

Knowledge Inference

- 1. Entity Disambiguation**
- 2. Label Propagation**
- 3. Link Prediction**
- 4. Data Integration**

Basic Questions in Knowledge Inference

Who are the entities?

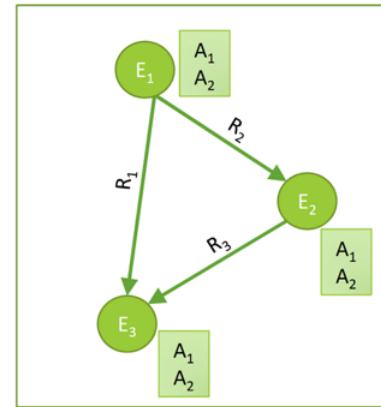
- Entity Linking / Disambiguation
- Data integration

What are their attributes?

- Collective Classification

How are they related?

- Link prediction



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 2

Knowledge inference concerns a wide number of problems that have been studied in many different contexts. Some of the basic examples are:

Entity linking and disambiguation, which concerns the problem of identifying which entity names represent the same real-world entity, respectively which entity is referred to in case of ambiguous entity names.

Schema integration, which concerns the problem which classes, attributes and relationships in one knowledge bases correspond to which in another one.

Collective classification, which concerns the problem of learning unknown attribute values from the available knowledge in a knowledge base.

Link prediction, which concerns the problem of learning unknown relationships from the available knowledge in a knowledge base.

1. ENTITY DISAMBIGUATION

Entity Disambiguation

Task: Link a text mention in a document to an entry in a knowledge base (e.g. Wikipedia or WikiData)

- Also called entity resolution and linking

Example: “**Schindler** is a Swiss industrial company. One of its main competitors is the American producer **Otis**.”

Schindler Group

From Wikipedia, the free encyclopedia

The Schindler Group is a manufacturer of elevators, moving walkways and escalators worldwide, founded in Switzerland in 1874. Schindler produces, installs, maintains and modernizes elevators and escalators in many types of buildings including residential, commercial and high-rise buildings. The company is present in more than 140 countries and employs more than 58,000.



Type	Public
Treated as	SIX: SCLG
Industry	Vertical transportation
Genre	Corporate histories
Founded	Lucerne, Switzerland (1874)

Otis Elevator Company

From Wikipedia, the free encyclopedia

This article may be in need of reorganization to comply with Wikipedia's layout guidelines. Please help by editing the article to make improvements to its overall structure. (January 2019) (Learn how and when to remove this template message)

The Otis Elevator Company is an American company that develops, manufactures and markets elevators, escalators, moving walkways, and related equipment. Based in Farmington, Connecticut.



Industry	Vertical transport systems
Founded	1853; 166 years ago (acquired in 1976)

©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 4

Once entities in a text have been recognized, one would link them to their corresponding counter-parts in a knowledgebase.

Challenge

Two problems

- **Homonyms**: entities with the same name
- **Synonyms**: different names for the same entity

Schindler's List

From Wikipedia, the free encyclopedia

This article is about the film. For the book that inspired this film (published in the U.S. as *Schindler's List*), see *Schindler's Ark*.

Schindler's List is a 1993 American epic historical period drama film directed and co-produced by Steven Spielberg and written by Steven Zaillian. It is based on the novel *Schindler's Ark* by Australian novelist Thomas Keneally. The film follows Oskar Schindler, a Sudeten German businessman, who saved



Otis, Colorado

From Wikipedia, the free encyclopedia Coordinates: 40°9'2"N 102°57'45"W

Otis is a Statutory Town in Washington County, Colorado, United States. The population was 534 at the 2000 census.

Town of Otis, Colorado

Town



Contents [hide]

- 1 History
- 2 Geography
- 3 Demographics
- 4 Climate
- 5 See also
- 6 References

©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 5

Entity disambiguation can however be quite challenging to the homonymy and synonymy problem. Handling these problems is essential for every text analytics tasks. Not being able to handle **homonymy** usually results in the **introduction of noise into the results (poor precision)** whereas not properly handling **synonymy** risks to **miss relevant documents (poor recall)**.

Sources of Information

Local information: similarity of text mention of the entity and the data in the knowledge base

Example: “Schindler” ≈ “Schindler’s list”

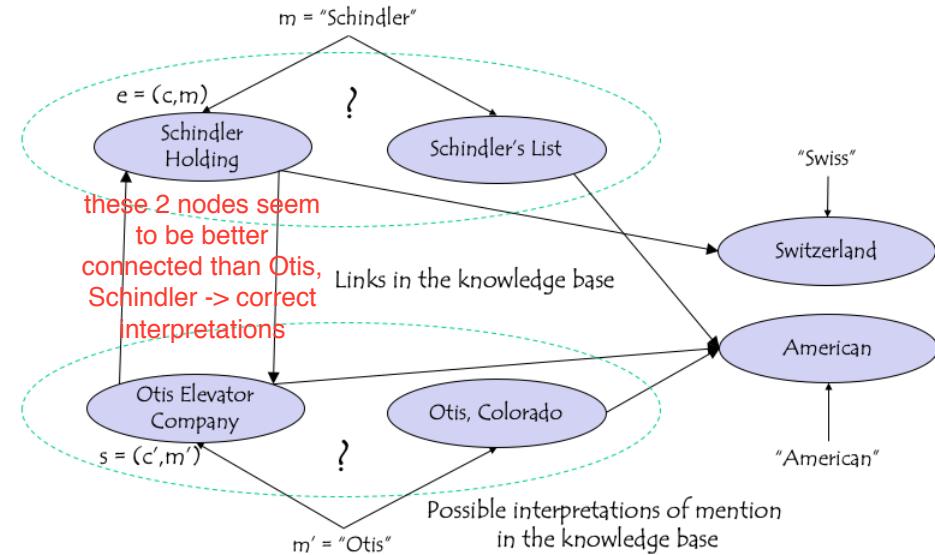
“Schindler” ≈ “Schindler Group”

Global information: coherence of the potential entities corresponding to the different text mentions within a document

→ entity graph

For performing entity disambiguation one can exploit two different sources of information. On the one hand local information extracted from the text mention, or its vicinity, can be exploited. This can be used to compare the text mention and its features with the representation in the knowledge base, in order to obtain evidence which entities in the knowledge base are potentially matched. A second source of information is the knowledge base itself, when multiple entities are mentioned in the same text. Since these entities from the same text are likely to be in some form of relationship among each other, it is likely that such relationships are also discovered in the knowledge base. This can help in entity disambiguation.

Entity Graph



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 7

Here we illustrate the basic model of how a knowledge base can be employed for entity disambiguation:

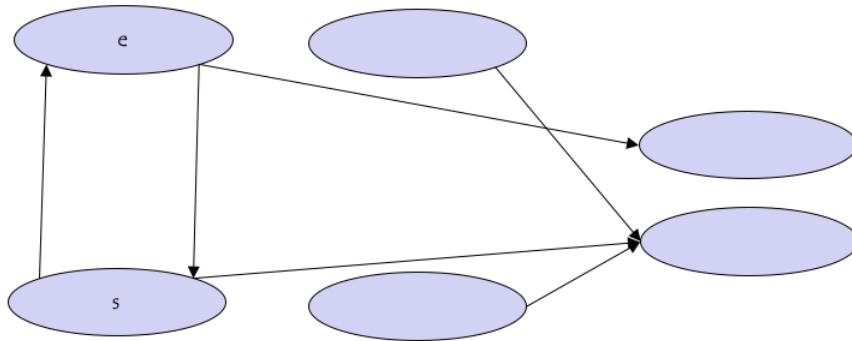
- m are text mentions of entities (extracted using NER)
- For each mention there are candidates in KB; these can be identified using local information from the text mention
- Each is a graph node; we can also associate a similarity measure to each node
- Edges between mention-candidate pairs are included if a link exists in the KB among the respective candidates

Since the same text mention can relate to several candidates, the problem of entity disambiguation is to determine which among the multiple candidates is the most likely to be correct.

