

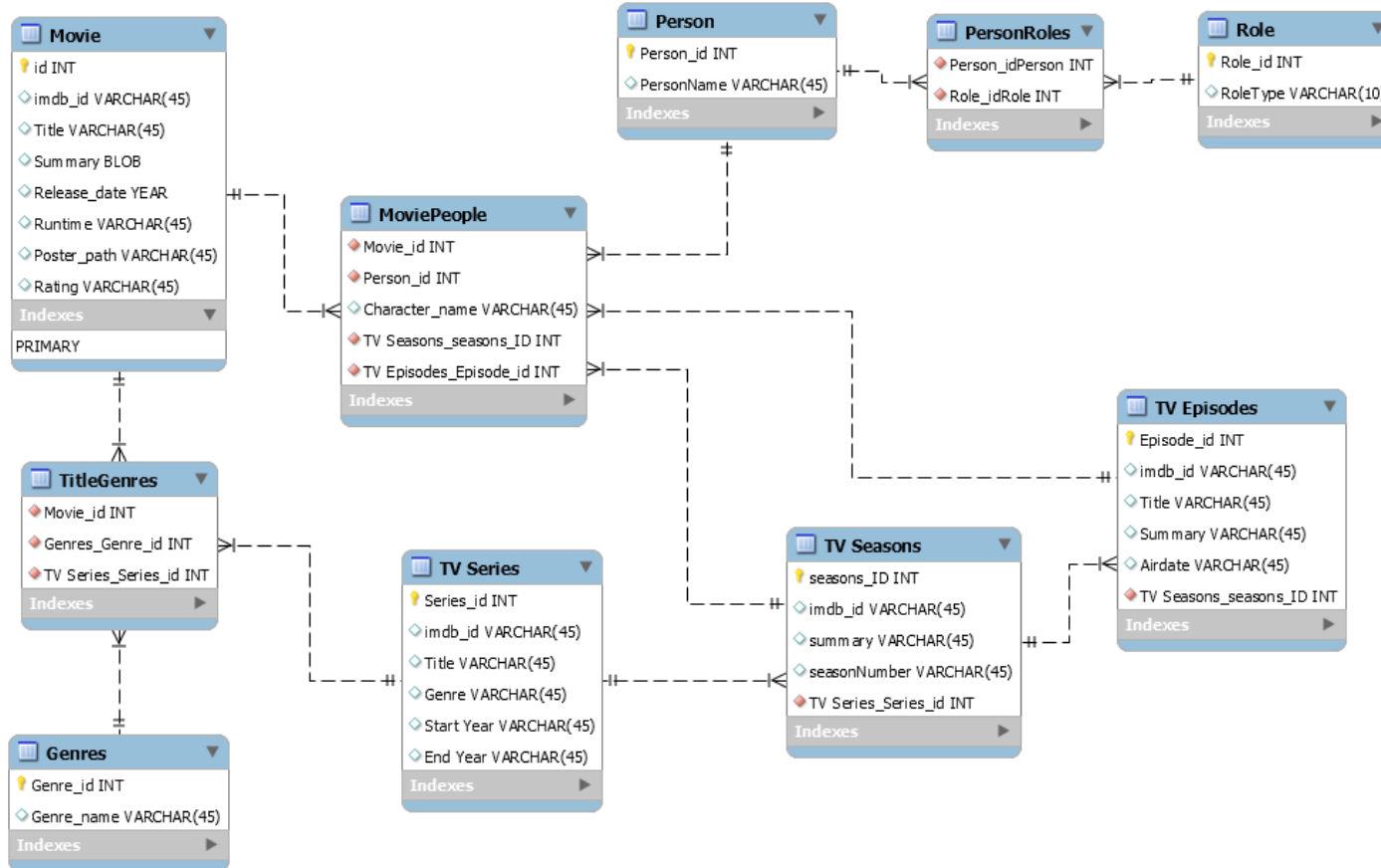


Evaluating Graph Embeddings and Graph Similarity Measurements for Database Systems

Josephine Rehak

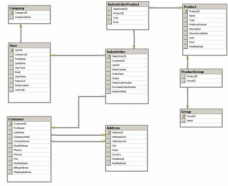
Profilprojekt Anwendungsforschung in der Informatik

