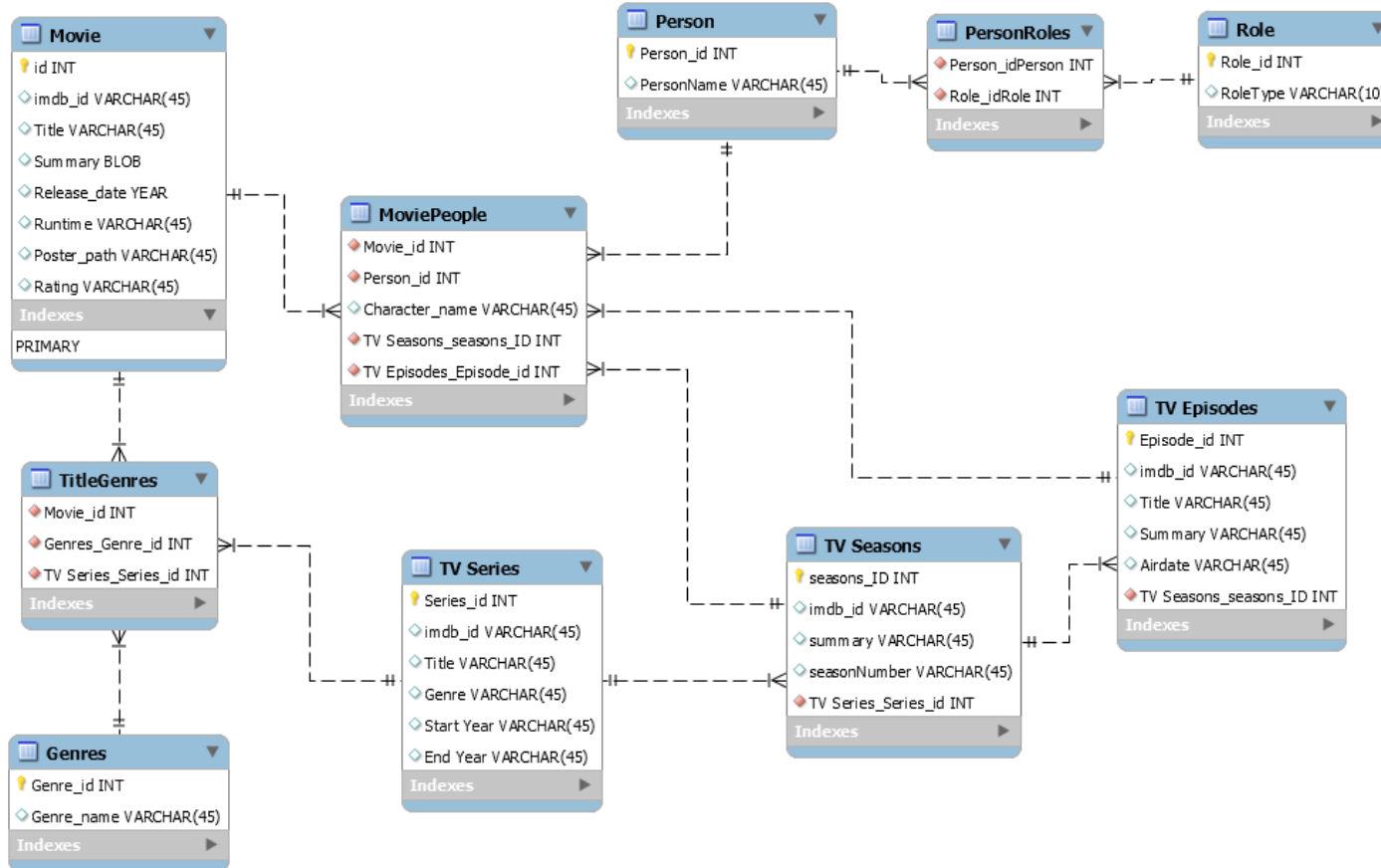


Evaluating Graph Embeddings and Graph Similarity Measurements for Database Systems