Coherence

How well does the choice of node s support the presence of node e ?

Other formulation: How relevant is node e for node s ?



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 8

To perform entity disambiguation, we pose the following question: given a node s in the entity graph, how well does it support the presence of another node e in the graph?

This question is actually related to a question that has been investigated in the context of personalized Web search: given a URL from the personal bookmark list of a user, how relevant is a page in the Web for this user. To answer this question a variation of the PageRank algorithm, which determines general relevance, has been proposed. It is called Personalized PageRank.

pagerank only answer how important is a node in a graph; not w.r.t. another node

Personalized PageRank

Same as PageRank, except that random jumps are always back to the same node (or same set of nodes)

- Original motivation: use personal bookmark list as source of rank

$$\vec{p} = c(qR \cdot \vec{p} + (1 - q)\vec{e})$$
$$\vec{e} = (1, 0, 0, \dots, 0)$$

Can be computed using Monte Carlo method

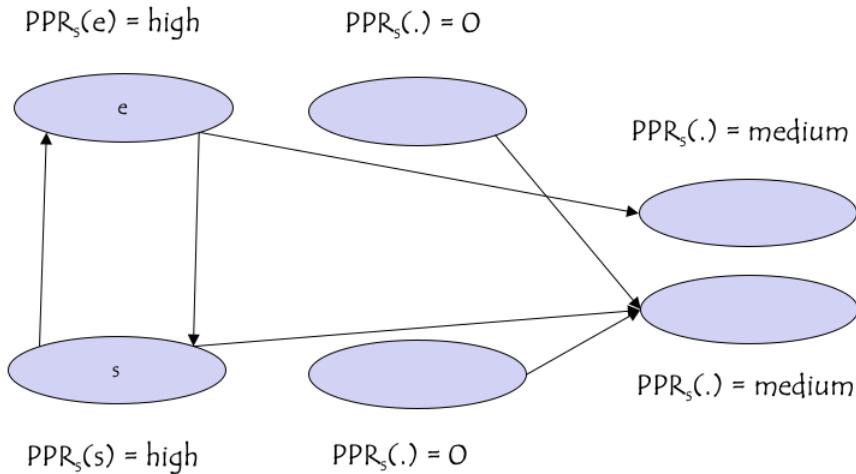
- Perform multiple independent random walks
- Compute distribution of end points of random walks

Personalized PageRank works almost the same as the original PageRank. The difference is that random jumps are not performed uniformly random to any other node of the graph, but to a selected subset of nodes. In the case of Web search this subset would be the personal bookmark list, and by jumping back to this subset nodes from that list will have a large influence on the ranking of a page, such that nodes that are well connected to the bookmarks will receive more ranking value.

PPR can be computed either like standard PageRank, or using a Monte Carlo method, by starting random walks independently and aggregate the distribution of the end points of those walks.

PPR on Entity Graph

Source node s



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 10

Now applying PPR, by considering a source node s as the selected subset of preferred nodes, will generate a distribution of ranking values for all other nodes. Nodes that are well connected to s will naturally receive higher ranking values. Intuitively it is clear, that in our example node e will be receive higher ranking when starting from s, and this be probably the preferred interpretation for the entity.

Approach

Finding the concept candidate linked to a mention m that is most likely to be valid

1. For all concept candidates c compute total support received from other nodes

$$e = (c, m), s = (c', m')$$
$$score_e = \sum_{s \in Contributors_e} PPR_s(e)$$

e-> very well connected with other sources s => high score

1. Select the candidate with highest score

The procedure for selecting the best concept candidate for a text mention is straightforward. One computes the total support that the candidate receives from other nodes s, by computing the PPR of e with the source node s a selected node. Which source nodes s are contributing to this computation we need still to determine. At the end the candidate with the highest score is selected.

Contributing Nodes

Only one interpretation c for a mention m is valid

- Competing nodes $s = (c', m)$ that have the same entity mention as $e = (c, m)$ cannot support e
1. Schidler's list cannot contribute as source to interpretation of Schindler group!!
- For multiple nodes s that have the same entity mention m , only the one with highest contribution is considered
Multiple src nodes that mention Schindler's list in the same way, only pick the highest

Thus

$$\text{Contributors}_e = \{(m', \underset{c}{\operatorname{argmax}} PPR_{(m', c)}(e), m' \neq m)\}$$

Not every node should contribute as source node for a mention m , if we assume that only one interpretation can be correct. First, we exclude nodes that are “competing” for the interpretation of the text mention. Second for multiple nodes related to the same other entity mention, we select only the contribution of the one that has the highest value. In this way we favor a unique interpretation for mentions.

Considering Popularity

The method can furthermore consider popularity measures for nodes

- If information is insufficient, favor popular nodes

$$\begin{aligned} coherence_s(e) &= PPR_s(e) \text{popularity}(s) && \text{popularity of src !} \\ coherence(e) &= \sum_{s \in \text{Contributors}_e} coherence_s(e) \\ score(e) &= coherence(e) + PPR_{avg} \text{popularity}(e) && \text{PPRavg * popularity of node e} \end{aligned}$$

To further improve the method, one can add a general popularity measure to weight the contributions of source nodes. This will favor popular nodes, which is beneficial if little information is available for disambiguation. Then it is better to choose a popular candidate since chances that this is correct are higher. One of possible choice of a popularity score could be the number of links a node has in the knowledge base.

Some Results

Models	Cucerzan	Kulkarni	Hoffart	Shirakawa	Alhelbawy	iSim	PPR	Without popularity	
								Other methods	Uses pageRank
Micro	51.03	72.87	81.82	82.29	87.59	62.61	85.56	91.77	
Macro	43.74	76.74	81.91	83.02	84.19	72.21	85.86	89.89	

Micro = fraction of correctly disambiguated entities

Macro = textual mentions, correctly disambiguated per entity, averaged over all entities

Experimental results show that the method works relatively well, with around 90% of entities that are correctly disambiguated. One can observe that the use of popularity helps to slightly improve the results.

2. LABEL PROPAGATION

Inferring Attributes

Example problem: Which users on Twitter have positive or negative emotion towards a topic?

- Users are nodes in a graph (follower network)
- Emotion is an attribute of the node (or a class)

Potential source of information

- Emoticons in tweets: indicate stance of user towards the topic
- Only a (small) fraction of the users is using emoticons

A second inference task in knowledge bases, after entities have been disambiguated, is to assign to the entities correct attribute values. For discrete ^{+ve/-ve} attributes this problem can be understood as a classification problem. This question has, for example, been studied for classifying users in a social network with respect to their stance or emotion towards a specific topic.

