



Data Quality Issues in Constructing Knowledge Graph 知识图谱构建的质量控制





Outline

- Introduction to DQ
- Computational DQ Problems
- Data Quality Issues in Constructing KG
 - Data Cleaning in KG
 - Entity Linking in KG
 - Data Imputation in KG

Conclusions





Knowledge Explosion?



DQ Problems in DBLP					Wei Wang:16Tao Wang:18Jun Zhang:21	
 Polyseme: 10+ a Synonyms: "Pei l 	Wei Li: 27 Lei Wang: 30 Michael Wagner: 5 Jim Smith: 3					
<pre>computer science bibliography (*) Search dblp > Hore </pre>	(+) Pei Li > Home > Pers (+) Other person (-) Journal Art 2015	Sons ns with a similar name ticles	(+) P (+) P (+) Of	Pei Lee me > Persor ther persons	with a similar name @	
Exact matches Wei Wang Wei Wang 000 National University of Singapore	■ [j19] 🗎 ይ 0	 Teng Li, Jian Mao, ' Rating cloud stor Jinxin Zhang, Chac 3-D simulation statistics 	2014	E & ¢	Pei Lee, Laks V. S. Lakshmanan, Evangelos E. Milios: CAST: A Context-Aware Story-Teller for Streaming Social	Co
 Wei Wang 000: College of Nanoscale Science, University at Albany / Purdue University Wei Wang 000 School of Life Science, Fudan University, China Wei Wang 000 Center for Engineering and Scientific Computation, Zhejiang University 	r ■ [j17] 🗎 ይ 0 ∋ 2014	Reliability 55(8): 11 Reliabin Duan, Pei L Interactive Learn (2015)	[c4]	8 L ¢	Pei Lee, Laks V. S. Lakshmanan, Evangelos E. Milios: Incremental cluster evolution tracking from highly dyna Pei Lee, Laks V. S. Lakshmanan, Mitul Tiwari, Sam Shah: Modeling Impression discounting in large-scale recomm	m
 show all ikely matches Wei Wang 0010 UCLA / University of North Carolina at Chapel Hill Wei Wang 0009 Fudan University, Shanghai, China Weidong Wang Wei-Fan Wang aka: Weifan Wang 	■ [j16] 🖹 ይ 🤇	 Yingwen Chen, Mir Empirical study c 2014: 180 (2014) Pei Li, Yunchuan S Modeling and per Communication St 	2013 [c2] 2012 [c1]	■ & ¢	Pei Lee, Laks V. S. Lakshmanan, Evangelos E. Milios: KeySee: supporting keyword search on evolving events i Pei Lee, Laks V. S. Lakshmanan, Jeffrey Xu Yu: On Top-k Structural Similarity Search. ICDE 2012: 774-785	ns

show all 351 matches

Difficult Names in Google Search

6

488941 britney spears 40134 brittany spears 36315 brittney spears 24342 britany spears 7331 britny spears 6633 briteny spears 2696 britteny spears 1807 brinev spears 1635 brittny spears 1479 brintey spears 1479 britanny spears 1338 britiny spears 1211 britnet spears 1096 britiney spears 991 britaney spears 991 britnay spears 811 brithney spears 811 brtiney spears 664 birtney spears 664 brintney spears 664 briteney spears 601 bitney spears 601 brinty spears 544 brittaney spears 544 brittnav spears 364 britev spears 364 brittiny spears 329 brtnev spears 269 bretney spears 269 britneys spears 244 britne spears 244 brytney spears 220 breatney spears 220 britiany spears 199 britnney spears 163 britnry spears 147 breatny spears 147 brittiney spears 147 britty spears 147 brotney spears 147 brutney spears 133 britteney spears 133 briyney spears 121 bittany spears

