

# SEMashup: Making Use of Linked Data for Generating Enhanced Snippets

Mazen Alsarem<sup>1</sup>, Pierre-Édouard Portier<sup>1</sup>, Sylvie Calabretto<sup>1</sup>, and Harald Kosch<sup>2</sup>

Université de Lyon, CNRS

<sup>1</sup>INSA-Lyon, LIRIS, UMR5205, F-69621, France

<sup>2</sup>Universität Passau, Innstr. 43, 94032 Passau, Germany

{firstname.lastname}@insa-lyon.fr

{firstname.lastname}@uni-passau.de

**Abstract.** We enhance an existing search engine’s snippet (i.e. excerpt from a web page determined at query-time in order to efficiently express how the web page may be relevant to the query) with linked data (LD) in order to highlight non trivial relationships between the information need of the user and LD resources related to the result page. Given a query, we first retrieve the top ranked web pages from the search engine results page (SERP). For each result, we build a RDF graph by combining DBpedia Spotlight [7] and a RDF endpoint connected to the DBpedia dataset. To each resource of the graph we associate the text of its DBpedia’s abstract. Given the initial result from the SERP and this textually enhanced graph, we introduce an iterative co-clustering approach in order to discover additional qualified relationships between the resources. Then, we apply a first PARAFAC tensor decomposition [6] to the graph in order to select the most promising nodes for a 1-hop extension from a DBpedia SPARQL endpoint. Finally, we compute a second tensor decomposition for finding hubs and authorities for the most relevant types of predicates. From this graph analysis, we build the enhanced snippet.

**Keywords:** Linked Data, Information Retrieval, Snippets, Co-Clustering, Tensor Decomposition

## 1 Introduction

A snippet is an excerpt from a web page determined at query-time and used to express, as efficiently as possible, the way a web page may be relevant to the query. By making explicit the reason why a document seems to meet, at least partially, the information need of a user, the snippet alleviates the semantic mismatch between relevance and aboutness.

In this work, we show that the LoD graph can be combined with analysis of the original textual content of a search engine’s results to create efficient enhanced snippets. Contrary to other approaches (such as [5]), we don’t rely exclusively on explicit semantic annotations (e.g. RDFa, ...) but we are able

to enhance any result even if the corresponding web page doesn't contain annotations. To attain this objective, we combine in a web mashup three existing services: the Google search engine API, the DBpedia Spotlight web service<sup>1</sup>, and a DBpedia SPARQL endpoint<sup>2</sup>.

The paper is organized as follows: we start with a brief state of the art for the domains of snippets' enhancement, and snippet generation for corpus of RDF documents (Section 2). Then we describe our approach (Section 3) before concluding (Section 4).

## 2 Related Works

Many approaches tried to enhance SERPs' snippets, but none of them seem to have used the LoD graph. First, Haas *et al.* [5] employed structured metadata (i.e. RDFa and several microformats) and information extraction techniques (i.e. handwritten or machine-learned wrappers designed for the top host names e.g., `en.wikipedia.org`, `youtube.com`,...) to enhance the SERP with multimedia elements, key-value pairs and interactive features. They chose not to use the LoD graph to avoid the problem of the transfer of trust between the Web of documents and the Web of Data.

Secondly, Google Rich Snippet (GRS) [12] is a similar initiative that relies exclusively on structured metadata authored by the web pages' publishers.

Thirdly, we should mention the works of Ge *et al.* [4], and Penin *et al.* [10] focused on the generation of snippets for ontology search, and also the work of Bai *et al.* [2] about the generation of snippets for native RDF documents.

In summation, we agree with Ge *et al.* [4] that the main benefit of possessing highly structured data from an RDF graph is the possibility to find non-trivial relations among the query terms themselves, and also between the query terms and the main concepts of the document. Moreover, we agree with Penin *et al.* [10] and Bai *et al.* [2] about the necessity to design a ranking algorithm for RDF statements that considers both the structure of the RDF graph and lexical properties of the textual data. However, we find ourselves in an inverted situation with genuine textual data from classical web pages, and RDF graphs generated from these web pages by using the LoD graph. We will now show that by combining textual and graphical features, we can make use of the LoD graph and still provide users with high added value snippets by avoiding introducing noise.