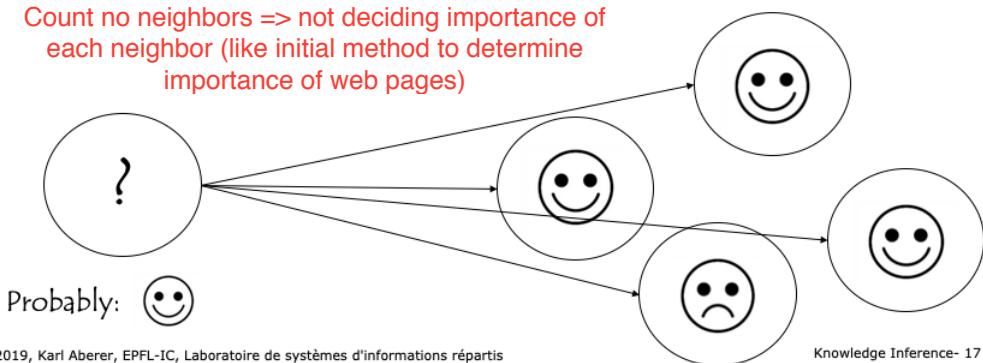
In the case of emotion analysis there exists typically indications of emotions, e.g., in the form of the use of emoticons or specific hashtags. However, only few users are using those.

Propagating Attribute Values

Assume that nodes that are connected by an edge, have a higher propensity of sharing the attribute of interest

- Twitter users following each other probably also are more likely to share the same emotion towards a topic

Count no neighbors => not deciding importance of each neighbor (like initial method to determine importance of web pages)



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 17

In many cases nodes that are connected by edges in a knowledge graph or as well in a social network, share properties.

In social networks this is quite apparent. People that are connected through social links (e.g follower, friend, retweet, reply etc) are in general more likely to share opinions than those that are not. Is it possible to exploit this property to predict the attributes (respectively class labels) for those users that have none?

Model

Graph $G = (V, E)$ with vertices V and Edges E

- Vertices V have a label from a set $L \cup \{\text{unknown}\}$ of size $n + 1$
- Edges are undirected and unweighted

Probability distributions for vertices

- $\mathbf{l}_{\text{apriori}}(v)$ is a vector of size $n + 1$
assigns weight 1 for label if known for $v \in V$ know apriori the label (ex users that used the smiley)
- $\mathbf{l}_{\text{inferred}}(v)$ is a vector of size $n + 1$
assigns a probability distribution to vertices true prob distrib
- $\mathbf{l}_{\text{unknown}}$ is a vector of size $n + 1$
assigns weight 1 for label unknown

©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 18

We consider the following model. A graph with vertices that can have either a label from a set L, or a label unknown. The edges are all undirected and unweighted. The model and approach can be extended for directed graphs with weighted edges.

We associate with vertices probability distributions that represent our knowledge about the assignment of a label. For some vertices we have an apriori assignment of labels (these are the vertices for which the label is known). For all vertices we will compute an inferred probability distribution. For vertices that have no apriori label assigned we have a vector to represent the unknown state.

Note that the apriori distribution can also be a true probability distribution, if we are initially not sure about the label, but have some knowledge.

Label Inference

We assume that all neighbors exert the same influence on a node

Thus we would require that

$$l_{inferred}(v) = \frac{1}{\deg(v)} \sum_{(v,w) \in E} l_{inferred}(w)$$

average of inferred weight of neighbors

Recursive equation resp. random walk model (like with PageRank)

In this model we assume that all neighbors of a node in the graph are equally influencing it. Thus, we can compute its expected label distribution by averaging the distributions of the neighbors. The resulting equation is analogous to the formulation of the PageRank model. Thus the Label Inference can be interpreted as a random walk model.

The model can be extended to weighted graphs, in which case the edge weight would be considered in the aggregation of the distributions of the neighboring vertices.

Label Propagation Algorithm

$l_{inferred}(v) = l_{apriori}(v)$ for nodes with known labels, otherwise $l_{unknown}$
while not converged

$$l_{inferred}(v) = \frac{1}{\deg(v)} \sum_{(v,w) \in E} l_{inferred}(w) \quad \text{continuously feeding apriori knowledge}$$
$$l_{inferred}(v) = p_v^{inj} l_{apriori}(v) + p_v^{continue} l_{inferred}(v) + p_v^{abandon} l_{unknown} \quad \begin{array}{ll} // \text{ inject apriori knowledge} \\ // \text{ infer from neighbors} \\ // \text{ stop} \end{array}$$

prob continue=> $p_v^{continue}$
prob abandon=> $p_v^{abandon}$

The probabilities p_v^{inj} , $p_v^{continue}$ and $p_v^{abandon}$ can be interpreted as decisions in a random walk

The full label propagation model adds additional aspects to the propagation of labels to neighbors.

- For vertices with apriori labels, at every step the apriori distribution is injected with a certain probability p_v^{inj} that depends on the vertex
- For all vertices the propagation process can also be abandoned with a certain probability $p_v^{abandon}$
- The propagation of the label distribution to neighbors occurs then with a (remaining) probability of $p_v^{continue}$

As in PageRank the process is iterated till convergence occurs.

Determining the Probabilities

Entropy of transition probabilities $H(v) = -\log \frac{1}{\deg(v)}$

$$c_v = \frac{\log 2}{\log(2+\deg(v))}$$

$d_v = (1 - c_v)\sqrt{H(v)}$ if v is labelled, 0 otherwise

$$z_v = \max(c_v + d_v, 1)$$

$$p_v^{cont} = \frac{c_v}{z_v}, \quad p_v^{inj} = \frac{d_v}{z_v} \text{ for labelled nodes, 0 otherwise}$$

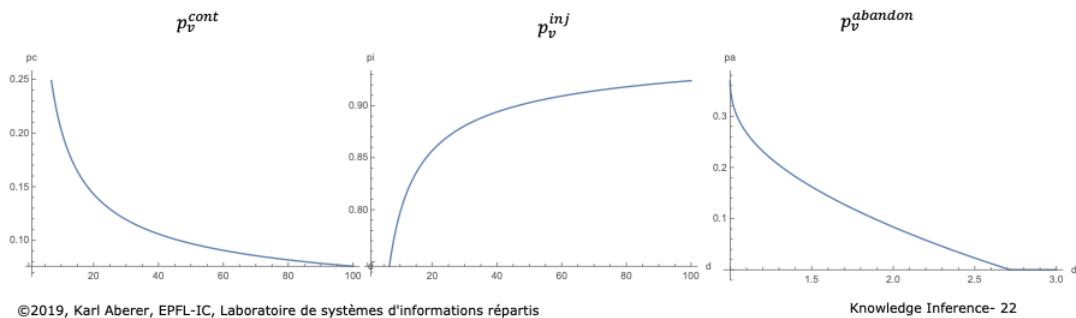
$$p_v^{abandon} = 1 - p_v^{cont} - p_v^{inj}$$

The transition probabilities depend on the properties of the vertices. In our model, the only relevant property is its degree.