29 britent spears 29 brittnany spears 29 britttany spears 29 btiney spears 26 birttney spears 26 breitney spears 26 brinity spears 26 britenay spears 26 britnevt spears 26 brittan spears 26 brittne spears 26 btittany spears 24 beitney spears 24 birtenv spears 24 brightney spears 24 brintiny spears 24 britanty spears 24 britenny spears 24 britini spears 24 britnwy spears 24 brittni spears 24 brittnie spears 21 biritney spears 21 birtany spears 21 bitenv spears 21 bratney spears 21 britani spears 21 britanie spears 21 briteany spears 21 brittay spears 21 brittinay spears 21 brtany spears 21 brtiany spears 19 birney spears 19 brirtney spears 19 britnaey spears 19 britnee spears 19 britony spears 19 brittanty spears 19 britttnev spears 17 birtny spears 17 brieny spears 17 brintty spears 17 brithy spears

9 brinttany spears 9 britanay spears 9 britinany spears 9 britn spears 9 britnew spears 9 britneyn spears 9 britrney spears 9 brtiny spears 9 brtittnev spears 9 brtny spears 9 brytny spears 9 rbitney spears 8 birtiny spears 8 bithney spears 8 brattany spears 8 breitny spears 8 breteny spears 8 brightny spears 8 brintav spears 8 brinttey spears 8 briotney spears 8 britanys spears 8 britley spears 8 britneyb spears 8 britnrey spears 8 britnty spears 8 brittner spears 8 brottany spears 7 baritney spears 7 birntey spears 7 biteney spears 7 bitiny spears 7 breateny spears 7 brianty spears 7 brintye spears 7 britianny spears 7 britly spears 7 britnej spears 7 britneyu spears 7 britniev spears 7 britnnay spears 7 brittian spears 7 briyny spears

7 brrittany spears

5 brney spears 5 broitney spears 5 brotny spears 5 bruteny spears 5 btivnev spears 5 btrittney spears 5 gritney spears 5 spritney spears 4 bittny spears 4 bnritney spears 4 brandy spears 4 brbritney spears 4 breatiny spears 4 breetney spears 4 bretiney spears 4 brfitnev spears 4 briattany spears 4 brieteny spears 4 briety spears 4 briitny spears 4 briittany spears 4 brinie spears 4 brinteney spears 4 brintne spears 4 britaby spears 4 britaev spears 4 britainey spears 4 britinie spears 4 britinney spears 4 britmney spears 4 britnear spears 4 britnel spears 4 britneuv spears 4 britnewy spears 4 britnmey spears 4 brittaby spears 4 brittery spears 4 britthey spears 4 brittnaey spears 4 brittnat spears 4 brittneny spears 4 brittnye spears 4 brittteny spears

4 briutney spears

3 britiy spears 3 britmeny spea 3 britneeev spears 3 britnehy spears 3 britnely spears 3 britnesy spears 3 britnetty spears 3 britnex spears 3 britnevxxx spears 3 britnity spears 3 brithtey spears 3 britnyey spears 3 britterny spears 3 brittneev spears 3 brittnney spears 3 brittnyey spears 3 brityen spears 3 brivtney spears 3 brltnev spears 3 broteny spears 3 brtaney spears 3 brtiiany spears 3 brtinav spears 3 brtinney spears 3 brtitany spears 3 brtiteny spears 3 brtnet spears 3 brytiny spears 3 btnev spears 3 drittney spears 3 pretney spears 3 rbritney spears 2 barittany spears 2 bbbritney spears 2 bbitney spears 2 bbritny spears 2 bbrittany spears 2 beitany spears 2 beitny spears 2 bertney spears 2 bertny spears 2 betney spears 2 betny spears 2 bhriney spears

spears spears 2 brirttany spears 2 brirttney spears 2 britain spears 2 britane spears 2 britaneny spears 2 britania spears 2 britann spears 2 britanna spears 2 britannie spears 2 britannt spears 2 britannu spears 2 britanvl spears 2 britanyt spears 2 briteenv spears 2 britenany spears 2 britenet spears 2 briteniv spears 2 britenys spears 2 britianey spears 2 britin spears 2 britinary spears 2 britmy spears 2 britnaney spears 2 britnat spears 2 britnbey spears 2 britndy spears 2 britneh spears 2 britneney spears 2 britney6 spears 2 britneve spears 2 britnevh spears 2 britneym spears 2 britneyyy spears 2 britnhey spears 2 britnjev spears 2 britnne spears 2 britnu spears 2 britonev spears 2 britrany spears

