

From hyperlinks to semantic links

A Probabilistic Approach

ABSTRACT

In this paper, we propose an approach for detecting the implicit relations between 2 entities.

1. INTRODUCTION

Semantifying web documents is a huge task. Recently, more and more hyperlinked documents emerge. There are almost semi-structured information everywhere, but unfortunately the links lack an explanation of why they are connected, and why human may consider them related. The types of hyperlink are article to article, article to knowledge base (or encyclopedia-like databases). Each link has two ends: the original one (the web page/entity page contains it) and the oriented one (the web page/entity page it points to), and each end has 2 possible types, an anchor url or an anchor entity. Hence, in all there are 4 possible combinations, which consists the 4 type of links. An anchor entity can be named entities such as movies, products, people names, locations etc. In this case we explain the relation between the entity and the original article. As for an anchor url, an automatic target entity extraction process [2011] can be performed, and then it becomes another semantifying problem between entities problem.

In this paper we focus on the problem of semantifying the first type of hyperlinks since other kinds of links can finally be deduced to this problem.

We argue that the context of an entity in an article is informative and deserves a higher occurrence and can produce fresh and latest relation of an entity, which can be useful in updating the KB [reverb]. In our approach the relations are not necessarily to exist, however the cluster of the relation [relation clustering] it belongs to should be conceptually right. For instance, the relationship `artist of` will be correct in the tuple of `(Leonardo Da Vinci, artist of, Mona Lisa)` [`Human, artist of, Painting`] and will be never correct in the tuple of `(Automobile, artist of, Paint-`

`ing)` [e.g. `(BMW i8, artist of, Mona Lisa)`] With rich context and large knowledge base, we can easily derive the fresh context relations and the concept of each entity,

On the other hand, we consider the co-occurrence of the 2 entities, based on the assumption of important relationship will be observed in various of documents.

For these reasons, we propose a relation explanation method leveraging the concept and co-occurrence of an entity, to explain the relation between the target entity and the related entity, thus semantifying the hyperlink.

2. RELATED WORKS

2.1 Link entities to Database

Linking entities to database, especially to Wikipedia, has been widely studied. Entity linking to Wikipedia [10, 9, 5] exploits Wikipedia as thesaurus and link web documents to it. In our work, instead of linking entities to the corresponding one in KB, we extract the target entity [1] and explain the semantic relation of the entity towards target entity.

2.2 Relation Discovery

to be read [12, 11, 7]

Recently, different efforts are devoted to relation Discovery [3, 13, 8] are studied on graph based approaches and text based approaches [6] [reverb]. However, these approaches will be limited largely by the incompleteness of Knowledge Base, and cannot discover new type of relations. In our work, we focus on semantic relations of entity by leveraging the concepts of entity extracted from knowledge bases to link entities with a probability.

2.3 Information Extraction

Attribute acquisition methods Among domain-dependent approaches, we can mention approaches that focus on products. In this domain, attributes have been used to improve product search and recommendation [18, 22], but also to enable data mining [27]

Attribute retrieval provides another granularity in Web search. This can interest communities that propose a more focused access to information or communities that envision aggregating pieces of information such as aggregated search [19, 15]. Wong et al. [27] combine tags and textual features in a Conditional Random Fields model to learn attribute extraction rules, but they need a seed of relevant documents manually fed.

2.4 Short Text Conceptualization

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

3. FRAMEWORK

this section need to be re-written

In this section, we present the problem and our solution to it.

3.1 Problem Definition

First, we do conceptualization.

Next, Judge whether the 2 entities are conceptually same. Then, there are 2 cases of the CanBeExplained function:

- Explain 2 conceptually similar entity

Table 1: conceptually similar entity

entity	concept
Steve jobs	Person
Bill Gates	Person

- Explain 2 conceptually different entity

Table 2: Add caption

entity	concept
Mona Lisa	Painting
Renaissance	Period

Note that the concept here are not unique.

