# SEED: 基于RDF知识库的实体扩展与知识纠错系统





研究背景
问题定义
问题定义
关键技术
实验结果
系统展示
总结与展望



问题定义



Problem: given some seeds, find relevant entities having some implicit common features with those seeds

## **Semantic Patterns**



$$\pi_1 = s_1 : p_{11}$$

$$\pi_2 = s_2 : \underline{p_{21}} \circ \underline{p_{12}}$$

$$\Phi(e_1, e_2, e_3) = \{\pi_1, \pi_2\}$$

关键技术

### □ 引入semantic pattern的概念

- Step 1: 根据种子节点找到相关的 semantic patterns;
- Step 2: 找到满足已找到的 semantic patterns的候选实体;
- Step 3: 根据semantic patterns对候 选实体排序(ranking model)



□ 实体相关度 
$$r(e) = \sum_{\pi \in \tilde{\Phi}(Q) \land e \models \pi} d(\pi) * r(\pi, Q)$$

$$\square \text{ Discriminability} \quad d(\pi) = \prod_{i=1}^{n} \min(\frac{1}{|E(\pi_i)|}, \frac{|E(\pi_i)|}{\sum_{e \in E(\pi_i)} |E(e:l_i)|})$$

**Relevance**  $r(\pi, Q) = \prod_{e \in Q} p(e, \pi)$ 

实验设置

□ RDF 知识库

DBPedia Version 3.9, 4M entities

□ 测试集

□ QALD-2/3/4, INEX-XER 2009

□ 评价指标

Precision@N, R-Precision, MRR

□ 对比对象

BBR, SEAL, LDSD, ARM and QBESS

主要结果

Table 1: Comparison on the QALD dataset						Table 2: Comparison on the INEX dataset							
SLT	seeds	p@5	p@10	p@20	MRR	R-pre	SLT	seeds	p@5	p@10	p@20	MRR	R-pre
SEAL	2	.377	.290	.208	.550	.269	SEAL	2	.412*	.388*	.331*	.542*	.327*
BBR	2	.350	.307	.253	.478	.283	BBR	2	.312	.265	.200	.501	.208
LDSD	2	.133	.100	.078	.273	.100	LDSD	2	.092	.081	.063	.237	.072
QBESS	2	.400	.353	.267	.545	.338	QBESS	2	.246	.223	.186	.330	.202
ARM	2	.503•	.420•	.322•	.661•	.375•	ARM	2	.223	.217	.188	.367	.199
ESER	2	.547 <sup>•</sup>	<b>.460</b> <sup>●</sup>	<b>.372</b> *	.699 <b>•</b>	<b>.457</b> <sup>●</sup>	ESER	2	.400*	.383*	.287*	<b>.551</b> *	.304*
SEAL	3	.453	.363	.267	.591	.340	SEAL	3	.462*	.433*	.354*	.547*	<b>.377</b> *
BBR	3	.467	.363	.287	.604	.337	BBR	3	.312	.288	.217	.506	.242
LDSD	3	.103	.082	.063	.194	.084	LDSD	3	.092	.065	.047	.247	.055
QBESS	3	.410	.343	.266	.536	.359	QBESS	3	.219	.194	.169	.310	.190
ARM	3	.550•	.468•	.372•	.665•	.446•	ARM	3	.188	.185	.168	.262	.187
ESER	3	.613 <b>•</b>	<b>.498</b> <sup>●</sup>	<b>.387</b> *	<b>.773</b> *	.501 <sup>•</sup>	ESER	3	<b>.500</b> *	.415*	.311*	<b>.684</b> *	.340*
SEAL	4	.420	.350	.270	.539	.354	SEAL	4	.423*	.383*	.319*	.530*	.339*
BBR	4	.467	.365	.290	.654	.345	BBR	4	.312	.290	.239	.527	.252
LDSD	4	.097	.077	.055	.252	.075	LDSD	4	.088	.069	.057	.193	.067
QBESS	4	.373	.290	.225	.465	.340	QBESS	4	.181	.163	.144	.222	.161
ARM	4	.527•	.430 <b>•</b>	.348•	.716•	.420	ARM	4	.254	.233	.186	.345	.208
ESER	4	.613 <sup>•</sup>	.502 <b>•</b>	. <b>392</b> *	<b>.801</b> *	<b>.525</b> *	ESER	4	.504.	.446*	.341*	<b>.633</b> *	<b>.376</b> *
SEAL	5	.410	.317	.247	.535	.352	SEAL	5	.377	.340	.284	.418	.311
BBR	5	.453	.362	.288	.627	.336	BBR	5	.346	.327	.257	.593	.282
LDSD	5	.130	.083	.060	.284	.090	LDSD	5	.092	.075	.059	.179	.069
QBESS	5	.370	.285	.227	.488	.353	QBESS	5	.173	.154	.135	.205	.153
ARM	5	.503•	.418 <b>•</b>	.342•	.665•	.426	ARM	5	.227	.227	.187	.322	.219
ESER	5	.563 <sup>•</sup>	.465 <b>•</b>	.381*	<b>.726</b> <sup>●</sup> <sub>*</sub>	<b>.515</b> *	ESER	5	.492 <b>•</b>	.433 <b>•</b>	.336 <sup>•</sup>	.629 <b></b> ∗	.381 <u>*</u>
SEAL	mix	.447	.347	.249	.592	.335	SEAL	mix	.473*	.398*	.305	.644	.330
BBR	mix	.447	.370	.292	.570	.347	BBR	mix	.358	.315	.251	.597	.279
LDSD	mix	.130	.100	.069	.268	.100	LDSD	mix	.112	.087	.062	.271	.063
QBESS	mix	.510	.417	.332	.630	.420	QBESS	mix	.427	.367	.273	.579	.300
ARM	mix	.537•	.443•	.337•	.646•	.433 <b>•</b>	ARM	mix	.350	.304	.239	.535	.273
ESER	mix	.633 <b>•</b>	.510 <sup>•</sup>	.403 <b>•</b>	<b>.799</b> <sup>●</sup>	<b>.</b> 559∗	ESER	mix	.515 <sup>•</sup>	<b>.440</b> <sup>●</sup>	.350*	<b>.701</b> *	<b>.409</b> <sup>●</sup>