## 3 Proposal

Our main purpose is to highlight on a practical and convincing use case the benefits of a conjoint use of the web of documents and the web of data. Thus, for each result of the SERP, we want to build a RDF graph and combine this graph with a textual analysis of the document in order to obtain features from

<sup>1</sup> <http://spotlight.dbpedia.org/>

<sup>2</sup> <http://dbpedia.org/sparql>

which to build an enhanced snippet. Meanwhile, our main concern is to limit the amount of noise introduced by this process. Our approach involves four steps that will be described in subsequent subsections.

Firstly, we use DBpedia Spotlight [7] to extract LoD resources from the content of the SERP result. Next, we use a SPARQL endpoint connected to the DBpedia dataset to introduce RDF statements between the resources we just found. At that time we have a first graph associated to the SERP result. We should note that if in the current implementation we used the Google search API, the DBpedia Spotlight text annotation service, and a DBpedia SPARQL endpoint, our system is modular enough to easily adapt to other web search APIs (e.g. BING<sup>3</sup>, Yahoo<sup>4</sup>,...), other annotation services (e.g. AlchemyAPI<sup>5</sup>,...), and other knowledge bases (e.g. ckan<sup>6</sup>, GeoNames<sup>7</sup>, umbel<sup>8</sup>, Ontobee<sup>9</sup>,...).

Secondly, we want to extend the graph in order to introduce resources relevant to the information need of the user but not present in the result. However, a systematic 1-hop extension of all the nodes would produce a graph too large and too noisy. To address this problem we introduce a new multi-step unsupervised co-clustering algorithm (see Section 3.1).

In the third step of our approach, we propose to combine the graph and the result of the clustering in order to decide which nodes of the graph are suitable for a 1-hop extension. To do this, we use a 3-way tensor to represent both the graph and the relations found thanks to the multi-step unsupervised co-clustering. Next, by interpreting the results of a PARAFAC tensor decomposition [6], we decide which nodes to extend based on their importance as hubs or authorities with respect to each *type of relation* (i.e. a LoD predicate or a set of terms associated with the resources found in a cluster). See Section 3.2 for details.

Finally, we do a new tensor decomposition on the extended graph in order to group the triples of the graph in topical factors that we can use for the construction of an enhanced snippet (see Section 3.3).

### 3.1 Multi-step Unsupervised Co-clustering

Since the graph obtained from entity extraction alone is poorly connected, our main objective is to discover relevant relations between the resources based on an analysis of the textual data associated with the resources. Indeed, we need to build a more informed graph in order to efficiently select, in a subsequent step, a relevant subset of nodes to extend. Thus, to each node of the graph obtained from a SERP result — by entity extraction and SPARQL query of the LoD — we associate the text of its DBpedia abstract and the fragments of text from the

<sup>3</sup> <http://datamarket.azure.com/dataset/bing/search>

<sup>4</sup> <https://developer.yahoo.com/boss/search/>

<sup>5</sup> <http://www.alchemyapi.com/>

<sup>6</sup> <http://ckan.org/>

<sup>7</sup> <http://www.geonames.org/>

<sup>8</sup> <http://umbel.org/>

<sup>9</sup> <http://www.ontobee.org/>

original document centered around the places where the entity was discovered. Following this first step, we have several options. For example, we could simply compute the cosine distances between the abstracts of the nodes. However, with this approach we would introduce a unique and unqualified proximity relation between resources (contrary to the typed relations obtained from the LoD). This point is crucial since we seek to explain to the user, as precisely as possible, how her information need is related to a SERP result.

