

# A SEMANTIC-BASED APPROACH FOR ARTIST SIMILARITY

Sergio Oramas<sup>1</sup>, Mohamed Sordo<sup>2</sup>, Luis Espinosa-Anke<sup>3</sup>, Xavier Serra<sup>1</sup>

<sup>1</sup>Music Technology Group, Universitat Pompeu Fabra

<sup>2</sup>Center for Computational Science, University of Miami

<sup>3</sup>TALN Group, Universitat Pompeu Fabra

{sergio.oramas, luis.espinosa, xavier.serra}@upf.edu, msordo@miami.edu

## ABSTRACT

This paper describes and evaluates a method for computing artist similarity from a set of artist biographies. The proposed method aims at leveraging semantic information present in these biographies, and can be divided in three main steps, namely: (1) entity linking, i.e. detecting mentions to named entities in the text and linking them to an external knowledge base; (2) deriving a knowledge representation from these mentions in the form of a semantic graph or a mapping to a vector-space model; and (3) computing semantic similarity between documents. We test this approach on a corpus of 188 artist biographies and a slightly larger dataset of 2,336 artists, both gathered from Last.fm. The former is mapped to the MIREX Audio and Music Similarity evaluation dataset, so that its similarity judgments can be used as ground truth. For the latter dataset we use the similarity between artists as provided by the Last.fm API. Our evaluation results show that an approach that computes similarity over a graph of entities and semantic categories clearly outperforms a baseline that exploits word co-occurrences and latent factors.

## 1. INTRODUCTION

Artist biographies are a big source of musical context information and have been previously used for computing artist similarity. However, only shallow approaches have been applied by computing word co-occurrences and thus the semantics implicit in text have been barely exploited. To do so, semantic technologies, and more specifically Entity Linking tools may play a key role to annotate unstructured texts. These are able to identify named entities in text and disambiguate them with their corresponding entry in a knowledge base (e.g. Wikipedia, DBpedia or BabelNet).

This paper describes a method for computing semantic similarity at document-level, and presents evaluation results in the task of artist similarity. The cornerstone of this work is the intuition that semantifying and formaliz-

ing relations between entity mentions in documents (both at in-document and cross-document levels) can represent the relatedness of two documents. Specifically, in the task of artist similarity, this derives in a measure to quantify the degree of relatedness between two artists by looking at their biographies.

Our experiments start with a preprocessing step which involve Entity Linking over artist biographical texts. Then, a knowledge representation is derived from the detected entities in the form of a semantic graph or a mapping to a vector-space model. Finally, different similarity measures are applied to a benchmarking dataset. The evaluation results indicate that some approaches presented in this paper clearly outperform a baseline based on shallow word co-occurrence metrics. Source code and datasets are available online<sup>1</sup>.

The remainder of this article is structured as follows: Section 2 reviews prominent work in the fields and topic relevant to this paper; Section 3 details the different modules that integrate our approach; Section 4 describes the settings in which experiments were carried out together with the evaluation metrics used; Section 5 presents the evaluation results and discusses the performance of our method; and finally Section 6 summarizes the main topics covered in this article and suggests potential avenues for future work.

## 2. RELATED WORK

Music artist similarity has been studied from the score level, the acoustic level, and the cultural level [9]. This work is focused on the latter approach, and more specifically in text-based approaches. Literature on document similarity, and more specifically on the application of text-based approaches for artist similarity is discussed next.

The task of identifying similar text instances, either at sentence or document level, has applications in many areas of Artificial Intelligence and Natural Language Processing [17]. In general, document similarity can be computed according to the following approaches: surface-level representation like keywords or n-grams [6]; corpus representation using counts [28], e.g. word-level correlation, jaccard or cosine models; Latent factor models, such as Latent Semantic Analysis [8]; or methods exploiting external



© Sergio Oramas<sup>1</sup>, Mohamed Sordo<sup>2</sup>, Luis Espinosa-Anke<sup>3</sup>, Xavier Serra<sup>1</sup>.

Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Sergio Oramas<sup>1</sup>, Mohamed Sordo<sup>2</sup>, Luis Espinosa-Anke<sup>3</sup>, Xavier Serra<sup>1</sup>. “A Semantic-based Approach for Artist Similarity”, 16th International Society for Music Information Retrieval Conference, 2015.