Josephine Rehak

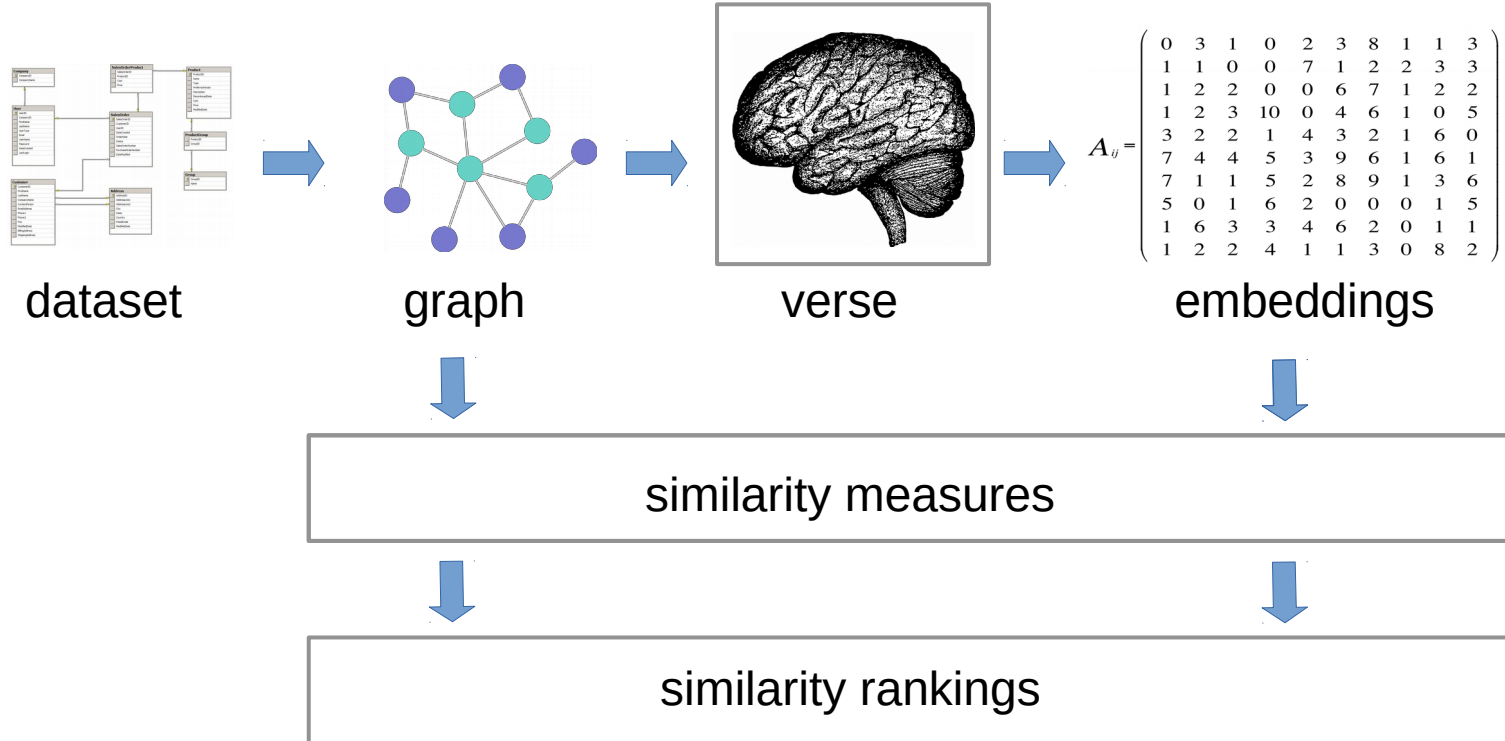
Profilprojekt Anwendungsforschung in der Informatik

Motivation



Motivation I

Overview



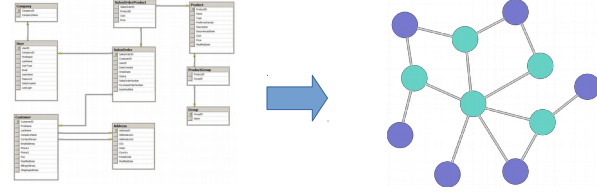


Fundamentals

Fundamentals I

Creating a directed graph from database

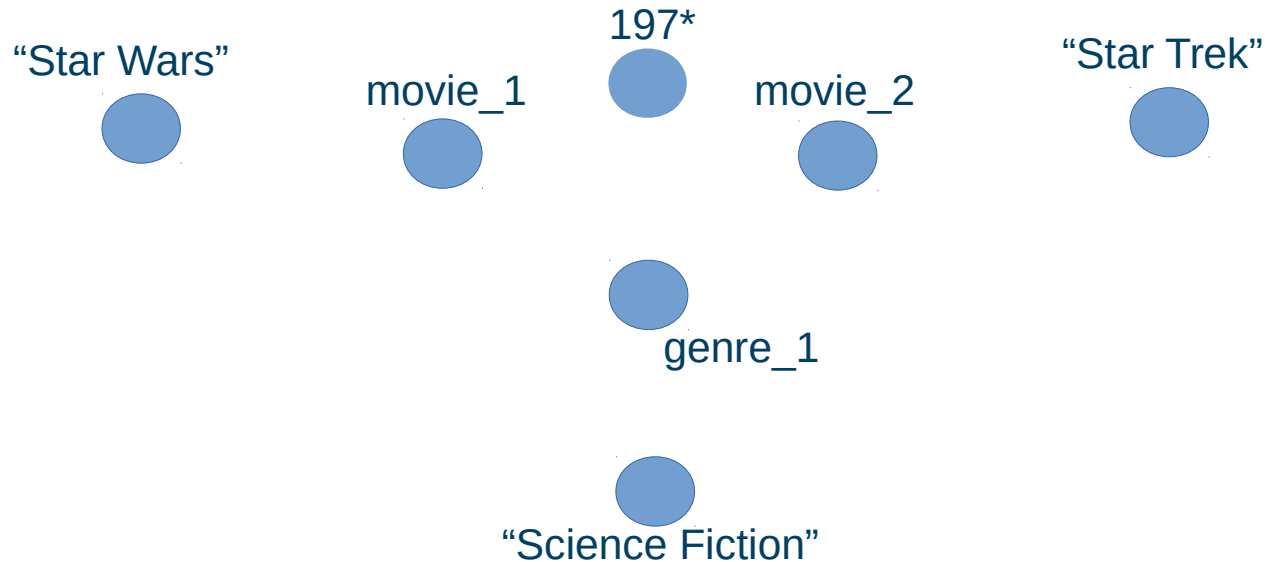
- Vertices: IDs + Table entries
 - Except: foreign keys, mn-tables



Fundamentals I



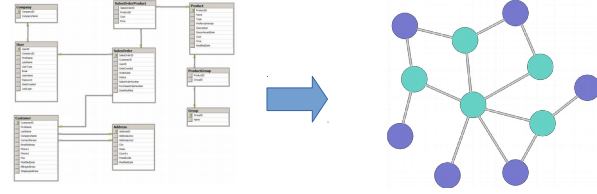
id	Name	Year	mid	gid	id	Genre
1	Star Wars	1977	1	1	1	Science Fiction
2	Star Trek	1979	2	1		