2 britreny spears

2 britsany spears

2 britry spears

Another Example with KBs

CID + ²	Name∗	Address*?	City₊⊃	Sex ^{∉∂}
11+2	张三↩	邯郸路 220 号计算机楼 527 室中	上海↩	0+2
24₽	李四↩	<u>覲秦路</u> 978 号 7 号楼 702 室↩	宁波↩	1+2
1				

CNO43	Name↩	Gender₊⊃	Address ↔	Phone/Fax¢
24+2	王五↩	F₊⊃	杭州市朝晖二区 555 号 2-308 室 310012+2	0571-88480666/+/
				0571-87074789₽
493+	李四↩	M↔	宁波市 <mark>鄭秦路</mark> 978 号 7 号楼 702 室 315012₽	0574-87074789+

	-		

NO⇔	Name⇔	Gender↔	Address≠	¢ity.₀	zip + [⊃]	Pone↔	Fax₊⊃	CID+	<u>Cno</u> +∂
1⇔	张三↩	F⇔	邯郸路 220 号	上	¢	4	¢	11₽	¢
			计算机楼 527	海↩					
			室↩						
2₽	李四↩	M+⊃	<u>鄭秦路 9</u> 78 号	宁	315012+	0574-87074789+2	¢	24₽	493¢ ²
			7702 室₽	波↩					
3₽	王五↩	F₄⊃	覲╦፩,555号	杭	310012+	1571-88480666+2	0571-+/	÷	24₽
			2-308 室↩	₩₽			88480667+2		

Different Schemas: e.g., "Sex"-"Gender", "Phone/Fax"-"Phone"+"Fax"

- Inconsistency values: e.g., "0/1"-"F/M"
- Missing values

Six DQ Dimensions



The Taxonomy of DQ Problems



Computational Data Quality Problems

- Data Integration
 - Schema Mapping
 - Record Matching
- Data Cleaning

- Data Imputation
- Data Provenance
- Data Uncertainty
- Data Constraints



Introduction to DQ Computational DQ Problems

Data Quality Issues in Constructing KG Data Cleaning in KG Entity Linking in KG Data Imputation in KG

Conclusions



- Open IE -> Knowledge Graph
- Bootstrapping Mechanisms
 - e.g.: KnowItAll, SnowBall, ProBase ...
- However, the <u>accuracy decreases sharply</u> after several iterations.



A Major Reason - Semantic drift happens



S1="Animals such as dog, cat, pig and chicken, grow fast."

S2="Yoga Postures are named after animals such as camel, pigeon, lion and cat."

S3="Common food from animals **such as** pork, beef and chicken." S4="Animals from African countries **such as** Giraffe and Lion."

(a)Semantic-based bootstrapping mechanism



- P1: "... X is a kind of mammal ..."
- P2: "Sometime, X is as clever as human beings"

(b)Syntax-based bootstrapping mechanism

Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)

Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)

16

Mutual Exclusion Bootstrapping

Pros and Cons: High Precision, Low Recall



Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)

Type Checking

Checking types of relevant entities

Pros and Cons: High Precision, Low Recall

Pillar, San Jose OK

Type Checking Arguments:

...companies such as Pillar... ... cities like San Jose...

X ,which is based in Y

Inclined pillar , foundation plate **NO**

Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)



Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)

22

Pattern-Relation Duality

Idea: The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.

Cons: still can not reach high precision and recall



Mainstream approaches

- Mutual Exclusion Bootstrapping (PACLING'07)
 - Drop those instances belonging to mutually exclusive classes
- Type Checking (WSDM'10)
 - Check the type of an entity for correctness
- Random Walk Ranking (ICDM'06)
 - Construct a graph, do random walk ranking
- Pattern-Relation Duality Ranking (WSDM'11)
 - The quality of a pattern (tuple) can be determined by the tuples (patterns) it extracts.
- A Model based on Detected Drifting Points (EDBT'14)

 Cleaning Model based on Detected Drifting Points
 Intuition: Drifting Points (DPs) are the reasons of Semantic Drift.





Properties of DPs

For a target class, the distribution of instances triggered by a DP is different from the distribution of instances that truly belong to the target class.