Last, We rank all the explanations in each step.

4. PROBABILITY RECALCULATION AFTER CONCEPTUALIZATION

Given an entity e , from Probase, we can acquire its concepts' set C and for each $c_i \in C$, the frequency $n(c_i, e)$ can be accordingly derived, which means how many times the e isA c_i pattern can be observed from the original corpus.

However, the concepts here has various forms as illustrated in Example 1. For our task, we only need relatively general concepts. The number of entities can be very large, but the number of top concepts and the relationship between them are limited, literally, we can find all the possible relationship between concepts instead of store all the long-tailed entities and their relations, which indicates the rationality of doing conceptualization.

EXAMPLE 1 (VARIOUS FORMS OF CONCEPTS). *Take the entity Mona Lisa as example, its concepts includes painting, famous painting, world's most famous painting, with corresponding frequency 33, 8, 1*

4.1 Problem Definition

Given Probase, for each entity e , we can derive $P(c_i|e)$, where $c_i \in C_{probase}$,

$$P(c_i|e) = \frac{n(c_i, e)}{n(e)}$$

We divide the $C_{Probase}$ into 2 parts, C_{simple} and C_{long} , where C_{simple} only contains one word and C_{long} are the rest.

C_{simple} are generated by the head modifier detection. The problem here is to recalculate the probability $P(\gamma_i|e)$ where $\gamma_i \in C_{simple}$, literally, we should contribute all the counts of C_{long} to C_{simple} .

4.2 Problem Solution

After head modifier detection, we have a set of $\gamma_i \in C_{simple}$, among all the $c_{long_j} \in C_{long}$, there are 2 cases in the probase determined by whether the c_{long_j} has an isA edge towards γ_i or not. The intuition of doing so is illustrated in the example 2:

EXAMPLE 2 (CONTRIBUTING LONG CONCEPTS). *Assume that Mona Lisa is a painting and Mona Lisa is a famous painting are observed respectively 33 times and 8 times from different documents, we will get the knowledge that Mona Lisa is a painting occurs 41 times instead of 33 times. There are less chance of occurring Famous painting is a painting so that there won't be necessarily an isA edge from famous painting to painting.*

Therefore, to calculate $P(\gamma_i|e)$, there are three cases:

1. e isA γ_i . The entity has has an isA edge towards one or more simple concept, which gives the original $P_{original}(\gamma_i|e)$
2. e isA c_{long} , c_{long} isA γ_i , In this case, we need to calculate the following equation

$$P(\gamma_i|e) = \sum_{c_{long}^* \in C_{long}} P(\gamma_i|c_{long}^*, e) \times P(c_{long}^*|e)$$

, where $P(c_{long}^*|e)$ can be obtained from Probase and

$$P(\gamma_i|c_{long}, e) = \frac{n(\gamma_i, c_{long}, e)}{n(\gamma_i, e)} \quad (1)$$

We assume that the occurrence of e does not affect $P(\gamma_i|c_{long})$ equivalently speaking, $P(\gamma_i|c_{long})$ is independent from e , thus Eq. 1 can be simplified

$$P(\gamma_i|c_{long}, e) = P_{probase}(\gamma_i|c_{long}) = \frac{n(\gamma_i, c_{long})}{n(\gamma_i)}$$

which can be obtained from Probase.

