

GRAMS: A Graph-based Approach for Inferring Semantic Descriptions of Wikipedia Tables

Binh Vu, Craig Knoblock, Pedro Szekely, Jay Pujara, Minh Pham

> Information Sciences Institute University of Southern California



Information Sciences Institute

Motivating Example



• Wikipedia has 7.5 millions tables covering many domains

List of players won Walter Payton Award

Members of 56th New Brunswick Legislature						Year +	P	layer	\$	Sc	hool	\$	Position +		
Name 🗢				Party	\$		Ri	1987	Kenny	Gamble		Colgate			RB
Hédard Albert			Liberal			Caraque	t	1988	Dave	Meggett		Towson	State		RB
David Alward			Progressive Conservative			Woodsto	stock 2018 2019 usie-Rest		Devlin	evlin Hodges Sar ey Lance Nor		Samford North Dakota State			QB QB
Donald Arseneault			Liberal			Dalhousi			Trey L					e	
John Betts			Progressive Conservative			Moncton Crescent			nt pericarditis						
Dereham Dereha		am	m Inmazeb			Regeneron			Ebola Virus						
Aldeby Aldeby			Olorofim			F2G			invasive mold infections						
Ashwellthorpe Ashwe		llth	orpe	1881	1939		Great	at Eastern							
Hotel Diablo			School			Rail		She	Shellharbour		19	1959		Website &	
FLYGOD is an Aweson			esome	GUD	VV	estside G	unn	нір г	юр						<u> </u>

Source Modeling Problem



- Building semantic descriptions of tables
 - Describing data source using classes and properties in ontologies



Third Presidents of National Council (Austria)

Source Modeling Problem



- Building semantic descriptions of tables
 - Describing data source using classes and properties in ontologies



Main Idea



• Information of entities in KGs can help source modeling \Rightarrow need little training data

President of the National Council (Austria)

From Wikipedia, the free encyclopedia

List of third presidents [edit]



USC Viterbi School of Engineering

Information Sciences Institute





Approach





USC Viterbi School of Engineering

Construct Candidate Graph: Discovering Links

2002

2006

2008

contex

Create a graph of cells and co

homa

Eva

Martin





Construct Candidate Graph: Discovering Links



• Add links discovered from knowledge in Wikidata





• Group links of cells from same source & target







P582: end time





P39 : position he P580: start time P582: end time











• Final candidate graph



After Building Candidate Graph



- Candidate (n-ary) relationships from the candidate graph
- Candidate columns' types from entities in table columns
- \Rightarrow Need to select the most appropriate relationships and



Approach

Inputs

- A target knowledge graph: Wikidata
- A linked relational table T
- A set of contextual values C
- 1. Construct candidate graph
- 2. Infer semantic description

Outputs:

• A semantic description of (T, C)



Collective Reasoning Problem



• Probabilistic Soft Logic (PSL)

"A probabilistic graphical models framework using firstorder logic"

- Two main elements: predicates and rules
 - Predicates have "soft" value in [0, 1]
 - Rules converted to exponential function to approximate P(x)

PSL Predicates (examples)

- CorrectRel(N₁, N₂, P): if a relationship is correct
 - CorrectRel(Name, stmt₁, P39)
 - CorrectRel(stmt₁, Entered Office, P580)
 - CorrectRel(stmt₁, Third President, P39)
- CorrectType(N₁, T): if a column type assignment is correct
 - CorrectType(Party, Organization)
 - CorrectType(Party, Political Party)
 - CorrectType(Name, Human)
- ... and more

P39: position held P580: start time P582: end time



P39

Third President



P580

Enterec Office

PSL Rules (examples)



By default, relationships/types are incorrect
1a. ¬ CorrectRel(N₁, N₂, P)
1b. ¬ CorrectType(N₁, T)

2. Relationships/types are correct/incorrect based on evidence

2a. FreqMatch(N₁, N₂, P) \rightarrow CorrectRel(N₁, N₂, P) 2b. FreqDiff(N₁, N₂, P) $\rightarrow \neg$ CorrectRel(N₁, N₂, P) 2c. FreqTypeMatch(N₁, T) \rightarrow CorrectType(N₁, T) 2d. ...and more



PSL Rules (examples)



3. If a statement value is incorrect, then the statement's qualifiers are also incorrect



4. We prefer fine-grain properties than high-level properties



5. ...and more

Information Sciences Institute



Post-Processing



School of Engineering

• PSL outputs probability of each relationships and types.