<sup>1</sup> <http://mtg.upf.edu/downloads/datasets/semantic-similarity>

knowledge bases like ontologies or encyclopedias [12].

The use of text-based approaches for artist and music similarity was first applied in [7], by computing co-occurrences of artist names in web page texts and building term vector representations. By contrast, in [30] term weights are extracted from search engine’s result counts. In [33] n-grams, part-of-speech tagging and noun phrases are used to build a term profile for artists, weighted by employing tf-idf. Term profiles are then compared and the sum of common terms weights gives the similarity measure. More approaches using term weight vectors have been developed over different text sources, such as music reviews [11], blog posts [4], or microblogs [29]. In [18] Latent Semantic Analysis is used to measure artist similarity from song lyrics. Domain specific ontologies have also been applied to the problem of music recommendation and similarity, such as in [5]. In [16], paths on an ontological graph extracted from DBpedia are exploited for recommending music web pages. However, to the best of our knowledge, there are scant approaches in the music domain that exploit implicit semantics and enhance term profiles with external knowledge bases.

### 3. METHODOLOGY

The method proposed in this paper can be divided in three main steps, as depicted in Fig 1. The first step performs entity linking, that is the detection of mentions to named entities in the text and their linking to an external knowledge base. The second step derives a semantically motivated knowledge representation from the named entity mentions. This can be achieved by exploiting natural language text as anchor between entities, or by incorporating semantic information from an external knowledge base. In the latter case, a document is represented either as a semantic graph or as a set of vectors projected on a vector space, which allows the use of well known vector similarity metrics. Finally, the third step computes semantic similarity between documents (artist biographies in our case). This step can take into consideration semantic similarity among entity mentions in document pairs, or only the structure and content of the semantic graph.

The following sections provide a more detailed description of each one of these steps, along with all the approaches we have considered in each step.

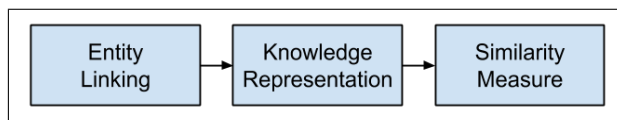


Figure 1. Workflow of the proposed method.

#### 3.1 Entity Linking

Entity linking is the task to associate, for a given candidate textual fragment, the most suitable entry in a reference Knowledge Base (KB) [23]. It encompasses similar subtasks such as Named Entity Disambiguation [2], which

is precisely linking mentions to entities to a KB, or Wikification [21], specifically using Wikipedia as KB.

We considered several state-of-the-art entity linking tools, including Babelify [23], TagMe [10], Agdistis [32] and DBpedia Spotlight [20]. However we opted to use the first one for consistency purposes, as in a later step we exploit *SensEmbed* [13], a vector space representation of concepts based on BabelNet [24]. Moreover, the use of a single tool across approaches guarantees that the evaluation will only reflect the appropriateness of each one of them, and in case of error propagation all the approaches will be affected the same.

Babelify [23] is a state-of-the-art system for entity linking and word sense disambiguation based on non-strict identification of candidate meanings (i.e. not necessarily exact string matching), together with a graph based algorithm that traverses the BabelNet graph and selects the most appropriate semantic interpretation for each candidate.

#### 3.2 Knowledge representation

##### 3.2.1 Relation graph

Relation extraction has been defined as the process of identifying and annotating relevant semantic relations between entities in text [15]. In order to exploit the semantic relations between entities present in artist biographies, we applied the method defined in [25] for relation extraction in the music domain. The method basically consists of three steps. First, entities are identified in the text by applying entity linking. Second, relations between pairs of entities occurring in the same sentence are identified and filtered by analyzing the structure of the sentence, which is obtained by running a syntactic parser based on the formalism of dependency grammar [1]. Finally, the identified entities and relations are modeled as a knowledge graph. This kind of extracted knowledge graphs may be useful for music recommendation [31], as recommendations can be conveyed to users by means of natural language. We apply this methodology to the problem of artist similarity, by creating a graph that connects the entities detected in every artist biography. We call this approach RG (relation graph). Figure 2 shows the output of this process for a single sentence.