Fundamentals I

Creating a directed graph from database

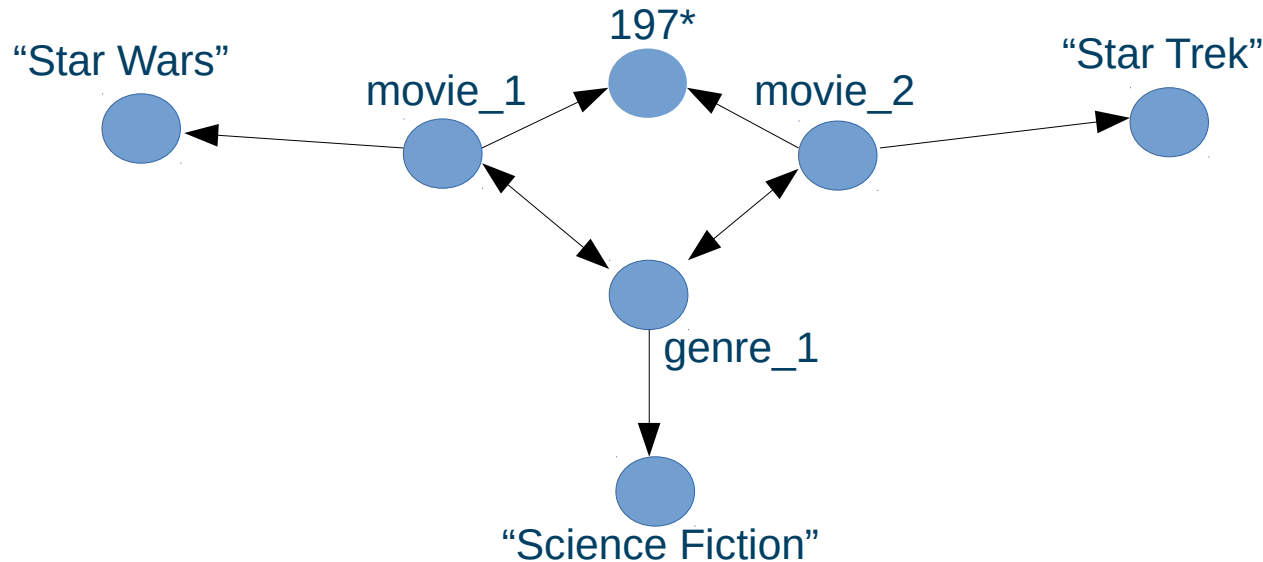
- Vertices: IDs + Table entries
 - Except: foreign keys, mn-tables
- Edges: Inner + Inter Table Relations
 - Directed edges: id to table values + Foreign keys
 - MN-Table-Relations to bidirected edges



Fundamentals I



id	Name	Year	mid	gid	id	Genre
1	Star Wars	1977	1	1	1	Science Fiction
2	Star Trek	1979	2	1		



Similarity Measures in Verse

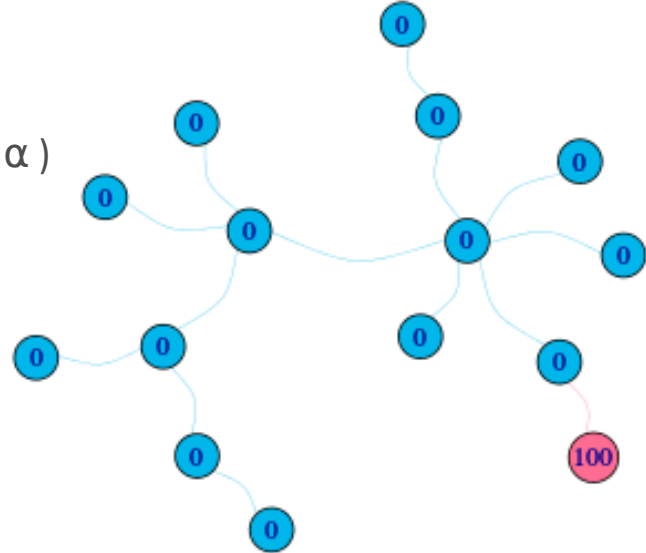
- Adjacency Similarity:
 - Similarity to neighbours only
 - Low complexity

$$sim_G^{ADJ}(u, v) = \begin{cases} 1/Out(u) & \text{if } (u, v) \in E \\ 0 & \text{otherwise} \end{cases}$$

Similarity Measures in Verse

- Personalized Page Rank:
 - Random walker with jump back probability (α)
 - Initial assignment, then recursive
 - Ranks converge to:

$$\pi_s = \alpha s + (1 - \alpha)\pi_s A$$



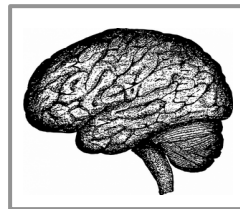
[3]

Similarity Measures in Verse

- SimRank:
 - Objects are similar when referenced by similar objects
 - Initial assignment, then recursive
 - importance of farther nodes : C
 - complexity of $O(n^4)$

$$\text{sim}_G^{\text{SR}}(u, v) = \frac{C}{|I(u)| |I(v)|} \sum_{i=1}^{|I(u)|} \sum_{j=1}^{|I(v)|} \text{sim}_G^{\text{SR}}(I_i(u), I_j(v))$$

Verse (*VER*tex Similarity Embeddings)



- Trains unsupervised 1-layer-NNs
- **Scalable:** Processes 10^6 nodes in less than a day
- **Global:** similarity on any pair of graph nodes
- **Versatile:** supports any similarity measure between nodes

$$\sum_{v \in V} \text{KL}(\text{sim}_G(v, \cdot) \parallel \text{sim}_E(v, \cdot))$$