In a more general model with edge weights the probabilities would depend on the distribution of edge weights of the edges connected to the vertex (Fan-out entropy heuristics)

Behavior of Probabilities: Labelled Nodes

- Injection probability increases with the degree of the nodes, while continuation probability decreases
- Abandoning probability positive for only very low degree nodes ($d = 1, 2$) nodes at the edges of graph should not inject too much



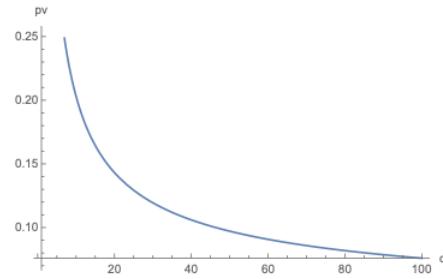
For labelled nodes the injection probability increases with the degree. Thus, well connected pre-labelled nodes have a lot of influence.

Behavior of Probabilities: Unlabelled Nodes

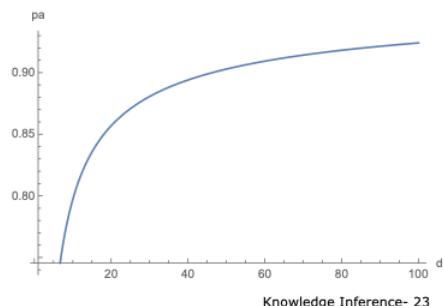
- Abandon probability increases with degree
- Prevents algorithm from propagating information through unlabeled, high-degree nodes

for high degree unconnected nodes we want less influence from them so they don't propagate lots of wrong knowledge

p_v^{cont}



$p_v^{abandon}$



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 23

For unlabeled nodes the behavior the abandon probability increases with the degree. Thus, high degree nodes have less influence.

Extensions

Label Propagation can be extended to

- A prior knowledge given as probability distribution
- Graphs with weighted edges
- Directed Graphs

Alternative algorithm: MAD (modified adsorption)

- Formulates an optimization problem and solves it directly
- Slightly better performance in practice

Semi-Supervised Learning

Label propagation is an example of a semi-supervised learning algorithm

- Exploit partial labelling
- Useful in cases where labels are sparse or labels can produced only for special cases using heuristics or background knowledge check in a cluster which is the dominating label => these nodes prob belong to this class
- Require that relationships among entities and their labels are correlated for some underlying principle

3. LINK PREDICTION

Knowledge Base Completion

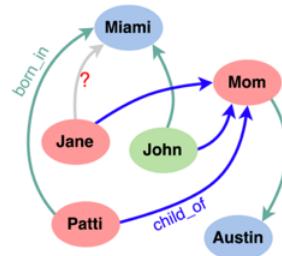
Large knowledge bases are usually incomplete

- DBpedia: 60% of persons miss place of birth
- FreeBase: 71% of persons miss place of birth etc

Try to predict missing links from existing data

Jane born in Miami?

Probably

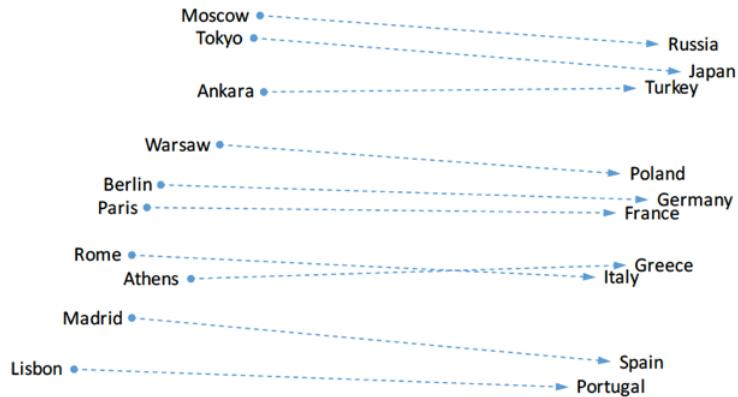


©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 27

In general, knowledge bases are incomplete. Thus there is a significant interest in completing in particular the relationships among entities. To do so one might exploit “patterns” that entities and relationships follow, and generalize them.

Observation on Word Embeddings



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 28

In order to tackle this problem, we come back to an observation that we had made for word embeddings. We found that (in some cases) relationships seem to be encoded as linear transformations.

Relations in Word Embeddings

$$\begin{aligned}\mathbf{v}_{Japan} - \mathbf{v}_{Tokyo} &\approx \mathbf{v}_{Germany} - \mathbf{v}_{Berlin} \\ \mathbf{v}_{Germany} - \mathbf{v}_{Berlin} &\approx \mathbf{v}_{Italy} - \mathbf{v}_{Rome} \\ \mathbf{v}_{Italy} - \mathbf{v}_{Rome} &\approx \mathbf{v}_{Portugal} - \mathbf{v}_{Lisbon}\end{aligned}$$

Idea: Find a vector
 $\mathbf{v}_{is_capital_of}$ such that

$$\begin{aligned}\mathbf{v}_{Tokyo} + \mathbf{v}_{is_capital_of} - \mathbf{v}_{Japan} &\approx \mathbf{0} \\ \mathbf{v}_{Berlin} + \mathbf{v}_{is_capital_of} - \mathbf{v}_{Germany} &\approx \mathbf{0} \\ \mathbf{v}_{Rome} + \mathbf{v}_{is_capital_of} - \mathbf{v}_{Italy} &\approx \mathbf{0} \\ \mathbf{v}_{Lisbon} + \mathbf{v}_{is_capital_of} - \mathbf{v}_{Portugal} &\approx \mathbf{0}\end{aligned}$$

relations are also represented as embedding vector!

This observation can be formally described in two different ways. We could say, that the differences of entity vectors that are in a relationship should result in similar values. Alternatively, we could also say, that there must be some vector that represents the relationship and that the sum of this vector with the difference vector of the entities should be approximately zero. This second formulation of the properties gives rise to an idea of how link prediction could be performed using an embedding technique.

Model

Knowledge graph G consists of (correct) triples (h, r, t)

where $h, t \in E$ and $r \in R$

E: nodes in graph

Define a plausibility score $f(h, r, t)$

$$f(h, r, t) > f(h', r', t')$$

If (h, r, t) is a plausible triple and (h', r', t') is a implausible triple

To formulate the model we assume that we have a knowledge graph consisting of triples (h, r, t) . h indicates the term “head” and t the term “tail”.

Then we introduce a function that measure how plausible a triple is (we can think of it as a probability). Clearly, the function should return higher values for plausible tuples than for implausible ones.

Learning the Model

Minimize the loss function

$$J(\theta) = \sum_{\substack{(h,r,t) \in G \\ (h',r',t') \in G' \\ \text{bad triple derived from good triple}}} \max(0, \gamma + f(h, r, t) - f(h', r', t'))$$

where $G'(h, r, t)$ is a set of incorrect triples, generated by corrupting the correct triple (h, r, t)

γ is a hyperparameter

Using the plausibility function, we can formulate a loss function that should be minimized. The minimization will be performed as usual using SGD.

TransE Model

One of the first embedding-based models for knowledge base completion

- Based on the intuition from text WE

$$f(h, r, t) = \|\mathbf{v}_h + \mathbf{v}_r - \mathbf{v}_t\|_{1/2}$$

tokyo - is_capital - japan = 0

Each entity and relationship is mapped to a low-dimensional vector, resulting in $\mathbf{v}_h, \mathbf{v}_r, \mathbf{v}_t$

With the generic model we can instantiate difference concrete approaches for link prediction, depending on the choice of f . One of the initial approaches that was strongly inspired by our initial observation on word embeddings from the beginning. It directly introduces a vector \mathbf{v}_r that represents the linear transformation corresponding to a relationship, and mapping the head and tail into the same (low-dimensional) vector space.