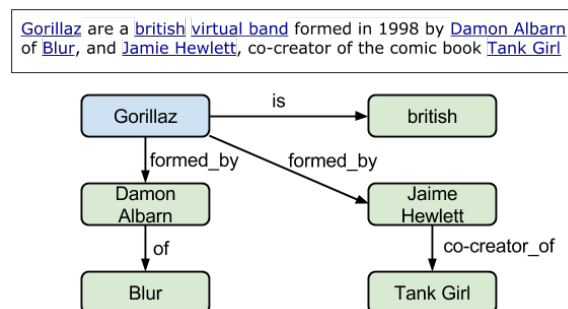


Figure 2. Relation graph of a single sentence

### 3.2.2 Semantically enriched graph

A second approach is proposed using the same set of linked entities. However, instead of exploiting natural language text, we use semantic information from the referenced knowledge base to enrich the semantics of the linked entities. We follow a semantic enrichment process similar to the one described in [27]. We use semantic information coming from DBpedia<sup>2</sup>. DBpedia resources are generally classified using the DBpedia Ontology, which is a shallow, cross-domain ontology based on the most common infoboxes of Wikipedia. DBpedia resources are categorized using this ontology among others (e.g. Yago, schema.org) through the `rdfs:type` property. In addition, each Wikipedia page may be associated with a set of Wikipedia categories, which link articles under a common topic. DBpedia resources are related to Wikipedia categories through the property `dcterms:subject`.

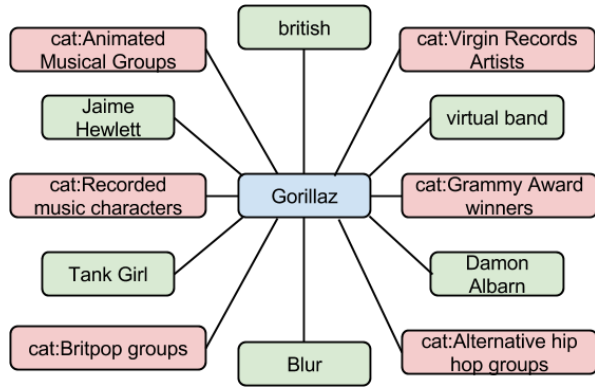
We take advantage of these two properties to build our semantically enriched graph. We consider three types of nodes for this graph: 1) artist entities obtained by matching the artist names to their corresponding DBpedia entry; 2) named entities detected by the entity linking step; and 3) Wikipedia categories associated to all the previous entities. Edges are then added between artist entities and the named entities detected in their biographies, and between entities and their corresponding Wikipedia categories. For the construction of the graph, we can select all the detected named entities, or we can filter them out according to the information related to their `rdfs:type` property. A set of six types was selected, including *artist*, *band*, *work*, *album*, *musicgenre*, and *person*, which we consider more appropriate to semantically define a musical artist.

From the previous description, we define five variants of this approach. The first variant, which we call AEC (Artists-Entities-Categories), considers all 3 types of nodes along with their relations (as depicted in Figure 3). The second variant, named AE (Artists-Entities) ignores the categories of the entities. The third and fourth variant, named AEC-FT and AE-FT, are similar to the first and second variant, respectively, except that the named entities are filtered using the above mentioned list of 6 entity types. Finally, the fifth variant, EC, ignores the artist entities of node type 1.

### 3.2.3 Sense embeddings

The semantic representation used in this approach is based on SensEmbed [13]. SensEmbed is a vector space semantic representation of words similar to word2vec [22], where each vector represents a BabelNet synset and its lexicalization. Let  $A$  be the set of artist biographies in our dataset. Each artist biography  $a \in A$  is converted to a set of disambiguated concepts  $Bf_{y_a}$  after running Babelify over it.