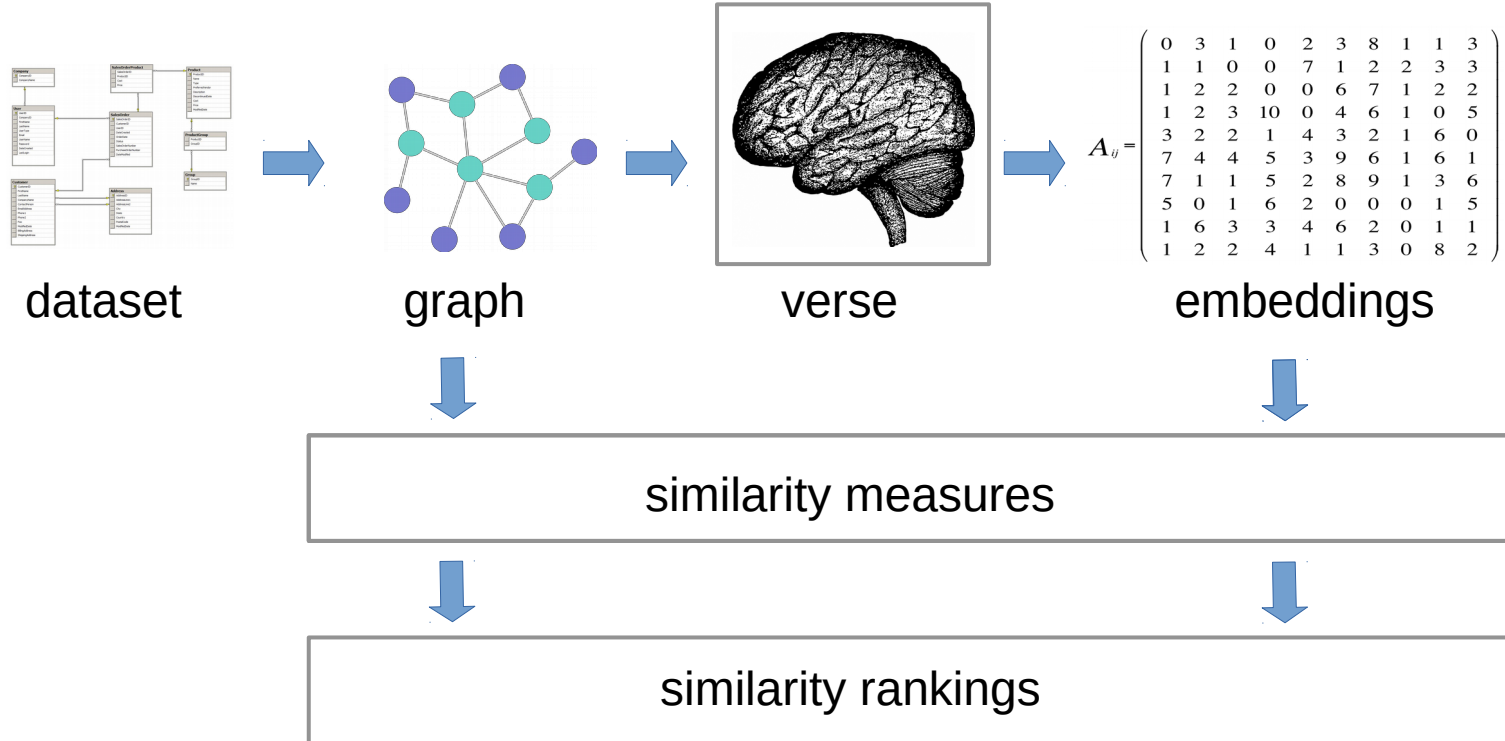
Verse Vertex Embeddings

- Mapping of graph vertex to vector space
 - Maintains: neighborhood, connections, graph topology
 - Efficiency for large graphs
 - individual dimensions have no meaning
- Usecases:
 - ML
 - similarity measures
 - NearestNeighbour
 - ...

$$A_{ij} = \begin{pmatrix} 0 & 3 & 1 & 0 & 2 & 3 & 8 & 1 & 1 & 3 \\ 1 & 1 & 0 & 0 & 7 & 1 & 2 & 2 & 3 & 3 \\ 1 & 2 & 2 & 0 & 0 & 6 & 7 & 1 & 2 & 2 \\ 1 & 2 & 3 & 10 & 0 & 4 & 6 & 1 & 0 & 5 \\ 3 & 2 & 2 & 1 & 4 & 3 & 2 & 1 & 6 & 0 \\ 7 & 4 & 4 & 5 & 3 & 9 & 6 & 1 & 6 & 1 \\ 7 & 1 & 1 & 5 & 2 & 8 & 9 & 1 & 3 & 6 \\ 5 & 0 & 1 & 6 & 2 & 0 & 0 & 0 & 1 & 5 \\ 1 & 6 & 3 & 3 & 4 & 6 & 2 & 0 & 1 & 1 \\ 1 & 2 & 2 & 4 & 1 & 1 & 3 & 0 & 8 & 2 \end{pmatrix}$$

Motivation I

Overview



Spearman coefficient



- Example, $n=3$

Movie id	r 1	r 2	d_i	d_i^2
1	2	1	1	1
2	1	3	-2	4
3	3	2	1	1

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}.$$

$$\rho = -0.5$$

- Values between -1 and 1
 - -1 inverted correspondence
 - 0 no correspondence
 - 1 strong correspondence



Experimental Setup

Experimental Setup

Data Set



The screenshot shows the Kaggle dataset page for 'The Movies Dataset'. The header includes a 'Dataset' icon and the title 'The Movies Dataset'. Below the title, it states 'Metadata on over 45,000 movies. 26 million ratings from over 270,000 users.' The creator is listed as 'Rounak Banik' with a note 'updated a year ago (Version 7)'. The page has tabs for 'Data', 'Kernels (74)', 'Discussion (9)', 'Activity', and 'Metadata'. On the right, there is a 'Download (228 MB)' link and a 'New Kernel' button. A small box in the top right corner of the image shows an upvote icon and the number '836'.