Thus, in order to address the weaknesses of a naive lexical approach, we introduce a new multi-step unsupervised co-clustering algorithm. First, we remove the stop-words from the texts associated to each node of the graph and we build a resource-term matrix (denoted  $A$ ). Then, our multi-step co-clustering process is a combination of the algorithm C of Listing 1 and of the algorithm O of Listing 2. Algorithm C is a classical co-clustering algorithm based on SVD [8] that we adapted for the X-means [9] algorithm to avoid having to estimate the number of clusters. We use a co-clustering approach because the relation between resources of a given cluster will be qualified by the terms present in this cluster.

However, the first time Algorithm C is executed, the internal quality of the found clusters is very poor — with a mean silhouette index [11] close to zero. Our approach is to decrease the size of the features' space (i.e. the number of terms and resources) until we find clusters of a good enough quality. When we find clusters of good internal quality and containing only terms (Listing 2 line 8) we remove those terms (Listing 2 line 13) and repeat the co-clustering (Listing 2 line 28). The underlying rationale for this strategy is that the terms that lie well in their own clusters but do not associate with resources will not help to discover and explain relationships between resources. However, we observed that this dimension reduction mechanism is insufficient since there can be clusterings with a poor average silhouette and with no clusters only made of terms. Therefore we also remove from the features' space the terms and the resources of the small clusters with a weak silhouette. Furthermore, when we can find no clusters only made of terms, no small clusters with a low silhouette, and when the average silhouette remains low, we heuristically re-apply Algorithm O on the clusters with a poor silhouette (Listing 2 line 39).

---

**Listing 1** Algorithm C
 

---

```

1 // given a resource-term matrix  $A$  of size  $m \times n$ 
2 // with  $m \neq 0$  and  $n \neq 0$ ,
3 // returns a co-clustering of the resources and the terms.
4 BEGIN
5    $A \leftarrow \text{tf-idf-weighting}(A)$ ;
6    $(U, \Sigma, V^T) \leftarrow \text{svd}(A)$ ;
7    $(U_k, \Sigma_k, V_k^T) \leftarrow \text{rank-k-approximation}(U, \Sigma, V^T)$ ;
8    $B \triangleq (b_{ij} \mid b_{ij} = (U_k \Sigma_k)_{ij} \text{ IF } i \in 1..m$ 
9        $\mid b_{ij} = (V_k \Sigma_k^T)_{(i-m)j} \text{ IF } i \in m+1..m+n$ ;
10  (mean-silhouette, clusters:SetOf((silhouette, terms-ids, >
11  resources-ids)))  $\leftarrow \text{X-Means}(B)$  with nb-max-clusters= $k$ 
12 END
```

