]} \sim KL \ (LM \ (q) | LM \ (e))$ 

LM (entity e) = prob. distr. of words seen in context of e



## LM's: from Entities to Facts

**Document / Entity LM's** 



LM for doc/entity: prob. distr. of words

- LM for query: (prob. distr. of) words
- LM's: rich for docs/entities, super-sparse for queries

richer query LM with query expansion, etc.



### **Triple LM's**

LM for facts: (degen. prob. distr. of) triple

LM for queries: (degen. prob. distr. of) triple pattern

LM's: apples and oranges

- expand query variables by S,P,O values from DB/KB
- enhance with witness statistics
- query LM then is prob. distr. of triples !

# LM's for Triples and Triple Patterns

[G. Kasneci et al.: ICDE'08; S. Elbassuoni et al.: CIKM'09, ESWC'11]



250/500

Cruyff coached FCBarca

witness statistics  $\Sigma$ : 2600

## LM's for Composite Queries

q: Select ?x,?c Where { France ml ?x . ?x p ?c . ?c in UK . }



queries q with subqueries  $q_1 \dots q_n$ results are n-tuples of triples  $t_1 \dots t_n$ LM(q): P[q<sub>1</sub>...q<sub>n</sub>] =  $\prod_i P[q_i]$ LM(answer): P[t<sub>1</sub>...t<sub>n</sub>] =  $\prod_i P[t_i]$ KL(LM(q)|LM(answer)) =  $\sum_i KL(LM(q_i)|LM(t_i))$ 





f31: ManU in UK 200 f32: Arsenal in UK 160 f33: Chelsea in UK 140

# LM's for Keyword-Augmented Queries

q: Select ?x, ?c Where {
 France ml ?x [goalgetter, "top scorer"].
 ?x p ?c .
 ?c in UK [champion, "cup winner", double].}



subqueries q<sub>i</sub> with keywords w<sub>1</sub> ... w<sub>m</sub> results are still n-tuples of triples t<sub>i</sub> LM(q<sub>i</sub>): P[triple t<sub>i</sub> | w<sub>1</sub> ... w<sub>m</sub>] =  $\prod_k \beta P[t_i | w_k] + (1-\beta) P[t_i]$ LM(answer f<sub>i</sub>) analogous KL(LM(q)|LM(answer f<sub>i</sub>)) =  $\sum_i KL (LM(q_i) | LM(f_i))$ 

#### result ranking prefers (n-tuples of) triples whose witnesses score high on the subquery keywords



- enhanced ranking
- efficiently computable
- plug into doc/entity/triples LM's

P[v=,,nineties" | x<sub>j</sub> = ,,mid 90s"] = 1/2 P[,,nineties" | ,,summer 1990"] = 1/30 P[,,nineties" | ,,last century"] = 1/10





Personal histories of queries & clicked facts → LM(user u): prob. distr. of triples !

#### **LM(q|u]** = $\mu$ **LM(q)** + (1- $\mu$ ) **LM(u)** then business as usual

u1 [classical music] $\rightarrow$  q: ?p from Europe . ?p hasWon AcademyAwardu2 [romantic comedy] $\rightarrow$  q: ?p from Europe . ?p hasWon AcademyAwardu3 [from Africa] $\rightarrow$  q: ?p isa SoccerPlayer . ?p hasWon ?a

**Open issue:** "insightful" results (new to the user)



# **Result Diversification**

[J. Carbonell, J. Goldstein: SIGIR'98]

q: Select ?p, ?c Where { ?p isa SoccerPlayer . ?p playedFor ?c . }

1 Beckham, ManchesterU

- 2 Beckham, RealMadrid
- 3 Beckham, LAGalaxy
- 4 Beckham, ACMilan
- 5 Zidane, RealMadrid
- 6 Kaka, RealMadrid