Agenda

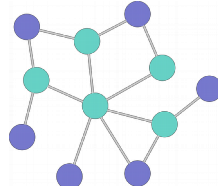


Agenda

Overview



dataset



graph



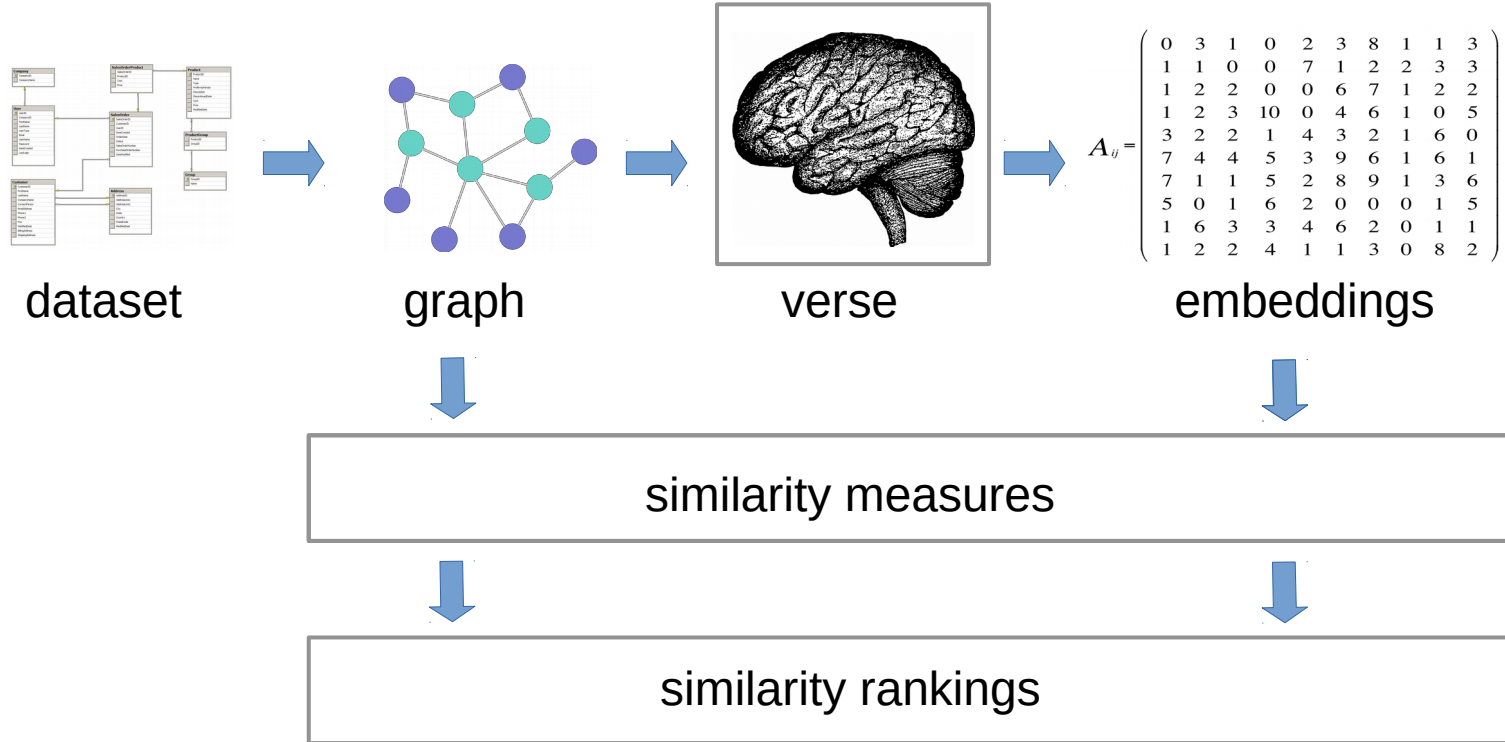
similarity measures



similarity rankings

Agenda

Overview



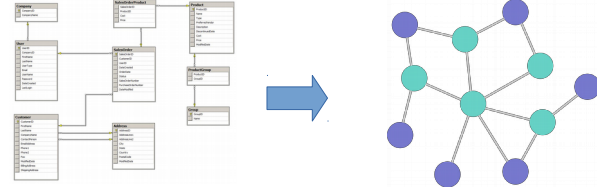


Fundamentals

Fundamentals I

Creating a directed graph from database

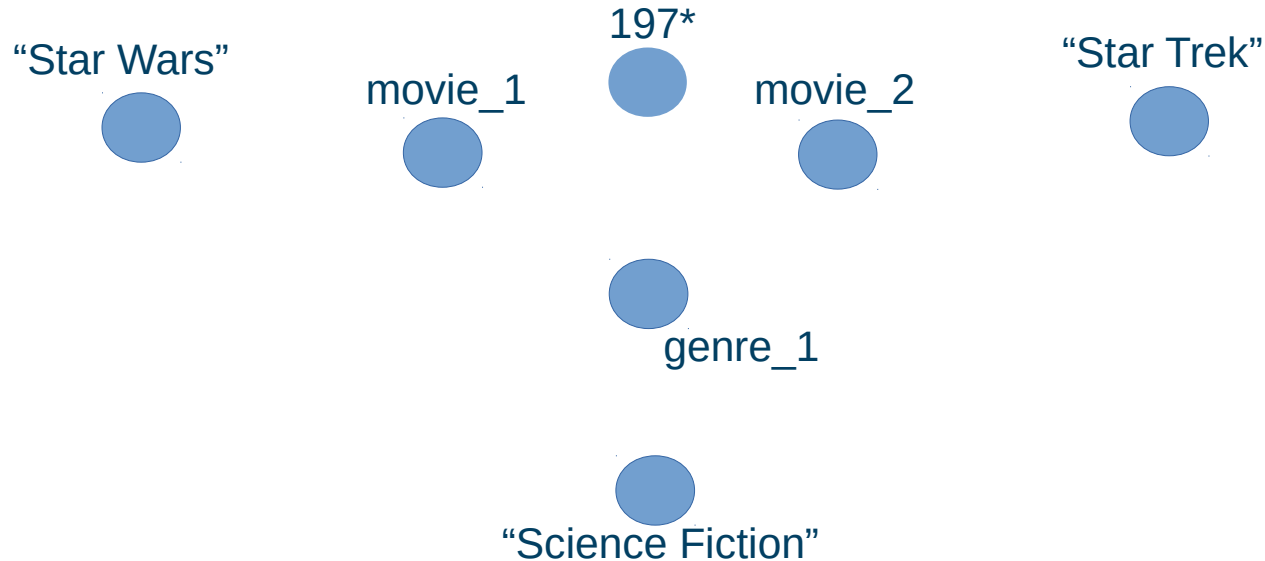
- Vertices: IDs + Table entries
 - Except: foreign keys, mn-tables



Fundamentals II



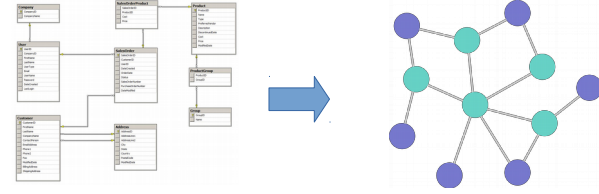
id	Name	Year	mid	gid	id	Genre
1	Star Wars	1977	1	1	1	Science Fiction
2	Star Trek	1979	2	1		