Z. Li et al., Overcoming Semantic Drift in Information Extraction, EDBT'14 Distributions of instances triggered by DPs and non-DPs

26

Finding Errors based on detected DPs



Z. Li et al., Overcoming Semantic Drift in Information Extraction, EDBT'14

Data Cleaning in KG – Experiments

Cleaning Method	p _{error}	r _{error}	p correct	r _{correct}
Before Cleaning	-	-	0.4305	1.0
MEx	0.9119	0.1570	0.4592	0.9832
TCh	0.9423	0.1451	0.4789	0.9724
RW-Rank	0.5753	0.5831	0.5636	0.6509
PRDual-Rank	0.5621	0.6545	0.5812	0.6940
DP Cleaning	0.9696	0.9145	0.8921	0.9393

(1) p_{error} : percentage of removed errors in all the removed instances; (2) r_{error} : percentage of removed errors in all the errors under each concept; (3) $p_{correct}$: percentage of remained correct instances in all the remained instance; (4) $r_{correct}$: percentage of remained correct instances in all the correct instances under each concept

Z. Li et al., Overcoming Semantic Drift in Information Extraction, EDBT'14

Outline

- Conventional Data Quality Problems
 - Introduction to DQ
 - Computational DQ Problems and Solutions
- Data Quality Issues in Knowledge Graph
 - Data Cleaning in KG
 - Entity Linking in KG
 - Data Imputation in KG
- □ Conclusions



Entity Linking in KG



- Also known as Entity Recognition and Disambiguation
- I. Polysemy (一词多义)
- E.g.: During his standout career at Bulls, Jordan also acts in the movies
 Space Jam.

Michael B. Jordan

(American Actor)

Michael I. Jordan

(Berkeley Professor)

Michael Jordan (NBA Player)

- D 2. Synonyms (多词一义)
- E.g.: Barack Hussein Obama(USA president)
 - m.02mjmr(Freebase)
 - Barack_Obama(Dbpedia)
 - 贝拉克·侯赛因·奥巴马(CN-Dbpedia)

Main Approaches for Solving Polysemy

- EL based on Local Compatibility (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)
- EL Based on Simple Relations (CIKM'08, AAAI'08)
- Pair-Wise Collective EL Approaches (ACL'10)
- Graph-Based Collective EL Approaches (SIGIR'11, 14)

Main Approaches for Solving Polysemy

EL based on Local Compatibility (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)

EL Based on Simple Relations (CIKM'08, AAAI'08)

Pair-Wise Collective EL Approaches (ACL'10)

Graph-Based Collective EL Approaches (SIGIR'11, 14)

- 32
- Local Compatibility Based Approaches (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)
 - Idea: Extract the discriminative features of an entity from its textual description, such as "NBA", "Basketball Player" to MJ.



During his standout career at **Bulls, Jordan** also acts in the movies **Space Jam**.

Main Approaches for Solving Polysemy

EL based on Local Compatibility (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)

EL Based on Simple Relations (CIKM'08, AAAI'08)

Pair-Wise Collective EL Approaches (ACL'10)

Graph-Based Collective EL Approaches (SIGIR'11, 14)

- Simple Relational Approaches (CIKM'08, AAAI'08)
 - Idea: the referent entity of a name mention should be coherent with its unambiguous contextual entities



Main Approaches for Solving Polysemy

EL based on Local Compatibility (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)

EL Based on Simple Relations (CIKM'08, AAAI'08)

Pair-Wise Collective EL Approaches (ACL'10)

Graph-Based Collective EL Approaches (SIGIR'11, 14)

Pair-Wise Collective Approaches (ACL'10)

Idea: Model and exploit the pair-wise interdependence between EL decisions (NP-HARD), and approximation

solutions are proposed.



Kulkarni, S. at al, Collective annotation of Wikipedia entities in web text. SIGKD'09

Main Approaches for Solving Polysemy

EL based on Local Compatibility (CIKM'07, EMNLP'07, IJCAI'09, COLING'10...)

EL Based on Simple Relations (CIKM'08, AAAI'08)

Pair-Wise Collective EL Approaches (ACL'10)

Graph-Based Collective EL Approaches (SIGIR'11, 14)

- Graph-Based Collective Approaches(SIGIR 11,14)
 - Idea: Model and exploit the global interdependence by graph-based collective EL method



Entity Linking in KG

- Also known as Entity Recognition and Disambiguation
- □ 1. Polysemy (一词多义)
- E.g.: During his standout career at Bulls, Jordan also acts in the movies
 Space Jam.

