

# 特定领域知识图谱构建初探

### 李涓子

清华大学计算机系知识工程研究室



# Outline



Knowledge graph and technologies

- Big scholar knowledge base Aminer II
- Knowledge graph building over enterprise data
- Conclusion





Increasing Connectivity





The Semantic Web. Tim Berners-Lee, James Hendler, and Ora Lassila. Scientific American, 2001.



### **Philosophy of ontology**



### Concept triangle



# [Ogden, Richards, 1923]

Ontology is the philosophical study of the nature of **being**, **becoming**, **existence**, or **reality**, as well as the **basic categories** of being and their relations. ---- Wikipedia







# **Our Knowledge graph definition**



### C – concepts

- A group of objects with same properties
- cars, students, professors

### I - instances

- A object which belongs to a concept
- Peter is a student

### T – ISA

- subConceptOf, instanceOf
- P properties
  - char本体中用于描述实例信息的其他语义关系
  - 如:instance-attribute-value (AVP)





## **Knowledge graph technologies**



### Manually KG building: Wordnet, Cyc, Hownet

### Taxonomy knowledge learning

- Learning from Wikipedia
- Learning beyond Wikipedia

### Factual knowledge learning

- Learning from Wikipedia
- Learning beyond Wikipedia



# Learning taxonomy knowledge from Wikipedia

### Category system in Wikipedia

Category system in Wikipedia as a conceptual network PHILOSOPHY and BELIEF (deals-with?) PHILOSOPHY and HUMANITIES (isa) PHILOSOPHY and SCIENCE (isa)

### Advantages:

- widely recognized concepts in human minds
- Large scale over millions of concepts and ten millions of instances

### Large coverage

### Problems:

- noise categories for different purposes
- inconsistence not well formally define



# Learning taxonomy knowledge from Wikipedia

- Using linguistic features of isa relationship
  - syntactic parsing:
    - head matching/modifier matching/Singular/plural forms
  - Lexico-patterns:
- Using structure of wikipedia
- Deriving a Large Scale Taxonc al. AAAI 07.
- Using external high quality is wordnet, Hownet, Cilin YAGO(WWW2007)
- isa relation validation using c Xlore (AAAI2014)

- 1. *NP2*,? (such as|like|, especially) *NP\* NP1 a stimulant such as caffeine*
- 2. such NP2 as NP\* NP1 such stimulants as caffeine
- 3. *NP1 NP*\* (and|or|,like) other *NP2 caffeine and other stimulants*
- 4. *NP1*, one of *det\_pl NP2 caffeine*, *one of the stimulants*
- 5. NP1, det\_sg NP2 rel\_pron caffeine, a stimulant which
- 6. NP2 like NP\* NP1 stimulants like caffeine



### Learning taxonomy knowledge beyond Wikipedia

### Using Web sources

Root concepts, search engine

- Hearst patterns
- Bootstrapping
- Taxonomy induction (structural learni domain specific taxonomy building EMNLP2010, ACL 2014

### Large scale taxonomy building

- Automatically generated from Web data
- 1.6 billion web pages
- Rich hierarchy of millions of concepted president
- Probabilistic knowledge base

SIGMOD2012

Probase: 2,653,872 concepts



Figure 1: Taxonomy Induction from Scratch.



George H. W. Bush, 0.021

George W. Bush, 0.019

20,757,54



### Factual knowledge learning



|                     | Supervised            | Semi-supervised   | Unsupervised                   |
|---------------------|-----------------------|---|--------------------------------|
| From<br>Wikipedia   |                       | Semantify Wikipdia-Kylin<br>Cross lingual IE-WikiCiKE                         |                                |
| Beyond<br>Wikipedia | Sematic<br>annotation | Distant<br>supervision(Stanford)<br>Coupled Semi-Supervised<br>Learning(NELL) | KnowItAll:<br>TextRuner<br>WOE |
|                     |                       |   |                                |



## **Automatic semantic annotation**



- Rule learning based approach Automatically learn annotation rules from the training data
- Classification based approach

Identify the boundary of tags in instances using classification models

Sequential labeling based approach

Consider the dependencies between tags

Constrained Hierarchical Conditional Random Fields

And Others ....



