

# A method to manage the difference of precision between Profiles and Items for Recommender System

## *Applied upon a news recommender system using SVM approach*

David Werner<sup>1</sup>, Christophe Cruz<sup>1</sup>

<sup>1</sup>LE2I Laboratory, UMR CNRS 6306, Université de Bourgogne, BP 47870, 21078 Dijon Cedex, France  
{david.werner, christophe.cruz}@u-bourgogne.fr

**Keywords:** recommender system, news, domain ontology, ontologies, knowledge base, indexing, recommendation, vector space model

**Abstract:** Contractors, commercial and business decision-makers need economical information to drive their decisions. The production and distribution of a press review about French regional economic actors represents a prospecting tool on partners and competitors for the businessman. Our goal is to propose a customized review for each user, thus reducing the overload of useless information. Some systems for recommending news items already exist. The usefulness of external knowledge to improve the process has already been explained in information retrieval. The system's knowledge base includes the domain knowledge used during the recommendation process. Our recommender system architecture is standard, but during the indexing task, the representations of content of each article and interests of users' profiles created are based on this domain knowledge. Articles and Profiles are semantically defined in the Knowledge base via concepts, instances and relations. This paper deals with the similarity measure, a critical subtask in recommendation systems. The Vector Space Model is a well-known model used for relevance ranking. The problematic exposed here is the utilization of the standard VSM method with our indexing method.

## 1 INTRODUCTION

The decision-making process in the economic field requires the centralization and intake of a large amount of information. The aim is to keep abreast with current market trends. Thus, contractors, business men and sales persons need to continuously be aware of the market conditions. This means to be up-to-date regarding ongoing information and projects undergoing development. With the help of economic monitoring, prospects can be easily identified, so as to establish new contracts. Our tool is specialized in the production and distribution of press reviews about French regional economic actors.

The reviews sent are the same for each user, but personalized according to a geographic area. All articles sent do not necessarily correspond to a person's needs, and can be a waste of time. To reduce the overload of useless information, we are moving towards a customized review for each user. Therefore, an opinion survey on magazine readers that covers a broad array of subjects, including news

services, was undertaken. Criteria for a relevant customization of the review were identified as a result of this survey as well as expert domain knowledge. These criteria are economic themes (i.e. main economic events), economic sectors, major transverse projects, temporal and localization data about each element underlined. Therefore, the complete production process of the review was redesigned to produce and to automatically distribute a customized review for each user. Another drawback in the existing process is the produced information storage. News articles are stored as PDF file reviews (i.e. the same format that is sent to customers, natural language), but this unstructured format is hard to handle and reuse.

The aim of the architecture is to manage all news produced, and provide the most relevant article for each customer, using our domain knowledge. The overload of news information is a particular case of information overload, which is a well-known problem, studied by Information Retrieval and Recommender Systems research fields. News recommender systems already exist (Middleton et al. 2004) (Getahun et al. 2009) (Liu et al. 2010)

(Tanev et al. 2008) SCENE (Li et al. 2011) NewsJunkie (Gabrilovich et al. 2004) Athena (Ijntema et al. 2010) GroupLens (Resnick et al. 1994) News Dude (Billsus, Pazzani. 1999) et YourNews (Brusilosky et al. 2007). Some of these systems use domain knowledge to improve the recommendation task (Ijntema et al. 2010) (Middleton et al. 2004).

To achieve this goal, a content-based recommender system is being developed. A recommender system is necessary for the item ranking. And content-base is required to analyze the content of each article to structure and preserve information content. The results of the analysis enable to link the domain knowledge to the articles. Because the domain knowledge can be reused to improve the recommendation task (Ijntema et al. 2010) (Middleton et al. 2004).

This paper is organized as follows: Section 2 presents a brief review of related work. In section 3, we outline the proposed solution of a recommender system. Section 4 presents the similarity measure task, and our implementation of VMS. Section 5 is the conclusion and future work.

