

A Query Refinement to Concept-based Information Retrieval

Silvia Calegari

DISCo
Università di Milano – Bicocca
V.le Sarca 336/14,
20126 Milano (Italia)
calegari@disco.unimib.it

Elie Sanchez

LIF, Faculte de Medecine
Universite Aix-Marseille II
27 Bd Jean Moulin,
13385 Marseille Cedex5, (France)
elie.sanchez@medecine.univ-mrs.fr

Abstract

This paper presents an approach to semantic Information Retrieval, based on the use of a fuzzy conceptual structure, called Object-Fuzzy Concept Network (O-FCN). It consists of a set of objects \mathbf{O}_{DB} stored in a database and of a Fuzzy Concept Network (FCN), which is a complete weighted graph, with edges among concept nodes, we called correlations. It was already introduced and described an Information Retrieval Algorithm involving an O-FCN. The algorithm is validated in a classical, crisp, case versus a fuzzy case. The users had the possibility to assign weights of importance to concepts modifiers in queries in order for the system to produce more relevant documents. For a better evaluation it has been finally introduced a new fuzzy accuracy measure that allowed to improve the results.

Keywords: Semantic Information Retrieval, Concept Modifiers, Object-Fuzzy Concept Network, Fuzzy Accuracy Measure.

1 Introduction

In the Semantic Web area of research, attention has recently been focused in two directions: understanding the requests from the users and finding documents that best satisfied their requests. An increasing number of

approaches have been presented in concept-based Information Retrieval [1]. We proposed to combine fuzzy ontologies to objects (stored in a data base) in order to search new documents semantically correlated to user's queries. The result is an Information Retrieval algorithm that is tested in the crisp versus the fuzzy case. Besides the fuzzy recall and fuzzy precision criteria, it is introduced a fuzzy accuracy measure to better evaluate the algorithm.

2 A new way of using concept modifiers

Humans typically use adverbs like *very* or *more or less* to formulate requests. For instance, in a e-commerce context, it is critical to distinguish between a customer who is interested in technical details and one who is *very* interested in these details. In [2] Zadeh introduced so-called linguistic hedges that could be viewed as operators acting on a fuzzy set representing the meaning of its operand. In detail, he has been defining precision concept modifiers. These linguistic hedge operations can be classified into two categories: concentration and dilation. The effect of dilation is opposite to that of concentration. The result of applying a concentration operator to a fuzzy set results in the reduction in the magnitude of the grades of membership, which is relatively large for elements with low membership. Technically, Zadeh achieves this by simply raising the degree of membership to the β -th power, where $\beta > 1$ (for concentration operator) is a constant. For example,

for *very* it is usually assigned $\beta = 2$ and so, if “Cabernet has a dry taste with value 0.8”, then “Cabernet has a *very* dry taste” will have value 0.64 (i.e., $0.8^2 = 0.64$).

Obviously, we obtain an opposite effect using dilation operator. Indeed, a dilation operator raises the degree of membership to the β -th power, where $0 < \beta < 1$, thereby increasing the degree of membership of elements with small value.

As previously said, humans use linguistic adverbs and adjectives to specify their interests and needs. For example, a user can be interested in finding in a web portal “a very fast car”, “a wine with a very astringent taste”, and so on. So that the necessity to handle the richness of natural languages used by humans emerges.

In this paper, we propose a new use of precision concept modifiers. The idea is to give the possibility to the user to attribute a weight to a concept written in a query. In this way, the weight assigns a different importance to each concept and so, during the calculation of relevance, the document having the concept with a major weight will be more relevant than another document identified with a concept associated with a lower weight. In the literature, one can find a lot of models in which a user can assign a weight to queries. But, in our approach we utilise the semantic features of concept modifiers (for an example, see SubSection 3.2).

2.1 Path Discovery Query

Path discovery query [3] is the most powerful and arguably the most interesting form of semantic queries. This type of query involves a number of entities (possibly just a pair of concepts) and attempts to return a set of paths (including relationships and intermediate concepts) that connect the concepts in the query. Each computed path represents a semantic association of the named concepts.