### Learning factual knowledge beyond Wikipedia-Knowledge Vault



# Current large scale knowledge graph is still not

### enough

| Relation       | % unknown   |  |
|----------------|-------------|--|
|                | in Freebase |  |
| Profession     | 68%         |  |
| Place of birth | 71%         |  |
| Nationality    | 75%         |  |
| Education      | 91%         |  |
| Spouse         | 92%         |  |
| Parents        | 94%         |  |
|                |             |  |

| Name                 | # Entity types | # Entity instances | # Relation types | # Confident facts (relation instances) |
|----------------------|----------------|--------------------|------------------|--|
| Knowledge Vault (KV) | 1100           | $45\mathrm{M}$     | 4469             | 271M                                   |
| DeepDive [30]        | 4              | $2.7\mathrm{M}$    | 34               | $7 M^a$                                |
| NELL [8]             | 271            | $5.19 \mathrm{M}$  | 306              | $0.435 \mathrm{M}^b$                   |
| PROSPERA [28]        | 11             | N/A                | 14               | 0.1M                                   |
| YAGO2 [18]           | 350,000        | 9.8M               | 100              | $4\mathrm{M}^{c}$                      |
| Freebase [4]         | 1,500          | 40M                | 35,000           | $637 \mathrm{M}^d$                     |
| Knowledge Graph (KG) | 1,500          | 570M               | 35,000           | $18,000\mathrm{M}^e$                   |

Table 1: Comparison of knowledge bases. KV, DeepDive, NELL, and PROSPERA rely solely on extraction, Freebase and KG rely on human curation and structured sources, and YAGO2 uses both strategies. Confident facts means with a probability of being true at or above 0.9.



### Learning factual knowledge beyond Wikipedia-Knowledge Vault



### Motivation

the new approach should automatically leverage alreadycataloged knowledge to build prior models of fact correctness
TXT: Distant supervision

### Framework



TXT: Distant supervisionDOM: DOM tree structure featuresTBL: Table informationANO: annotated tags in htmls



Priors: Path ranking algorithm Priors: Neural network method



# Learning factual knowledge beyond Wikipedia-Knowledge Vault



| System    | #     | # > 0.7 | # > 0.9           | Frac. $>0.9$ | AUC   |
|-----------|-------|---------|-------------------|--------------|-------|
| TBL       | 9.4M  | 3.8M    | 0.59M             | 0.06         | 0.856 |
| ANO       | 140M  | 2.4M    | $0.25 \mathrm{M}$ | 0.002        | 0.920 |
| TXT       | 330M  | 20M     | 7.1M              | 0.02         | 0.867 |
| DOM       | 1200M | 150M    | 94M               | 0.08         | 0.928 |
| FUSED-EX. | 1600M | 160M    | 100M              | 0.06         | 0.927 |

 Table 2: Performance of different extraction systems.





# Summary



- Various knowledge representation and learning methods
- What are the effective methods used for domain specific knowledge graph building?
- What are the proper representation for domain specific knowledge graph?



# Outline



- Knowledge graph and technologies
- Big scholar knowledge base Aminer II
- Knowledge graph building over enterprise data
- Conclusion









Xindong Wu 📀

Supervised Learning

Professor

③ 35 views

Similar

H-Index: 45 | #Paper: 331 | #Citation: 9644

P Department of Computer Science, University of Vermont

Machine Learning Information Extraction Bayesian Networks Data Mining

Grid computing Database marketing Parallel computing

#### 数据挖掘

2+ Follow

数据挖掘(Data Mining)是通过分析每个数据,从大量数据中寻找其规律的技术,主要有数 据准备、规律寻找和规律表示3个步骤。数据挖掘的任务有关联分析、聚类分析、分类分 析、异常分析、特异群组分析和演变分析等。

上位词:

资料分析 计算机科学基础理论 决策支持系统 信息管理术语 数据挖掘 形式科学

#### Miner

All

09

· Italy Germany
Denmark

India

 New Zealand Japan

Israel

**Computer Science** 

High Performance

Computer Network Net and Information Security

Computing

Whatever comes to your mind

🕂 Home | 🍈 🕞

| Rank            | Confere  | ence (Full Name)  | Short Name     | Impact Factor |
|-----------------|----------|---|----------------|---------------|
| 1               | Science  |   |                | 162.00        |
| 2               | Nucleic  | Acids Research  | NAR            | 128.00        |
| 3               | IEEE Co  | nference on Computer Vision and Pattern Recognition     | CVPR           | 112.00        |
| 4               | IEEE Tra | insactions on Pattern Analysis and Machine Intelligence | TPAMI          | 101.00        |
| 5               | NeuroIm  | lage  | NeuroImage     | 99.00         |
| 6               | IEEE Tra | insactions on Industrial Electronics                    | TIE            | 80.00         |
| 7 Bioinforme    |          | Tratics   | Bioinformatics | 79.00         |
|                 |          | munications Magazine                                    | ICM            | 74.00         |
|                 |          | sactions on Signal Processing                           | TSP            | 73.00         |
|                 |          | erence on Human Factors in Computing Systems            | CHI            | 71.00         |
|                 |          | tology ePrint Archive                                   |                | 71.00         |
|                 |          | sactions on Image Processing                            | TIP            | 69.00         |
|                 |          | a   | Automatica     | 69.00         |
|                 |          | national Conference on Computer Communications          | INFOCOM        | 69.00         |
| time ser<br>sam | les      | cations of the ACM                                      | Commun. ACM    | 69.00         |
|                 |          | ns on Automatic Control                                 | TAC            | 67.00         |
|                 |          | al World Wide Web Conferences                           | www            | 66.00         |
|                 |          | gs of the IEEE  | Proc. IEEE     | 64.00         |
|                 |          |   |                |               |

## Conference Ranking

ACM Knowledge Discovery and Data Mining

|   | semi-supervised learning                      |
|---|---|
| Search for Conference                   | unlabeled data proposed framework             |
| Conference/Journal                      | case study optimization problem topic model   |
| ACM Knowledge Discovery and Data Mining | feature selection data mining                 |
| From Year                               | efficient algorithm large graph use           |
| 2008                                    | synthetic data social media social network    |
| To Year                                 | search engine JUUICI IIUUVUIN                 |
| 2013                                    | different type proposed algorithm wate stream |
|   | real-world data set                           |
| Submit                                  |   |

Author Distribution







Top cited authors

| 2008                       | 2009                               | 2010                           | 2011                                       | 2012                               | 2013                             |
|----------------------------|------------------------------------|--------------------------------|--|------------------------------------|----------------------------------|
| Yehuda Koren:1/924         | Jure Leskovec:1/778 <sup>+</sup>   | Wei Chen:1/437↑                | Eunjoon Cho:1/541↑                         | Jing Yuan:1/137 <sup>+</sup>       | Chris Thornton:1/64              |
| Jure Leskovec:1/474        | Wei Chen:1/578                     | Yu Wang:1/168                  | Dashun Wang:1/208                          | Thanawin<br>Rakthanmanon:1/108↑    | Bin Liu:1/33↑                    |
| Jie Tang:1/440             | Jie Tang:1/384                     | Hongning Wang:1/162            | Jing Yuan:1/171↑                           | Alan Ritter:1/103                  | Yu Zheng:1/33↑                   |
| Yabo Xu:1/423              | Mohsen Jamali:1/312                | Deng Cai:1/135 <sup>+</sup>    | Rainer Gemulla:1/167                       | Isabelle Stanton:1/86 <sup>+</sup> | Hongzhi Yin:1/31                 |
| Victor S. Sheng:1/420      | Justin Ma:1/259                    | Maayan Roth:1/131↑             | Salvatore Scellato:1/164                   | Ashton Anderson:1/781              | Arjun Mukherjee:1/31             |
| David Crandall:1/404       | Albert Bifet:1/249                 | Liang Xiang:1/128 <sup>+</sup> | Marco Pennacchiotti:1/106 <sup>+</sup>     | Ling-Yin Wei:1/75                  | Charalampos<br>Tsourakakis:1/27↑ |
| Aris Anagnostopoulos:1/392 | Theodoros Lappas:2/235             | Arik Friedman:1/125            | Noman Mohammed:1/93                        | Rui Li:1/75                        | Quan Yuan:1/26↑                  |
| Huanhuan Cao:1/313         | Anna Monreale:1/229 <sup>+</sup>   | Min-Ling Zhang:1/113           | Mao Ye:1/90↑                               | Jie Tang:2/68 <sup>↑</sup>         | Reza Zafarani:1/24               |
| David J. Crandall:1/218    | Sayali Kulkarni:1/228 <sup>↑</sup> | Yong Ge:1/104                  | Wei Liu:1/89                               | Jiliang Tang:2/67                  | Madhav Jha:1/231                 |
| lan Porteous:1/216         | Prem Melville:1/207                | Michael Jahrer:1/103↑          | Robson Leonardo Ferreira<br>Cordeiro:1/74↑ | Xiwang Yang:1/65↑                  | Wook-Shin Han:1/22↑              |



