



Discovering Graph Patterns for Fact Checking in Knowledge Graphs

Peng Lin Qi

Qi Song

Jialiang Shen

Yinghui Wu

Washington State University

Beijing University of Posts and Telecommunications Washington State University Pacific Northwest National Laboratory







What is fact checking?

Knowledge Graph (KG): G = (V, E, L)



Fact: a triple predicate

Triple $< v_x, r, v_y >$

- v_x and v_y are two nodes;
- x and y are node labels;
- r is a relationship;

e.g., <Cicero, influencedBy, Plato>

- v_x = "Cicero", v_y = "Plato"
- x, y = "philosopher"
- r = "influencedBy"

Fact checking answers if a fact belongs to the missing part of KG.



Graph structure can be evidence for fact checking.



Graph structure can be evidence for fact checking.

Rule Model: Graph Fact Checking Rules (GFC)



A GFC rule contains two patterns connected by two anchored nodes.

- Given: G = (V, E, L)
- GFC $\varphi : P(x, y) \rightarrow r(x, y)$
- True facts Γ⁺:
 - sampled from the edges *E* in *G*.
- False facts Γ⁻:
 - sampled from node pairs (v_x, v_y) that have no *r* between them.
 - following partial closed world assumption (PCA)



Statistical measures are defined in terms of graph and a set of training facts.

Support and Confidence

GFC: $\varphi : \mathbf{P}(\mathbf{x}, \mathbf{y}) \rightarrow \mathbf{r}(\mathbf{x}, \mathbf{y})$

• supp
$$(\varphi) = \frac{|P(\Gamma^+) \cap r(\Gamma^+)|}{|r(\Gamma^+)|}$$

Ratio of facts can be covered out of r(x, y) triples.

$$r(vx, v_y) = r(vx, v_y)$$

supp = 2/3
$$P = P$$

$$(vx, v_y) = (vx, v_y)$$

$$(vx, v_y) = (vx, v_y)$$

$$(vx, v_y) = (vx, v_y)$$

$$(vx, v_y) = (vx, v_y)$$

Ρ

•
$$\operatorname{conf}(\varphi) = \frac{|P(\Gamma^+) \cap r(\Gamma^+)|}{|P(\Gamma^+)_N|}$$

Ratio of facts can be covered out of (x, y) pairs, under **PCA**.

Support and confidence are for pattern mining.

Ρ

GFC: $\varphi : P(x, y) \rightarrow r(x, y)$

G-Test score

$$\operatorname{sig}(\varphi, p, n) = 2|\Gamma^{+}|(p\ln\frac{p}{n} + (1 - p)\ln\frac{1 - p}{1 - n})$$

p and n are the supports of P(x, y) for positive and negative facts, respectively.

A "rounded up" score $\max\{sig(\varphi, p, \delta), sig(\varphi, \delta, n)\}$ is used in practice. where δ is a small positive to prevent infinities.

In our work, we also normalize it between 0 and 1 by a sigmoid function.

Significance is the ability to distinguish true and false facts.

Diversity

S is a set of GFCs.

$$\operatorname{div}(\boldsymbol{S}) = \frac{1}{|\Gamma^+|} \sum_{t \in \Gamma^+} \sqrt{\sum_{\varphi \in \Phi_t(\boldsymbol{S})} \operatorname{supp}(\varphi)}$$

 $\Phi_t(S)$ is the GFCs in S that cover a true fact t. E.g. $S_1 = \{P_1, P_2, P_3\}, S_2 = \{P_4, P_5, P_6\}$



Diversity is to measure the redundancy of a set of GFCs

To cope with diversity, the total significance $\operatorname{sig}(S) = \sqrt{\sum_{\varphi \in S} \operatorname{sig}(\varphi)}$.

Coverage function:
$$cov(S) = sig(S) + div(S)$$

Problem formulation:

Given graph G, support threshold σ and confidence threshold θ , and a set of true facts Γ^+ and a set of false facts Γ^- , and integer k, identify a size-k set of GFCs S, such that:

(a) For each GFC φ in S, supp $(\varphi) \ge \sigma$, conf $(\varphi) \ge \theta$. (b) cov(S) is maximized.

More significance, less redundancy.

• $\operatorname{cov}(S)$ is a set function. marginal gain: $\operatorname{mg}(S) = \operatorname{cov}(S \cup \{\varphi\}) - \operatorname{cov}(S)$

- cov(S) is monotone.
 Adding elements to S does not decrease cov(S).
- $\operatorname{cov}(S)$ is submodular. If $S_1 \subseteq S_2$ and $\varphi \notin S_2$, then $\operatorname{mg}(S_2) \leq \operatorname{mg}(S_1)$.

Submodularity is a good property for set optimization problem.

• OPT = max{cov(S)}

- Cannot afford to enumerate every size-k set of GFCs.
- cov(S) is a monotone submodular function.
- A greedy algorithm can have $(1 \frac{1}{\rho})$ approximation of OPT.

GFC_batch:

- 1. Mine all the patterns satisfying support and confidence. 2. $S = \emptyset$
- 3. While |S| < k, do
- 4. Select the pattern P with the largest marginal gain.

GFC_batch: mining in batch and selecting greedily

- GFC_batch is infeasible and slow.
 - Still, it requires mine all patterns first.
 - Can we do better?