7 Cristiano Ronaldo, RealMadrid

8 Raul, RealMadrid

9 van Nistelrooy, RealMadrid10 Casillas, RealMadrid

- 1 Beckham, ManchesterU
- 2 Beckham, RealMadrid
- 3 Zidane, RealMadrid
- 4 Kaka, ACMilan
- 5 Cristiano Ronaldo, ManchesterU
- 6 Messi, FCBarcelona
- 7 Henry, Arsenal
- 8 Ribery, BayernMunich
- 9 Drogba, Chelsea
- 10 Luis Figo, Sporting Lissabon



rank results  $f_1 \dots f_k$  by ascending  $\delta \text{KL}(\text{LM}(q) \mid \text{LM}(f_i)) - (1-\delta) \text{KL}(\text{LM}(f_i) \mid \text{LM}(\{f_1..f_k\} \setminus \{f_i\}))$ implemented by greedy re-ranking of  $f_i$ 's in candidate pool

# Take-Home Message (Ranking)

Don't re-invent the wheel:

LM's are elegant and expressive means for ranking consider both data & workload statistics

Extensions should be conceptually simple:

can capture informativeness, personalization,

**relaxation**, **diversity** – all in same framework

 Unified ranking model for complete query language: still work to do



# Outline

### ✓ Motivation

### ✓ Searching for Entities & Relations

### ✓ Informative Ranking

### **★** Entity-Name Disambiguation




## **Entity Diversity**

∑ sig.ma - Semantic Infor	mation MAshup +	
♦ ★ http://sig.r	ma/search?q=Pushkin	🟫 - ሮ 🚷 - Google 🛛 🔎 🍙 🗷
SIG, M	NTIC INFORMATION	Sources (9) Approved (0) Rejected (0) <
Pushkin	Add More Info Start New Options # Use it P	1 <u>Pushkin</u> 15 facts   2010-07-29 http:// <b>www.w3.org</b> /2006/03/wn/wn20/instances/synset
Pushkin	Order ∞	2 <u>Pushkin</u> 10 facts   2010-02-04 http://wordnet.rkbexplorer.com/id/wordsense-Pushkin
arg peri:	2.79203 [8]	3 <u>Pushkin</u> 16 facts   2010-02-04
aphelion:	3.6556025 [8]	
alt names:	1977 QL3 [8]	http://wordnet.rkbexplorer.com/id/word-Pushkin
abs magnitude:	10.96 [8]	5 Pushkin 7 facts   2011-05-27
albedo:	0.0497 [8]	http://www.w3.org/2006/03/wn/wn20/instances/wordsen
comment:	(Russian poet (1799-1837)) [1,3]	6 Pushkin 4 facts   2011-05-27
is contains word	Pushkin [2]	http://www.w3.org/2006/03/wn/wn20/instances/word-Pu
sense of:		7 <u>Pushkin</u> 3 facts   2010-05-21
contains word	Aleksandr Sergevevich Buchkin [4,2]	http://pushkin-mtg.livejournal.com/
sense:	Alexander Bushkin [1,3]	8 <u>Pushkin</u> 28 facts   2009-09-18
	Pushkin [1 3]	http://dbpedia.org/resource/2208_Pushkin
discoverer:	http://dbnedia.org/resource/NSChernykh_181	9 <u>Aleksandr Sergeyevich Pu</u> 155 facts   2010-06-24 http://dbnedia.org/resource/Alexander Pushkin
discovery site:	Crimean Astrophysical Observatory [8]	
designations:	ves [8]	reject all T approve all T
ennch:	2008-05-14 [8]	http://example.loc/document.rdf
eccentricity:	0.0460085 [8]	
gloss:	(Russian poet (1799-1837)) [1,3]	
hypernym:	poet [1,3]	http://sig.ma
is hyponym of:	poet [1 3]	

## **Entity-Name Ambiguity**

## <sameAs>

#### interlinking the Web of Data

#### About: 2208 Pushkin

An Entity of Type : <u>Thing</u>, from Named Graph : <u>http://dbpedia.org</u>, within Data Space : <u>dbpedia.org</u>