## 2 RELATED WORK

Our work is related to several works in News recommender systems, SCENE (Li et al. 2011) NewsJunkie (Gabrilovich et al. 2004) Athena (Ijntema et al. 2010) News Dude (Billsus, Pazzani. 1999) et YourNews (Brusilosky et al. 2007). The survey of K. Nageswara Rao (Nageswara, Talwar, 2008) proposed a general comparison of the main advantages and drawbacks of each kind of Recommender System (e.g. content based or collaborative filtering). The advantages of content-based recommender systems for news recommendation are also explained in (Liu et al. 2010) to improve the Google news platform. The main drawback of collaborative filtering systems is the recommendation of new items. Novelty is very important in the particular case of news; each new article must be quickly recommended. Waiting for enough users to have read it before recommending is a big waste of time. Furthermore, we need to be able to recommend for very particular user profiles, as some customer needs are unique.

There are many systems that work without Knowledge (Liu et al. 2010) (Billsus, Pazzani. 1999) (Resnick et al. 1994). The advantages of using exterior knowledge for enhancing the recommendation were exposed by Ijntema (Ijntema

et al. 2010). He uses the name semantic-based recommender system to distinguish standard content-based systems from the systems using external knowledge (e.g. domain ontologies or lexical knowledge as WordNet (Fellbaum 1998)). Lexical knowledge is used by (Getahun et al. 2009) and domain knowledge by (Middleton et al. 2004). Athena uses both (Ijntema et al. 2010). Ontologies used by these systems already exist or are created by hand, and maintained. Unlike previous systems, the Knowledge base containing domain Knowledge is used as an index for articles and profiles, as it is explained in section 3. To compare profiles and articles, classic VSM is not directly usable, so we have adapted it, as presented in section 4.

## 3 IMPLEMENTATION

Our system is an ontology drive content based recommender system. An ontology schema is created and populated with the help of company experts, in order to model the domain knowledge in a knowledge base. In a classic content-based recommender system, we distinguish two main tasks. The first is indexing. The task is to create a representation of the users' needs, and item content. The Knowledge base will be populated during this task. The quality of content analysis is important for the knowledge base population and for indexing, so our system is semi-supervised by an expert. The second task is comparison. This task is the comparison with item representation so as to measure the degree of relevance for each profile. Items are ranked with the help of the similarity measure, after being provided to the user. These subjects are developed as follows.

### 3.1 The Knowledge Base

The knowledge base  $\mathcal{K}$  used for this system is composed of several ontologies (Werner et al. 2012). An Upper level ontology is used to manage information shared by all application areas (in the case of an extension of the system to new fields of application). High level concepts like location, geospatial information, temporality, events, agents, etc. Domain ontology is used to manage domain-specific knowledge. Concepts of this ontology are mainly specialization of concepts from the upper level. Other ontologies are used to manage articles, profiles, and lexical resources used by a gazetteer. The lexical resource ontology is based on the ontology PROTON used on the KIM platform

(Popov et al. 2003). In this paper the knowledge base model defined by Ehrig et al. (Ehrig et al. 2004) and based on the Karlsruhe Ontology (SOURCE) Model is used.

**Definition 1** Ontology:

$$O = (C, T, \leq_C, \leq_T, R, A, \sigma_a, \sigma_R, \leq_R, \leq_A)$$

Wherein  $C, T, R, A$  are disjoint sets of concepts, data types, relations and attributes,  $\leq_C, \leq_T, \leq_R, \leq_A$  are hierarchy of classes, data types, relations and attributes, and  $\sigma_a, \sigma_R$  are function that provide the signature for each  $\sigma_a: A \rightarrow C \times T$  attribute and  $\sigma_R: R \rightarrow C \times C$  relation.

**Definition 2** Knowledge Base:

$$\mathcal{K} = (C_{KB}, T_{KB}, R_{KB}, A_{KB}, I, V, \iota_C, \iota_T, \iota_R, \iota_A)$$

Wherein  $C_{KB}, T_{KB}, R_{KB}, A_{KB}, I, V$  are disjoint sets of concepts, data types, relation attributes, instances and data values.  $\iota_C$  is the classes instantiation function  $\iota_C: C_{KB} \rightarrow 2^I$ .  $\iota_T$  is the data type instantiation function  $\iota_T: T_{KB} \rightarrow 2^V$ .  $\iota_R$  is the relation instantiation function  $\iota_R: R_{KB} \rightarrow 2^{I \times I}$ .  $\iota_A$  is the attribute instantiation function  $\iota_A: A_{KB} \rightarrow 2^{I \times V}$ .