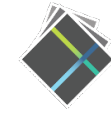
Fundamentals III

Creating a directed graph from database

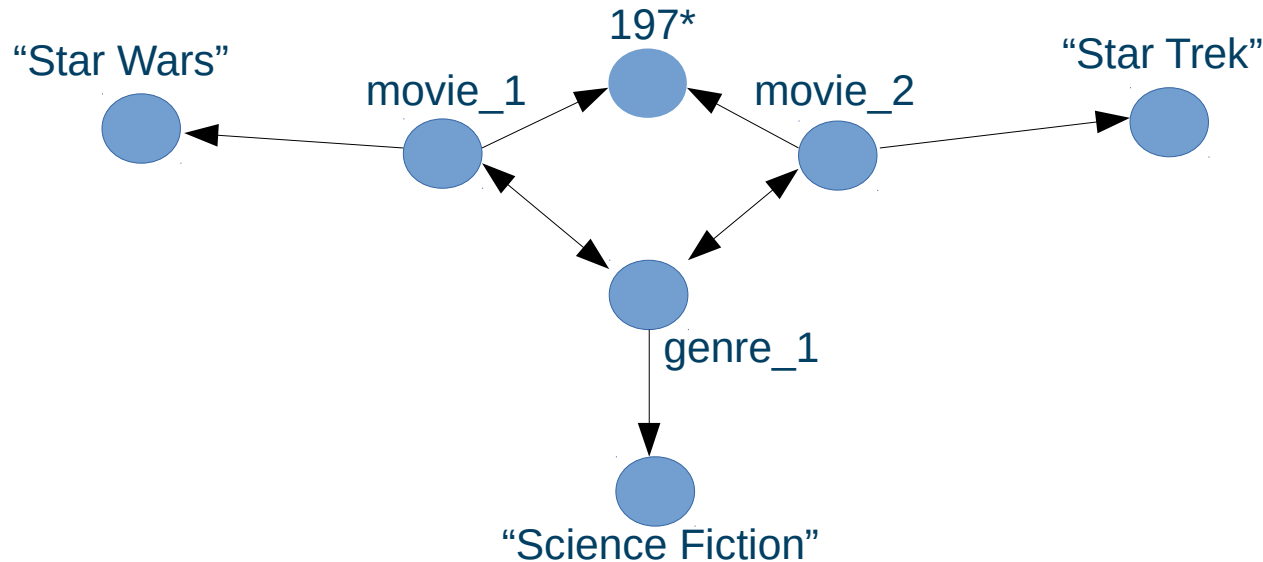
- Vertices: IDs + Table entries
 - Except: foreign keys, mn-tables
- Edges: Inner + Inter Table Relations
 - directed edges: id to table values
 - Bidirected edges: MN-Table-Relations + Foreign Keys



Fundamentals IV

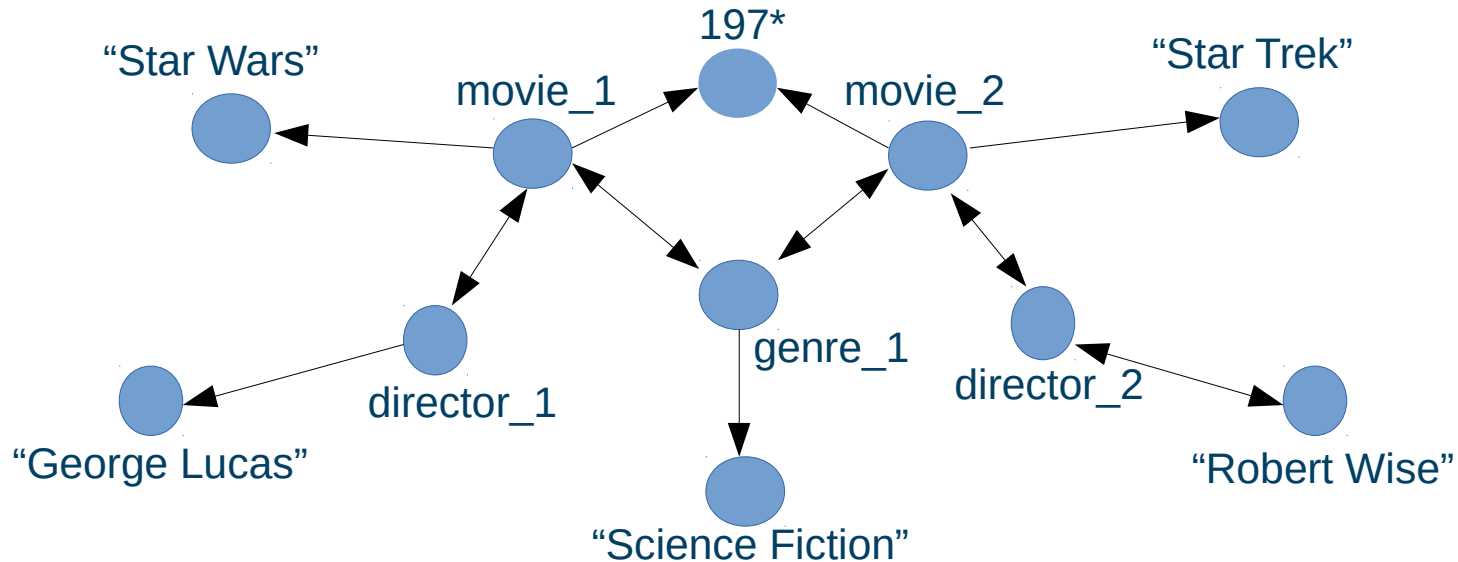


id	Name	Year	mid	gid	id	Genre
1	Star Wars	1977	1	1	1	Science Fiction
2	Star Trek	1979	2	1		