Formally, a query has the form $q = m_{\beta_1}C_1, m_{\beta_2}C_2, \dots, m_{\beta_n}C_n$ where $C_i \in \mathbf{C}$ are concepts of a fuzzy ontology and m_{β_i} defines a precision concept modifier. Using these short queries we utilised a simple model where it

is not necessary to have a common parser in order to understand and interpret the user’s queries.

In this work, we have adopted the Khang et al.’s [4] algorithm in order to give a semantic interpretation (i.e., a specific degree according to the context) to the concept modifiers. Indeed, this algorithm allows to define a concept modifier with a length not known a priori. For example, given a finite set of fuzzy modifiers like *little*, *very*, a possible set of combinations will be *veryverylittle*, *little*, *veryvery* ... In this way a dynamic set of modifiers, not predictable by the expert, can be obtained. For the new methodology proposed in this paper, we have to use precision concept modifiers in an opposite situation than the one presented in the literature (see SubSection 2.2).

2.2 O-FCN

In the Semantic Web domain, a crucial topic is to define a dynamic knowledge of a domain adapting itself to the context. A goal is to retrieve semantic information in order to satisfy the user’s query. Indeed, an increasing number of approaches to Information Retrieval have proposed models based on concepts rather than on keywords in the last years. In [5] it has been proposed a system that allows to achieve these objectives using a new non-taxonomic fuzzy relation, named *correlation* (here *corr*). It consists in the determination of a semantic correlation among the concepts that are searched together, for example, in a query or when a document has been inserted into a database. The notion of Fuzzy Concept Network (FCN) is properly suited for path discovery semantic query. Indeed navigating it, new non a priori predictable semantic associations among concepts are obtained. The FCN is extended incorporating Database Objects so that, concepts and information can similarly be represented in the network [6].

Definition

An Object-Fuzzy Concept Network (O-FCN) is a weighted graph $\mathcal{N}_{fo} = \{\mathbf{O}_{\mathbf{DB}}, \mathcal{N}_f\}$, where $\mathbf{O}_{\mathbf{DB}}$ is the set of the objects stored in the database and $\mathcal{N}_f = \{\mathbf{C}, F, m\}$ is a Fuzzy Con-

cept Network (FCN). Each object is described by the concepts of the FCN, i.e. $\forall o_i \in \mathbf{O}_{DB} o_i = \{c_1, \dots, c_n\}$ where $c_i, \dots, c_n \in \mathbf{C}$.

The set \mathbf{O}_{DB} identifies all the information that is contained into the database, such as documents, digital pictures, videos, and so on. In this work we don't have investigated the indexing problem. In a concept-based Information Retrieval model the meaning of a text (or document) depends on conceptual relationships among concepts. Thus, to determine which are the representative concepts of an object is a crucial topic. At the moment, in this approach, we have directly analysed the database where this phase has been just performed.

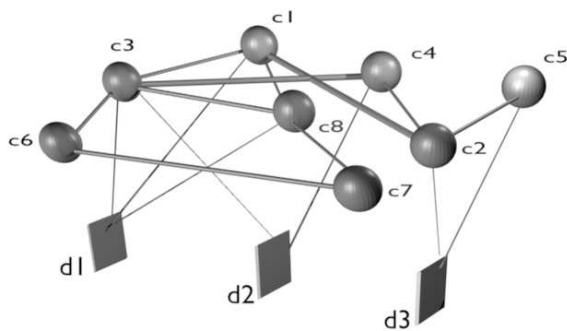


Figure 1: An example of an O-FCN.

In this case, the edges of the O-FCN define the correlation relationships between concepts in \mathbf{C} , i.e. $F := corr$, where $corr : \mathbf{C} \times \mathbf{C} \mapsto [0, 1]$. The different thickness of the links identifies how strongly the concepts are correlated. The thicker the link the more correlated are the two concepts (i.e. the closer to 1 is the fuzzy value).