# **Reviewer** Suggestion

Inerest matching COI avoiding Load balancing Forcast review quality

#### ArnetMiner: extraction and mining of academic social networks. Authors (Seperated by comma) Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, Zhong Su ICDM × PAKDD Abstracts This paper addresses several key issues in the ArnetMiner system, which aims at extracting and mining academic social networks. Specifically, the system focuses on: 1) Extracting researcher profiles automatically from the Web; 2) Integrating the publication data into the network from existing digital libraries; 3) Modeling the entire academic network; and 4) Providing search services for the academic + Relevant Publications: network. So far, 448,470 researcher profiles have

been extracted using a unified tagging approach. We integrate publications from online Web databases and propose a probabilistic framework to deal with the name ambiguity problem. Furthermore, we propose a unified modeling approach to

| Conference | (Journal) |
|------------|-----------|
|------------|-----------|



**ArnetMiner** 

Title

SEARCH



#### Yajun Wang Microsoft Research Asia

H-index: 14, #Papers: 49, #Citations: 677

social network, principal component analysis, Shortest Path

|   | Ming-Syan Chen ( ALIAS: Ming-Syan Syan Chen ) |
|---|---|
|   | National Taiwan University                    |
| Z | H-index: 42, #Papers: 278, #Citations: 9978   |
|   | Data Mining,Data Streams,Data Replication     |
|   | + Relevant Publications:                      |
|   | Michael R. Berthold (ALIAS: Michael Berthold) |









+ Relevant Publications:

German Software Development Lab, IBM

#### H-index: 4, #Papers: 12, #Citations: 35

KNIME.com, University of Konstanz

H-index: 17, #Papers: 77, #Citations: 1180

Knowledge-Based Methods, Proof Transformation, Der rechtliche Schutz von

International Symposium, Data Analysis, Second International Symposium

| + Relevant Publications: | 🚹 🗊 🧼    |
|--------------------------|----------|
|                          |          |
|                          | <u> </u> |



#### Zementis

#### H-index: 4, #Papers: 7, #Citations: 68

cloud computing, neural networks, open standard, predictive analytics, data mining, predictive model markup language, pmml

| + Relevant Publications: | ∩ 🖸 🤛 |
|--------------------------|-------|
|                          |       |



#### H-index: 4, #Papers: 12, #Citations: 124

social network, interactive recommendation, Efficient Confident Search

#### + Relevant Publications:



ET 🍑

🖪 🗊 🧼

**Account** 

Welcome jietang

# **Reviewer Suggestion**





# AMiner II (ArnetMiner)



- Academic Social Network Analysis and Mining system—AMiner (http:// aminer.org)
  - Online since 2006
  - >38 million researcher profiles
  - >76 million publication papers
  - >241 million requests
  - >12.35 Terabyte data
  - 100K IP access from 170 countries per month
- 10% increase of visits per month
   Deep analysis, mining, and search