The mappings from the entities and relationships to the latent space are performed (as in word embeddings) using embedding matrices, which constitute the model parameters that have to be learnt using SGD.

Performing SGD

Initialize vectors with random values

- From interval $[-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}}]$ where k is the embedding dimension
- In each iteration
 - Sample a correct triple or batch
 - Derive a corrupt triple from the correct one: replace h or t by a random entity
 - Update embeddings by minimizing loss function
 - Normalize all entity vectors to 1 (not relationship vectors!)

The SGD algorithm proceeds as follows. First the vectors are initialized with random values (the choice is motivated by empirical findings from neural network training). Then in every iteration a triple (or several triples are randomly chosen). Negative samples are generated by randomly replacing head or tail (not both). The update of the embeddings is performed as usual by computing the differential of the loss function. Entity vectors are normalized to 1 in every iteration (this avoids the model to find a trivial solution).

Qualitative Results

Bold: good answer

Italic: also in KB, not for this learning but also good

INPUT (HEAD AND LABEL)	PREDICTED TAILS
J. K. Rowling influenced by	<i>G. K. Chesterton, J. R. R. Tolkien, C. S. Lewis, Lloyd Alexander,</i> Terry Pratchett, Roald Dahl, Jorge Luis Borges, Stephen King , Ian Fleming
Anthony LaPaglia performed in	<i>Lantana, Summer of Sam, Happy Feet, The House of Mirth,</i> Unfaithful, Legend of the Guardians , Naked Lunch, X-Men, The Namesake
Camden County adjoins	Burlington County, Atlantic County, Gloucester County, Union County, Essex County, New Jersey, Passaic County, Ocean County, Bucks County
The 40-Year-Old Virgin nominated for	<i>MTV Movie Award for Best Comedic Performance,</i> BFCA Critics' Choice Award for Best Comedy, <i>MTV Movie Award for Best On-Screen Duo,</i> MTV Movie Award for Best Breakthrough Performance, MTV Movie Award for Best Movie , MTV Movie Award for Best Kiss, D. F. Zanuck Producer of the Year Award in Theatrical Motion Pictures, Screen Actors Guild Award for Best Actor - Motion Picture
Costa Rica football team has position	<i>Forward, Defender, Midfielder, Goalkeepers,</i> Pitchers, Infielder, Outfielder, Center, Defenseman
Lil Wayne born in	<i>New Orleans, Atlanta, Austin, St. Louis,</i> Toronto, New York City, Wellington, Dallas, Puerto Rico
WALL-E has the genre	Animations, Computer Animation, <i>Comedy film,</i> Adventure film, Science Fiction, Fantasy , Stop motion, Satire, Drama

The qualitative results show that the method works reasonably well. The bold phrases are the correct predictions. However, given that the knowledge bases used for evaluation are incomplete, it is not out of question that also other predictions are meaningful.

Alternative Models

Model	Score function $f(h, r, t)$	Opt.
Unstructured	$\ v_h - v_t\ _{\ell_{1/2}}$	SGD
SE	$\ \mathbf{W}_{r,1}v_h + \mathbf{W}_{r,2}v_t\ _{\ell_{1/2}} ; \mathbf{W}_{r,1}, \mathbf{W}_{r,2} \in \mathbb{R}^{k \times k}$	SGD
SME	$(\mathbf{W}_{1,1}v_h + \mathbf{W}_{1,2}v_r + \mathbf{b}_1)^\top (\mathbf{W}_{2,1}v_t + \mathbf{W}_{2,2}v_r + \mathbf{b}_2)$ $\mathbf{b}_1, \mathbf{b}_2 \in \mathbb{R}^n ; \mathbf{W}_{1,1}, \mathbf{W}_{1,2}, \mathbf{W}_{2,1}, \mathbf{W}_{2,2} \in \mathbb{R}^{n \times k}$	SGD
TransE	$\ v_h + v_r - v_t\ _{\ell_{1/2}} ; v_r \in \mathbb{R}^k$ $\ (\mathbf{I} - r_p r_p^\top)v_h + v_r - (\mathbf{I} - r_p r_p^\top)v_t\ _{\ell_{1/2}}$ $r_p, v_r \in \mathbb{R}^k ; \mathbf{I} : \text{Identity matrix size } k \times k$	SGD
TransH	$\ \mathbf{W}_r v_h + v_r - \mathbf{W}_r v_t\ _{\ell_{1/2}} ; \mathbf{W}_r \in \mathbb{R}^{n \times k} ; v_r \in \mathbb{R}^n$	SGD
TransR	$\ (\mathbf{I} + r_p r_p^\top)v_h + v_r - (\mathbf{I} + r_p r_p^\top)v_t\ _{\ell_{1/2}}$ $r_p, v_r \in \mathbb{R}^n ; h_p, t_p \in \mathbb{R}^k ; \mathbf{I} : \text{Identity matrix size } n \times k$	SGD
TransD	$\ (\mathbf{I} + r_p r_p^\top)v_h + v_r - (\mathbf{I} + r_p r_p^\top)v_t\ _{\ell_{1/2}}$ $r_p, v_r \in \mathbb{R}^n ; h_p, t_p \in \mathbb{R}^k ; \mathbf{I} : \text{Identity matrix size } n \times k$	AdaDelta
lppTransD	$\ (\mathbf{I} + r_{p,1} h_p^\top v_h + v_r - (\mathbf{I} + r_{p,2} t_p^\top)v_t\ _{\ell_{1/2}}$ $r_{p,1}, r_{p,2}, v_r \in \mathbb{R}^n ; h_p, t_p \in \mathbb{R}^k ; \mathbf{I} : \text{Identity matrix size } n \times k$	SGD
STransE	$\ \mathbf{W}_{r,1}v_h + v_r - \mathbf{W}_{r,2}v_t\ _{\ell_{1/2}} ; \mathbf{W}_{r,1}, \mathbf{W}_{r,2} \in \mathbb{R}^{k \times k} ; v_r \in \mathbb{R}^k$	SGD
TranSparse	$\ \mathbf{W}_r^\theta(\theta_r^\theta)v_h + v_r - \mathbf{W}_r^\theta(\theta_r^\theta)v_t\ _{\ell_{1/2}} ; \mathbf{W}_r^\theta, \theta_r^\theta \in \mathbb{R}^{n \times k} ; \theta_r^\theta, \theta_r^\theta \in \mathbb{R} ; v_r \in \mathbb{R}^n$	SGD
DISTMULT	$v_h^\top \mathbf{W}_r v_t ; \mathbf{W}_r \text{ is a diagonal matrix} \in \mathbb{R}^{k \times k}$	AdaGrad
NTN	$v_r^\top \tanh(v_h^\top \mathbf{M}_r v_t + \mathbf{W}_{r,1}v_h + \mathbf{W}_{r,2}v_t + \mathbf{b}_r)$ $v_r, \mathbf{b}_r \in \mathbb{R}^n ; \mathbf{M}_r \in \mathbb{R}^{k \times k \times n} ; \mathbf{W}_{r,1}, \mathbf{W}_{r,2} \in \mathbb{R}^{n \times k}$	L-BFGS
HolE	$\text{sigmoid}(v_r^\top (v_h \circ v_t)) ; v_r \in \mathbb{R}^k, \circ \text{ denotes circular correlation}$	AdaGrad
Bilinear-COMP	$v_h^\top \mathbf{W}_{r_1} \mathbf{W}_{r_2} \dots \mathbf{W}_{r_m} v_t ; \mathbf{W}_{r_1}, \mathbf{W}_{r_2}, \dots, \mathbf{W}_{r_m} \in \mathbb{R}^{k \times k}$	AdaGrad
TransE-COMP	$\ v_h + v_{r_1} + v_{r_2} + \dots + v_{r_m} - v_t\ _{\ell_{1/2}} ; v_{r_1}, v_{r_2}, \dots, v_{r_m} \in \mathbb{R}^k$	AdaGrad
ConvE	$v_t^\top g(\text{vec}(g(\text{concat}(\bar{v}_h, \bar{v}_r) * \Omega)) \mathbf{W}) ; g \text{ denotes a non-linear function}$	Adam
ConvKB	$w^\top \text{concat}(g([v_h, v_r, v_t]) * \Omega)) ; * \text{ denotes a convolution operator}$	Adam