In [6] it has been introduced and described an Information Retrieval algorithm using O-FCN. This algorithm allows to derive a unique path among the concepts involved in the query in order to obtain the maximum semantic associations in the knowledge domain. Now a very general step-by-step description of this new algorithm [7] is given below (see also Fig. 2):

The O-FCN has been involved in all the steps of the algorithm in order to semantically enrich the results that were obtained. The al-

| | |
|--|-------------------------------------|
| 'O-FCN'-IR Search (E_q : concept vector) | |
| 1: | 'O-FCN'-based E_q expansion |
| 2: | 'O-FCN'-based documents extraction |
| 3: | 'O-FCN'-based relevance calculation |
| return ranking of the documents | |

Figure 2: New Information Retrieval Algorithm using O-FCN

gorithm input is a vector E_q identifying the terms in the query. The O-FCN is used to calculate the relevance of the documents in order to sort them in decreasing order. In particular, thanks to the O-FCN characterising functions F and m , the weights for the concepts c_i in each selected document are determined according to the following equation:

$$w(c_i) = m(c_i)^{\frac{1}{\beta_{c_i}}} \cdot \frac{\sum_{c_j \in K, c_j \neq c_i} [F(c_i, c_j)]^{\frac{1}{\beta_{c_i, c_j}}}}{l_{c_i}^2} \quad (1)$$

where K is the set of the concepts identified by the specific document, $\beta_{c_i, c_j} \in \mathbf{R}$ is a precision concept modifier value used to alter the fuzzy value associated with the concepts. In this application $F(c_i, c_j) = corr(c_i, c_j)$, $m(c_i) = corr(c_i, c_i)$ and l_{c_i} is the level of the concept. The level stores the position of the concept into the graph itself: the greater this value the more semantically distant are the words written into the query. In this formula the precision concept modifiers are used in contrast with the classic approach presented in Section 2. For example, let us consider the query $q := car, very_red$ where red in the O-FCN has value of 0.2, i.e. $corr(red, red) := 0.2$. In the classical approach we obtain $0.2^2 := 0.04$, where *squaring* is the interpretation of *very*. But in this new approach a user might be interested in finding information about *car* having mainly color *red*. So, in the new system we have $0.2^{\frac{1}{2}} \simeq 0.45$ (as default, the concept *car* is assigned $\beta = 1$). In this way, the document will have a relevance value higher than the one from the usual methodology.

From another point of view, in our case, what is usually considered as a concentration mod-

ifier becomes a dilation modifier and vice versa.

3 Validation Test

This Section is divided as follows: in the first part it is introduced the environment for the experiment related with the approach proposed in the previous Section, whereas in the second part the analytic study is reported.

3.1 Description of the experiment

A creative learning environment is the context chosen to test the new Information Retrieval algorithm based on O-FCN. In particular, the ATELIER (Architecture and Technologies for Inspirational Learning Environments) project has been considered. ATELIER is an EU-funded project that was part of the Disappearing Computer initiative. The aim of this project was to build a digitally enhanced environment, supporting a creative learning process in architecture and interaction design education. The work of the students was supported by many kinds of devices (e.g., large displays, RFID technology, barcodes, ...) and a hyper-media database (HMDB) has been used to store all digital materials produced. Every day the students have created a very large amount of documents and artifacts and they collected a lot of material (e.g., digital pictures, notes, videos, and so on). In this context, it emerges that the evolution of the O-FCN is mainly given by the words of the documents inserted in a HMDB and from the concepts written during the definition of a query by the students.

We have studied the dynamic evolution of the O-FCN by examining 485 documents and 200 queries of the students (a history file has been used). For each query a user had the opportunity to include up to 5 different concepts and the possibility to semantically enrich his/her requests by using the following list of concept modifiers: *{little, enough, moderately, quite, very, totally}*.

The algorithm previously presented has been

tested in two different situations: classical and fuzzy approaches. In the first case, the crisp situation has been reported assigning value 1.0 to the correlations values and without considering the concept modifiers into the queries of the students. Instead, in the last case, all the parameters described in this paper have been considered.