### 7.32 million IP from 220 countries/regions



### **Top 10 countries**

- 1. USA
- 2. China
- 3. Germany
- 4. India
- 5. UK

- 6. Canada
- 7. Japan
- 8. Spain
- 9. France
  - 10. Italy



# **Motivating Example**





() 清華大学 Tsinghua University





| 70.60% of the researchers have<br>at least one homepage or an<br>introducing page |                                     |  |  |
|---|-------------------------------------|--|--|
| 85.6% from universities   | 14.4% from companies                |  |  |
| 71.9% are homepages   | 28.1% are<br>introducing<br>pages   |  |  |
| 40% are in lists and tables   | 60% are<br>natural<br>language text |  |  |



**CRFs** 



- Green nodes are hidden vars, - Purple nodes are observations





## **Processing Flow for Profiling**





NEG

# Profiling Results—5-fold cross validation

| Profiling Task | Unified | Unified_NT | SVM   | Amilcare |
|----------------|---------|------------|-------|----------|
| Photo          | 89.11   | 88.64      | 88.86 | 31.62    |
| Position       | 69.44   | 64.70      | 64.68 | 56.48    |
| Affiliation    | 83.52   | 72.16      | 73.86 | 46.65    |
| Phone          | 91.10   | 78.72      | 79.71 | 83.33    |
| Fax            | 90.83   | 64.28      | 64.17 | 86.88    |
| Email          | 80.35   | 75.47      | 79.37 | 78.70    |
| Address        | 86.34   | 75.15      | 77.04 | 66.24    |
| Bsuniv         | 67.38   | 57.56      | 59.54 | 47.17    |
| Bsmajor        | 64.20   | 59.18      | 60.75 | 58.67    |
| Bsdate         | 53.49   | 40.59      | 28.49 | 52.34    |
| Msuniv         | 57.55   | 47.49      | 49.78 | 45.00    |
| Msmajor        | 63.35   | 61.92      | 62.10 | 57.14    |
| Msdate         | 48.96   | 41.27      | 30.07 | 56.00    |
| Phduniv        | 63.73   | 53.11      | 57.01 | 59.42    |
| Phdmajor       | 67.92   | 59.30      | 59.67 | 57.93    |
| Phddate        | 57.75   | 42.49      | 41.44 | 61.19    |
| Overall        | 83.37   | 72.09      | 73.57 | 62.30    |



# Outline



- Knowledge graph and technologies
- Big scholar knowledge base Aminer II
- Knowledge graph building over enterprise data
- Conclusion





## Knowledge Graph over Enterprise Data

### Motivation

# The current constructions of the knowledge graph are mainly from two aspects: Web, Domains, Science



There is huge demand on knowledge graph building based on internal data of enterprise



### Building Knowledge Graph based on Mobile Customer Care Documents

- Document Parsing based Logical Structure Extraction
- Heuristic Table Extraction
- Hierarchical Concept Extraction
- Iterative Instance Identification & Property Extraction
- Evaluation Throughout Performance



# **Knowledge Graph over Enterprise Data**

### Evaluations

Document Parsing Evaluation

 $section\_precision = \frac{correct\_section}{sum\_ext\_section}$ 

block\_prec ison = correct\_bl ock sum\_ext\_bl ock → bloc
 Table Alignment Evaluation

- via manual evaluation
- Knowledge Graph Evaluation
  - Coverage



block\_recall = 
$$\frac{correct_block}{sum_org_block} \leftrightarrow$$

# Summary



- Domain data characteristics
- Problem to be solved
- Building pipeline
- Linking with external knowledge graphs
- Visualization and Evaluation
  - On each process
  - Knowledge base evaluation

Human interaction



### Conclusions



- There have a lot of available knowledge graphs and large scale linked data
- There have various knowledge representation and learning methods for different kinds of format sources
- knowledge graph building pipeline and toolkits within specific domains can facilitate the fast building of specific domain knowledge graph
- Making using existing knowledge graphs as possible as we can





# **Thanks & Questions**