©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

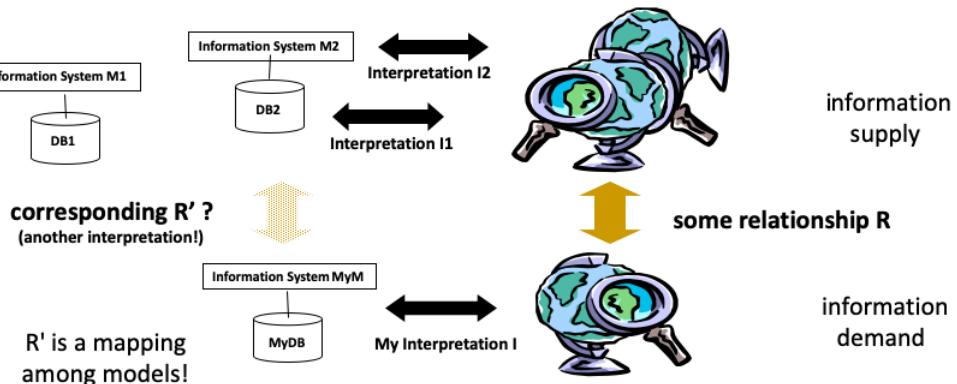
Knowledge Inference- 35

The TransE method is only one example of very many methods that have been in the meanwhile proposed to tackle the link prediction problem.

4. DATA INTEGRATION

Key Tasks in Distributed Information Management

More data! ... More models!? ... More useful information?



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 37

This figure illustrates the problem of information starvation: there exist many information systems supplying data, but each having their own view on the world, which does not necessarily match the needs or understanding of a specific consumer. Every information system is interpreting its model differently with respect to the real world and relating to different views on the real world. Though there exists some relationships among all these views on the real world (let's denote it as R), and it surely implies some relationship R' among the different models used in the different information systems, the consumers of the information cannot easily understand the relationship R , and thus can also not easily relate their models to the models of others via the relationship R' . From the viewpoint of the data providers, introducing R' introduces a new interpretation of their data with respect to the model used by the data consumer.

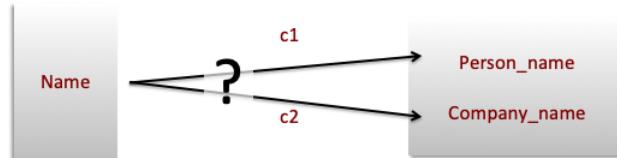
Schema Matching

Integration of heterogeneous data sources

- Every project on Big Data analysis first has to integrate data from different, heterogeneous data sources
- Different schemas, taxonomies, knowledge graphs
- One of the long-standing open problems in data management

How to find good “matches”?

How to choose the “best matches”?



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 38

Schema matching is a long standindg problem, both in industry and research.

Schema Matching Approaches

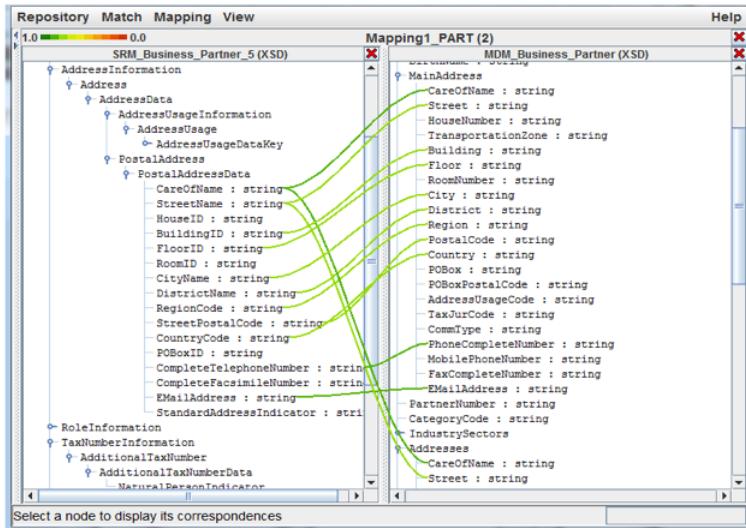
Manual matching

- still common practice today

Schema matching tools

- Based on structural and content features
 - names, domains, structure, values, ...
- Establish correspondences and rank according to quality
 - Errors are frequent and unavoidable
 - Works well for small schemas

Schema Matching Tools



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 40

The Problem

Given two (hierarchical) knowledge graphs D1, D2

Goal: find a 1:1 matching of nodes (entities, classes) that have the same or similar “meaning”

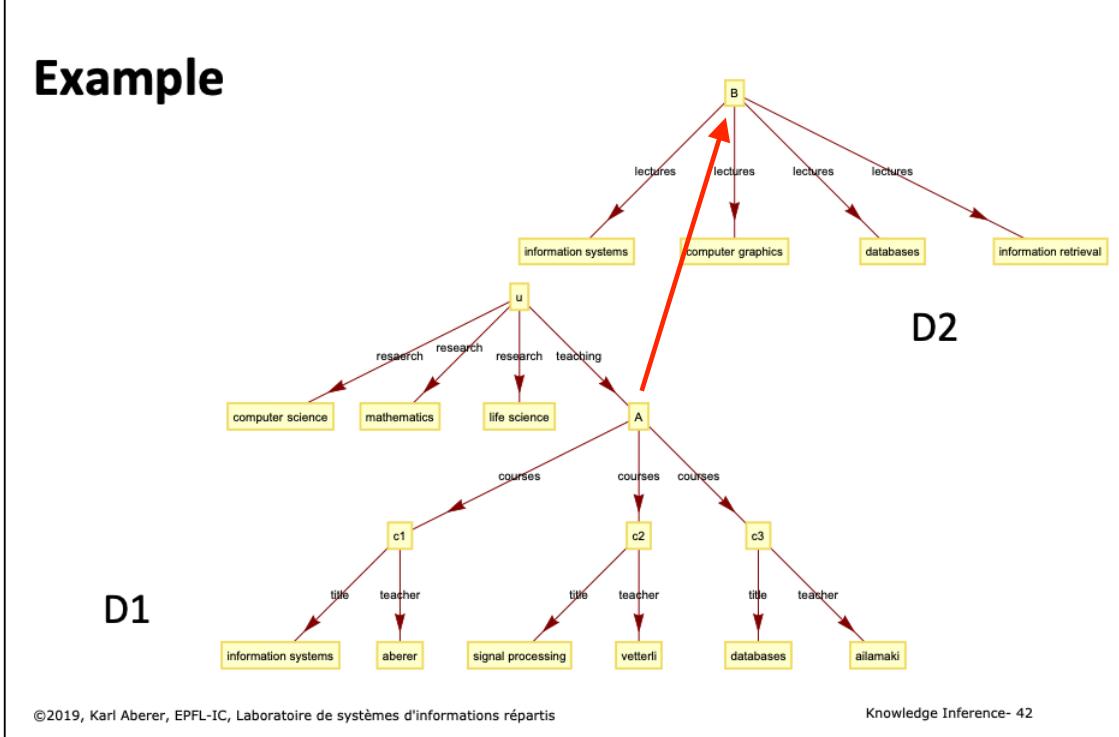
Assumption: two nodes (classes) are considered similar, if they have the same or similar instances (entities)