3.2 Analytic Considerations

Now, all the analytic considerations evaluated in our analysis are reported. Examining the whole path that we have followed it is possible to have a clearer vision of the final obtained results.

Recall and precision measures are the usual parameters used in Information Retrieval Systems (IRs) in order to evaluate information retrieval algorithms. In detail, *recall measure* is defined as $R := \frac{|R_T \cap R_L|}{|R_L|}$, i.e. it is the proportion of relevant documents (R_L) that are retrieved (R_T), and *precision measure* is defined as $P := \frac{|R_T \cap R_L|}{|R_T|}$, i.e. it is the proportion of retrieved documents that are relevant. These evaluations are mainly used to compare algorithms involving different techniques. Here, we investigated the same algorithm in two different situations: crisp and fuzzy cases. For this reason, *fuzzy recall and fuzzy precision* [8] have been applied in our validations [6]. In detail,

$$R_F := \frac{\sum_{d_i \in Q_\theta} \mu_Q(d_i)}{\sum_{d_i \in \mathcal{D}} \mu_Q(d_i)} \quad (2)$$

where \mathcal{D} is the set of all documents, i.e. $\mathcal{D} := \{d_1, d_2, \dots, d_{n-1}, d_n\}$, and Q_θ is the θ -cut of Q defined as $Q_\theta := \{d_i \in \mathcal{D} \text{ s.t. } \mu_Q(d_i) \geq \theta\}$. Note that Q_θ could be rewritten as $R_T(\theta)$, meaning that it is the (crisp) set of documents that are retrieved, above a threshold θ . Moreover, all documents are relevant, but *Relevance* is Fuzzy, so documents are relevant with a degree (of course a degree equal to zero for a document means that it is not relevant at all).

$$P_F := \frac{\sum_{d_i \in Q_\theta} \mu_Q(d_i)}{|Q_\theta|} \quad (3)$$

Table 1 reports the average values of fuzzy precision and fuzzy recall for the 200 queries performed in the two approaches. Retrieved documents are ranked up to a theta threshold (θ). In particular, we have chosen three values of θ (0.35, 0.50 and 0.75) to validate the algorithm in different situations. Ideally, high precision and high recall values (for crisp and fuzzy formulae) are both desired [9].

Table 1: Average values of Fuzzy Precision and Fuzzy Recall in the fuzzy and crisp cases.

| Fuzzy Case | | |
|----------------|--------------|-----------|
| θ value | F. Precision | F. Recall |
| 0.35 | 0.573 | 0.612 |
| 0.50 | 0.602 | 0.523 |
| 0.75 | 0.912 | 0.221 |
| Crisp Case | | |
| θ value | F. Precision | F. Recall |
| 0.35 | 0.590 | 0.622 |
| 0.50 | 0.604 | 0.593 |
| 0.75 | 0.942 | 0.234 |

Table 1 reports non-significant differences, although in the crisp case the values for both of the two measures are always higher. So, apparently, it can seem that by using the crisp methodology better results than fuzzy ones are obtained.

During this work, other exams were needed in order to show that, by using a fuzzy approach, accurate results were obtained. Indeed, in the fuzzy case it has been observed a better *variability* of relevant documents. This result was derived from the analysis of *coefficient variance* based on fuzzy precision measure (here CV_P). In detail,

$$CV_P := \left(\frac{\sigma}{P_F}\right) \cdot 100 \quad (4)$$

where σ is the standard deviation calculated on the relevance of the documents and P_F is the fuzzy precision. Thus, it is a useful statistic for comparing the degree of variation from one data series to another. In general, the larger this number, the greater the variability in data.

Analysing all the queries, it emerged that the fuzzy methodology exhibits higher CV_P values than the crisp one. This means that the fuzzy approach identifies refined results. For instance, we can examine the query number 69, i.e., $q_{69} = \text{"city, castle"}$. In this case a user was interested in finding documents where a "castle" is located into a "city". For $\theta = 35$ we obtained 14 documents. Table 2 reports the two results for the crisp and fuzzy cases.