Fundamentals V

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



Similarity Measures in Verse

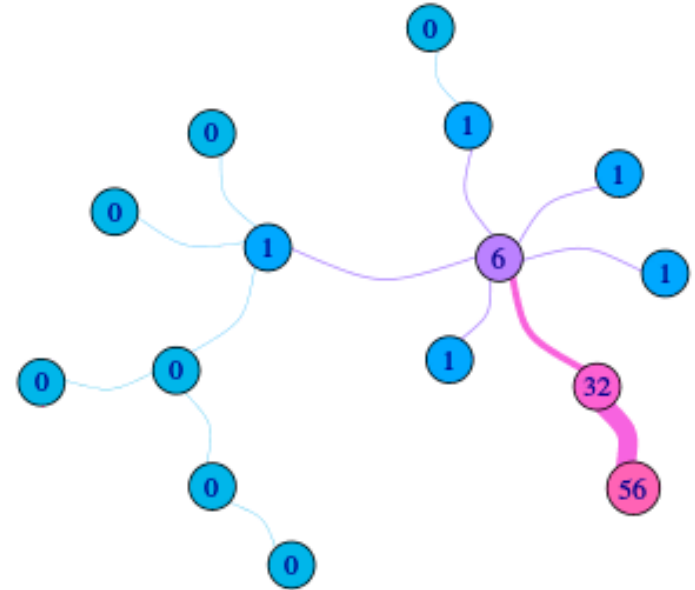
- Adjacency Similarity:
 - Similarity to neighbours only
 - Low complexity

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

Similarity Measures in Verse

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

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



[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
- Exceeds previous embedding methods
- **Scalable:** Processes 10^6 nodes in less than a day
- **Versatile:** supports any similarity measure between nodes



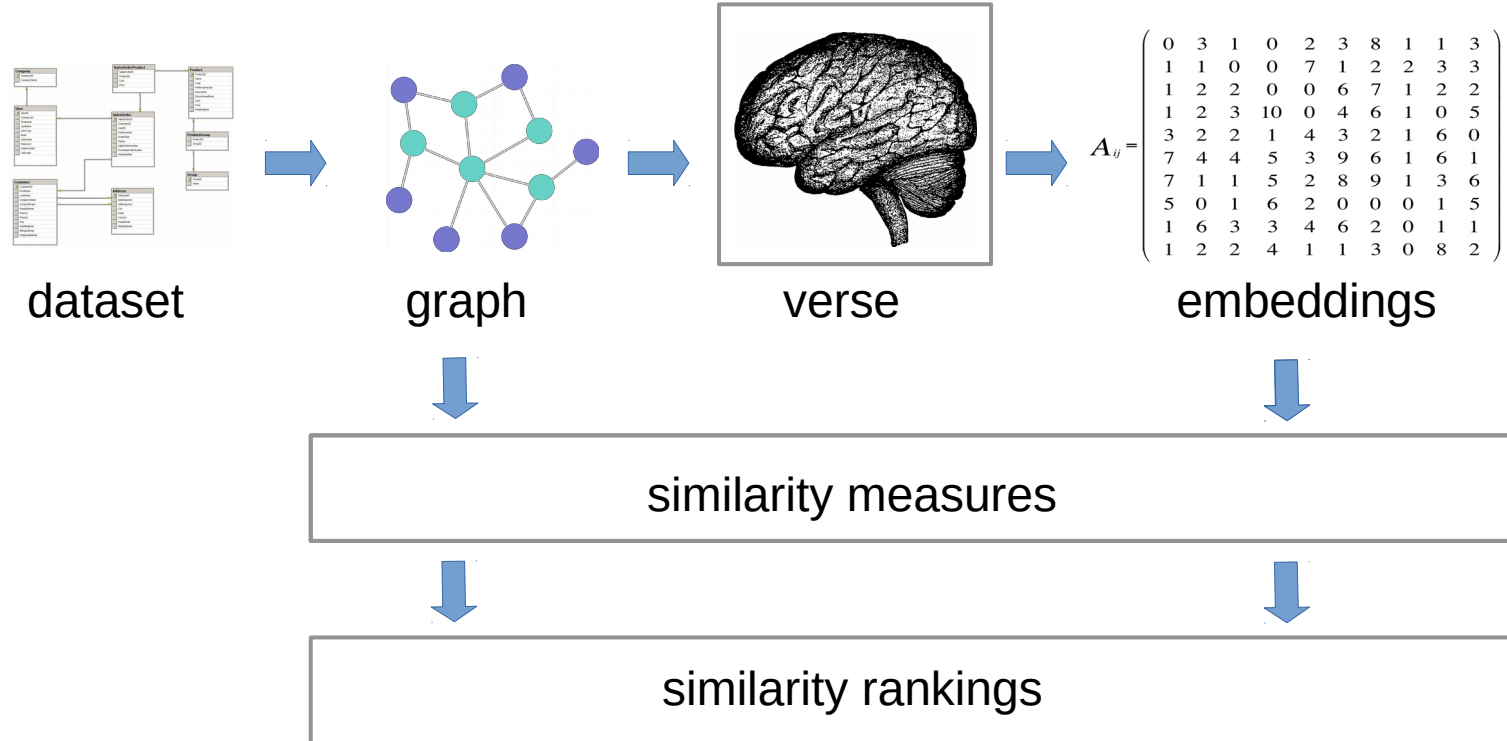
$$\sum_{v \in V} \text{KL}(\text{sim}_G(v, \cdot) \parallel \text{sim}_E(v, \cdot))$$

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
 - ...