In the following we will study the problem of integrating heterogeneous databases respectively knowledge graphs. We assume that we want to integrate data that are available in two knowledge graphs that are semantically related to each other.

We will make a few assumptions:

- the knowledge graph is hierarchical. It could also be a taxonomy.
- Some entities represent classes. We are in particular interested in matching those.
- Entities that represent classes share similar kinds of instances, which are entities.

Example



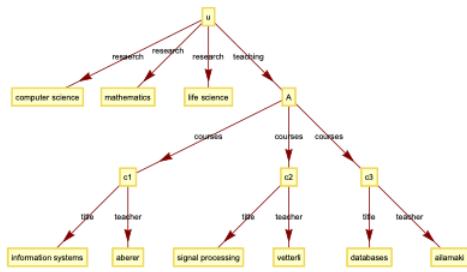
Let us assume that we have two example knowledge graphs. By inspecting the databases, a human can easily recognize that class A in graph D1 corresponds to class B in graph D2. The problem is how to perform the task of identifying this correspondence automatically.

Universe of a Database

Universe U is a finite set of possible instances

Example:

$U = \{u, \text{computer science}, \text{mathematics}, \text{life sciences}, A, c1, c2, c3, \text{information systems}, \text{aberer}, \text{signal processing}, \text{vetterli}, \text{databases}, \text{ailamaki}, \dots\}$



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 43

One basic approach to assess whether two classes are similar, is to measure their level of similarity at the content level, i.e., to test to which extent the elements (=instances) of the two classes overlap.

Similarity of Classes

A, B are classes, classes are subsets of U

Similarity measure (Jaccard similarity)

$$sim(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{P(x \in A \text{ and } x \in B)}{P(x \in A \text{ or } x \in B)} = \frac{P(A, B)}{P(A, B) + P(\bar{A}, B) + P(A, \bar{B})}$$

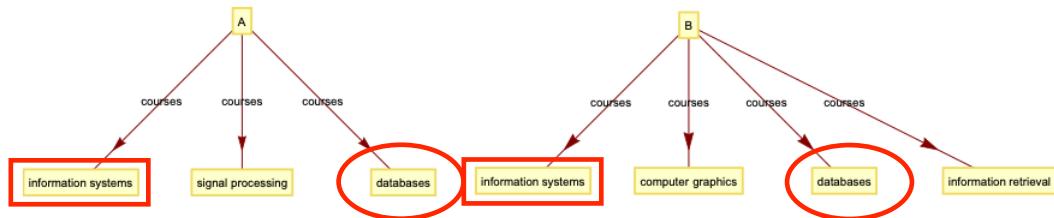
where $P(A, B) = P(x \in A \text{ and } x \in B)$ and

$P(A, \bar{B}) = P(x \in A \text{ and } x \notin B)$ etc

The Jaccard measure is a standard approach to measure such a set similarity.

Example

$$sim(A, B) = \frac{2}{5}$$



Note: instances of A are
“information systems”, “signal processing”, “databases”

In this example the two classes have 2 common elements, and the number of all elements in their union is 5, thus the Jaccard similarity is 2/5.

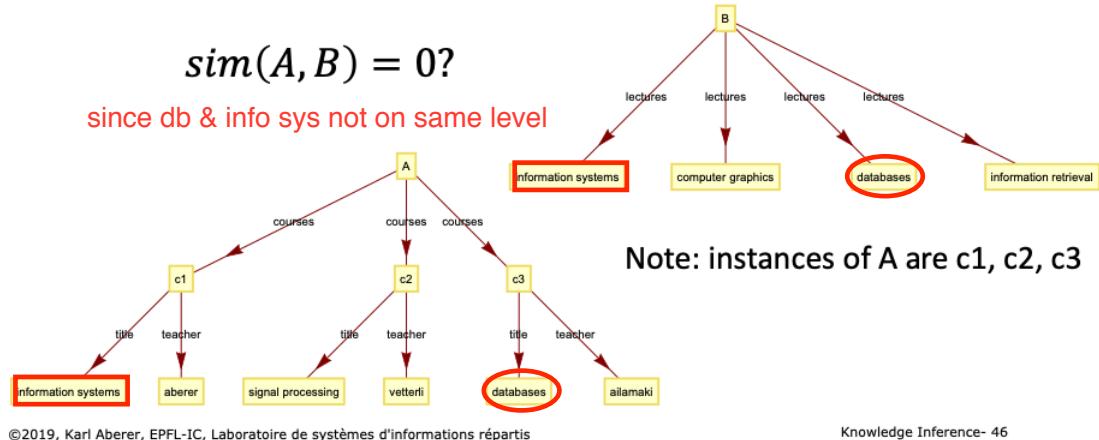
Problem 1

A similar to B?

Same information has structurally different representation

$$sim(A, B) = 0?$$

since db & info sys not on same level



In the previous example, the instances of the two classes have been leaf nodes in the knowledge graph. Therefore it was straightforward to compare their instances directly. In a general knowledge graph, the correspondences are not always that simple. In this example, we see that classes A and B correspond to each other, even if A has as instances complex data (nodes ci) and B atomic data (Strings). More complex features of classes are required. For the example shown, this is easy to resolve, by considering the atomic data found at leaf level for an inner node, such as A.

Problem 2

Two databases might have

- **Classes** with similar meaning (e.g. Course)
- **Different Instances for those classes** (e.g. different courses at different universities)

Intension is similar, but extension is different

Due to different database extensions and different naming conventions even classes that have a strong semantic similarity often have no common instances in two different databases.

Complex Features of Classes

Considering the instances (direct children) of classes is only an example of a simple class feature

Many complex features can be exploited (and have been exploited in schema matching)

- Names of attributes and relationships
- Structural relationships (data types)
- Distributional features (data values)
- Content features (e.g., text)

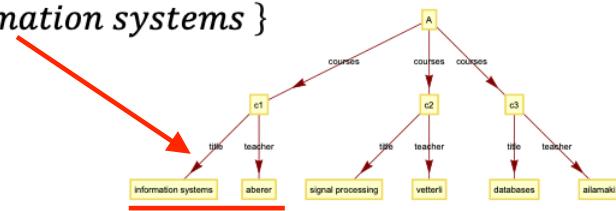
In database integration, many different features that can be associated with structural elements of a database have been considered for establishing semantic correspondences. For illustration, we will use in the following a simple feature: the set of atomic data values (strings) that can be reached from a node in the knowledge graph.