Table 2: Relevant documents for $\theta = 35$ and query number 69.

| Fuzzy Case | |
|------------|-------------|
| Relevance | # Documents |
| 100% | 1 |
| 85% | 1 |
| 50% | 11 |
| 40% | 1 |
| Crisp Case | |
| Relevance | # Documents |
| 100% | 2 |
| 50% | 12 |

Thus, by using a fuzzy methodology it appeared a wider variability than in the crisp case. This means that in the fuzzy approach more accurate classifications of the information are obtained. In fact, Figure 3 depicts the difference of CV_P between fuzzy and crisp approaches with respect to the θ values assumed (i.e. $\theta = 0.35$, $\theta = 0.50$ and $\theta = 0.75$). In the fuzzy case we can observe higher CV_P values for the fuzzy case, for all the queries analysed. Thanks to this Figure the considerations above are consolidated.

In this paper we now define and investigate a new fuzzy measure for IRs, named *fuzzy accuracy*. Indeed, precision and recall are measures from the domain of information retrieval and they are focussed only on the class relevant of the ranked information. There is no interest in measuring the degree in which the system does not retrieve irrelevant instances. In the literature, Weiss and Kulikowski [10] use a basic accuracy measure that is simply defined as the ratio of correctly assigned

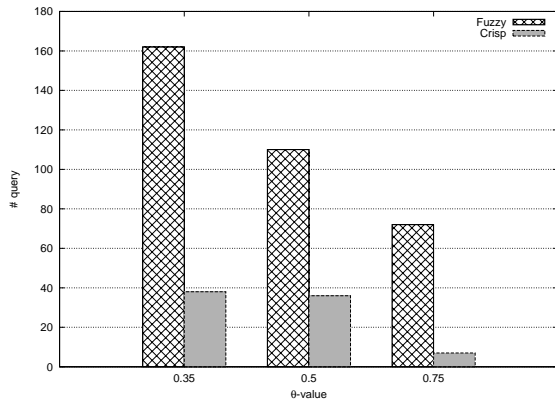


Figure 3: Trend of CV_P value for each query.

items over the total of items. Namely, in the crisp case, *accuracy measure* is defined as $A := \frac{(|R_L \cap R_T|) \cup (|\bar{R}_L \cap \bar{R}_T|)}{|\mathcal{D}|}$, i.e. the fraction of its classifications that are correct.

Thus, having proposed fuzzy precision and fuzzy recall, we define now a fuzzy accuracy measure as:

$$A_F := \frac{\sum_{d_i \in Q_\theta} \mu_Q(d_i) + \sum_{d_i \in \bar{Q}_\theta} (1 - \mu_Q(d_i))}{|\mathcal{D}|} \quad (5)$$

Table 3 reports the average values of fuzzy accuracy measures, for the different θ , as follows:

Table 3: Average values of Fuzzy Accuracy in the fuzzy and crisp cases.

| θ value | Fuzzy Accuracy | |
|----------------|----------------|------------|
| | Fuzzy Case | Crisp Case |
| 0.35 | 0.330 | 0.344 |
| 0.50 | 0.326 | 0.492 |
| 0.75 | 0.134 | 0.157 |

Also in this case, we obtained higher values in the crisp evaluation than in the fuzzy one. This result is coherent with the analysis of Table 1 where, apparently, better results are produced in the crisp situation.

A next step has been made and some new empirical evaluations have been proposed. Indeed, analysing CV_P factors it is shown how a major variability of relevance is obtained in the fuzzy case. From another point of view,

this means that in the crisp case there are groups of clusters of documents. Namely, in the non-fuzzy situation higher relevance values than fuzzy ones are obtained. In Table 2 an example is reported. In this way, it is given a possible explanation of the fact that we have lower fuzzy precision and fuzzy recall values in the fuzzy case, although a major identification (or variability) of the documents relevance is obtained. Thus, we have analysed the relevance values in order to confirm this fact. In Table 4 the average values of the weighted arithmetic mean calculated on the documents relevance for each query are reported. As expected in the crisp case higher relevance values than fuzzy ones are obtained.