Dataset

The Movies Dataset

Metadata on over 45,000 movies. 26 million ratings from over 270,000 users.

Rounak Banik • updated a year ago (Version 7)

Data Kernels (74) Discussion (9) Activity Metadata

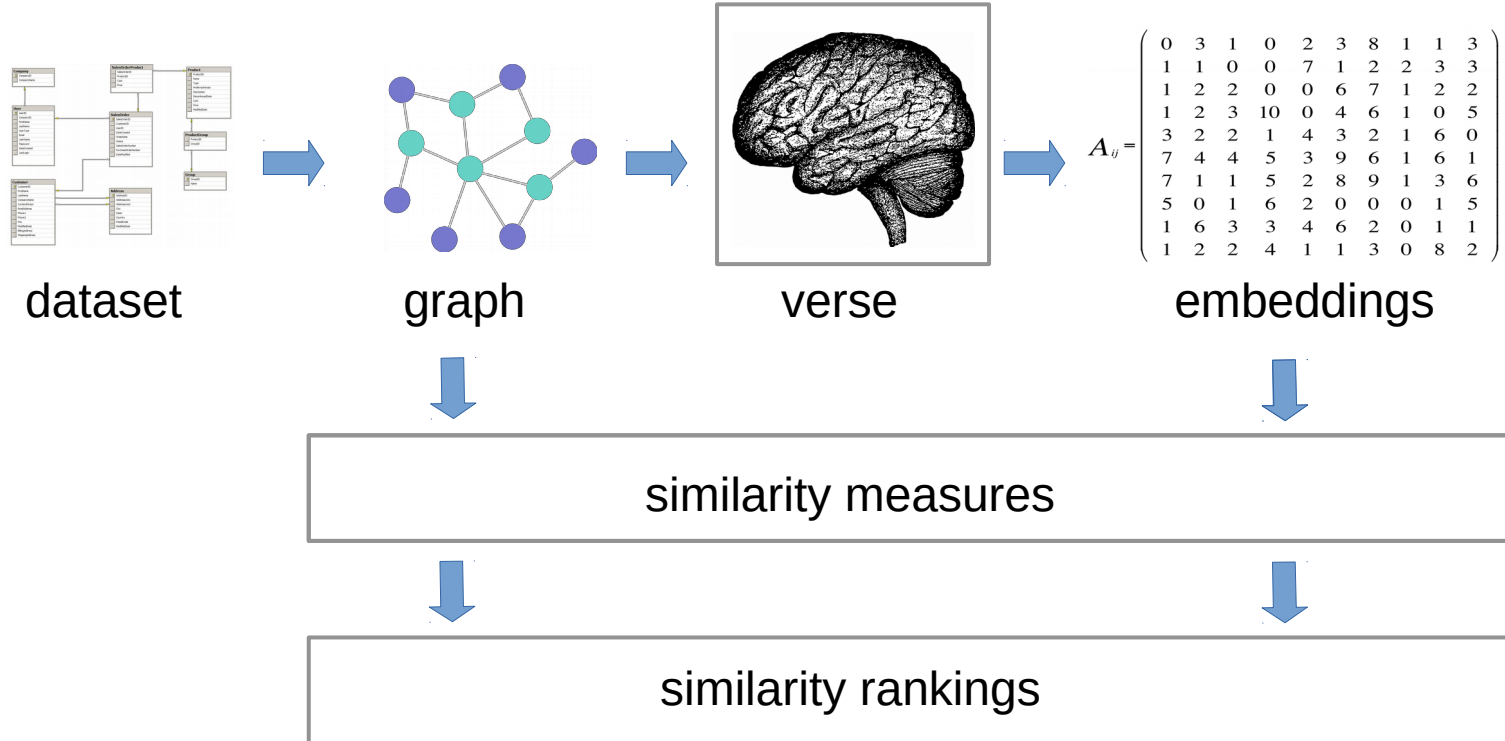
Download (228 MB) New Kernel

Used excerpt has: 45433 movies, 18001 directors, 20 genres