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



**Figure 3.** Semantically enriched subgraph of the same sentence from Figure 2, variant AEC with  $h=1$

## 3.3 Similarity approaches

### 3.3.1 SimRank

SimRank is a similarity measure based on a simple graph-theoretic model [14]. The intuition is that two nodes are similar if they are referenced by similar nodes. In particular we use the definition of bipartite SimRank [14]. We build a bipartite graph with named entities and their corresponding Wikipedia categories (the EC variant from Section 3.2.2). The similarity between two named entities (say  $p$  and  $q$ ) is computed with the following recursive equation:

$$s(p, q) = \frac{C}{|O(p)||O(q)|} \sum_{i=1}^{|O(p)|} \sum_{j=1}^{|O(q)|} s(O_i(p), O_j(q)) \quad (1)$$

where  $O$  denotes the out-neighboring nodes of a given node and  $C$  is a constant between 0 and 1. For  $p = q$ ,  $s(p, q)$  is automatically set up to 1. Once the similarity between all pairs of entities is obtained, we proceed to calculate the similarity between pairs of artists (say  $a$  and  $b$ ) by aggregating the similarities between the named entities identified in their biographies, as shown in the following formula:

$$sim(a, b) = Q(a, b) \frac{1}{N} \sum_{e_a \in a} \sum_{e_b \in b} s(e_a, e_b) \quad \text{if } s(e_a, e_b) \geq 0.1 \quad (2)$$

where  $s$  denotes the SimRank of entities  $e_a$  and  $e_b$  and  $N$  is the number of  $(e_a, e_b)$  pairs with  $s(e_a, e_b) \geq 0.1$ . This is done to filter out less similar pairs. Finally,  $Q(a, b)$  is a normalizing factor that accounts for the pairs of artists with more similar entity pairs than others.

### 3.3.2 Maximal common subgraph

Maximal common subgraph (MCS) is a common distance measure on graphs. It is based on the maximal common subgraph of two graphs. MCS is a symmetric distance metric, thus  $d(A, B) = d(B, A)$ . It takes structure as well as content into account. According to [3], the distance between two non empty graphs  $G_1$  and  $G_2$  is defined as

$$d(G_1, G_2) = 1 - \frac{|mcs(G_1, G_2)|}{\max(|G_1|, |G_2|)} \quad (3)$$

It can also be seen as a similarity measure  $s$ , assuming that  $s = 1 - d$ , as applied in [19]. To compute this similarity measure we need to have a graph for each artist. This can be achieved by finding subgraphs in the graph approaches defined in Section 3.2. A subgraph will include an artist entity node and its neighboring nodes. Furthermore, we apply the notion of h-hop item neighborhood graph defined in [26] to a semantic graph. Let  $G = (E, P)$  be an undirected graph where  $E$  represent the nodes (entities), and  $P$  the set of edges with  $P \subseteq E \times E$ . For an artist item  $a$  in  $G$ , its h-hop neighborhood subgraph  $G^h(a) = (E^h(a), P^h(a))$  is the subgraph of  $G$  formed by the set of entities that are reachable from  $a$  in at most h hops, according to the shortest path. Following this approach, we obtain an h-hop item neighborhood graph for each artist node of the semantic graph. Then, maximal common subgraph is computed between each pair of h-hop item neighborhood graphs. For each artist, the list of all similar artists ordered from the most similar to the less one is finally obtained.

### 3.3.3 Cumulative cosine similarity

For each pair of concepts  $c \in \text{Bfy}_a$  and  $c' \in \text{Bfy}'_a$  (as defined in Section 3.2.3), we are interested in obtaining the similarity of their closest senses. This is achieved by first deriving the set of associated vectors  $V_c$  and  $V_{c'}$  for each pair of concepts  $c, c'$ , and then optimizing

$$\max_{v_c \in V_c, v_{c'} \in V_{c'}} \left( \frac{v_c \times v_{c'}}{\|v_c\| \|v_{c'}\|} \right) \quad (4)$$

i.e. computing cosine similarity between all possible senses (each sense represented as a vector) in an all-against-all fashion and keeping the highest scoring similarity score for each pair. Finally, the semantic similarity between two artist biographies is simply the average among all the cosine similarities between each concept pair.

## 4. EXPERIMENTAL SETUP

To evaluate the accuracy of the proposed approaches we designed an experimental evaluation over two datasets. The first dataset contains 2,336 artists and it is evaluated using the list of similar artists provided by the Last.fm API as a ground truth. The second dataset contains 188 artists, and it is evaluated against user similarity judgements from the MIREX Audio Music Similarity and Retrieval task. Apart from the defined approaches, a pure text-based approach for document similarity is added to act as a reference for the obtained results.