• GFC_stream:

- Interleave pattern generation and rule selection.
- Find the top-*k* GFCs *on-the-fly*.
- One pass of pattern mining.
- $(\frac{1}{2} \epsilon)$ approximation of OPT

GFC_stream: mining and selecting on-the-fly!

Discovery Algorithms

PGen: pattern generation

- Generates patterns in a stream way.
- Pass the patterns for selection
- Can be in any order, e.g., Apriori, DFS, or random.
 pattern

PSel: pattern selection

- Selects and constructs GFCs on-the-fly.
- Based on a "sieve" strategy, $\left(\frac{1}{2} \epsilon\right)$ OPT Fast compute!
 - 1. Estimate the range of OPT by $max{cov(P)}$
 - 2. Each one is a size-k sieve with an estimation m for OPT.

PGen

PSel

stream

decision

- 3. While the sieves are not full
- 4. if $mg(P, S) \ge (\frac{m}{2} cov(S))/(k |S|)$, add P to sieve S.
- 5. Signal **PGen** to stop and output the sieve with largest cov.

GFC_stream: mining and selecting on-the-fly!

➢GFact_R: Using GFCs as rules:

- Invokes GFC_stream to find top-k GFCs.
- "Hit and miss"
 - True if a fact is covered by one GFC.
 - False If no GFC can cover the fact.
- A typical rule model to compare with: AMIE+

➢GFact: Using GFCs in supervised link prediction:

- A feature vector of size k.
- Each entry encodes the presence of one GFC.
- Build a classifier, by default, Logistic Regression.
- A typical rule models to compare with: PRA

Experiment settings

Dataset	category	V	E	# node labels	# edge labels	# < <i>x</i> , <i>r</i> , <i>y</i> >	
Yago	Knowledge base	2.1 M	4.0 M	2273	33	15.5 K	
DBpedia	Knowledge base	2.2 M	7.4 M	73	584	8240	
Wikidata	Knowledge base	10.8 M	41.4 M	18383	693	209 К	
MAG	Academic network	0.6 M	1.71 M	8665	6	11742	
Offshore	Social network	1.0 M	3.3 M	356	274	633	

Tasks	Rule Mining	Fact Checking				
Our methods	GFC_batch, GFC_stream	GFact, GFact _R				
Baselines	AMIE+, PRA	AMIE+, PRA, KGMiner				
Evaluation Metrics	running time vs. $ E $, $ \Gamma^+ $	prediction rate, precision, recall, F1				

Experiment: efficiency

Overview

- GFC_stream takes 25.7 seconds to discover 200 GFCs over Wikidata with 41.4 million edges and 6000 training facts.
- On average, GFC_stream is 3.2 times faster than AMIE+ over DBpedia.



Experiment: effectiveness

Compared with AMIE+, PRA and KGMiner, respectively, on average:

- GFact achieves additional 30%, 20%, and 5% gains of precision over DBpedia.
- GFact achieves additional 20%, 15%, and 16% gains of F1-score over Wikidata.



Case study: are two anonymous companies same? (Offshore)



Low accuracy.

If an officer is both a shareholder of company u_x and a beneficiary of company u_y, and u_x has an address and is registered through a jurisdiction in a place, and u_y is active in the same place, then they are likely to be the same anonymous company.

Conclusions and future work

- Graph Fact Checking Rules (GFCs)
- Top-k GFCs discovery problem Maximize a submodular cov function.
- A stream-based rule discovery algorithm
 - One pass, $\left(\frac{1}{2} \epsilon\right)$ OPT
- Evaluation of GFCs-based techniques
 - Rule models, fact checking (2 methods), efficiency, and case studies.
- Our future work: scalable GFC-based methods
 - Parallel mining, Distributed learning





Discovering Graph Patterns for Fact Checking in Knowledge Graphs

Thank you!

Related work: Gstream (IEEE BigData 2017)

Event Pattern Discovery by Keywords in Graph Streams Mohammad Hossein Namaki, Peng Lin, Yinghui Wu https://ieeexplore.ieee.org/abstract/document/8258019/



 Table 1. Effectiveness: average accuracy.

	YAGO				DBpedia			Wikidata				MAG				
Model	Pred	Prec	Rec	F_1	Pred	Prec	Rec	F_1	Pred	Prec	Rec	F_1	Pred	Prec	Rec	F_1
GFact	0.89	0.81	0.60	0.66	0.91	0.80	0.55	0.63	0.92	0.82	0.63	0.68	0.90	0.86	0.62	0.71
$GFact_R$	0.73	0.40	0.75	0.50	0.70	0.43	0.72	0.52	0.85	0.55	0.64	0.55	0.86	0.78	0.55	0.64
AMIE+	0.71	0.44	0.76	0.51	0.69	0.50	0.85	0.58	0.64	0.42	0.78	0.48	0.70	0.53	0.62	0.52
PRA	0.87	0.69	0.34	0.37	0.88	0.60	0.41	0.45	0.90	0.65	0.51	0.53	0.77	0.88	0.21	0.32
KGMiner	0.87	0.62	0.36	0.40	0.88	0.75	0.60	0.63	0.90	0.63	0.49	0.52	0.76	0.74	0.17	0.27

More effectiveness results



Time vs. k

Time vs. support

More efficiency results

More thorough experiments to compare various methods:

http://eecs.wsu.edu/~plin1/pdfs/2017-Preprint-Factchecking-experiments.pdf