#### 2208 Pushkin (1977 QL3) is an outer main-belt asteroid discovered on August 22, 1977 by N. S. $^{\circ}$

	Property	Value
The Web of Data has many equivalent LIPIs	http://dbpedia.org/meta/editlink	http://en.wikipedia.org/w/index.php?title=2208_Pushkin&action=edit
The web of Data has many equivalent on is.	http://dbpedia.org/meta/pageid	<ul> <li>16477278 (xsd:integer)</li> </ul>
This convine holes you to find as references between	http://dbpedia.org/meta/revision	http://en.wikipedia.org/w/index.php?title=2208_Pushkin&oldid=438390754
different data sets. Enter a known URI, or use Sindice to search first.	<ul> <li>data sets.</li> <li>common URI, or use Sindice to search first.</li> <li>data sets.</li> </ul>	
	dbpprop:absMagnitude	<ul> <li>10.96 (xsd:double)</li> </ul>
	dbpprop:abstract_live	<ul> <li>2208 Pushkin (1977 QL3) is an outer main-belt asteroid discovered on August 22, 1977 by N. S. Chernykh at the Crimean Astrophysical Observatory. It was named after the leading writer of Russia, Pushkin.</li> </ul>
	dbpprop:albedo	<ul> <li>0.0497 (xsd:double)</li> </ul>
Search results from Sindice, with co-references applic	dbpprop:altNames	■ 1977 QL3
- oodron rosalis ir on <u>omaroo</u> , mili oo ronoronoos appin	dbpprop:aphelion	<ul> <li>3.6556025 (xsd:double)</li> </ul>
	dbpprop:argPeri	<ul> <li>2.79203 (xsd:double)</li> </ul>
	dbpprop:ascNode	<ul> <li>79.48005 (xsd:double)</li> </ul>
	dbpprop:bgcolour	FFFFC0
Sindle "2208 Pushkin" ୟ	dbpprop:comment_live	2208 Pushkin (1977 QL3) is an outer main-belt asteroid discovered on August 22, 1977 by N. S.
1 http://dbpedia.org/resource/2208 Pus	hkin	

http://dbpedia.org/resource/2208\_Pushkin

2. http://rdf.freebase.com/ns/guid.9202a8c04000641f8000000007d408a8

rdf+xml · n3 · json · text

#### 蓤 ndce "Pushkin" 🔍

<sameAs>

<sameAs> { 1. http://wordnet.rkbexplorer.com/id/synset-Pushkin-noun-1

2. http://www.w3.org/2006/03/wn/wn20/instances/synset-Pushkin-noun-1





## **Named-Entity Disambiguation**



#### Three NLP tasks:

- 1) named-entity detection: segment & label by HMM or CRF (e.g. Stanford NER tagger)
- 2) co-reference resolution: link to preceding NP (trained classifier over linguistic features)
- named-entity disambiguation: map each mention (name) to canonical entity (entry in KB)

## Mentions, Meanings, Mappings



weighted undirected graph with two types of nodes



Popularity (m,e):

- freq(m,e|m)
- length(e)
- #links(e)

Similarity (m,e):

 cos/Dice/KL (context(m), context(e))



weighted undirected graph with two types of nodes



Popularity (m,e):

- freq(m,e|m)
- length(e)
- #links(e)

Similarity (m,e):

 • cos/Dice/KL (context(m), context(e))



- Coherence (e,e'): • dist(types) • overlap(links)
- overlap
  - (anchor words)

weighted undirected graph with two types of nodes



max planck institut informatik

(anchor words)

weighted undirected graph with two types of nodes



Popularity (m,e):

- freq(m,e|m)
- length(e)
- #links(e)

Similarity (m,e):

 • cos/Dice/KL (context(m), context(e))



- Coherence (e,e'): • dist(types) • overlap(links) • overlap
  - (anchor words)

weighted undirected graph with two types of nodes



Popularity (m,e):