Content Feature of Classes

We consider as content feature, the set of terms associated with the paths of an node i to the leaves

$T_i = \{t_1, \dots, t_n\}$ with repeated occurrences (bag of words)

- $T_A = \{ aberer, information\ systems, signal\ processing, vetterli, \dots \}$
- $T_{c1} = \{ aberer, information\ systems \}$



©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 49

Another feature that can be exploited for this case would be the labels used along the paths.

Finding Corresponding Classes

If U_1 and U_2 are the universes of DB1 and DB2 we may assume

$$\underline{U_1 \cap U_2 = \emptyset}$$

Thus no way to compute $P(A,B)$ directly!

Given $i \in A$, the question is whether it would be likely that considering its features also $i \in B$, even if i is not part of U_2

Even if we have identified features that can help to spot correspondences between structurally different, but semantically equivalent elements of two databases, it might be the case that the coverage of a real-world aspect in two databases is different (e.g. courses in two different universities). Thus, directly comparing the features (e.g. the instances of a class) does not help to detect the correspondences.

Probabilistic Approach

We want to construct a function that gives the probability for a given instance i with feature T_i to belong to a class A

$$P(A|T_i) = P(T_i|A)P(A)P(d) \propto P(T_i|A)P(A) \text{ (Bayes law)}$$

$P(A|T_i)$ is a Naïve Bayes classifier to determine whether i belongs to A

Objective: determine $P(T_i|A)$ and $P(A)$

To tackle the problem, we will construct a model that determines (probabilistically) whether a given data instance with certain features is semantically related to a class.

Naïve Bayes Classifier

We know $P(A) = \frac{|A|}{|U_1|}$

Independence assumption: $P(T_i|A) = P(t_1|A) \dots P(t_n|A)$

With T_A being the bag of all terms occurring in all instances of A we have

$$P(t|A) = \frac{|t \in T_A|}{\sum_{t'} |t' \in T_A|}$$

Computing $P(A)$ is straightforward. We have just to determine the relative frequency of instances of A (how big A is) in the set of all possible instances.

For computing $P(T_i | A)$ we first make an independence assumption: we assume that different terms occurring in the instance i of a class A are independent of each other. In practice this is not the case, but is a generally accepted assumption for simplicity. Then we have to compute the probability of a single term t in class A to occur.

By computing $P(A)$ and $P(T_i | A)$ we obtain (up to a constant) the probability of for a (new) instance I to belong to class A.

Example

If T_{cn} is better matching in class U or A, figure out A is better

Classifier for class $A = \{c1, c2, c3\}$

$T_A = \{information\ systems, aberer, signal\ processing, \dots\}$

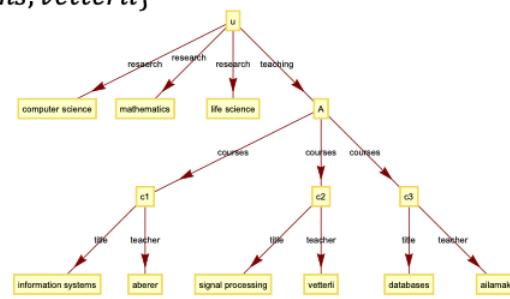
New instance cn : $T_{cn} = \{information\ systems, vetterli\}$

$$P(vetterli|A) = \frac{1}{6}, P(vetterli|U) = \frac{1}{9}$$

$$P(T_{cn}|A) = \frac{1}{36}, P(T_{cn}|U) = \frac{1}{81}$$

$$P(A) = \frac{3}{13}, P(U) = \frac{4}{13} \text{ (13 instances total)}$$

$$P(A|T_{cn}) \propto \frac{1}{36} \frac{3}{13}, P(U|T_{cn}) \propto \frac{1}{81} \frac{4}{13}$$



Thus instance cn is considered to correspond more likely to A than to U

We show here a complete example of how for a new instance cn , the most likely corresponding class can be determined in the database. The Naïve Bayes classifier assigns a higher probability for cn to belong to A than to U (which coincides with our intuition).

Computing Similarity between Classes A and B

1. Take all instances of U_1 and train a classifier to decide whether an instance belongs to A or not
2. Select all instances in U_2 belonging to B : U_2^B
3. Apply the classifier trained with U_1 to all instances in U_2^B to produce set U_2^{AB}
4. Do the same with roles of A and B exchanged
5. Compute $P(A, B) = \frac{|U_1^{AB}| + |U_2^{AB}|}{|U_1| + |U_2|}$
6. Compute similarly $P(A, \bar{B})$ etc
7. Then compute $sim(A, B) = \frac{P(A, B)}{(P(A, B) + P(\bar{A}, B) + P(A, \bar{B}))}$

©2019, Karl Aberer, EPFL-IC, Laboratoire de systèmes d'informations répartis

Knowledge Inference- 54

For computing Jaccard similarity between two classes A and B we need to know how many elements are contained in the intersection of A and B and the union of A and B . To that end we can first separate in each of the two databases the elements that belong to the class present in the database (e.g. A in database D_1 with universe U_1), and then apply the classifier learnt from the other database, whether an element in each of those two sets belongs also to the other class, or not. We consider U_2^{AB} as the set of elements of B that are likely to belong to A , if they were part of database D_1 . In this way we determine (or better estimate) the sizes of the potential intersection of A and B and the union of A and B and can based on this compute an estimate for Jaccard similarity.

Node Mapping

With the similarity values, alternative class mappings are possible

Naïve approach

- Order matchings by probability
- Choose the most probable matching and produce a mapping among the classes
- Remove the mapped classes
- Choose the next most probable matching and repeat

Drawback

- Obvious consistency constraints may be violated, e.g., if all children of a class are mapped, then also their parent class should be mapped
- Solution: more sophisticated mapping algorithms that incorporate such general and domain-specific constraints (research problem)

Once the similarity values are computed the process of matching for establishing correspondences is completed. Now the problem of mapping remains, i.e. determine which classes should be considered as equivalent. This is in general not uniquely determined, since a class can be similar to several others. On the other hand we can assume that one class should only be mapped to exactly another class in the other database. This imposes a constraint on the mapping.

A simple approach to establish a mapping that observes this constraint is by proceeding in a greedy fashion. First the best matches are used to produce mappings, then the the mapped classes are removed (to assure that each class is mapped to a unique other class) and then the process continues.

References

Course material based on

- Pershina, Maria, Yifan He, and Ralph Grishman. "Personalized page rank for named entity disambiguation." *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2015.
- Talukdar, Partha Pratim, and Koby Crammer. "New regularized algorithms for transductive learning." *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, Berlin, Heidelberg, 2009.
- Bordes, Antoine, et al. "Translating embeddings for modeling multi-relational data." *Advances in neural information processing systems*. 2013.
- Nguyen, Dat Quoc. "An overview of embedding models of entities and relationships for knowledge base completion." arXiv preprint arXiv:1703.08098 (2017).
- Doan, AnHai, et al. "Learning to map between ontologies on the semantic web." *Proceedings of the 11th international conference on World Wide Web*. ACM, 2002.