---

**Listing 2** Algorithm O

---

```

1 // given a resource-term matrix  $A$ , recursively applies
2 // co-clustering Algorithm C in order to
3 // find good clusters described by relevant terms.
4 GLOBAL VAR results initialized at creation with the value  $\emptyset$ ;
5 BEGIN
6 (mean-silhouette, clusters)  $\leftarrow$  AlgoC( $A$ );
7 cs1  $\leftarrow$  {c | c  $\in$  clusters,  $\triangleright$ 
8 c.silhouette > s-min, c.resources-ids =  $\emptyset$ };
9 cs2  $\leftarrow$  {c | c  $\in$  clusters, |c.resources-ids| < c-min};
10 VAR poor-small-clusters  $\leftarrow$  false;
11 IF | cs1  $\neq$   $\emptyset$   $\rightarrow$ 
12     FOR c in cs1  $\rightarrow$  remove columns of  $A$  with ids in  $\triangleright$ 
13         c.terms-ids
14     ROF;
15 | cs2  $\neq$   $\emptyset$   $\rightarrow$ 
16     IF | mean-silhouette > mean-s-min  $\rightarrow$ 
17         SKIP;
18     | otherwise  $\rightarrow$ 
19         poor-small-clusters  $\leftarrow$  true;
20         FOR c in cs2  $\rightarrow$  remove columns of  $A$  with ids in  $\triangleright$ 
21             c.terms-ids;
22             remove rows of  $A$  with ids in c.resources-ids;
23         ROF;
24     FI;
25 | otherwise  $\rightarrow$ 
26     SKIP;
27 FI;
28 IF | (cs1  $\neq$   $\emptyset$   $\vee$  poor-small-clusters)  $\wedge$   $\neg$ null( $A$ )  $\rightarrow$  AlgoO( $A$ );
29 | (cs1  $\neq$   $\emptyset$   $\vee$  poor-small-clusters)  $\wedge$  null( $A$ )  $\rightarrow$  RETURN;
30 | (cs1 =  $\emptyset$   $\wedge$   $\neg$ poor-small-clusters)  $\wedge$  cs2  $\neq$   $\emptyset$   $\rightarrow$ 
31     FOR c in {c | c  $\in$  clusters, c.silhouette > s-min}  $\rightarrow$ 
32         addResult(c);
33     ROF;
34 | (cs1 =  $\emptyset$   $\wedge$   $\neg$ poor-small-clusters)  $\wedge$  cs2 =  $\emptyset$   $\rightarrow$ 
35     IF | mean-silhouette > mean-s-min  $\rightarrow$ 
36         FOR c in {c | c  $\in$  clusters, c.silhouette > s-min}  $\rightarrow$ 
37             addResult(c);
38         ROF;
39     | otherwise  $\rightarrow$ 
40         FOR c in clusters  $\rightarrow$ 
41             IF | c.silhouette > s-min  $\rightarrow$ 
42                 addResult(c);
43             | otherwise  $\rightarrow$  applyAlgoO(c);
44         FI;
45     ROF;
46     FI;
47 FI;
48 END

```

```

49
50 PROC addResult(c) →
51 results ← results ∪ {c};
52 IF |c.resources-ids| > c-max → applyAlgoO(c);
53   | otherwise → RETURN;
54 CORP
55
56 PROC applyAlgoO(c) →
57 B ← A;
58 remove columns of B with ids not in ▷
59 c.terms-ids;
60 remove rows of B with ids not in ▷
61 c.resources-ids;
62 IF |¬null(B) → AlgoO(B);
63   | null(B) → RETURN;
64 CORP

```

---

### 3.2 Graph Extension by Tensor Decomposition

We would not benefit much from the LoD if we didn't extend the graph in order to find facts not present in the original document. However, we must be careful not to introduce noise during the extension. To do this, we propose to guide this process with the knowledge obtained from the previous multi-step co-clustering. Therefore, we have to combine the connectivity information of the graph with the labeled clusters. We propose to represent these two kinds of information in one structure: a 3-way tensor (denoted  $\mathcal{T}$ ).

As in [3], the three modes of the tensor are associated respectively to the subject, the object and the predicate of the triples that constitute the graph. Thus, for each predicate, there is one horizontal slice that represents the adjacency matrix of the subgraph only made of the triples that link resources thanks to that predicate. We propose to add a horizontal slice for each cluster returned by the multi-step co-clustering algorithm. Such a slice represents the clique of all the resources present in the cluster. To each edge of the clique, we assign a weight linearly proportional to the silhouette index of the cluster. This coefficient of proportionality is chosen so as to make all the slices of the tensor comparable.

In order to select the nodes to extend, we start by applying a PARAFAC [6] tensor decomposition algorithm to the tensor. This decomposition computes a representation of the tensor as a sum of rank-one tensor (a rank-one three-way tensor is the outer product of three vectors), i.e.,

$$\mathcal{T} = \sum_{r=1}^R \mathbf{s}_r \circ \mathbf{o}_r \circ \mathbf{p}_r. \quad (1)$$

If there are  $n_r$  resources and  $n_p$  predicates, the lengths of each  $\mathbf{s}_r$ ,  $\mathbf{o}_r$  and  $\mathbf{p}_r$  vectors are respectively  $n_r$ ,  $n_r$  and  $n_p$ . The components of the vectors  $\mathbf{s}_r$  and  $\mathbf{o}_r$

represent, for the number  $r$  factor, the importance of each resource as it plays respectively the part of subject and object in relations involving the high-scored predicates of  $\mathbf{p}_r$ . Our strategy is to select for each factor  $r$ , with at least one high-scored predicate found in  $\mathbf{p}_r$ , the high scored resources of  $\mathbf{s}_r$  and  $\mathbf{o}_r$  as the nodes to extend (i.e. we query a DBpedia SPARQL endpoint to find new triples for which these nodes are subject).