- freq(m,e|m)
- length(e)
- #links(e)

Similarity (m,e):

 • cos/Dice/KL (context(m), context(e))



- Coherence (e,e'): • dist(types) • overlap(links) • overlap
  - (anchor words)

## **Different Approaches**

Combine Popularity, Similarity, and Coherence Features (Cucerzan: EMNLP'07, Milne/Witten: CIKM'08):

- for sim (context(m), context(e)): consider surrounding mentions and their candidate entities
- use their types, links, anchors as features of context(m)
- set m-e edge weights accordingly
- use greedy methods for solution



**Collective Learning with Prob. Factor Graphs** (Chakrabarti et al.: KDD'09):

- model P[m|e] by similarity and P[e1|e2] by coherence
- consider likelihood of P[m1 ... mk | e1 ... ek]
- factorize by all m-e pairs and e1-e2 pairs
- use hill-climbing, LP, etc. for solution

# **Joint Mapping**



- Build mention-entity graph or joint-inference factor graph from knowledge and statistics in KB
- Compute high-likelihood mapping (ML or MAP) or dense subgraph such that: each m is connected to exactly one e (or at most one e)

K. Kulkarni et al.: Collective Annotation of Wikipedia Entities in Web Text, KDD'09 J. Hoffart et al.: Robust Disambiguation of Named Entities in Text, EMNLP'11

## **Mention-Entity Similarity Edges**

Precompute characteristic keyphrases q for each entity e: anchor texts or noun phrases in e page with high PMI:

weight  $(q, e) = \log \frac{freq (q, e)}{freq (q) freq (e)}$ 

"Eurovision song contest"

Match keyphrase q of candidate e in context of mention m



score  $(e \mid m) \sim \sum_{\substack{q \in keyphrases \ in \ context \ (m)}} score (q) dist (cover(q), m)^{-\alpha}$ 

## **Entity-Entity Coherence Edges**

#### Precompute overlap of incoming links for entities e1 and e2

$$mw - coh(e1, e2) \sim 1 - \frac{\log \max(in(e1, e2)) - \log(in(e1) \cap in(e2))}{\log |E| - \log \min(in(e1), in(e2))}$$

Alternatively compute overlap of anchor texts for e1 and e2

 $ngram - coh(e1, e2) \sim \frac{|ngrams (e1) \cap ngrams (e2)|}{|ngrams (e1) \cup ngrams (e2)|}$ 

or overlap of keyphrases, or similarity of bag-of-words, or ... optionally filtered by words or n-grams in entire input text

Optionally combine with type distance of e1 and e2 (e.g., Jaccard index for type instances) and other (precomputed) measures

max planck institut informatik



Compute dense subgraph to

maximize min weighted degree among entity nodes such that:

each m is connected to exactly one e (or at most one e)

Greedy approximation:

iteratively remove weakest entity and its edges



Compute dense subgraph to

maximize min weighted degree among entity nodes such that:

each m is connected to exactly one e (or at most one e)

Greedy approximation:

iteratively remove weakest entity and its edges



Compute dense subgraph to

maximize min weighted degree among entity nodes such that:

each m is connected to exactly one e (or at most one e)

Greedy approximation:

iteratively remove weakest entity and its edges

[J. Hoffart et al.: EMNLP'11]



 Compute dense subgraph to maximize min weighted degree among entity nodes such that:

each m is connected to exactly one e (or at most one e)

Greedy approximation:

iteratively remove weakest entity and its edges

## **AIDA Accurate Online Disambiguation**

#### http://www.mpi-inf.mpg.de/yago-naga/aida/



Agnetha, Björn, Benny, and Anr successful pop music group. The and SOS.