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

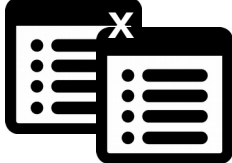
Overview



Spearman coefficient

- Example,

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 I

Data Set



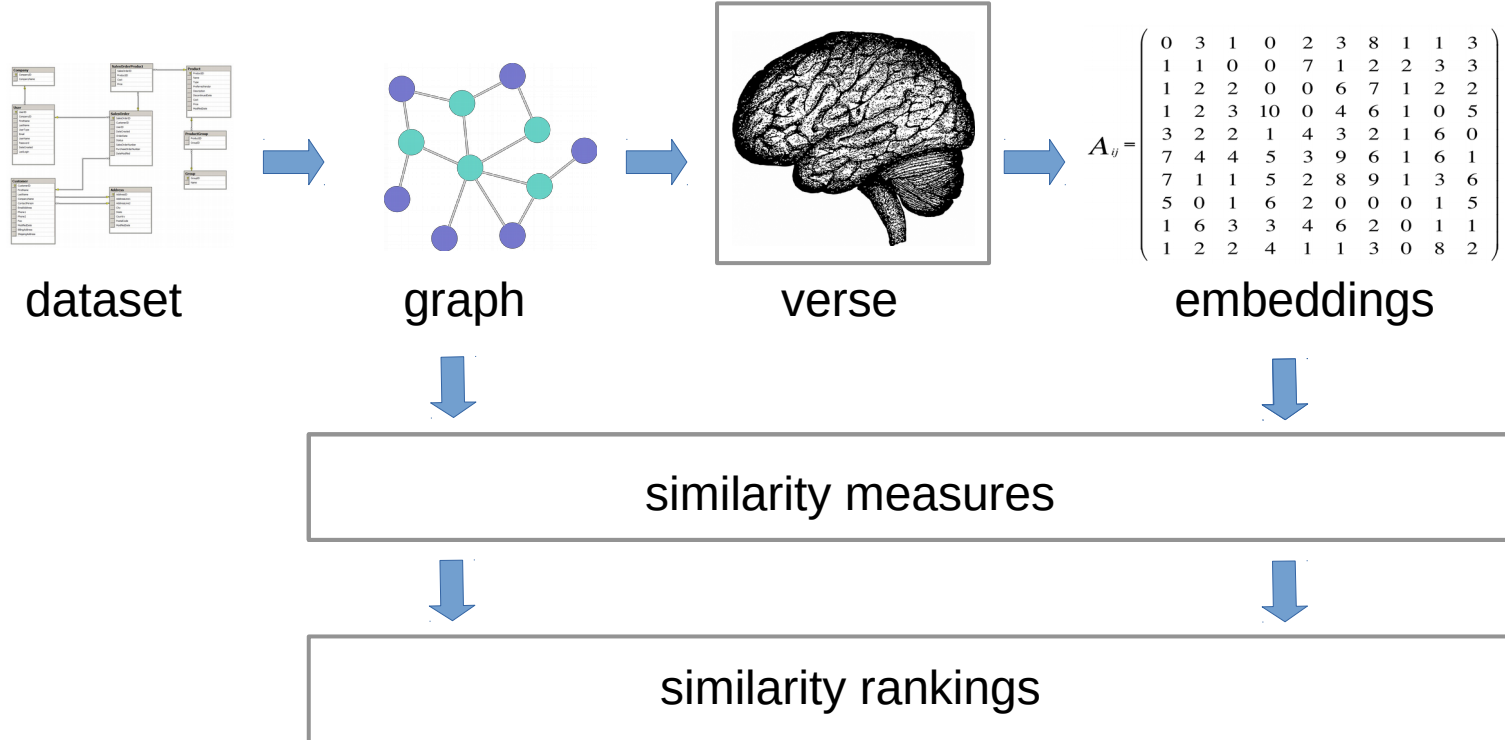
The screenshot shows the Kaggle dataset page for 'The Movies Dataset'. The header includes a 'Dataset' icon and the title 'The Movies Dataset'. Below the title, it states: 'Metadata on over 45,000 movies. 26 million ratings from over 270,000 users.' The creator is listed as 'Rounak Banik · updated a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