<https://www.kaggle.com/rounakbanik/the-movies-dataset/home>

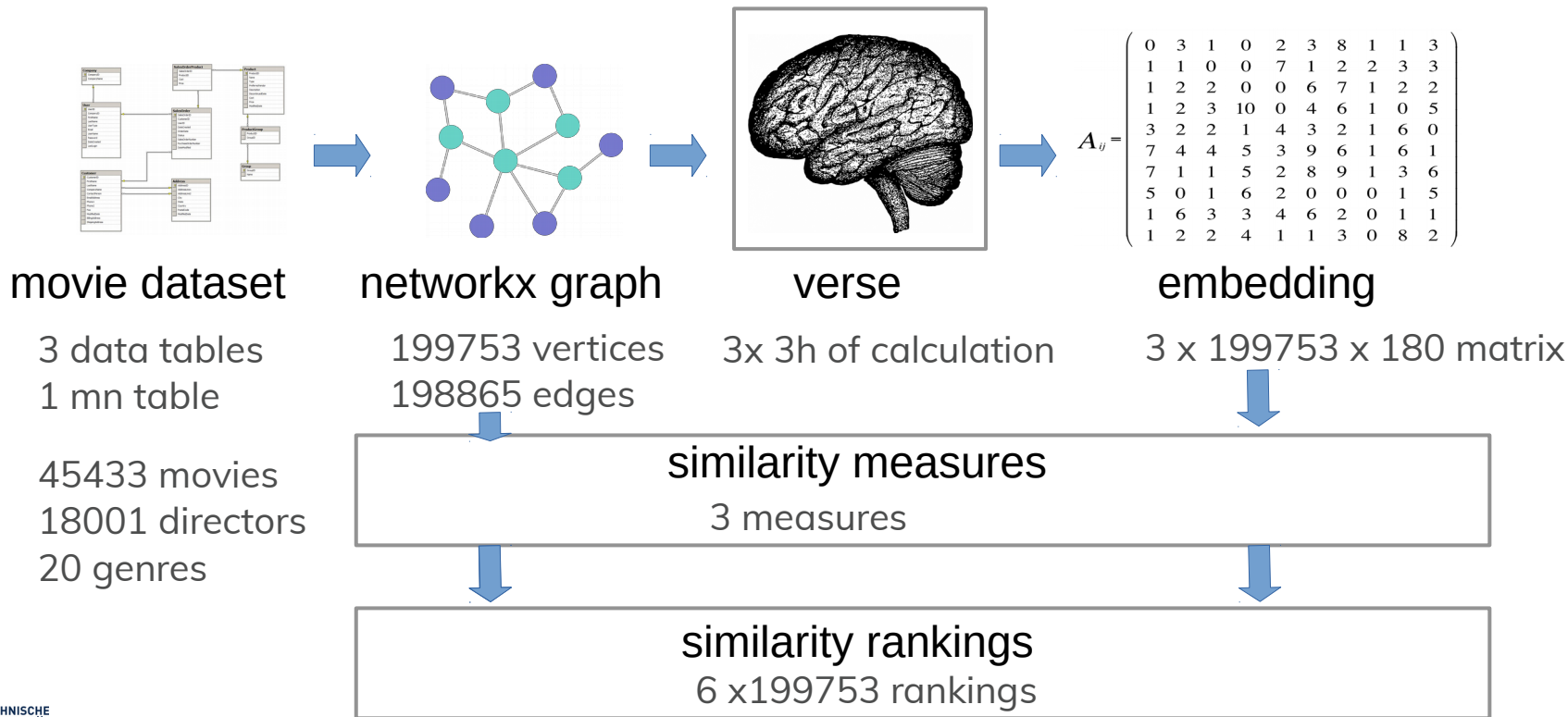
Motivation I

Overview



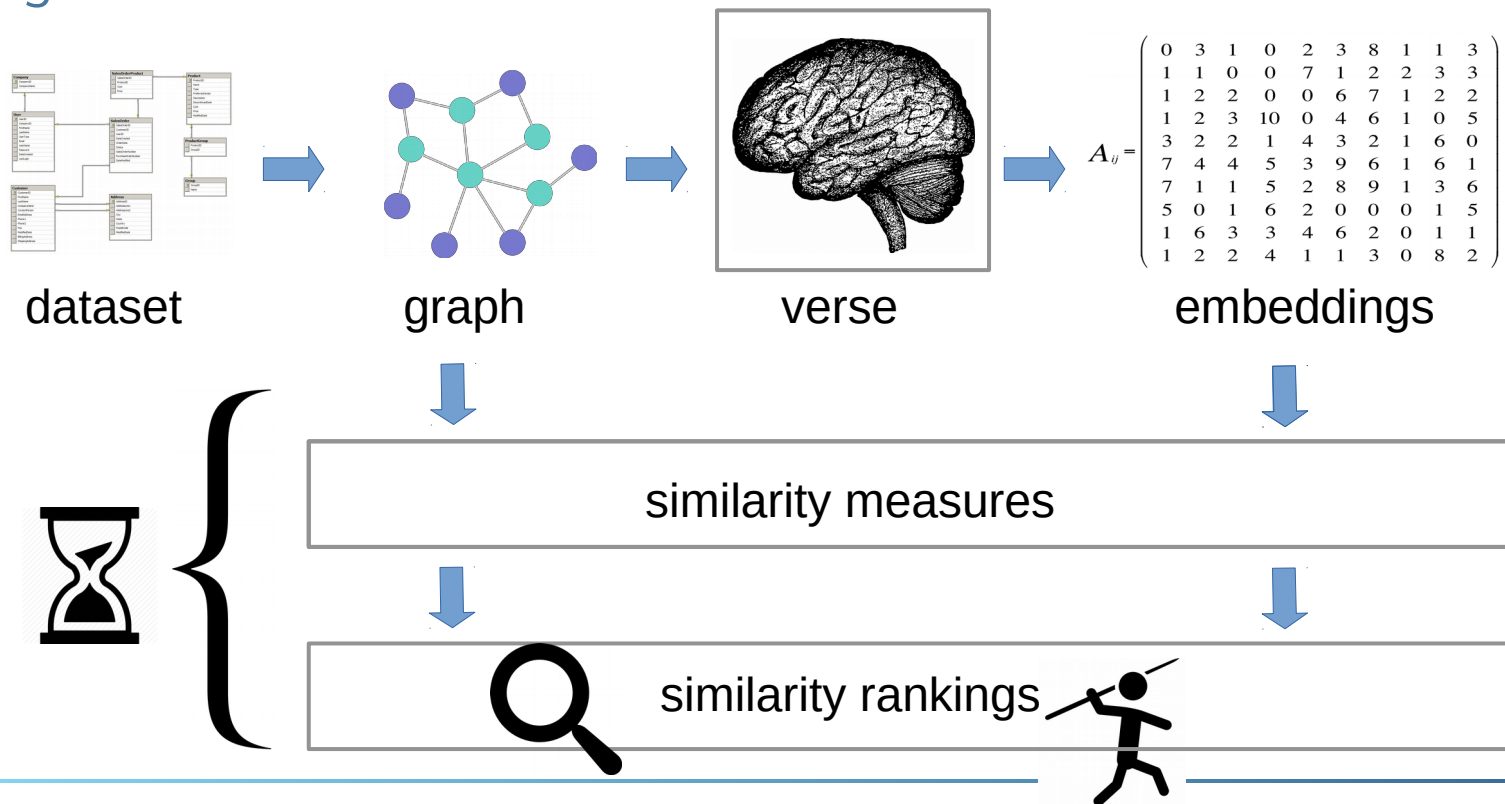
Experimental Setup

Overview



Motivation I

Investigation





Results

Results I

PPageRank Graph

	Movie Names
1	Star_Wars
2	Nana_Neul
3	Aki_Kaurismäki
4	The_Circle
5	Cross_Creek
6	Jean_Mach