- □ 2. Synonyms (多词一义)
- E.g.: Barack Hussein Obama(USA president)
 - m.02mjmr(Freebase)
 - Barack_Obama(Dbpedia)
 - 贝拉克·侯赛因·奥巴马(CN-Dbpedia)

Entity Linking in KG – Synonyms

Approaches for Solving Synonym Problems

- String-matching based methods (CITISIA'09)
 - Edit Distance, Jaccard, Cosine, Hybrid Metrics...
- Collective alignment methods (VLDB'11, SIGKDD'13)
 Use various information of entities such as *Properties*,
 - Relations, Instances to construct a probabilistic matching model
- Based on structure similarity only (CCKS'16)
 Whole Knowledge Base Embedding

Entity Linking in KG – Synonyms

Based on structure similarity only(CCKS 16)

Idea: (1)give some initial alignments(seed entity alignments); (2) learn the embedding of the two KBs in a uniform embedding vector space connected by the seed entities "bridge"



Fig. 2. Selecting seed entities in two KBs.

Yanchao Hao at al. A Joint Embedding Method for Entity Alignment of Knowledge Bases. CCKS 2016

Outline

Conventional Data Quality Problems

- Introduction to DQ
- Computational DQ Problems and Solutions

Data Quality Issues in Knowledge Graph

- Data Cleaning in KG
- Entity Linking in KG
- Data Imputation in KG
- □ Conclusions



Data Imputation in KG

Data Imputation in KG aims at increasing the coverage of KG

Tasks

- Missing entities
- Missing types for entities (known as classification)
- Missing relations that hold between entities

Data Imputation in KG – Approaches

Type Assertions

- Internal Knowledge-based
 - SDType (ISWC'13); and some other methods
- External Knowledge-based
 - Tipola (ISWC'12); Classifier based on Wiki Links (LDOW'12)

Relation Prediction

- Internal Knowledge-based
 - Neural Tensor Network (NIPS'13) ; Mining Association Rules(ISWC'15)
- External Knowledge-based
 - Matching HTML Tables to DBpedia(WIMS'15); and some other methods

Data Imputation in KG – Approaches

□ Type Assertions

- Internal Knowledge-based
 - SDType (ISWC'13); and some other methods
- External Knowledge-based
 - Tipola (ISWC'12); Classifier based on Wiki Links (LDOW'12)

Relation Prediction

- Internal Knowledge-based
 - Neural Tensor Network (NIPS'13) ; Mining Association Rules(ISWC'15)
- External Knowledge-based
 - Matching HTML Tables to DBpedia(WIMS'15); and some other methods

SDType: using Statistical Distribution of types in the subject and object positions for predicting the instance's types.

Туре	Subject (%) Object (%)
owl:Thing	100.0 88.6
dbpedia-owl:Place	69.8 87.6
dbpedia-owl:PopulatedPlace	0.0 84.7
dbpedia-owl:ArchitecturalSt	ructure 50.7 0.0
dbpedia-owl:Settlement	0.0 50.6
dbpedia-owl:Building	34.0 0.0
dbpedia-owl:Organization	29.1 0.0
dbpedia-owl:City	0.0 24.2

Table 1. Type distribution of the property dbpedia-owl:location in DBpedia



H. Paulheim et al. Type Inference on Noisy RDF Data, ISWC'2013

Internal Methods for Type Assertions

Other Internal methods

- Training a Classification Model (e.g., SVMs)
 - E.g., Exploiting interlinks between the knowledge graphs to classify instances in one knowledge graph based on properties present in the other.
- Association Rule Mining for predict missing information.
 - Exploit association rules to predict missing types in DBpedia based on such redundancies.
- Using Topic Modeling for type prediction
 - E.g., LDA is applied to find topics for documents of entities.