## 3.2 Indexing

To archive the recommendation of articles to customers, the system needs are a representation of the content of each article, and representation of the needs of each customer. The index used in our system is the same for articles and profiles. The knowledge base used has an index. Articles and profiles are represented by instances in our knowledge base. Some relations in the system ontologies are used to model of article content, and users' interests.

### 3.2.1 Article Indexing

The ambition is to create a machine understandable representation of the content of each article, so as to compare with profiles. The unstructured information contained in articles is analyzed. Two kinds of information can be distinguished: explicit pieces of information (e.g. places, persons, organizations and so on) and implicit pieces of information (e.g. the theme of each article). In the system, the theme is one criterion, corresponding to the main event related by the article. Company experts have predefined a hierarchical list of themes wherein each article must be classified. The tasks of information extraction, annotation, indexing are done with the help of the GATE platform (Cunningham 2002). A

web interface was developed. It enables the writers to create articles. Results of the analysis are presented in it. They can be validated/corrected by the writer.

### Analysis

Some post processes are applied into articles, such as tokenizers, sentence splitters, POS taggers, gazetteers (which use the knowledge base lexical resources as a dictionary) before JAPE patterns matching engine. In the first prototype, we used handmade lexico-semantic patterns. The aim is to extract important entities (explicit pieces of information like Persons, Organizations, places and dates). Results of analysis are hand-checked, corrected and validated. Implicit pieces of information are specifically handmade.

### Population

For each article analyzed, an instance of the concept article is created in the knowledge base, so as to represent the article. Automatic analysis, correction done by hand and specifications of not automatically funded criteria are used to characterize the content of each article. Relations are created in the knowledge base between the article's instance, and criterion's instances (result of the analysis). Instances of criteria and relations with the instances of articles permit to index the article and create a semantic representation understandable by the machine.

### 3.2.2 Profiles Indexing

In the company, sellers are in charge of understanding the needs of each customer. Several phone calls are necessary so as to acquire future customers. During the phone call the seller proposes a free trial period. This helps to create a first handmade profile for each customer by an expert, and avoids the problem of a cold start, common to all content-based recommender systems.

The profile indexing process is the same as the articles. A profile instance is created in the knowledge base. Relations are created between the profile instance and criteria instances. A web interface was developed to enable sellers to define / change the user profile. The choices are reflected in the Knowledge base that permits to index the profile, and create a semantic representation understandable by the machine.

## 3.3 Recommendation

The recommendation task is mainly based on the comparison between the profile and available items. The knowledge base is used by the system like an index, and profiles and articles are presented by a set of instances and relations. For each article validated by writers, the full content is inserted in the database and a representation of the article is created in the knowledge base. An instance of the concept article and an instance of each relation “isAbout” is created to link the article with its criteria. For each profile made by the sellers, the representation is created in the knowledge base. An instance of the concept profile is created, and each “isInterestedIn” relation between the profile instance and its criteria are instanced. We present the comparison method used in the following section.

## 4 COMPARISON

The comparison task enables to evaluate the pertinence of an article for a profile, via a similarity measure between them.

**Definition 3** similarity:  $SIM(x, y): I \times I \rightarrow (0,1)$  is a function to measure the degree of similarity between  $x$  and  $y$ . The similarity function can satisfy some properties:

$\forall x, y \in I \quad SIM(x, y) \geq 0$  Positiveness

$\forall x, y \in I \quad SIM(x, x) = 1$  Reflexivity

The last one is the symmetry,  $\forall x, y \in I \quad SIM(x, y) = SIM(y, x)$  but in our context we want an asymmetric function, because we consider that comparing profiles and articles is not the same as comparing two articles.