[Agnetha Fältskog] Agnetha ,
[Björn Ulvaeus] Björn , [Benny
Andersson]Benny, and
[Anni-Frid Lyngstad] Anni-Frid
formed [ <mark>Sweden</mark> ]Sweden 'S
most successful pop music
group. Their greatest hits
were [ <u>Waterloo (ABBA</u>

song)]Waterloo and SOS.

al Steps	
only)	
ME Similarity	Weighted Degree
497821536934663	0.052519420551120015
278548264326E-5	0.011433304988143484
37274091523E-5	0.009133432457122746
	0.006144100802016364
410256151456E-4	0.005857037672735628
37580795959E-4	0.005835433432846912
192167752377E-5	0.005348033055157968
	0.0047467918338561935
	0.004242218418100741
	0.00398109454783811
	0.002440125447239848
179001724556E-4	0.0022059134686564985
784286215732E-5	0.002197047514610515
	0.002174127922480215
251094038047E-4	0.0021561290904151646
27315240582E-5	0.0020782134411012295
	0.002051145658234978
	0.002051145658234978
909796136458E-5	0.0018885971232344612
	0.0018776092471261214
	0.0017481163638816684



## **AIDA Accurate Online Disambiguation**

## http://www.mpi-inf.mpg.de/yago-naga/aida/

Disambiy	juanon meniou	•	
prior	prior+sim	prior+sim+coherence (graph)	
	Pa	rameters: (default should be OK)	
Similarity	Impact: <mark>0.9</mark>		
Ambiguity	degree 5		
Coherenc	e threshold: 0.9		

#### **Mention Extraction:**

Stanford NER Manual

You can manually tag the mentions by putting them between [[ and manual mode.

		B	I	U	ABC	≣	≣	≣		Sty	les	
Ж	6	2	T	1	A	A.A.	:=		-7	(4	HTML	A
			13	3		17	1 m	l Y	I		-	2

Tottenham	Crouch
Bayern	Robben
Shakhtar	Adriano
ManU	Beckham
Chelsea	Ballack
Real	Raul
Milano	Basten

[Tottenham Hotspur F.C.] Tottenham [Peter Crouch] Crouch [FC Bayern Munich] Bayern [Arjen Robben] Robben [FC Shakhtar Donetsk]Shakhtar [Adriano Leite Ribeiro]Adriano [Manchester United F.C.]ManU [David Beckham] Beckham [Chelsea F.C.] Chelsea [Michael Ballack] Ballack [Real Madrid C.F.]Real [Raúl González] Raul [A.C. Milan] Milano [John Basten] Basten

Gr	aph	Removal S	teps						
n									
solve	d by lo	ocal sim. only	)						
ity	МІ	E Similarity	Weighted Degree	Weighted De remove					
	8.91260	07921453174E-4	0.282339625632226	0.156253808					
			0.060602238345761804	-1.0					
		39009157388E-4	0.0459465069323101	-1.0					
			0.02171347034201928	-1.0					
n the			0.0020525462869665713	-1.0					
n me		)459592857E-5	3.075469550958973E-5	3.075469550					
			0.0	0.0					
			0.0	0.0					
3	<b>*</b>		0.0	0.0					
				•					
			A						
		car sim. only	y)						
		local sim. or	ıly)						
			, /						
		irsim. only)							
		local sim. or	1ly)						
		local sim. or	 וע)						
	=								
		pcal sim. only)							



## **Record Linkage (Entity Resolution)**



Find equivalence classes of entities, and records, based on:

- similarity of values (edit distance, n-gram overlap, etc.)
- joint agreement of linkage

 $\rightarrow$  similarity joins, grouping/clustering, collective learning, etc.

Halbert L. Dunn: Record Linkage. American Journal of Public Health. 1946 H.B. Newcombe et al.: Automatic Linkage of Vital Records. Science, 1959.

# **Record Linkage (Entity Resolution)**



Halbert L. Dunn: Record Linkage. American Journal of Public Health. 1946 H.B. Newcombe et al.: Automatic Linkage of Vital Records. Science, 1959.