用户界面

::	(-			ler v					c
) Tin	_gray_(c	e: 5733 ms	michael_stonebra		Cr	iarie:	S_bachman x	vith New Triples	s: 🛛 Y
ntiti	es			Sen	nanti	c Pa	tterns		
1¦	ତ ତ	Entity	Up	1;	¢	©	Anchor Entity	Predicate	Filte
1	۲	michael_stonebraker	t	1	$\bigcirc$	0	category:database_researchers	subject	
2		jim_gray_(computer_scientist)	t	2		0	category:fellows_of_the_association_for_computing_machi	<u>subject</u>	
3	$\bigcirc$	edgar_fcodd	t	3	$\bigcirc$	0	computer_science	field	
4	$\bigcirc$	charles_bachman	t	4	$\bigcirc$	0	category:american_computer_scientists	<u>subject</u>	
5	$\bigcirc$	david_dewitt	t	5	$\bigcirc$	+	category:turing_award_laureates	subject	
6	0	peter_chen	Ť	6	$\bigcirc$	0	category:sigmod_edgar_fcodd_innovations_award_winners	subject	
7	$\bigcirc$	hvjagadish	Ť	7	$\bigcirc$	+	turing_award	award	
8	$\bigcirc$	stanley_zdonik	Ť	8	$\bigcirc$	0	university_of_michigan	almamater	
9	$\bigcirc$	joseph_mhellerstein	Ť	9	$\bigcirc$	0	category:living_people	subject	

ICDE'16 Demo

知识纠错

- □ Semantic patterns能够帮助我们预测知识缺陷
- □ Association rule mining方法的引入
  - Top相关实体作为transactions
  - ■每个实体对应谓词集合作为一个transaction的items
  - ■找到所有的频繁项X
  - ■如果存在一个频繁项X能够imply一个谓词p,则预测 p可能是e的一个谓词
  - □同理,可以预测另一个实体e'是不是的e关联实体
  - 进而预测e是否满足一个SP: e':p, 即<e, p, e'>是否 成立

系统展示

- □ 前端:
  - Html5+CSS3+Javascript+Ajax
    - Bootsrap
    - Echarts



□ URL:

http://202.112.114.205:8080/SEED

## 总结与展望

#### □ 总结

- ■老问题上的新方法
- ■性能优良
- ■扩展性强

#### □ 展望

- □ 加强纠错功能
  - 多版本控制
  - 内部错误挖掘
- ■提升查询性能
  - ■速度
  - ■准确度

#### □谢谢关注!

#### □ 陈跃国

□ chenyueguo@ruc.edu.cn