### 4.1 VSM

An approach based on the vector space model (Salton 1970) was used in the prototype. Articles and profiles are represented by vectors on a space wherein each dimension is a potential instance of criteria. Several methods can be used to compare vectors; the most common is the cosine similarity.

$$SIM(\vec{a}, \vec{p}) = \cos \theta = \frac{\vec{a} \cdot \vec{p}}{|\vec{a}| \times |\vec{p}|} \quad (1)$$

An article can be defined like a vector of instances of Entities and criteria. For the recommendation task in the prototype only instances of criteria are used.  $\vec{a} = \{i_1, i_2, \dots, i_t\}$

Wherein, an article is represented by a vector  $\vec{a}$  composed by a set of instances  $i_x$ .  $i_x \in I'$

Instances  $i_x$  are instances of concepts belonging to the set of concepts  $C'$  defined by indexing criteria.

The set  $I'$  contains all instances of all concepts in the set  $C'$ .  $I' \subseteq I$

The set  $I'$  is a subset of the set of all instances  $I$  of the Knowledge Base  $\mathcal{K}$ .

A profile can be defined as a vector of instances of criteria.  $\vec{p} = \{i_1, i_2, \dots, i_t\}$

Where, a profile is represented by a vector  $\vec{p}$  composed of a set of instances  $i_x$ .

In our implementation, one vector for each criterion is used. This enables to weigh or use different similarity measures (e.g. cosine, jacquard, Euclidian) for each criterion. For example, location, theme and sectors are much more important than project and organization and so highly weighed.

$$SIMF(\vec{a}, \vec{p}) = \frac{\sum w_c SIM_c(\vec{a}_c, \vec{p}_c)}{\sum w_c} \quad (2)$$

$SIM_c(\vec{a}_c, \vec{p}_c)$  Is the similarity between profile  $\vec{p}_x$  and article  $\vec{a}_x$  and  $\mathcal{W}_x$  the coefficient for a specific criterion  $c$ .  $\forall i_{x,c} \in I'_c \quad \vec{p}_c = \{i_{1,c}, i_{2,c}, \dots, i_{t,c}\}$  And  $\vec{a}_c = \{i_{1,c}, i_{2,c}, \dots, i_{t,c}\}$ . One or more concepts are defined for each criterion.

Methods from the information retrieval fields can be used to enhance the recommendation. One of the first systems using external knowledge to improve the understanding of user needs is (Voorhees 1994). The Voorhees approach used WordNet (Fellbaum 1998) to provide a query expansion. We can translate this kind of method to recommender systems, instead of query expansion; we can name this method ‘profile expansion’. Middleton (Middleton et al. 2004) uses this method without naming it. Ijntema (Ijntema et al. 2010) also uses it, but unlike Middleton, he uses more powerful ontology relations (not just is\_a) to expand the user profile. In our system, the profile expansion takes the form of adding instances, e.g. if the user’ U profile shows an interest for company Co and in the knowledge base a symmetric relation like is\_aSubsidiaryOf is instanced with another company Co’, it is possible to add the other company to his profile. The Expanded VSM is developed in the following section.

### 4.2 Expanded Vectors

The 4.1 section presents an implementation of the similarity measure between two instances, by the creation of vectors of related instances. It was explained by Voorhees (Voorhees 1994) in the VSM that all dimensions are orthogonal and so, all elements of each vector are considered as independent. That is not really the case for the lexical used by Voorhees, and instances used in our system. In Voorhees’s method, the metronomic and

synonymic relations defined in WorldNet are used to add related lexical to the vector. In our case, relations between individuals exist, because we used a Knowledge base to manage the domain knowledge, and our criteria are hierarchically defined in it.

Meronomic relations exist between instances of Location. With Voorhees's method, it seems logical to expand the profile (the query for Voorhees) with all sub instances enclosed by the most general instance. For example, if the profile is interested in Bourgogne (Region, biggest Administrative division), it seems logical to add Cote d'Or and Yonne (Departement, sub administration division of Region) which are two departments within the Bourgogne region, and Dijon, Beaune, Chenove, Auxerre, which are cities within these departments. But, in the real case with this method, it is necessary to add four departments and its 2047 cities to the vector. The similarity between a profile interested in Bourgogne and an article about Dijon will be very low with this method.

So our method is to expand profile and articles, not only the profile and limiting the size of vector, instances added are including instances, not included (by meronomic relations). So our method is analogue to graph-based methods which search the common ancestor. For example, if the profile is interested by Dijon, Cote d'Or and Bourgogne, these are added to the vector. But if the profile is interested just by Bourgogne, nothing is added.