### Linked Data: Record Linkage at Web Scale



Source: Christian Bizer, Tom Heath, Tim Berners-Lee, Michael Hausenblas, WWW 2010 Workshop on Linked Data on the Web





### Linked Data: Record Linkage at Web Scale



## **Open Problems**

- More efficient graph algorithms (multicore, etc.)
- Allow mentions of unknown entities, mapped to null
- Short and difficult texts:
  - tweets, headlines, etc.
  - fictional texts: novels, song lyrics, etc.
  - incoherent texts
- Disambiguation beyond entity names:
  - coreferences: pronouns, paraphrases, etc.
  - common nouns, verbal phrases (general WSD)



# Outline

✓ Motivation

#### ✓ Searching for Entities & Relations

## ✓ Informative Ranking

### ✓ Entity-Name Disambiguation





#### KB Applications: Achievements & Challenges Search for Entities and Relations

good success story on entities, problems left:

- coverage beyond celebrities (long tail)
- complex queries → Sparql awareness & extensions
- reconcile Web with Web of Data

Ranking

good progress on ER language models, challenges left:

- better understanding of time & space
- mapping names & phrases to entities & relations
- efficiency & scalability

**Entity-Name Disambiguation** 

good progress, challenges left:

- entities in the long tail, newly emerging entities
- robustness on short & difficult inputs
- apply at Web scale: tables, lists, text+RDFa, LOD, etc.

### UI: Structured Keyword Search [Ilyas et al. Sigmod'10]

Need to map (groups of) keywords onto entities & relationships based on name-entity similarities/probabilities



q: German football clubs that won (a match) against Real



#### **UI: Natural Language Questions**

translate question into Sparql query:

dependency parsing to decompose question

mapping of question units onto entities, classes, relations



### **UI: Natural Language Questions**

translate question into Sparql query:

- dependency parsing to decompose question
- mapping of question units onto entities, classes, relations



## **Performance/Scalability: Open Issues**

- Many joins between triples and inverted lists
- Query relaxation very expensive
- Entity disambiguation within QP
- Top-k processing with early termination
- Everything distributed (LOD cloud)



## **Efficient& Scalable RDF Query Processing**

**RDF-3X engine** [T. Neumann et al.: VLDB'08, SIGMOD'09, VLDBJ'10]:

- no-tuning RISC, versioning, online updates, transactions
- aggressive indexing (all SPO permutations & projections)
- fast DP-based join-order optimization for 20-30 joins

aggressive sideways information passing



## **Overall Take-Home**

#### Web of Data

→ RDF: entities, relationships, structure, no schema

- **Querying**
- $\rightarrow$  extended SPARQL:

W3C, tripleX patterns, schema-free joins

- **Ranking**
- → Language Models: from docs to triples, composable, relaxable, temporal, individual

### **Disambiguation**

max planck institut informatik

→ coherence graph: powerful model, harness KB and statistics, robustness & efficiency challenge