### 4.1 Datasets

#### 4.1.1 Last.fm dataset

A dataset of 2,336 artist biographies was gathered from Last.fm. The artists in this dataset share a set of restrictions. Their biography has at least 500 characters and is

written in English. All of the artists have a correspondent Wikipedia page, and we have been able to mapped it automatically, obtaining the DBpedia URI of every artist. For every artist, we queried the getSimilar method of the Last.fm API and obtained an ordered list of similar artists. Every artist in the dataset fulfills the requirement of having at least 10 similar artists within the dataset. We used these lists of similar artists as the ground truth for our evaluation.

#### 4.1.2 MIREX dataset

To build this dataset, the gathered artists from Last.fm were mapped to the MIREX Audio Music Similarity task dataset. The AMS dataset (7,000 songs from 602 unique artists) contains human judgments of song similarity. According to [29], the similarity between two artists can be roughly estimated as the average similarity between their songs. We used the same approach in [29], that is, two artists were considered similar if the average similarity score between their songs was at least 25 (on a fine scale between 0 and 100).

After the mapping, we obtained an overlap of 268 artists. As we want to evaluate Top-10 similarity, every artist in the ground truth dataset should have information of at least 10 similar artists. However, not every artist in the MIREX evaluation dataset fulfills this requirement. Therefore, after removing the artists with less than 10 similars, we obtained a final dataset of 188 artists, and used it for the evaluation.

## 4.2 Baseline

In order to assess the goodness of our approaches, we need to define a baseline approach with which to compare to. The baseline used in this paper is a classic vector-based model approach used in many Information Retrieval systems. A text document is represented as a vector of word frequencies (after removing English stopwords and words with less than 2 characters), and a matrix is formed by aggregating all the vectors. The word frequencies in the matrix are then re-weighted using TF-IDF, and finally latent semantic analysis (LSA) [8] is used to produce a vector of concepts for each document. The similarity between two documents can be obtained by using a cosine similarity over their corresponding vectors.

## 4.3 Evaluated approaches

From all possible combinations of knowledge representations, similarity measures and parameters, we selected a set of 10 different approach variants. The prefixes AEC, RG and AE refer to the graph representations (see Sections 3.2.1 and 3.2.2). SE refers to the sense embeddings approach, and LSA to the latent semantic analysis baseline approach. When these prefixes are followed by FT, it means that the entities in the graph have been filtered by type. The second term in the name refers to the similarity measure. MCS refers to maximal common subgraph, and SimRank and Cosine to SimRank and cumulative cosine similarity measures. MCS approaches are further followed by a number indicating the number of h-hops of the neighborhood subgraph.

Approach variants	Genres							
	Blues	Country	Edance	Jazz	Metal	Rap	Rocknroll	Overall
Ground Truth	5.78	5.46	6.88	7.04	7.10	8.68	5.17	6.53
LSA	4.43	4.12	3.80	4.64	5.79	5.08	4.74	4.69
RG MCS 1-hop	2.63	3.50	1.50	2.95	4.00	2.54	1.70	2.68
RG MCS 2-hop	4.14	4.92	1.69	2.80	3.78	3.06	2.77	3.27
AE MCS	5.52	5.15	4.36	7.00	4.34	5.36	4.46	5.11
AE-FT MCS	5.43	6.12	4.16	6.20	6.32	5.36	3.77	5.26
AEC MCS 1-hop	<b>7.22</b>	5.92	5.24	7.12	5.48	6.92	4.86	6.02
AEC MCS 2-hop	4.22	3.69	4.56	6.20	4.55	4.64	4.09	4.54
AEC-FT MCS 1-hop	6.91	<b>6.80</b>	<b>6.04</b>	<b>7.60</b>	<b>6.79</b>	<b>7.12</b>	<b>5.37</b>	<b>6.59</b>
AEC-FT MCS 2-hop	4.09	4.36	5.56	6.72	4.39	4.16	3.77	4.67
EC SimRank	6.74	5.38	3.16	6.40	4.59	4.44	3.80	4.85
SE Cosine	3.39	5.50	5.32	5.16	4.31	5.36	4.31	4.75

**Table 3.** Average genre distribution of the top-10 similar artists using the MIREX dataset. In other words, on average, how many of the top-10 similar artists are from the same genre as the query artist. LSA stands for Latent Semantic Analysis, RG for Relation Graph, SE for Sense Embeddings, and AE, AEC and EC represent the semantically enriched graphs with Artists-Entities, Artist-Entities-Categories, and Entities-Categories nodes, respectively. As for the similarity approaches, MCS stands for Maximum Common Subgraph.