Table 4: Average values of the weighted arithmetic mean in the fuzzy and crisp cases.

| θ value | Weighted arithmetic mean | |
|----------------|--------------------------|------------|
| | Fuzzy Case | Crisp Case |
| 0.35 | 44 % | 46 % |
| 0.50 | 48 % | 49 % |
| 0.75 | 51 % | 53 % |

Thus, we have made again the same analysis (i.e., Fuzzy Precision, Fuzzy Recall and Fuzzy Accuracy) considering new aspects. In particular, we have analysed queries where more or less the same relevance for the fuzzy and crisp case was obtained. Namely, we have calculated the difference of the weighted arithmetic mean of the relevance on these cases and only the queries having this value close to 0 have been considered.

Accuracy is a measure often used for evaluating machine learning classification problems. In this case, we have an O-FCN that evolves in time according to the content of the queries. The goal of this dynamic behaviour is to adapt the O-FCN's knowledge to the context. So, we have also included, in this new evaluation, the queries with a higher fuzzy accuracy value in the fuzzy case than in the crisp one. Table 5 reports fuzzy precision and fuzzy recall given for this new analysis.

Now, the fuzzy case presents a better result than the previous one (see Table 1). Thus,

Table 5: Average values of Fuzzy Precision and Fuzzy Recall in the fuzzy and crisp cases.

| Fuzzy Case | | |
|----------------|--------------|-----------|
| θ value | F. Precision | F. Recall |
| 0.35 | 0.585 | 0.717 |
| 0.50 | 0.621 | 0.701 |
| 0.75 | 0.925 | 0.224 |
| Crisp Case | | |
| θ value | F. Precision | F. Recall |
| 0.35 | 0.581 | 0.696 |
| 0.50 | 0.619 | 0.753 |
| 0.75 | 0.920 | 0.211 |

considering an analogous situation, the fuzzy methodology reports some advantages. Finally, we have calculated fuzzy accuracy measures (see Table 6). It is evident that a major correctness of relevance in the fuzzy case is established. In particular, analysing $\theta = 0.50$ almost all of the queries are considered in the new evaluation. This means that the fuzzy methodology allows to identify more precisely the relevance of the documents. In particular, it allows to better classify the information that is relevant for the users.

Table 6: Average values of Fuzzy Accuracy in the fuzzy and crisp cases.

| Fuzzy Accuracy | | |
|----------------|------------|------------|
| θ value | Fuzzy Case | Crisp Case |
| 0.35 | 0.410 | 0.405 |
| 0.50 | 0.386 | 0.290 |
| 0.75 | 0.153 | 0.149 |

For example, examining the query number 142, i.e. $q_{142} = \text{“veryNapoleon, portrait”}$. In this case a user was interested in finding documents about “Napoleon” and possibly having some information on its “portrait”. Let us note that this statement is correct using the new semantics proposed in this paper. Indeed, the user wants to find information where the concept “Napoleon” is more important than “portrait”. Using $\theta = 35$ we obtained 254 documents. When a lot of documents is obtained, a suited classification of their rele-

vance is needed. Indeed, in the fuzzy case a higher fuzzy accuracy than in the crisp case is obtained, i.e. $A_F := 0.482$ and $A_F := 0.479$ respectively.

Table 7 reports the two results for the crisp and fuzzy cases.

Table 7: Relevant documents for $\theta = 35$ and query number 142.

| Fuzzy Case | |
|------------|-------------|
| Relevance | # Documents |
| 71% | 1 |
| 50% | 156 |
| 49% | 43 |
| 48% | 13 |
| 45% | 1 |
| 44% | 2 |
| 43% | 2 |
| 42% | 12 |
| 41% | 4 |
| 40% | 6 |
| 39% | 4 |
| 36% | 8 |
| 35% | 2 |
| Crisp Case | |
| Relevance | # Documents |
| 85% | 1 |
| 50% | 182 |
| 42% | 71 |

In this case, not only a greater variability of relevance is given in the fuzzy approach, but also a major identification of relevance is obtained. Let us consider the first document where a different relevance value in the crisp case and in the fuzzy case has been calculated. In the non-fuzzy methodology this document has been over estimated.