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

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

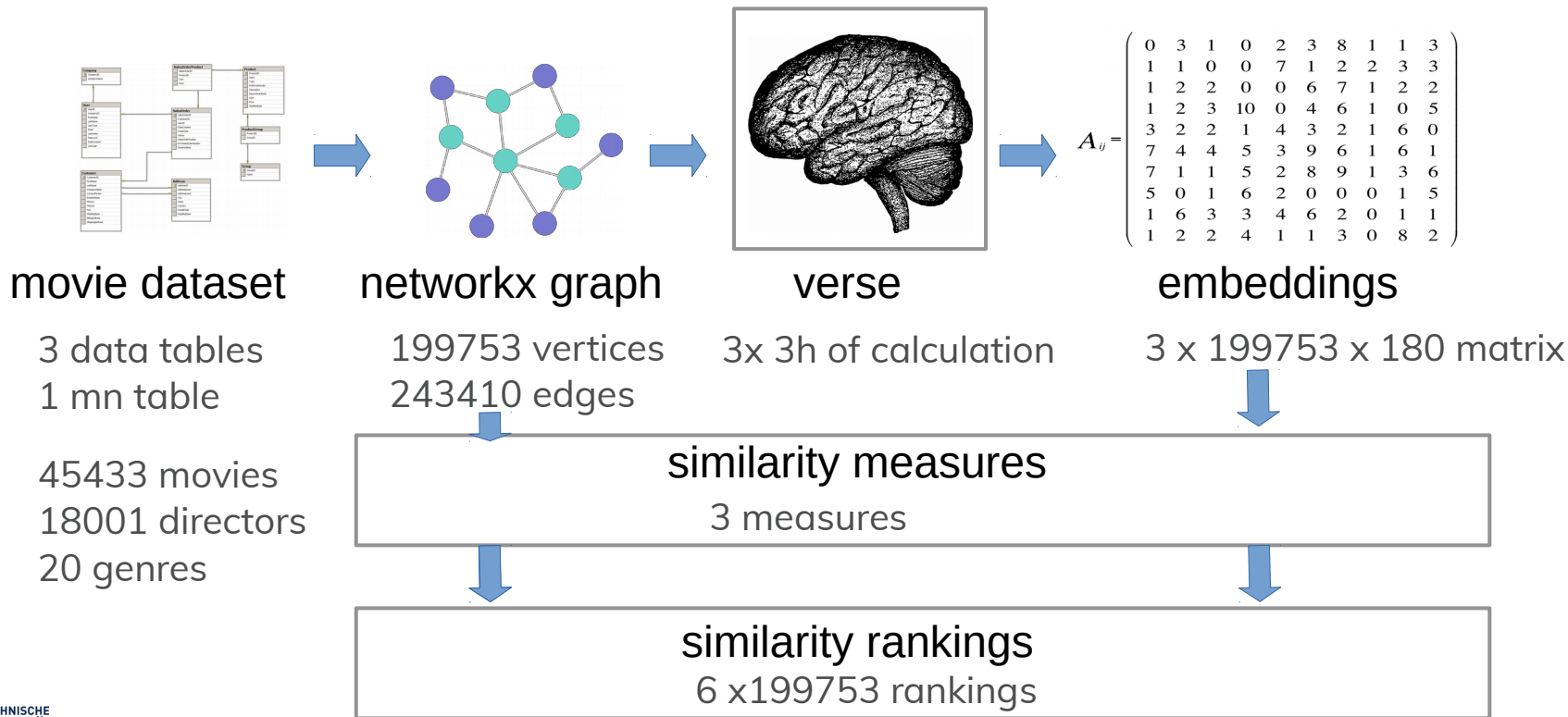
Experimental Setup II

Overview



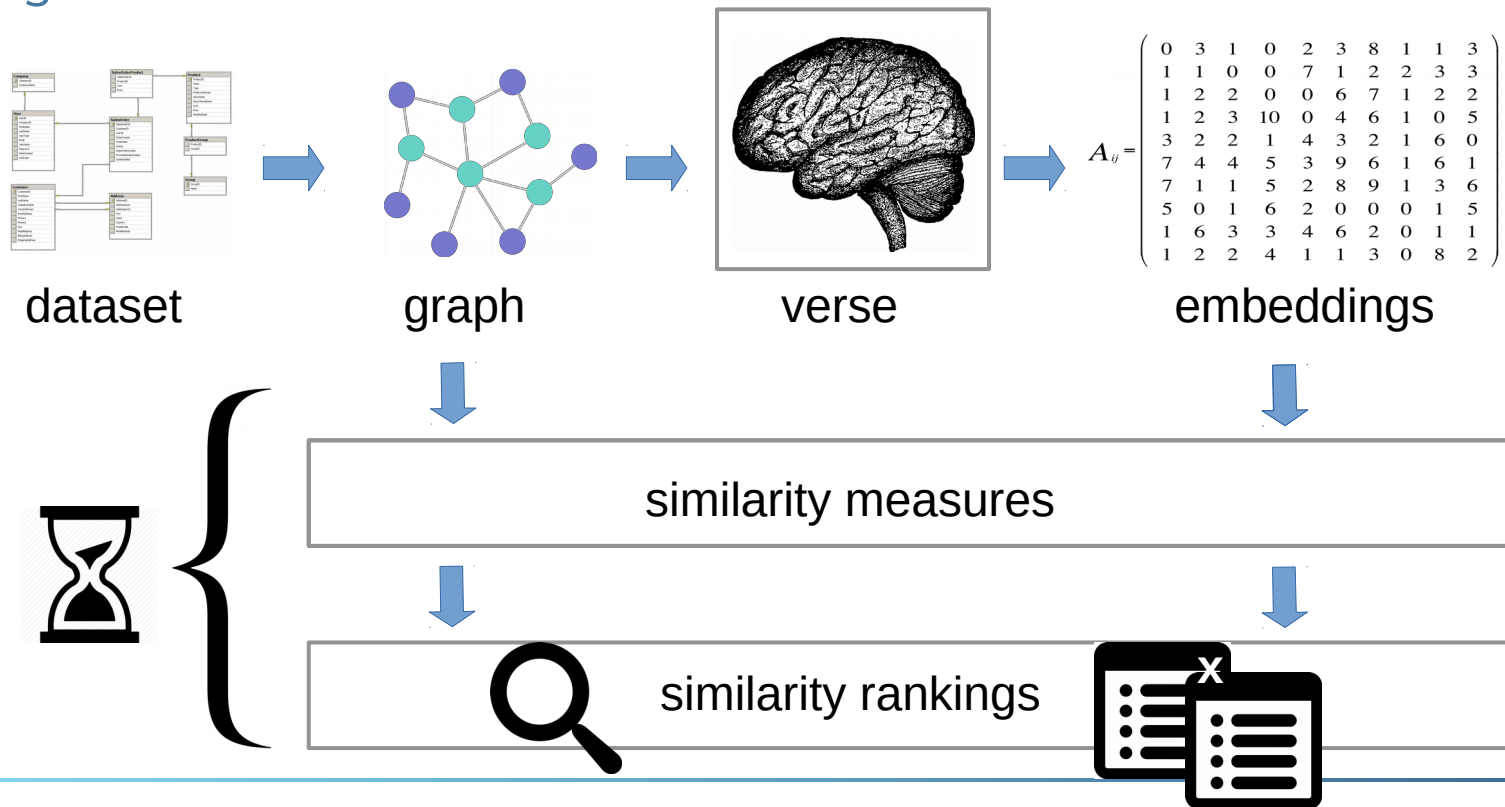
Experimental Setup III

Overview



Experimental Setup IV

Investigation





Results

Results I

PPageRank Graph

	Movie Names
1	Star Wars
2	Star Wars Episode II: Attack of the Clones
3	American Grafitti
4	THX1138
5	Star Wars Episode III: Revenge of the Sith