3. e isA c_{long} , c_{long} has no edge towards γ_i . The edge here refers to the isA relationship in Probase. Example 2 pointed out that there won't be necessarily an isA edge from famous painting(c_{long}) to painting(γ_i), however c_{long} obviously belongs to γ_i . In this case, since it's detected by the head modifier method, we assume

$$P_{head}(\gamma_i|c_{long_j}) = 1$$

Another reason why we do head modifier detection here is that even if the long concept c_{long} has an isA edge towards a certain concept γ_i' , it still sometimes not include the head concept of the long concept which is very plausible. The tradeoff of the 2 method is described in Example 3

Notice that the boundary between case 2 and case 3 are not strict, there are such edges that have low observation in

Example 3. So that if we consider them as a whole, we can derive:

$$P(\gamma_i|c_{long}) = \lambda P_{head}(\gamma_i|c_{long}) + (1 - \lambda) P_{probase}(\gamma_i|c_{long}) \quad (2)$$

where λ is a parameter **principle: related to plausibility, number of occurrence, varies for different c_{long} should it be derived from learning ?** since we assume $P_{head}(\gamma_i|c_{long})$ to be 1, Eq. 2 is simplified to:

$$P(\gamma_i|c_{long}) = \lambda + (1 - \lambda) P_{probase}(\gamma_i|c_{long})$$

EXAMPLE 3 (HEAD CONCEPTS VS ORIGINAL CONCEPTS). Again take **famous painting** as example, whose concepts **image**, **treasure** are reasonable but implausible, since their occurrence are twice and once respectively. However, the most plausible concept **painting** is not among the concepts. On the other hand, there exists several concepts that have also reasonable. For example **topaz**(a kind of yellow gemstone) has the concept **precious stones**, and **precious stones** has an edge towards **material** which is reasonable.

Finally $P(\gamma_i|e)$ is calculated using the following equation:

$$P(\gamma_i|e) = P_{original}(\gamma_i|e) + \sum_{c_{long}^* \in C_{long}} [\lambda_i^* + (1 - \lambda_i^*) P(\gamma_i|c_{long}^*)] \times P(c_{long}^*|e) \quad (3)$$

The process of calculation is illustrated in the example 4

EXAMPLE 4 (CALCULATING $P(\gamma_i|e)$). As illustrated in Fig. 4.2, the process of calculating the typicality a concept is as follows, where **painting** is γ_i and **Mona Lisa** is e . Then $P(\text{painting}|\text{MonaLisa})$ consists of 2 parts, the direct edge $P_{original}(\gamma_i|e) = 0.23$, and the second part

$$\sum_{c_{long}^* \in C_{long}} [\lambda_i^* + (\alpha_i^*) P(\gamma_i|c_{long}^*)] \times P(c_{long}^*|e)$$

($\alpha_i^* + \gamma_i^* = 1$) Thus we get

$$P = 0.007 \times \lambda_{i2} + 0.05 \times \lambda_{i1} + 0.04 \times (\lambda_{i3} + 0.65\alpha_{i3})$$

For **piece**, it is the similar process. The relation here is only part of the whole graph.

We consider only 2 layers of isA relationship for 2 reasons. The first one is that more layers will lead to noisy concepts such as **issue**, **factor**, **element**, which are concepts for almost everything, Secondly, discussing the transitive relation between concepts is beyond the scope of this paper.

5. FIND ALIAS FOR ATTRIBUTES

For a pair (Sherlock holmes, United Kindom), **country** is a merely-ok attribute, on the contrary, **residence**, **deathPlace** are better since they are more specific and more seemingly plausible to be an attribute. We argue that for each pair of entity, there is a selectional preference for attribute.

5.1 Problem Definition

Given a set of concept pairs (γ_1, γ_2) , where $\gamma_1 \in C_1$ and $\gamma_2 \in C_2$, we want to find a set of attributes A , where for each $a \in A$: we can form a (γ_1, a, γ_2) pair which best describe the relationship between γ_1 and γ_2 .