Approach variants	Precision@N		nDCG@N	
	N=5	N=10	N=5	N=10
LSA	0.100	0.169	0.496	0.526
RG MCS 1-hop	0.059	0.087	0.465	0.476
RG MCS 2-hop	0.056	0.101	0.433	0.468
AE MCS	0.106	0.178	0.503	0.517
AE-FT MCS	0.123	0.183	0.552	0.562
AEC MCS 1-hop	0.120	0.209	0.573	0.562
AEC MCS 2-hop	0.086	0.160	0.550	0.539
AEC-FT MCS 1-hop	<b>0.140</b>	<b>0.218</b>	<b>0.588</b>	<b>0.578</b>
AEC-FT MCS 2-hop	0.100	0.160	0.527	0.534
EC SimRank	0.097	0.171	0.509	0.534
SE Cosine	0.095	0.163	0.454	0.484

**Table 1.** Precision and normalized discounted cumulative gain for Top-N artist similarity using the MIREX dataset ( $N=\{5, 10\}$ )

#### 4.4 Evaluation measures

To measure the accuracy of the artist similarity we adopt two standard performance metrics such as Precision@N, and nDCG@N (normalized discounted cumulative gain). Precision@N is computed as the number of relevant items (i.e., true positives) among the top-N items divided by N, when compared to a ground truth. Precision considers only the relevance of items, whilst nDCG takes into account both relevance and rank position. Denoting with  $s_{ak}$  the relevance of the item in position  $k$  in the Top-N list for the artist  $a$ , then nDCG@N for  $a$  can be defined as:

$$\text{nDCG@N} = \frac{1}{\text{IDCG@N}} \sum_{k=1}^N \frac{2^{s_{ak}} - 1}{\log_2(1 + k)} \quad (5)$$

where IDCG@N indicates the score obtained by an ideal or perfect Top-N ranking and acts as a normalization factor. We run our experiments for  $N = 5$  and  $N = 10$ .

Approach variants	Precision@N		nDCG@N	
	N=5	N=10	N=5	N=10
LSA	0.090	0.088	0.233	0.269
RG MCS 1-hop	0.055	0.083	0.126	0.149
AE MCS	0.124	0.200	0.184	0.216
AE-FT MCS	0.136	0.201	0.224	0.260
AEC MCS 1-hop	0.152	0.224	0.277	0.297
AEC-FT MCS 1-hop	<b>0.160</b>	<b>0.242</b>	<b>0.288</b>	<b>0.317</b>

**Table 2.** Precision and normalized discounted cumulative gain for Top-N artist similarity using the Last.fm dataset ( $N=\{5, 10\}$ )

## 5. RESULTS AND DISCUSSION

We evaluated all the approach variants described in Section 4.3 on the MIREX dataset, but only a subset of them on the Last.fm dataset, due to the high computational cost of some of the approaches.

Table 1 shows the Precision@N and nDCG@N results of the evaluated approaches using the MIREX dataset, while Table 2 shows the same results for the Last.fm dataset. We obtained very similar results in both datasets. The approach that gets best performance for every metric, dataset and value of N is the combination of the Artists-Entities-Categories graph filtered by types, with the maximal common subgraph similarity measure using a value of  $h = 1$  for obtaining the h-hop item neighborhood graphs.

Furthermore, given that the MIREX AMS dataset also provides genre data, we analyzed the distribution of genres in the top-10 similar artists for each artist, and averaged them by genres. The idea is that an artist’s most similar artists should be from the same genre as the seed artist. Table 3 presents the results. Again, the best results are obtained with the approach that combines the Artists-Entities-Categories graph filtered by types, with the maxi-

mal common subgraph similarity measure using a value of  $h = 1$  for the h-hop item neighborhood graphs.