- Structure, No Schema
- SPARQL Extensions
- Language Models
- Entity Coherence

## **Recommended Readings: Search & Ranking**

- G. Kasneci, F. Suchanek, G. Ifrim: NAGA: Searching and Ranking Knowledge. ICDE 2008
- S. Elbassuoni, M. Ramanath, G. Weikum: Query Relaxation for Entity-Relationship Search. ESWC 2011
- S. Elbassuoni, M. Ramanath, et al.: Language-model-based ranking for queries on RDF-graphs. CIKM 2009
- Z. Nie, Y. Ma, S. Shi, J.-R. Wen, W.-Y. Ma: Web Object Retrieval. WWW 2007
- H. Bast, A. Chitea, F. Suchanek, I. Weber: ESTER: efficient search on text, entities, and relations. SIGIR 2007
- J. Pound, P. Mika, H. Zaragoza: Ad-hoc object retrieval in the web of data. WWW 2010
- J. Pound, I.F. Ilyas, G.E. Weddell: Expressive and flexible access to web-extracted data: a keyword-based structured query language. SIGMOD 2010
- O. Udrea, D. Recupero, V.S. Subrahmanian: Annotated RDF. ACM Trans. Comput. Log. 11(2), 2010
- V. Hristidis, H. Hwang, Y. Papakonstantinou: Authority-based keyword search in databases. TODS 33(1), 2008
- •T. Cheng, X. Yan, K. Chang: EntityRank: Searching Entities Directly and Holistically. VLDB 2007
- D. Vallet, H. Zaragoza: Inferring the Most Important Types of a Query: a Semantic Approach. SIGIR 2008
- R. Blanco, H. Zaragoza: Finding support sentences for entities. SIGIR 2010
- P. Serdyukov, D. Hiemstra: Modeling Documents as Mixtures of Persons for Expert Finding. ECIR 2008
- R. Kaptein, P. Serdyukov, A.P. de Vries, J. Kamps: Entity ranking using Wikipedia as a pivot. CIKM 2010
- H. Fang, C. Zhai: Probabilistic Models for Expert Finding. ECIR 2007
- D. Petkova, W.B. Croft: Hierarchical Language Models for Expert Finding in Enterprise Corpora. ICTAI 2006
- M. Pasca: Towards Temporal Web Search. SAC 2008
- K. Berberich, S.J. Bedathur, O. Alonso, G. Weikum: A Language Modeling Approach for Temporal Information Needs. ECIR 2010: 13-25
- C. Zhai: Statistical Language Models for Information Retrieval. Morgan&Claypool, 2008
- B.C. Ooi (Ed.): Special Issue on Keyword Search, IEEE Data Eng. Bull. 33(1), 2010
- T. Neumann, G. Weikum: The RDF-3X engine for scalable management of RDF data. VLDB J. 19(1), 2010
- S. Ceri, M. Brambilla (Eds.): Search Computing: Challenges and Directions, Springer, 2010

## **Recommended Readings: Disambiguation**

- J. Hoffart, M. A. Yosef, I. Bordino, et al.: Robust Disambiguation of Named Entities in Text. EMNLP 2011
- R.C. Bunescu, M. Pasca: Using Encyclopedic Knowledge for Named entity Disambiguation. EACL 2006
- S. Cucerzan: Large-Scale Named Entity Disambiguation Based on Wikipedia Data. EMNLP 2007
- D.N. Milne, I.H. Witten: Learning to link with wikipedia. CIKM 2008
- S. Kulkarni, A. Singh, G. Ramakrishnan, S. Chakrabarti: Collective annotation of Wikipedia entities in web text. KDD 2009
- G. Limaye, S. Sarawagi, S. Chakrabarti: Annotating and Searching Web Tables Using Entities, Types and Relationships. PVLDB 3(1), 2010
- S. Rüd, M. Ciaramita, J. Müller, H. Schütze: Piggyback: Using Search Engines for Robust Cross-Domain Named Entity Recognition. ACL 2011
- S. Singh, A. Subramanya, F.C.N. Pereira, A. McCallum: Large-Scale Cross-Document Coreference Using Distributed Inference and Hierarchical Models. ACL 2011
- L. Ratinov, D. Roth, D. Downey, M. Anderson: Local and Global Algorithms for Disambiguation to Wikipedia. ACL 2011
- R. Navigli: Word sense disambiguation: A survey. ACM Comput. Surv. 41(2), 2009
- T. Heath, C. Bizer: Linked Data: Evolving the Web into a Global Data Space. Morgan&Claypool, 2011
- F. Naumann, M. Herschel: An Introduction to Duplicate Detection. Morgan&Claypool, 2010
- H. Köpcke, A. Thor, E. Rahm: Learning-Based Approaches for Matching Web Data Entities. IEEE Internet Computing 14(4): 23-31 (2010)



# **Thank You!**



# **Thank You!**





max planck institut informatik



DFG Deutsche Forschungsgemeinschaft