AdjacencySimilarity Graph

	Movie Names
1	Star Wars
2	Hellboy
3	Get Carter
4	Der Totmacher
5	Begotten
6	Laws of Gravity

Results II

PageRank Embeddings

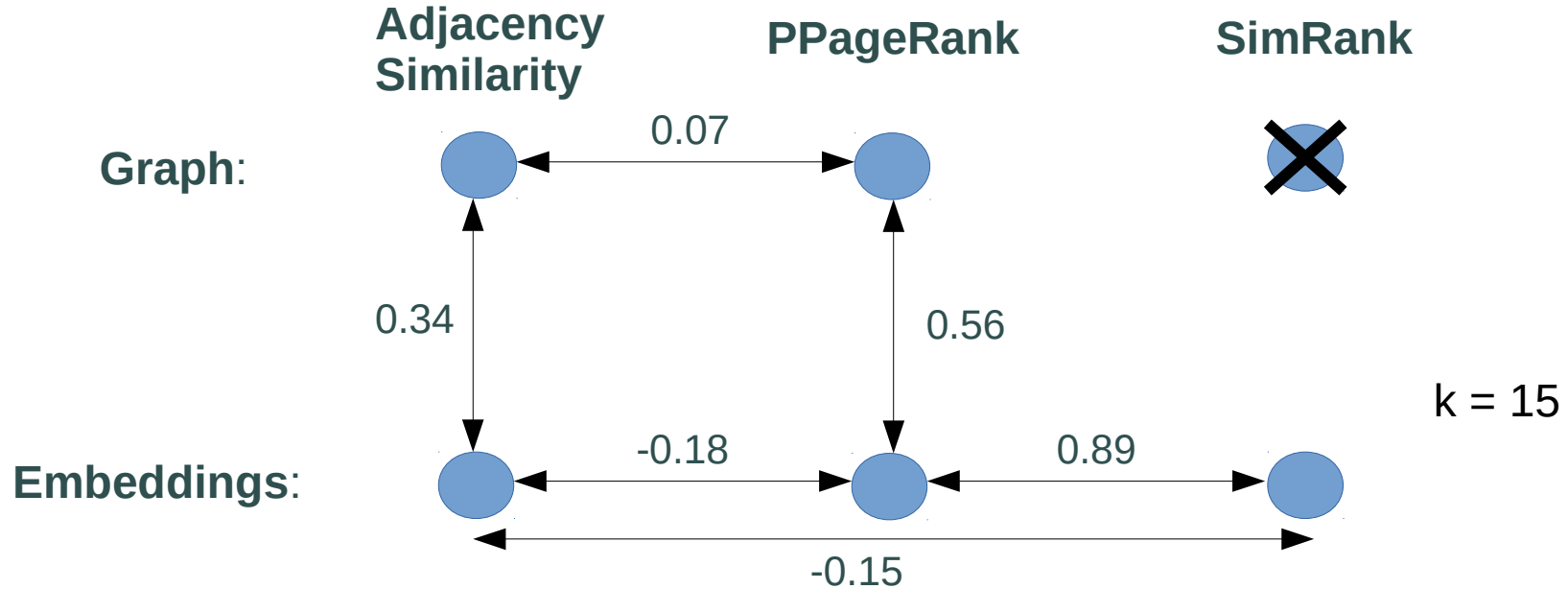
	Movie Names
1	Star Wars
2	Star Wars: Episode III - Revenge of the Sith
3	THX 1138
4	Star Wars: Episode II Attack of the Clones
5	Star Wars: Episode I The Phantom Menace

SimRank Embeddings

	Movie Names
1	Star Wars
2	Star Wars: Episode II Attack of the Clones
3	Electronic Labyrinth THX 1138 4EB
4	THX 1138
5	Star Wars: Episode I The Phantom Menace

Results III

Spearman on Similarity Rankings



Results IV

Computation Time

**Adjacency
Similarity**

PPageRank

SimRank

Graph:

2.77s

12.82s

∞

Embeddings:

0.32s

0.36s

0.21s



Excursus

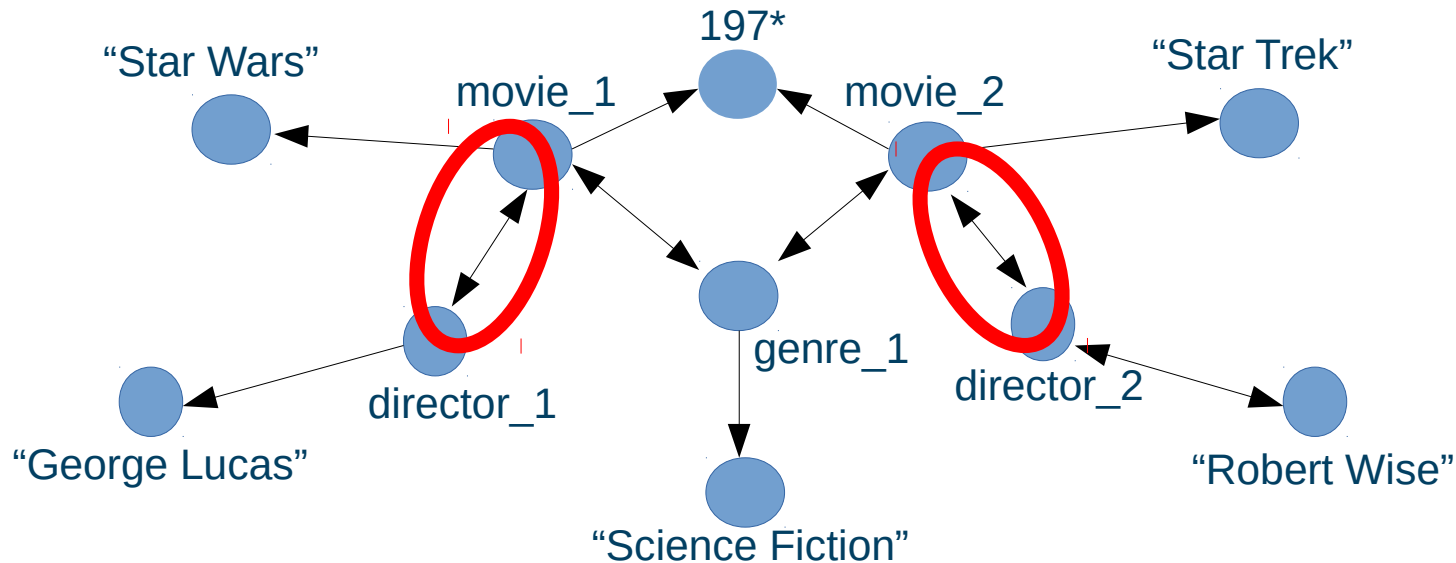
Excursus I



id	Name	Year	d_id
1	Star Wars	1977	1
2	Star Trek	1979	2

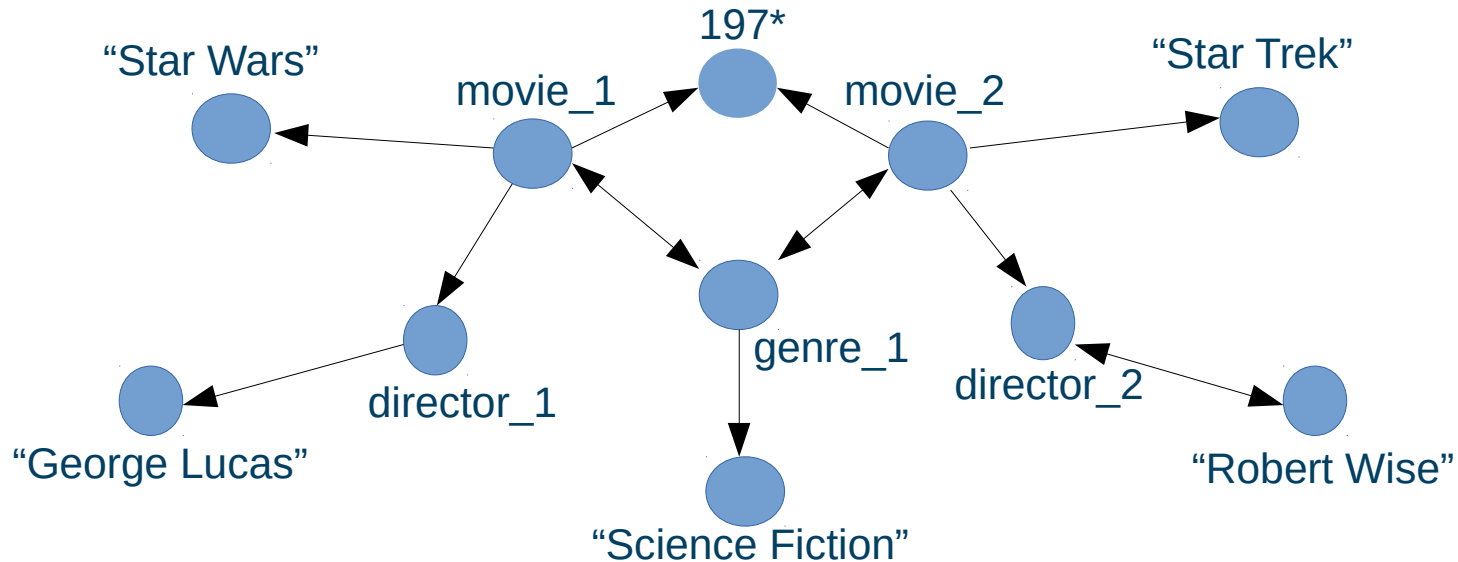
M 1

id	Director
1	George Lucas
2	Robert Wise



Excursus II

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



Excursus III

Directed: SimRank



	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

Bidirected: SimRank



	Movie Names
1	Star Wars
2	Star Wars: Episode II Attack of the Clones
3	Electronic Labyrinth THX 1138 4EB
4	THX 1138
5	Star Wars: Episode I The Phantom Menace



Conclusion

Conclusion I

Computation Time

**Adjacency
Similarity**

PPageRank

SimRank

Graph:

slow

slow

∞

Embeddings:

bad ranking

*efficient
good ranking*

*efficient
good ranking*

Conclusion I



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



- Verse needs preparation time
- Similarity measure and graph construction coupled
- Needs parameter tuning



Fragen?



Extra Slides

Kullback-Leibler-Divergenz

- Also called „*Information Gain*“
- Measure for divergence of probability distributions

$$D(P\|Q) = KL(P, Q) = \sum_{x \in X} P(x) \cdot \log \frac{P(x)}{Q(x)}$$

Verse Steps

- Sampling: NCE-based
- *Dimensionality reduction: (linear)*
- *Normalisation: softmax*
- *Optimisation: KL + Stochastic gradient descent*
- *ML: 1-layer NN + input/output layer*

Vertice Similarity Measures



closeness

Structure-based

- Random Walks
- PpageRank
- SimRank

indirect neighborhood

- Cosine similarity
- Euclidean distance
- Jaccard coeficient
- PPageRank
- SimRank

Direct neighborhood

- Adjacency Similarity

complexity