Finding paths in large data graphs

Bruno Guillon Mikaela Keller Charles Paperman

UCA

LIMOS – UCA, Clermont-Ferrand

June 4, 2020

- nondeterministic two-way transducers and weighted automata
- rational transductions and beyond
- unary transducers
- descriptional complexity of automata

- extension of Monadic Second Order Logic on words
- regular transducers with origin semantics
- 2017-2018 Postdoc in Milan (with G.Pighizzini):
 - reversibility in automata and transducers
 - descriptional complexity of automata

- nondeterministic two-way transducers and weighted automata
- rational transductions and beyond
- unary transducers
- descriptional complexity of automata

- extension of Monadic Second Order Logic on words
- regular transducers with origin semantics
- 2017-2018 Postdoc in Milan (with G.Pighizzini):
 - reversibility in automata and transducers
 - descriptional complexity of automata

- nondeterministic two-way transducers and weighted automata
- rational transductions and beyond
- unary transducers
- descriptional complexity of automata

- extension of Monadic Second Order Logic on words
- regular transducers with origin semantics
- 2017-2018 Postdoc in Milan (with G.Pighizzini):
 - reversibility in automata and transducers
 - descriptional complexity of automata

- nondeterministic two-way transducers and weighted automata
- rational transductions and beyond
- unary transducers
- descriptional complexity of automata

- extension of Monadic Second Order Logic on words
- regular transducers with origin semantics
- 2017-2018 Postdoc in Milan (with G.Pighizzini):
 - reversibility in automata and transducers
 - descriptional complexity of automata

- nondeterministic two-way transducers and weighted automata
- rational transductions and beyond
- unary transducers
- descriptional complexity of automata

2016-2017 Postdoc in Warsaw (with M.Bojańczyk):

- extension of Monadic Second Order Logic on words
- regular transducers with origin semantics
- 2017-2018 Postdoc in Milan (with G.Pighizzini):
 - reversibility in automata and transducers
 - descriptional complexity of automata

November 2018-* Postdoc at INRIA Lille (with C.Paperman):

algorithms for Regular Path Queries on large graphs

- work in progress

- Graph databases
- ▶ Renewal with the noSQL trend, *e.g.*, neo4j (SPARQL),

back-end for the rdf and semantic web

- more flexible than relational databases
- easier to distribute than relational databases
 - notion of data locality, e.g., TitanDB, JanusGraph (gremlin)

- Graph databases
- ► Renewal with the noSQL trend, *e.g.*, neo4j (SPARQL),

back-end for the rdf and semantic web

- more flexible than relational databases
- easier to distribute than relational databases
 - notion of data locality, e.g., TitanDB, JanusGraph (gremlin)

▶ [Sun *et al.*'15 – SQLGraph: an efficient relational-based property graph]



Paths

Definition (Regular Path Queries (RPQ))

Given

• a property graph $G = \langle V, E \rangle$ with ppt : $V \mapsto \mathbb{D}$;

•
$$s, t \in V$$
;

• a regular expression e on \mathbb{D} :

find a path $\pi = s \cdot v_1 \cdots v_\ell \cdot t$ such that $\pi \in [|e|]$.

Paths

Definition (Regular Path Queries (RPQ))

Given

• a property graph $G = \langle V, E \rangle$ with ppt : $V \mapsto \mathbb{D}$;

•
$$s, t \in V$$
;

• a regular expression e on \mathbb{D} :

find a path $\pi = s \cdot v_1 \cdots v_\ell \cdot t$ such that $\pi \in [|e|]$.

Examples (on DBLP property graph)

find path from given source to given target visiting

- only papers published after 2016
- only journal papers/not arxiv papers
- only French authors
- only paper whose abstract does not contain the word "experimental"

Paths

Definition (Regular Path Queries (RPQ))

Given

• a property graph $G = \langle V, E \rangle$ with ppt : $V \mapsto \mathbb{D}$;

•
$$s, t \in V$$
;

• a regular expression e on \mathbb{D} :

find a path $\pi = s \cdot v_1 \cdots v_\ell \cdot t$ such that $\pi \in [|e|]$.

Examples (on DBLP property graph)

find path from given source to given target visiting

- only papers published after 2016
- only journal papers/not arxiv papers
- only French authors
- only paper whose abstract does not contain the word "experimental"
- papers in chronological order along the path



from given source to given target

How to find regular paths in large property graphs efficiently? $\bigcup_{\substack{G = \langle V, E \rangle \\ \text{ppt} : V \mapsto D}}^{G = \langle V, E \rangle}$







goal: find path from given source to given target by visiting few nodes



goal: find path from given source to given target by visiting few nodes **allowed:**

computation time when operating in RAM is "free"



goal: find path from given source to given target by visiting few nodes **allowed:**

- computation time when operating in RAM is "free"
- ▶ precomputed information (on disk): $O(n \cdot \text{polylog}(n))$ space