With this method the first drawback is solved, synonymic and meronomic relations between instances are managed, but the second still is not. Our similarity function is symmetric, but we want to compare a profile and an article to different things. The precision of an article must not have the same consequences as the precision of a profile. This problem is the subject of the following section.

### 4.3 Expanded Vectors

The previous section looked at how to take into account the relations between instances in the VSM. This section is about the managing of the difference of precision between profiles and articles. To solve this problem we used an intermediate vector for each criterion, a sub vector  $\vec{s}_c$  composed of common instances of the article  $\vec{a}_c$  and profile  $\vec{p}_c$  vectors for the criteria.

**Definition 4** Precision: In the hierarchy of concepts, more general concepts enclosed more specific one. In the hierarchy of instances it is the same. Instances from the top of the hierarchy are less specific than instances from the bottom.

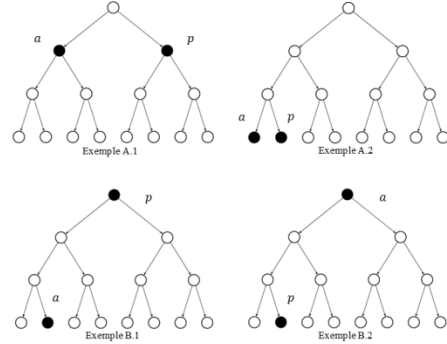


Figure 1: Examples of profiles and articles for one criterion.

If the Article is about an instance from the bottom and the Profile is interested in top instance (of the same branch), the similarity must be higher than if the Article is about a top instance and the Profile is inserted in an instance from the bottom of the hierarchy because there is a loss of precision if the profile is more specific than the article, so the article is less relevant.  $S_c = I'_{p,c} \cap I'_{a,c}$

$S_c$  is the sub set of common elements of the set of instances related to the profile  $I'_{p,c}$  and the article  $I'_{a,c}$ .  $\forall i_{x,c} \in S_c \vec{s}_c = \{i_{1,c}, i_{2,c}, \dots, i_{t,c}\}$

The vector  $\vec{s}_c$  is composed by elements of the set  $S_c$ .

$$simGlob_c(\vec{a}_c, \vec{p}_c) = \frac{w_{1,c} \times sim_c(\vec{a}_c, \vec{s}_c) + w_{2,c} \times sim_c(\vec{p}_c, \vec{s}_c)}{w_{1,c} + w_{2,c}} \quad (3)$$

It is possible to weigh differently the precision of the profile and the precision of the article with this method. In our implementation we used  $w'_{1,c} = 1$  and  $w'_{2,c} = 4$ , because we consider that the difference of precision of the profile mustn't influence the final note over 20%. However, it is possible to change the value, and it is also possible to use different values according to the criterion.

$$SIMF(\vec{a}, \vec{p}) = \frac{\sum w_c \times simGlob_c(\vec{a}_c, \vec{p}_c)}{\sum w_c} \quad (4)$$

The final similarity  $SIMF(\vec{a}, \vec{p})$  value is the sum of the similarity measure for each criterion. This measure is used for the ranking of articles proposed to the user according to his profile.

With the method the similarity value for the case A.2 (figure 1.) is higher than the case A.1, because in the case A.2, a and p have more common ancestors than in the case A.1. Moreover, the cases B.1 and B.2 figure the problem of precision. With our asymmetric method the value of similarity between a and p in the case B.1 is higher than in the case B.2 because the user needs a specific and the

article information (about this criterion) very general relative to the user needs. So the article is less relevant for the user. The following section presents the conclusion and future work.

## 5 CONCLUSION

In this paper we have presented the adaptation of a standard VSM recommender system to our specific method of indexing (e.g. articles and profiles are semantically defined in the knowledge base via relations with the domain knowledge already defined in it). We first presented the context, our goal, and the existing approaches, then we explained the architecture in detail. Finally we explained the specific task of comparison that we adapted to our case.

## ACKNOWLEDGEMENTS

This project is founded by the company Actualis SARL and the financing CIFRE research grant from the French agency ANRT.