We extract some insights from these results. First, semantic approaches are able to improve pure text-based approaches. Second, using knowledge from an external knowledge base provides better results than exploiting the relations inside the text. Third, using a similarity measure that exploits the structure and content of a graph, such as maximal common subgraph, overcomes other similarity measures based on semantic similarity among entity mentions in document pairs.

## 6. CONCLUSION

In this paper we presented a methodology that exploits semantic technologies for computing artist similarity, which can be divided in three main steps: First, named entity mentions are identified in the text and linked to a knowledge base. Then, these entity mentions are used to construct a semantically motivated knowledge representation. Finally a similarity function is defined on top of the knowledge representation to compute the similarity between artists. For each one of these steps we explored several approaches, and evaluated them against a small dataset of 188 artist biographies, and a larger dataset of 2,336 artists, both obtained from Last.fm.

Results showed that a combination of the Artists-Entity-Categories graph filtered by types, and a maximal common subgraph similarity measure using a value of  $h = 1$  for obtaining the h-hop item neighborhood graphs, clearly outperforms a baseline approach that exploits word co-occurrences and latent factors. In the light of these results, the following conclusions can be drawn: First, semantic approaches may outperform pure text-based approaches. Second, we observe that knowledge leveraged from external ontologies may improve the accuracy of the similarity measure. Third, reducing noise by filtering linked entities by type is a rewarding step that contributes to an improved performance. Finally, we show that similarity measures that take into consideration the structure and content of a graph representation may achieve much higher performance.

There are still many avenues for future work. We would like to compare our semantic-based approach with acoustic and collaborative filtering approaches. In addition, the use of text sources different from artist biographies can be studied. Finally, in order to improve the results obtained by our semantic approach, different state-of-the-art entity linking tools can be applied, or a specific entity linking tool for the music domain could be created for this purpose.

## 7. REFERENCES

- [1] Bernd Bohnet. Very high accuracy and fast dependency parsing is not a contradiction. In *Proceedings of the 23rd International Conference on Computational Linguistics, COLING '10*, pages 89–97, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [2] Razvan Bunescu and Marius Pasca. Using Encyclopedic Knowledge for Named Entity Disambiguation. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics (EACL-06)*, pages 9–16, Trento, Italy, 2006.
- [3] Horst Bunke and Kim Shearer. A graph distance metric based on the maximal common subgraph. *Pattern Recognition Letters*, 19(3-4):255–259, March 1998.
- [4] Òscar Celma, Pedro Cano, and Perfecto Herrera. Search Sounds An audio crawler focused on weblogs. In *7th International Conference on Music Information Retrieval (ISMIR)*, 2006.
- [5] Òscar Celma and Xavier Serra. FOAFing the music: Bridging the semantic gap in music recommendation. *Web Semantics*, 6:250–256, 2008.
- [6] Hung Chim and Xiaotie Deng. Efficient phrase-based document similarity for clustering. *Knowledge and Data Engineering, IEEE Transactions on*, 20(9):1217–1229, 2008.
- [7] William W. Cohen and Wei Fan. Web-collaborative filtering: recommending music by crawling the Web. *Computer Networks*, 33:685–698, 2000.
- [8] Scott C. Deerwester, Susan T Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. Indexing by latent semantic analysis. *JASIS*, 41(6):391–407, 1990.
- [9] Daniel P. W. Ellis, Brian Whitman, Adam Berenzweig, and Steve Lawrence. The quest for ground truth in musical artist similarity. In *Proc. International Symposium on Music Information Retrieval (ISMIR 2002)*, pages 170–177, 2002.
- [10] Paolo Ferragina and Ugo Scaiella. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 1625–1628. ACM, 2010.
- [11] X Hu, JS Downie, Kris West, and AF Ehmann. Mining Music Reviews: Promising Preliminary Results. In *ISMIR*, pages 536–539, 2005.
- [12] Xiaohua Hu, Xiaodan Zhang, Caimei Lu, Eun K Park, and Xiaohua Zhou. Exploiting Wikipedia as external knowledge for document clustering. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 389–396. ACM, 2009.
- [13] Ignacio Iacobacci, Mohammad Taher Pilehvar, and Roberto Navigli. Sensembled: Enhancing word embeddings for semantic similarity and relatedness. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Beijing, China, July 2015. Association for Computational Linguistics.