AdjacencySimilarity Graph

	Movie Names
1	Star_Wars
2	Hellboy
3	Get_Carter
4	Der_Totmacher
5	Begotten
6	Laws_of_Gravity

PageRank Embeddings

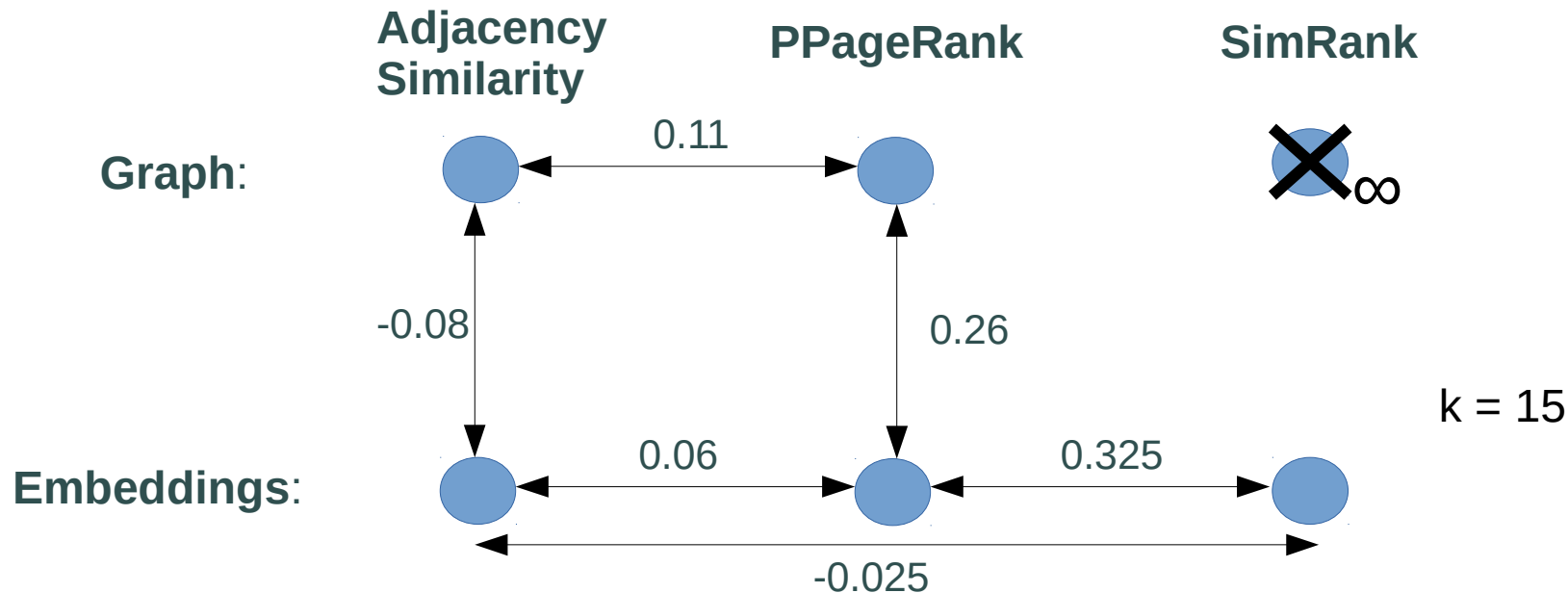
	Movie Names
1	Star_Wars
2	Soldier's Girl
3	Незнайка_на_Луне
4	Lignes_de_Front
5	Best_Night_Ever
6	London_Boulevard

SimRank Embeddings

	Movie Names
1	Star_Wars
2	Roswell
3	THE_DARK_SIDE_OF_DIMENSIONS
4	Perfect
5	Love_and_Distrust
6	National_Lampoon's_Last_Resort

Results I

Spearman on Similarity Rankings



Results I

Computation Time

**Adjacency
Similarity**

PPageRank

SimRank

Graph:

2.77s

12.82s

∞

Embeddings:

0.32s

0.36s

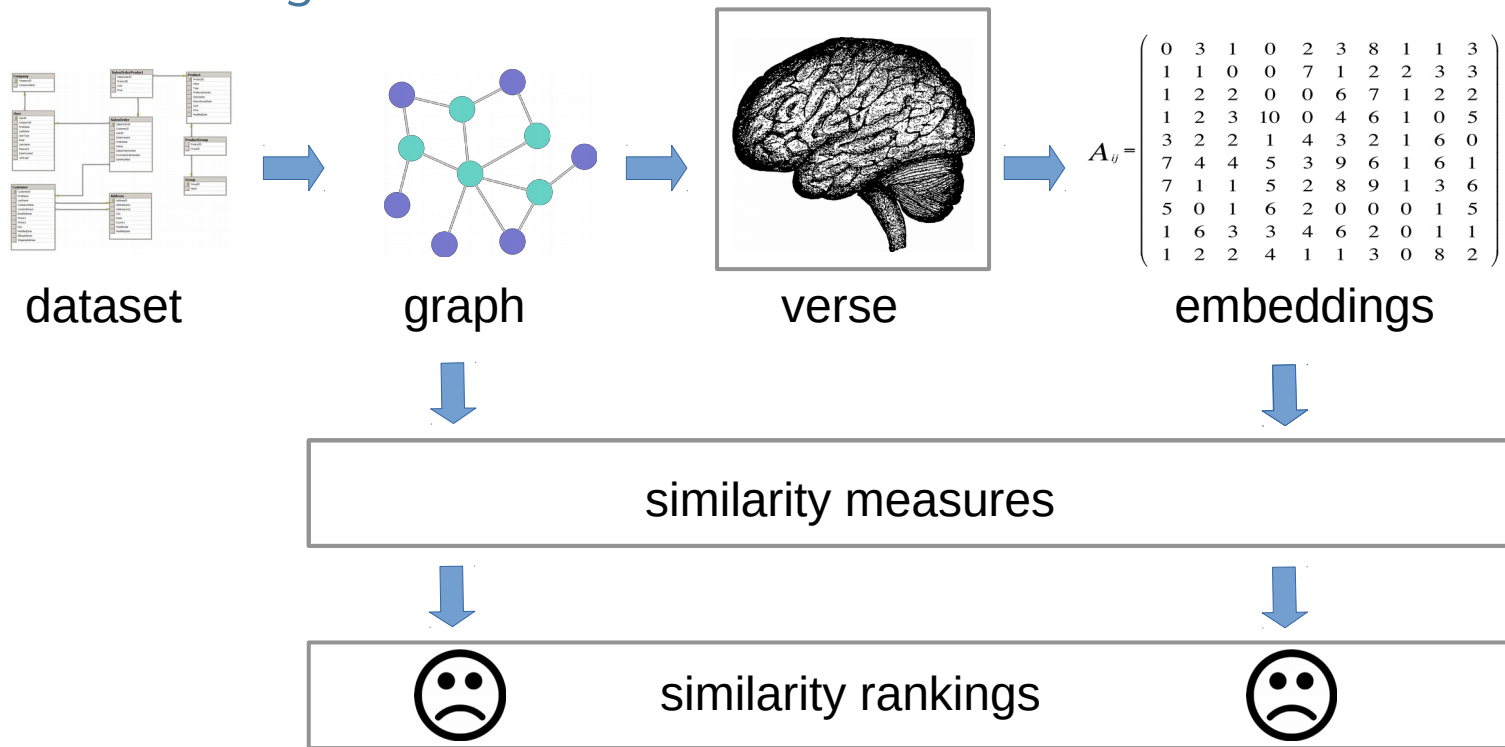
0.21s



Discussion

Discussion

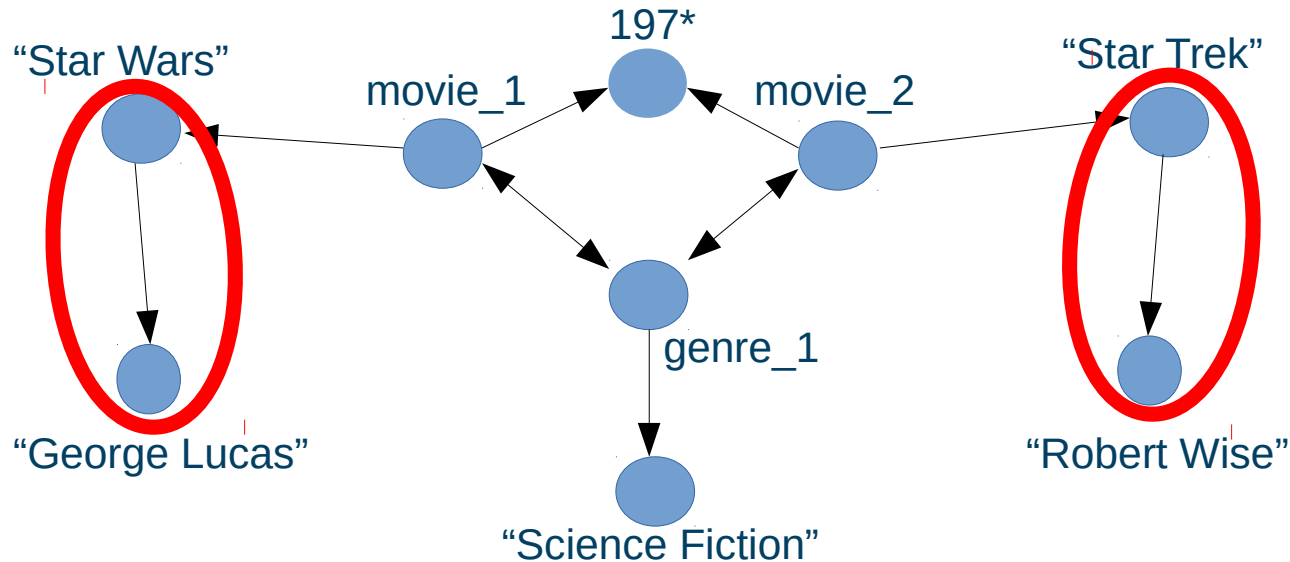
What went wrong?



Discussion



id	Name	Year	d_id		id	Director
1	Star Wars	1977	1	<u>M</u> 1	1	George Lucas
2	Star Trek	1979	2		2	Robert Wise





Conclusion

Conclusion I



- Translates across tasks/databases
- Speeding up similarity calculation
- Supports many similarity measures
- ‚Simple‘ solution for complex db problem



- Verse needs preparation time
- Good parameters / similarity measure required
- Dependant of graph structure



Fragen?