Outline





Outline

1. Search paths

- Bilateral Best First Search
- Tradeoff between short path and efficiency
- 2. Learning graph embeddings
- 3. Experimental work
 - What to measure
 - Results

4. Future work

Given an graph exploration algorithm

G

Given an graph exploration algorithm we extend it to a graph regular exploration algorithm

G

Given an graph exploration algorithm we extend it to a graph regular exploration algorithm,

by exploring the direct product of the graph and an automaton

 $G \times A$

Given an graph exploration algorithm we extend it to a graph regular exploration algorithm,

by exploring the direct product of the graph and an automaton

 $\underline{G \times A}$ not materialized

▶ Breadth-First-Search: very bad (social networks have small diameter)

- ► Breadth-First-Search: very bad (social networks have small diameter)
- ► Bilateral BFS: far better

- ► Breadth-First-Search: very bad (social networks have small diameter)
- Bilateral BFS: far better
- ▶ Best First Search algorithms: use heuristic information to guide the search

- ► Breadth-First-Search: very bad (social networks have small diameter)
- Bilateral BFS: far better
- Best First Search algorithms: use heuristic information to guide the search
 A*: if admissible heuristics (shortest path found)

- ► Breadth-First-Search: very bad (social networks have small diameter)
- Bilateral BFS: far better
- Best First Search algorithms: use heuristic information to guide the search
 A*: if admissible heuristics (shortest path found)
 - ▶ otherwise, a path may be found, not necessarily a shortest one



- $\mu : \{ \mathsf{nodes} \} \mapsto \mathbb{R}^d$
- distance: $||\mu(v) \mu(u)||$



- μ : {nodes} $\mapsto \mathbb{R}^d$
- distance: $||\mu(v) \mu(u)||$

What is a good embedding?



- μ : {nodes} $\mapsto \mathbb{R}^d$
- distance: $||\mu(v) \mu(u)||$

What is a good embedding? not clear.



- μ : {nodes} $\mapsto \mathbb{R}^d$
- distance: $||\mu(v) \mu(u)||$

What is a good embedding? not clear.

Intuitively, we want that being close in the embedding means being close in the graph topology

We assume an embedding $\mu : \{ \text{nodes} \} \mapsto \mathbb{R}^d$.

```
We assume an embedding \mu : {nodes} \mapsto \mathbb{R}^d.
```

```
Algorithme : Best-First-Search
input : graph G, embedding \mu, source s, target t
output : a path from s to t in G
Visited \leftarrow \{s\};
while t not in Visited do
   v \leftarrow select best node in adherence of Visited ;
   add v to Visited :
return a path from s to t in the subgraph induced by Visited;
```

```
We assume an embedding \mu : \{ \text{nodes} \} \mapsto \mathbb{R}^d.
```

```
Algorithme : Best-First-Search
input : graph G, embedding \mu, source s, target t
output : a path from s to t in G
Visited \leftarrow \{s\};
while t not in Visited do
   \rightarrow minimizing ||\mu(t) - \mu(v)||
   add v to Visited ;
return a path from s to t in the subgraph induced by Visited;
```

```
We assume an embedding \mu : \{ \text{nodes} \} \mapsto \mathbb{R}^d.
```

```
Algorithme : Best-First-Search
input : graph G, embedding \mu, source s, target t, parameter \alpha
output : a path from s to t in G
Visited \leftarrow \{s\};
while t not in Visited do
    v \leftarrow \texttt{select} \ best \ node in adherence of Visited ;
                                 \rightarrow minimizing ||\mu(t) - \mu(v)|| \cdot \text{depth}(v)^{\alpha}
    add v to Visited :
return a path from s to t in the subgraph induced by Visited;
```

depth(v): the shortest path length from s to v in the explored subgraph

```
We assume an embedding \mu : \{ \text{nodes} \} \mapsto \mathbb{R}^d.
```

```
Algorithme : Best-First-Search
input : graph G, embedding \mu, source s, target t, parameter \alpha
output : a path from s to t in G
Visited \leftarrow \{s\};
while t not in Visited do
    v \leftarrow select best node in adherence of Visited ;
                                \rightarrow minimizing ||\mu(t) - \mu(v)|| \cdot \text{depth}(v)^{\alpha}
    add v to Visited :
return a path from s to t in the subgraph induced by Visited;
```

depth(v): the shortest path length from s to v in the explored subgraph
 when α = 0: deep search - when α → ∞: breadth first search

```
We assume an embedding \mu : \{\text{nodes}\} \mapsto \mathbb{R}^d.
```

```
Algorithme : Best-First-Search
input : graph G, embedding \mu, source s, target t, parameter \alpha
output : a path from s to t in G
Visited \leftarrow \{s\};
while t not in Visited do
   v \leftarrow select best node in adherence of Visited ;
                               \rightarrow minimizing ||\mu(t) - \mu(v)|| \cdot depth(v)^{\alpha}
    maintain a layered DAG structure (depth);
    add v to Visited :
return a path from s to t extracted from the layered DAG;
```

depth(v): the shortest path length from s to v in the explored subgraph
 when α = 0: deep search - when α → ∞: breadth first search