## REFERENCES

- Billus, D., Pazzani, M.J., 1999. A Personal News Agent that Talks, Learns and Explains. In: The Third Annual Conference on Autonomous Agents, ACM, pp. 268–275.
- Ahn, J., Brusilovsky, P., Grady, J., He, D., Syn, S.Y., 2007. Open User Profiles for Adaptive News Systems: Help or Harm? In: 16th international conference on World Wide Web, ACM, pp. 11–20
- Cunningham, H., 2002 GATE, A General Architecture for Text Engineering. *Computers and the Humanities* 36 pp. 223–254
- Ehrig, M., Haase, P., Stojanovic, N., Hefke, M., 2005 Similarity for Ontologies A Comprehensive Framework. ECIS. Regensburg, Germany.
- Fellbaum, C., ed., 1998. *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, MA
- Gabrilovich, E., Dumais, S., Horvitz, E., 2004. Newsjunkie: providing personalized newsfeeds via analysis of information novelty. In *WWW*, pp. 482–490.
- Getahun, F., Tekli, J., Richard, C., Viviani, M., Yetongnon, K., 2009. Relating RSS News/Items. In: 9th International Conference on Web Engineering, Springer, pp. 442–452
- IIntema, W., Goossen, F., Frasinca, F., Hogenboom, F., 2010 Ontology-Based News Recommendation. In: *International Workshop on Business intelligence and the Web, BEWEB, EDBT/ICDT Workshops*. pp. 16:1–16:6. ACM, New York, USA
- Li, L., Wang, D., Li, T., Knox, D., Padmanabhan, B., 2011. SCENE: A scalable two-stage personalized news recommendation system. In *Proc. the 34th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, Beijing, China, Jul. 25–29, 2011, pp.124–134.
- Liu, J., Dolan, P., Pedersen, E.R., 2010. Personalized news recommendation based on click behavior. In *Proceedings of the 15th International Conference on Intelligent User Interfaces*, pp. 31–40. ACM.
- Middleton, S.E., Shadbolt, N.R., Roure, D.C.D., 2004. Ontological User Profiling in Recommender Systems. *ACM Transactions on Information Systems* 22, pp. 54–88
- Nageswara Rao, K., and Talwar, V.G., 2008. Application Domain and Functional Classification of Recommender Systems A Survey, *Journal of Library & Information Technology*, Vol. 28, No. 3: 17–35.
- Piskorski, J., Tanev, H., Wennerberg, P.O., 2007. Extracting Violent Events From OnLine News for Ontology Population. In: 10th International Conference on Business Information Systems, BIS. *Lecture Notes in Computer Science*, vol. 4439, pp. 287–300. Springer-Verlag Berlin Heidelberg.
- Popov, B., Kiryakov, A., Kirilov, A., Manov, D., Ognyanoff, D., Goranov, M., 2003. KIM Semantic Annotation Platform.
- Resnick, P., Iacovou, N., Suchak, M., Bergstrom, P., Riedl, J., 1994. GroupLens: An open architecture for collaborative filtering of netnews. In *Proceedings of the 1994 ACM Conference on Computer Supported Cooperative Work*, Chapel Hill, North Carolina, United States. ACM Press, New York, pp. 175–186.
- Salton, G., 1970. *The SMART retrieval system : Experiments in automatic document processing*. Prentice Hall.
- Stumme, G., Ehrig, M., Handschuh, S., Hotho, A., Maedche, A., Motik, B., Oberle, D., Schmitz, C., Staab, S., Stojanovic, L., Stojanovic, N., Studer, R., Sure, Y., Volz, R., Zacharias, V., 2003. *The Karlsruhe view on ontologies*. Technical report, University of Karlsruhe, Institute AIFB.
- Tanev, H., Piskorski, J., Atkinson, M., 2008. Real-Time News Event Extraction for Global Crisis Monitoring. In: 13th International Conference on Natural Language and Information Systems: Applications of Natural Language to Information Systems, NLDB. *Lecture Notes in Computer Science*, vol. 5039, pp. 207–218. Springer-Verlag Berlin Heidelberg
- Voorhees, E., 1994. *Query Expansion using Lexical-Semantic Relations*.
- Werner, D., Cruz, C., Nicolle, C., 2012. Ontology-based Recommender system of economic articles; In *Proceedings of the 2012 Webist Conference*, Porto, Portugal.