We should also note that there is no effective way to compute the number of factors for a tensor decomposition. In [3] the authors used an implementation of the CORCONDIA heuristic-based algorithm [1] to find a suitable number of factors. However, although we generate tensors in a way similar to [3], we discovered that the CORCONDIA heuristic didn't apply well on our tensors. Thus, we found by experiment that, with the data we generate, the optimal number of factors is close to 10. Therefore, we test multiple decompositions with a number of factors varying around 10 and we keep the decomposition that offers the best fit (i.e. for which the recomposed tensor is closer to the original one).

### 3.3 Snippet Construction

We apply a new tensor decomposition on the extended graph and we select only the factors for which there are values in  $\mathbf{p}_r$  larger than an experimentally fixed threshold (in our case 0.1). For those factors, we then select the resources of  $\mathbf{s}_r$  and  $\mathbf{o}_r$  with a weight greater than this same threshold. Thus, for each factor that describes an important sub-graph (i.e. for which there exist high-scored predicates) we isolate the main-resources that act either as hubs or authorities.

To each main-resource we associate a set of explanatory triples and we select a sentence of the SERP's result for its capacity to contextualize the main-resource. The triples associated to a main-resource are trivially derived from the factors where this resource appears either as an important subject or as an important object: we select the triples for which at least the subject and the predicate or the predicate and the object have good scores. The sentence associated to a main-resource is found by ranking the sentences of the document according to seven factors: (1) the presence of a DBpedia Spotlight annotation for the main-resource, (2) the presence of a DBpedia Spotlight annotation for a resource that appears with a high score in a group where the main-resource also has a high score, (3) the presence in the text of the local name of a predicate that belongs to one of the main-resource's groups, (4) the presence of query's terms, (5) the presence in the text of the local name of the resource, (6) the presence of a DBpedia Spotlight annotation for a resource that appears in the query, and (7) the presence of a DBpedia Spotlight annotation for a resource adjacent to the main-resource in the graph.

Finally, we apply a similar ranking strategy to find a sentence that is both close to the query and a place of co-occurrence for as many main-resources as possible. This sentence is shown in addition to the original excerpt chosen by the search engine. Figure 1 gives an example of an enhanced snippet, and Figure 2

shows the effect of a user's click on a main-resource. A more detailed version of this work, screenshots and a demonstration are available online<sup>1011</sup>.

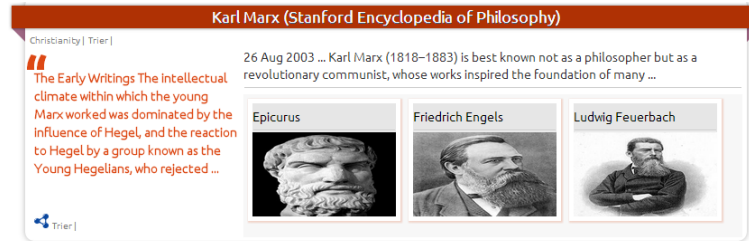


Fig. 1: Screenshot of an enhanced snippet for the query "Karl Marx"