Outline

1. Search paths

- Bilateral Best First Search
- Tradeoff between short path and efficiency

2. Learning graph embeddings

3. Experimental work

- What to measure
- Results

4. Future work





find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix





find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix



SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix

Random walks + neural networks: use learning techniques of natural language processing



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 Uniform random walks: deepWalk (2011)



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 - ► Uniform random walks: deepWalk (2011)
 - ► Biased random walks: node2vec (2017)



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 Uniform and deep walkey deep Walk (2011)
 - Uniform random walks: deepWalk (2011)
 Biased random walks: node2vec (2017)
- Landmarks select a small set of nodes, called *landmarks* compute distances from any node to any landmarks



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 Uniform random walks: deepWalk (2011)
 - Biased random walks: node2vec (2017)
- Landmarks select a small set of nodes, called *landmarks* compute distances from any node to any landmarks
 - ▶ already used to guide path searching on road networks (2007)



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 - ► Uniform random walks: deepWalk (2011)
 - Biased random walks: node2vec (2017)
- Landmarks select a small set of nodes, called *landmarks* compute distances from any node to any landmarks
 - ▶ already used to guide path searching on road networks (2007)
 - ▶ and to estimate node distance in social networks (2009)



- SVD: find eigenvectors of maximum eigenvalues of the inverse Laplacian matrix
- Random walks + neural networks: use learning techniques of natural language processing
 Italian and a star Walk (2011)
 - ► Uniform random walks: deepWalk (2011)
 - Biased random walks: node2vec (2017)
- Landmarks select a small set of nodes, called *landmarks* compute distances from any node to any landmarks
 - ▶ already used to guide path searching on road networks (2007)
 - ▶ and to estimate node distance in social networks (2009)
- many other embeddings to test...

What to measure?

efficiency: how many nodes have been visited with respect to the standard bilateral BFS?

quality: are the found paths much longer than the shortest paths?

Experimental results

method	SVD	node2vec	deepWalk	BFS	
--------	-----	----------	----------	-----	--

Powerlaw cluster graph (#nodes: 2000, #edges: 7979)

mean visited	65.06	76.36	64.13	86.76
score	74.99%	88.01%	73.92%	
mean error	0.03%	0.05%	0.01%	

Regular graph (#nodes: 2000, #edges: 7000)

mean visited	53.55	53.02	56.16	98.43
score	54.41%	53.86%	57.05%	
mean error	0.02%	0.03%	0.0%	

DBLP Graph (#nodes: 5 745K, #edges: 12 205K)

mean visited		1652.14	11632.02
score		14.2%	
mean error		0.26%	
median error		0.22%	



► improve embeddings



- ► improve embeddings
- maintain embeddings under updates
 - force algorithms, continuous learning, landmarks



- improve embeddings
- maintain embeddings under updates
 - force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths



- improve embeddings
- maintain embeddings under updates

 force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths
- subgraph extraction using the embedding convex closure



- improve embeddings
- maintain embeddings under updates

 force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths
- subgraph extraction using the embedding convex closure
- index for searching nodes close to a point in the topology

 quadtree, octrees, k-d-tree (good for small dimensions)
 space-filling curves



- improve embeddings
- maintain embeddings under updates

 force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths
- subgraph extraction using the embedding convex closure
- index for searching nodes close to a point in the topology

 quadtree, octrees, k-d-tree (good for small dimensions)
 space-filling curves





- improve embeddings
- maintain embeddings under updates

 force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths
- subgraph extraction using the embedding convex closure
- index for searching nodes close to a point in the topology

 quadtree, octrees, k-d-tree (good for small dimensions)
 space-filling curves





- improve embeddings
- maintain embeddings under updates

 force algorithms, continuous learning, landmarks
- from filtered paths to data regular paths
- subgraph extraction using the embedding convex closure
- index for searching nodes close to a point in the topology

 quadtree, octrees, k-d-tree (good for small dimensions)
 space-filling curves
- graph clustering

Thank you for your attention!

More on DBLP graph

- ► sources: https://dblp.uni-trier.de/, January 16th, 2019 (2.4Gb xml file)
- ▶ graph: #nodes = 5746803 (2152092 authors) + json for each node, #edges = 12205333, (199Mb for edges (txt); 1.2Gb for nodes + jsons)
- connected comp.: 57 186 one large (5 500 304) and others small (< 101)
- ▶ degrees: avg=5.24769..., med=3, stdev=10, max=1626,



- distances (main cc, evaluated with sample > 110 000 pairs of nodes): avg=10.9, med=11, stdev=2, max=27
- alternative encoding "nodes are authors, edges are journals" yields much more edges (19972747)