Data Imputation in KG – Approaches

Type Assertions

- Internal Knowledge-based
 - SDType (ISWC'13); and some other methods
- External Knowledge-based
 - Tipola (ISWC'12); Classifier based on Wiki Links(LDOW'12)

Relation Prediction

- Internal Knowledge-based
 - Neural Tensor Network (NIPS'13) ; Mining Association Rules(ISWC'15)
- External Knowledge-based
 - Matching HTML Tables to DBpedia(WIMS'15); and some other methods

External Methods for Type Assertions

Tipalo Algorithm: identifies the most appropriate types for an entity by interpreting its natural language definition.



External Methods for Type Assertions

Classifier based on wiki Links

- using Wikipedia link graph to predict types in a KG
- interlinks between Wikipedia pages are exploited to create feature vectors, e.g., based on the categories of the related pages.



Nuzzolese et al. Type inference through the analysis of Wikipedia links, LDOW'12

Data Imputation in KG – Approaches

Type Assertions

- Internal Knowledge-based
 - SDType (ISWC'13); and some other methods
- External Knowledge-based
 - Tipola (ISWC'12); Classifier based on Wiki Links(LDOW'12)

Relation Prediction

- Internal Knowledge-based
 - Neural Tensor Network (NIPS'13) ; Mining Association Rules(ISWC'15)
- External Knowledge-based
 - Matching HTML Tables to DBpedia(WIMS'15); and some other methods

Internal Methods for Relation Prediction

Neural tensor network is suitable for reasoning over relationships between two entities.



R. Socher et al. Reasoning with neural tensor networks for knowledge base completion, NIPS'13

Internal Methods for Relation Prediction

- □ Mining Association Rules for predicting relations.
 - Mining of association rules which predict relations between entities in DBpedia from Wikipedia categories is proposed.



Kim et al. The association rule mining system for acquiring knowledge of DBpedia from Wikipedia categories, ISWC'15

Data Imputation in KG – Approaches

Type Assertions Internal Knowledge-based SDType (ISWC'13); and some other methods External Knowledge-based Tipola (ISWC'12); Classifier based on Wiki Links(LDOW'12)

Relation Prediction

- Internal Knowledge-based
 - Neural Tensor Network (NIPS'13) ; Mining Association Rules(ISWC'15)
- External Knowledge-based
 - Matching HTML Tables to DBpedia(WIMS'15); and some other methods

External Methods for Relation Prediction

Matching HTML Tables to Dbpedia

- Challenges:
 - pairs of table columns have to be matched to properties in the DBpedia ontology
 - rows in the table need to be matched to entities in Dbpedia
- Solution:
 - evaluated on a gold standard mapping for a sample of HTML tables from the WebDataCommons Web Table corpus

University	Present President		<university andrew="" d.="" hamilton="" of="" oxford,="" present="" president,=""></university>
University of Oxford	Andrew D. Hamilton		
University of Cambridge	Leszek Krzysztof Borysiewicz		 University of Oxford, present_president, Andrew D. Hamilton >
University College London	Michael Arthur		<university andrew="" d.="" hamilton="" of="" oxford,="" present_president,=""></university>

Ritze et al. Matching HTML Tables to Dbpedia, WIMS'15

External Methods for Relation Prediction

Distant supervision with a large text corpora;

- Step 1: Seed Entities in the knowledge graph are linked to the text corpus by means of Named Entity Recognition
- Step 2: Seek for text pattern which correspond to relation types
- Step 3: Apply those patterns to find additional relations in the text corpus
- A Bootstrapping way with starting seeds in KG.
- Based on web search engines:
 - Discover frequent context terms for relations
 - Use those frequent context terms to formulate search engine queries for filling missing relation values.
- Based on another KG
 - Using Interlinks between KGs to fill gaps and do knowledge transfer

Outline

Conventional Data Quality Problems

- Introduction to DQ
- Computational DQ Problems and Solutions
- Data Quality Issues in Knowledge Graph
 - Data Cleaning in KG
 - Entity Linking in KG
 - Data Imputation in KG
- Conclusions



Conclusions

- Big Data -> Big Dirty Data
 - More Challenges ...
 - More Opportunities...
- What can we do?
 - Use the rich knowledge
 - Better Precision and Recall
 - Pay Attention to Efficiency
 - Pay Attention to Cost



Thanks!