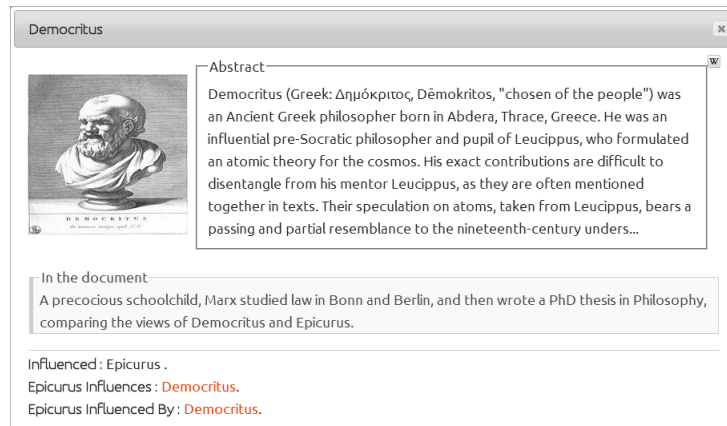


Fig. 2: Screenshot of a main-resource's description

## 4 Conclusion

We proposed a way of enhancing a traditional snippet with structured data coming from the LoD. Our intention was firstly to identify relationships between the main concepts of a document relevant to the user's information need, and secondly to link these concepts with either relevant concepts not present in the document, or facts about them. Our main challenge was to succeed in this

<sup>10</sup> <http://demo.ensen-insa.org>

<sup>11</sup> <http://blog.ensen-insa.org>



task while introducing as little noise as possible. In this respect, we identified the keystone as the graph extension step for which more relational information was needed in order to efficiently choose the resources that would benefit most from an extension. Thus, we proposed a multi-step unsupervised co-clustering algorithm in order to discover additional relationships between the resources by taking into account the textual data coming from the result of the search engine and also from textual nodes in the LoD (e.g. a DBpedia abstract). Next, we represented both the structures coming from the LoD and those added by our co-clustering process with a 3-way tensor before running a tensor decomposition so as to identify the main resources we should extend. Our proposal has been implemented and we are now in the process of thoroughly evaluating its usability and efficiency.

## 5 Acknowledgment

We would like to thank Adrian Delmarre and Sylvain Gourgeon for contributing the analysis of the CORCONDIA algorithm [1] on the graphs we generate.

## References

1. Andersson, C.A., Bro, R.: The n-way toolbox for matlab. *Chemometrics and Intelligent Laboratory Systems* 52(1), 1–4 (2000)
2. Bai, X., Delbru, R., Tummarello, G.: Rdf snippets for semantic web search engines. In: *On the Move to Meaningful Internet Systems: OTM 2008*, pp. 1304–1318. Springer (2008)
3. Franz, T., Schultz, A., Sizov, S., Staab, S.: Triplerank: Ranking semantic web data by tensor decomposition. In: *The Semantic Web-ISWC 2009*, pp. 213–228. Springer (2009)
4. Ge, W., Cheng, G., Li, H., Qu, Y.: Incorporating compactness to generate term-association view snippets for ontology search. *Information Processing & Management* pp. 513–528 (2012)
5. Haas, K., Mika, P., Tarjan, P., Blanco, R.: Enhanced results for web search. In: *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. pp. 725–734. ACM (2011)
6. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM review* 51(3), 455–500 (2009)
7. Mendes, P.N., Jakob, M., García-Silva, A., Bizer, C.: Dbpedia spotlight: shedding light on the web of documents. In: *Proceedings of the 7th International Conference on Semantic Systems*. pp. 1–8. I-Semantics '11, ACM (2011)
8. Park, L.A., Leckie, C.A., Ramamohanarao, K., Bezdek, J.C.: Adapting spectral co-clustering to documents and terms using latent semantic analysis. In: *AI 2009: Advances in Artificial Intelligence*. pp. 301–311. Springer (2009)
9. Pelleg, D., Moore, A.W., et al.: X-means: Extending k-means with efficient estimation of the number of clusters. In: *ICML*. pp. 727–734 (2000)
10. Penin, T., Wang, H., Tran, T., Yu, Y.: Snippet generation for semantic web search engines. In: *The Semantic Web*, pp. 493–507. Springer (2008)

11. Rousseeuw, P.J.: Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics* 20, 53–65 (1987)
12. Steiner, T., Troncy, R., Hausenblas, M.: How google is using linked data today and vision for tomorrow. *Proceedings of Linked Data in the Future Internet* 700 (2010)