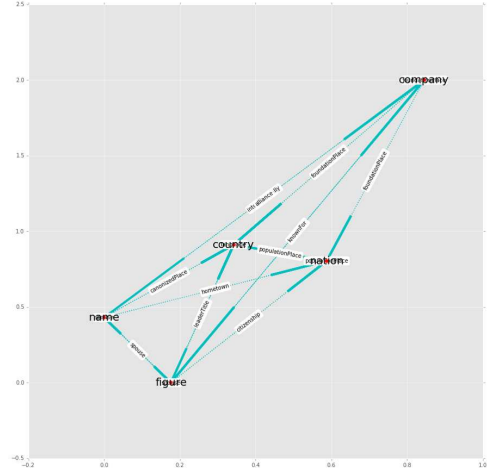


Figure 2: Subgraph of Entity Attribute Graph

5.2 Problem Solution

5.2.1 Entity Attribute Graph Construction

For any $(entity_1, attribute, entity_2)$ tuple, later denoted as (e_1, a, e_2) , where e_1 and e_2 are also referred to as **domain** and **range** of the attribute. We can conceptualize e_1 and e_2 using the method in section 4, and get a set of concept C_1, C_2 , accompanied with a set of probabilities $P(\gamma_1|e_{1i}), P(\gamma_2|e_{2j})$, where $\gamma_1 \in C_1, \gamma_2 \in C_2$.

To construct the Entity Attribute Graph, we only need topK concepts to form (γ_1, γ_2) pair, K through **case study** is around 5, so we here set $K=10$.

Thus for any attribute a , given a pair of entity (e_{1i}, e_{2j}) , we can define: **should i use joint ratio here?**

$$P_{(e_{1i}, e_{2j})}((\gamma_1, \gamma_2)|a) = P_{before}(\gamma_1|a) \times P_{after}(\gamma_2|a) = P(\gamma_1|e_{1i})P(e_{1i}|a) \times P(\gamma_2|e_{2j})P(e_{2j}|a) \quad (4)$$

where we use $P_{(e_{1i}, e_{2j})}((\gamma_1, \gamma_2)|a)$ to denote observing a single pair (e_{1i}, e_{2j}) , how likely is a combination of (γ_1, a, γ_2) to occur.

Consequently,

$$P((\gamma_1, \gamma_2)|a) = \sum_{e_{1i} \in E_1, e_{2j} \in E_2} P_{(e_{1i}, e_{2j})}((\gamma_1, \gamma_2)|a) \quad (5)$$

where E_1, E_2 denoting the whole set of domain entity and range entity, The $P(e_{1i}|a)$ and $P(e_{2j}|a)$ here has only 2 values 1 and 0, depending on whether e_{1i} occurs before a or e_{2j} occurs after a . Apparently, only (e_{1i}, a, e_{2j}) occurs will give the equation a non-zero value, therefore, Eq. 5 is finally equal to Eq. 6.

$$P((\gamma_1, \gamma_2)|a) = \sum_{(e_{1i}, a, e_{2j}) \in KB} P_{(e_{1i}, e_{2j})}((\gamma_1, \gamma_2)|a) = \sum_{(e_{1i}, a, e_{2j}) \in KB} P(\gamma_1|e_{1i}) \times P(\gamma_2|e_{2j}) \quad (6)$$