- Use BANK algorithm to choose the most probable relationships
 - Avoid unnecessary loops
 - Prefer tree structure if possible



Information Sciences Institute

Evaluation of GRAMS



- Collective reasoning is beneficial
 - Avoid cascading errors from subject column detection phase
 - Handle complex schema: multiple entities' types and n-ary relationshima

	Dataset	Method	(CPA		CTA			
	2 414200		Precision	Recall	$\mathbf{F_1}$	Precision	Recall	$\mathbf{F_1}$	
Wilkingdig		MantisTable	0.535	0.442	0.484	0.928	0.331	0.488	
wikipedia		$MantisTable^*$	0.559	0.569	0.564	0.940	0.394	0.556	
lables	950WT	BBW	0.796	0.123	0.214	0.850	0.233	0.367	
	250 W 1	BBW*	0.740	0.559	0.638	0.759	0.777	0.768	
		GRAMS-ST	0.526	0.681	0.594	-	-	-	
		GRAMS	0.824	0.650	0.726	0.819	0.813	0.816	
Synthetic		MantisTable	0.985	0.976	0.981	0.977	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	0.880	
lables	SamTab2020	GRAMS-ST 0.526 0.681 0.594 GRAMS 0.824 0.650 0.726 MantisTable 0.985 0.976 0.981 b2020 BBW 0.996 0.995 0.995 GRAMS-ST 0.990 0.989 0.990	0.995	0.980	0.980	0.980			
	Sem 1a02020	GRAMS-ST	0.990	0.989	0.990	-	0.331 0 0.394 0 0.233 0 0.777 0 - 0.813 0 0.800 0 0.980 0 - 2 0.981 0	-	
		GRAMS	0.996	0.994	0.995	0.982	0.981	0.982	

MantisTable* and BBW* are modified to retrieve correct subject

Related Work



				Modeling Capabilities					
	Method		Data Hungry	Handle Literal Columns	Handle Qualifiers	Denormalize d Tables			
Custom	Taheriyan et al. 2016		Υ	Y	Υ	Υ			
s	Vu et al. 2	2019	Υ	Y	Υ	Υ			
KG Ontologie s		Ritze et al. 2015	_	Y	Ν	Ν			
	lterative Method	Zhang et al. 2017	_	Y	Ν	Ν			
		SemTab systems	_	Y	Ν	Ν			
	Graphic al Models	Limaye et al. 2010	_	Ν	Ν	Υ			
		Mulward et al. 2013	_	N	Ν	Υ			
		GRAMS	_	Y	Υ	Y			

Discussion and Future work



- Contribution: A novel graph-based approach, GRAMS, for building semantic descriptions of Wikipedia Tables.
 - The candidate graph makes it easy to represent and discover n-ary relationships.
 - Using PSL to collectively infer correct relationships and types.
- Future work:
 - Handle unlinked tables

Entity Linking		GRAMS			human (Q5) 1 party (P102) political party (07278) 1				
Name	Year	Name	Year		rdfs:label rdfs:label				
Willi Brauneder	1996	Willi Brauneder	1996	Name position held (P39) resident of the National Council of Austria (0223	Name wikibase:Statement 1 Party end time (P582)				
Thomas Prinzhorn	2002	Thomas Prinzhorn	2002		position held (P39) start time (P580) resident of the National Council of Austria (Q22328268)				
Eva Glawischnig-Piesczek	2006	Eva Glawischnig-Piesczek	2006		Entered Office				

Generate large labeled dataset from Wikipedia tables to train semantic modeling systems