As previously said, in this new evaluation some empirical considerations are given. But it needs a deeper analysis: for instance, to find an analogy between the fuzzy accuracy measure and CV_P analysis.

4 Concluding Remarks

In this paper it has been proposed a new use of precision concept modifiers. Users could assign weights of importance in concepts involved in queries, with the effect to better fit their needs. In order to improve the evaluation of our Information Retrieval algorithm, we have introduced and defined a fuzzy accuracy measure. A wider use of this criteria will be expected in future applications. In particular when comparing our algorithm to other ones from the literature [11, 12], but this is a future work.

Acknowledgements

The work presented in this paper had been partially supported by the ATELIER project (IST-2001-33064).

References

- [1] Croft, W.B., Das, R. (1990). Experiments with query acquisition and use in document retrieval systems. In: SIGIR '90: Proceedings of the 13th annual international ACM SIGIR conference on Research and development in information retrieval, New York, NY, USA, ACM (1990) 349–368
- [2] L. A. Zadeh (1972). A fuzzy-set-theoretic interpretation of linguistic hedges. *Journal of Cybernetics*, 2(3):4–34.
- [3] B. Arpinar A. Sheth and V. Kashyap (2002). Relationships at the Heart of Semantic Web: Modeling, Discovering, and Exploiting Complex Semantic Relationships. Technical report, LSDIS Lab, Computer Science, University of Georgia, Athens GA 30622.
- [4] T. D. Khang, H.P. Störr, and S. Hölldobler (2002). A fuzzy description logic with hedges as concept modifiers. In *Third International Conference on Intelligent Technologies and Third Vietnam-Japan Symposium on Fuzzy Systems and Applications*, pages 25–34.
- [5] S. Calegari and F. Farina (2007). Fuzzy Ontologies and Scale-free Networks Analysis. *International Journal of Computer Science and Applications*, IV(II):125–144.
- [6] S. Calegari and E. Sanchez (2007). A Fuzzy Ontology-Approach to improve Semantic Information Retrieval. Proceedings of the Third ISWC Workshop on Uncertainty Reasoning for the Semantic Web - URSW'07 (F. Bobillo, P. C. G. da Costa, C. D'Amato, N. Fanizzi, F. Fung, T. Lukasiewicz, T. Martin, M. Nickles, Y. Peng, M. Pool, P. Smrz, and P. Vojtas, eds.), CEUR Workshop Proceedings, vol. 327, CEUR-WS.org, 2007.
- [7] S. Calegari and E. Sanchez (2008). Object-Fuzzy Concept Network: an enrichment of Ontologies in Semantic Information Retrieval. *International Journal of the American Society for Information Science and Technology*. Submitted.
- [8] E. Sanchez and P. Pierre (1994). Fuzzy Logic and Genetic Algorithms in Information Retrieval. In T. Yamakawa, editor, *Proceedings of the 3rd Int. Conf. on Fuzzy Logic, Neural Nets and Soft Computing*, pages 29–35. Jono Printing Co..
- [9] G. Salton and M. J. McGill (1986). Introduction to Modern Information Retrieval. McGraw-Hill, Inc., New York, NY, USA.
- [10] S. M. Weiss and C. A. Kulikowski (1991). Computer Systems that Learn. Morgan Kaufmann.
- [11] P. Resnik (1995). Using Information Content to Evaluate Semantic Similarity in a Taxonomy. *IJCAI*, pages 448-453.
- [12] J. Jiang and D. Conrath (1997). Semantic similarity based on corpus statistics and lexical taxonomy. In *Proceedings of the 10th international conference on research in computational linguistics*, Taipei, Taiwan pages 19-33.