The process of calculating is demonstrated in Example. 5

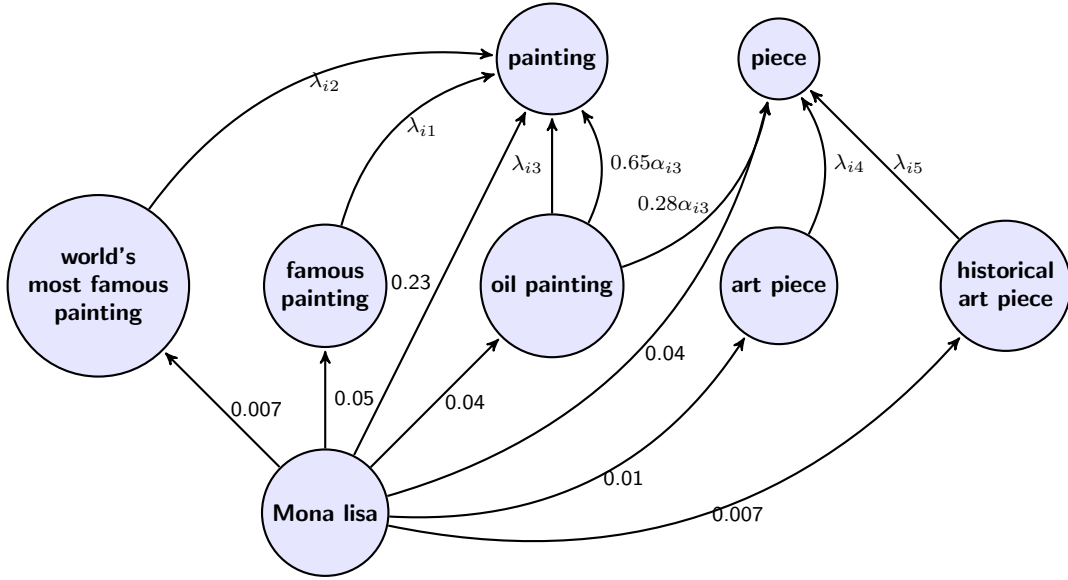


Figure 1: calculating $P(\gamma_i | \text{MonaLisa})$

EXAMPLE 5 (CALCULATING $P((\gamma_1, \gamma_2) | a)$). As illustrated in Fig. 3, the process of calculating $P((\gamma_1, \gamma_2) | a)$ is as follows

insert a graph

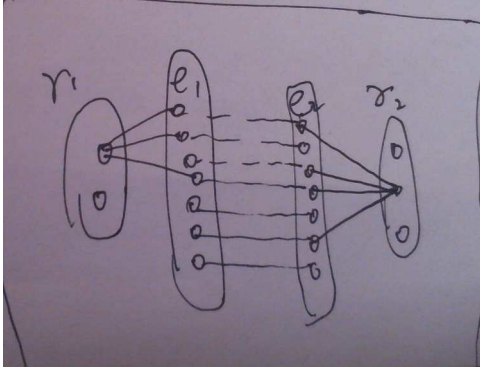


Figure 3: Calculating $P((\gamma_1, \gamma_2) | a)$ for

To Construct the Entity Attribute Graph, we calculate $P((\gamma_1, \gamma_2) | a)$ for each attribute.

Note that we only consider the attributes whose range is an entity, and ignore those numerical values or date-and-time values such as (*MonaLisa*, *Year*, 1503).

For each (γ_1, a, γ_2) tuple, we can calculate $P((\gamma_1, \gamma_2) | a)$ for each

5.3 For multiple hops

So far, we have tackled with the relations and generated the edges in the relationship graph. This problem is similar to **hierarchy ranking problems in a directed graph** [4]. Originally, it was a minimum feedback arc set problem on a weighted network which is a classic NP-hard problem [2]. A few approaches [14] have been proposed on

unweighted directed graphs, for weighted graphs, the extended agony[hierarchies in directed network(unpublished kdd15)] algorithm can be utilized to generate hierarchy results in this approach the K (number of hierarchies) is fixed, maybe we can make it adaptive to data here?

In this section, we first formulate the problem of finding semantic link into a maximum flow problem on the concept network with multiple-sources and multiple-sinks, and then, we cut out the subgraph and perform **improved agony** to derive the concept of the middle entities. Last we use co-occurrence to verify the validness of the relation.

5.3.1 Improved agony

5.4 Find the best alias

We then Use an arg max model use KL divergence? to minimize D_{KL} to solve the problem.