- [14] Glen Jeh and Jennifer Widom. Simrank: a measure of structural-context similarity. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 538–543. ACM, 2002.
- [15] J. Jiang and C. Zhai. A systematic exploration of the feature space for relation extraction. In *HLT-NAACL*, pages 113–120, 2007.
- [16] José Paulo Leal, Vânia Rodrigues, and Ricardo Queirós. Computing Semantic Relatedness using DBpedia. *1st Symposium on Languages, Applications and Technologies, SLATE 2012*, 2012.
- [17] Hongzhe Liu and Pengfei Wang. Assessing Text Semantic Similarity Using Ontology. *Journal of Software*, 9(2):490–497, 2014.
- [18] Beth Logan and Daniel P W Ellis. Toward Evaluation Techniques for Music Similarity. *SIGIR 2003: Workshop on the Evaluation of Music Information Retrieval Systems*, pages 7–11, 2003.
- [19] Mathias Lux and Michael Granitzer. A Fast and Simple Path Index Based Retrieval Approach for Graph Based Semantic Descriptions. In *Proceedings of the Second International Workshop on Text-Based Information Retrieval*, 2005.
- [20] Pablo N Mendes, Max Jakob, Andrés García-Silva, and Christian Bizer. DBpedia spotlight: shedding light on the web of documents. In *Proceedings of the 7th International Conference on Semantic Systems*, pages 1–8. ACM, 2011.
- [21] Rada Mihalcea and Andras 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.
- [22] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic Regularities in Continuous Space Word Representations. In *HLT-NAACL*, pages 746–751, 2013.
- [23] Andrea Moro, Francesco Cecconi, and Roberto Navigli. Multilingual word sense disambiguation and entity linking for everybody. In *Proceedings of the 13th International Conference on Semantic Web (P&D)*, 2014.
- [24] Roberto Navigli and Simone Paolo Ponzetto. BabelNet: Building a very large multilingual semantic network. In *Proceedings of the 48th annual meeting of the association for computational linguistics*, pages 216–225. Association for Computational Linguistics, 2010.
- [25] Sergio Oramas, Mohamed Sordo, and Luis Espinosa-anke. A Rule-Based Approach to Extracting Relations from Music Tidbits. In *2nd Workshop in Knowledge Extraction from Text, WWW’15*, 2015.
- [26] Vito Claudio Ostuni, Tommaso Di Noia, Roberto Mirizzi, and Eugenio Di Sciascio. A Linked Data Recommender System using a Neighborhood-based Graph Kernel. *15th International Conference on Electronic Commerce and Web Technologies*, pages 1–12, 2014.
- [27] Vito Claudio Ostuni, Sergio Oramas, Tommaso Di Noia, Xavier Serra, and Eugenio Di Sciascio. A Semantic Hybrid Approach for Sound Recommendation. *24th International World Wide Web Conference (WWW 2015)*, pages 3–4, 2015.
- [28] Mark Rorvig. Images of similarity: A visual exploration of optimal similarity metrics and scaling properties of TREC topic-document sets. *Journal of the American Society for Information Science*, 50(8):639–651, 1999.
- [29] Markus Schedl, David Hauger, and Julián Urbano. Harvesting microblogs for contextual music similarity estimation: a co-occurrence-based framework. *Multimedia Systems*, 20(6):693–705, 2013.
- [30] Markus Schedl, Peter Knees, and Gerhard Widmer. A Web-Based Approach to Assessing Artist Similarity using Co-Occurrences. In *Proceedings of the 4th International Workshop on Content-Based Multimedia Indexing {(CBMI’05)}*, 2005.
- [31] Mohamed Sordo, Sergio Oramas, and Luis Espinosa. Extracting Relations from Unstructured Text Sources for Music Recommendation. In *20th International Conference on Applications of Natural Language to Information Systems*, pages 1–14, 2015.
- [32] Ricardo Usbeck, Axel-Cyrille Ngonga Ngomo, Sören Auer, Daniel Gerber, and Andreas Both. Agdistis-agnostic disambiguation of named entities using linked open data. In *International Semantic Web Conference*, page 2, 2014.
- [33] B. Whitman and S. Lawrence. Inferring descriptions and similarity for music from community metadata. In *Proceedings of the 2002 International Computer Music Conference*, pages 591–598, 2002.