Given (e_1, e_2) , our goal is to find the best attribute for it. We denote it as:

$$\arg \max P((e_1, e_2) | a)$$

where

$$P((e_1, e_2) | a) =$$

6. EXPERIMENT

6.1 Head Concept Vs Original Concept

6.2 Find alias

6.2.1 compare

Compare $P((\gamma_{1i}, \gamma_{2i} | a))P(\gamma_{1i} | a) \times P(\gamma_{2i} | a)$

6.2.2 Sense Disambiguation

We can solve the problem of sense disambiguation problem well by applying this method since there are many entities belongs to the same concept and we only consider topK (γ_1, γ_2) pairs that has high typicality $P((\gamma_1, \gamma_2) | a)$, so that

the weird (γ_1, γ_2) patterns as manifest in Example. 6 can be easily filtered.

cut the figure smaller



Figure 4: (γ_1, γ_2) plot for attribute Manufacturer

EXAMPLE 6 (SENSE DISAMBIGUATION). Consider the following (e_1, a, e_2) tuple (iphone, manufacturer, apple). Suppose it is our query, where apple's sense can either be a kind of fruit or a company. Fig. 4 is a heatmap for all the concepts pairs (γ_1, γ_2) of attributes manufacturer. The horizontal axis represents the e_1 and the vertical axis stands for e_2 . The darker the blue is, the higher typicality it will be. In Fig. 4, We can observe that the top concepts of e_2 in the heatmap are company, manufacturer, ... and top 10 pairs also does not include fruit. The intuition for this is that there exists thousands of (e_1, a, e_2) tuple such as (BMW_Z4, manufacturer, BMW), (PlayStation_4, manufacturer, Sony) other than (iphone, manufacturer, apple) tuple, which results in a reasonable distribution.

6.3 Selectional Preference

6.4 Evaluation

7. CONCLUSION

8. REFERENCES

- [1] N. Dalvi, R. Kumar, and M. Soliman. Automatic wrappers for large scale web extraction. *Proceedings of the VLDB Endowment*, 4(4):219–230, 2011.
- [2] I. Dinur and S. Safra. On the hardness of approximating minimum vertex cover. *Annals of mathematics*, pages 439–485, 2005.
- [3] L. Fang, A. D. Sarma, C. Yu, and P. Bohannon. Rex: explaining relationships between entity pairs. *Proceedings of the VLDB Endowment*, 5(3):241–252, 2011.
- [4] M. Gupte, P. Shankar, J. Li, S. Muthukrishnan, and L. Iftode. Finding hierarchy in directed online social networks. In *Proceedings of the 20th international conference on World wide web*, pages 557–566. ACM, 2011.
- [5] X. Han, L. Sun, and J. Zhao. Collective entity linking in web text: a graph-based method. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*, pages 765–774. ACM, 2011.
- [6] T. Hasegawa, S. Sekine, and R. Grishman. Discovering relations among named entities from large corpora. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 415. Association for Computational Linguistics, 2004.
- [7] N. Konstantinova. Review of relation extraction methods: What is new out there? In *Analysis of Images, Social Networks and Texts*, pages 15–28. Springer, 2014.
- [8] G. Luo, C. Tang, and Y.-l. Tian. Answering relationship queries on the web. In *Proceedings of the 16th international conference on World Wide Web*, pages 561–570. ACM, 2007.
- [9] R. Mihalcea and A. Csomai. Wikify!: linking documents to encyclopedic knowledge. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 233–242. ACM, 2007.
- [10] D. Milne and I. H. Witten. Learning to link with wikipedia. In *Proceedings of the 17th ACM conference on Information and knowledge management*, pages 509–518. ACM, 2008.
- [11] N. Nakashole, G. Weikum, and F. Suchanek. Discovering and exploring relations on the web. *Proceedings of the VLDB Endowment*, 5(12):1982–1985, 2012.
- [12] N. Nakashole, G. Weikum, and F. Suchanek. Discovering semantic relations from the web and organizing them with patty. *ACM SIGMOD Record*, 42(2):29–34, 2013.
- [13] D. Shahaf and C. Guestrin. Connecting the dots between news articles. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 623–632. ACM, 2010.
- [14] N. Tatti. Faster way to agony. In *Machine Learning and Knowledge Discovery in Databases*, pages 163–178. Springer, 2014.