

## Model of an intellectual search engine represented by the fuzzy sets theory

Hristina Moneva

**Abstract:** *In this report an attempt to explain the role and place of the fuzzy set theory with its typical techniques and methods in the process of the knowledge representation is proffered. A usage of the fuzzy sets and fuzzy logic for organizing a knowledge base is described. Their role with their typical characteristics in the process of designing a system for searching and classifying documents based on knowledge is examined. A generalized model of an intellectual system for searching and classifying documents in the Web is proposed.*

**Key words:** *fuzzy sets, intellectual systems, systems based on knowledge, search engines, classification.*

### INTRODUCTION

When information is searched in Internet or LAN a request is send and a response is received. This is a list of documents and sometimes for the user it is hard to locate the desired knowledge. The received documents are not thematically grouped, not ordered and also it is not possible to achieve the level of adequacy that is needed.

To solve these problems a new approach for designing a search engine is proffered – an intellectual search engine based on knowledge. In its knowledge base will have data about documents, keywords that describe them and the domains that they belong to (basically). Systems, which have a kernel that is knowledge base or model of domain, are described with some high level language close to natural language. They are called intellectual [2]. With this approach, a person who publishes a few documents has to give them certain characteristics before adding them to the knowledge base. In this case they are the experts whose knowledge is stored in the intellectual system. In contrast to standard systems is that knowledge base is not static, but is getting larger all the time.

Attempting to formalize the common knowledge, the researchers faced the problem to represent adequately the vagueness in a conventional mathematical way.

The attempt to formalize the common knowledge met the problem to represent it by using conventional mathematical apparatus to solve it. In case we are dealing with a non-numerical quantity that cannot be measured against a numerical scale, than we cannot use a numerical universe. Than it is said that the elements were taken from a psychological continuum. An example of such a universe could be {large, medium, small, very small}. These characteristics usually are fuzzy and cannot be synonymously interpreted.

For that reason tasks that are solved in the intellectual systems often have to work with an uncertain knowledge which cannot be interpreted as a pure true or false (logical true/false or 0/1). There is a knowledge which reliability can be expressed with an intermediate number, for example 0.7.

How to represent such knowledge formally keeping the fuzziness and uncertainty? In order to solve these problems, the American mathematician Lotfi Zadeh proffered a formal apparatus of fuzzy logic and fuzzy arithmetic in the beginning of 70's. Later this trend marked the beginning for one of the areas of AI –soft computing [2].

Zadeh introduced one of the basic concepts in fuzzy logic – the concept of a linguistic variable. *LINGUISTIC VARIABLE* – this is a variable which values are words or sentences. The set of values it can take is called its term set. Each value in the term set is a fuzzy variable defined over a base variable. The base variable defines the universe of discourse for all the fuzzy variables in the term set. The fuzzy set is defined with some base scale B and membership function  $\mu(x)$ ,  $x \in B$  which values are into  $[0,1]$ . Consequently the fuzzy set B is a collection of ordered pairs  $(x, \mu(x))$ , where  $x \in B$ .

When the membership function gives the subjective expert level, the concrete value of base scale is the definition of a fuzzy set. This function is not a probabilistic that has an objective character and is not following any other mathematical dependencies.

The fuzzy sets give possibilities to define the subjective opinion of different individuals.

### FORMULATION OF THE TASK

Let **U** be the universe of discourse (for short just universe) containing all sets **A<sub>1</sub>, A<sub>2</sub>...A<sub>n</sub>** as elements. Elements of these sets are domains and sub domains that are organized in hierarchies. The set **U** is called universe, because the system pretends to be able to consist of all domains and sub domains that can come into consideration.

Let the elements of set **D** be the documents **d<sub>1</sub>, d<sub>2</sub>...d<sub>m</sub>**.

Let the set **R** be the binary relation from a set **U** to a set **D**, which gives information about classification of documents into domains and sub domains. The relation **R** is described by characteristic function of set **R** -  $\|\mu_{Rij}(A, d)\|$ ,  $i=1,2...n$ ,  $j=1,2...m$ , define as

$$\mu_{Rij}(A, d) = \begin{cases} 1, & ako(A_i, d_j) \in R \\ 0, & ako(A_i, d_j) \notin R \end{cases} [1].$$

Let **K** be the set with elements **k<sub>1</sub>, k<sub>2</sub>...k<sub>l</sub>** – keywords describing some document of the set **D** with the fuzzy relation **S** defined with  $\mu_S(d, k) \in [0,1]$ .

Now we have described all elements of the knowledge base, that is part of an intellectual classifying and searching system.

To define how these keywords are describing all documents in a domain, specify this domain too. To define this, we have to find max-min and min-max composition of **R** and **S** (or optimistic and pessimistic expectations [4]). Now we have the fuzzy relations **T<sub>1</sub>** and **T<sub>2</sub>**,

where for **T<sub>1</sub>**:  $\mu_{R \circ S}(A, k) = \sup_{d \in D} [\min(\mu_R(A, d), \mu_S(d, k))]$

and for **T<sub>2</sub>**:  $\mu_{R \circ S}(A, k) = \inf_{d \in D} [\max(\mu_R(A, d), \mu_S(d, k))]$  [3].

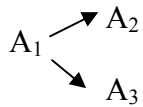
The relation **T<sub>2</sub>** is always set of zero elements because normally always will exist a keyword that does not appear in all domains and sub domains, otherwise this keyword is not important for us.

Now all documents are classified in different domains / sub domains with the set of keywords using the fuzzy sets and the fuzzy relations theory. The domains and sub domains are characterized with keywords too.

All this can be used for searching on keywords and letting us represent the result classified in different domains / sub domains consisting of a list of documents for each of them.

#### For example:

The sets **A<sub>2</sub>** and **A<sub>3</sub>** are subsets of **A<sub>1</sub>**:  $A_2 \subset A_1, A_3 \subset A_1$  or presented as oriented graph:



R	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>
A <sub>1</sub>	1	1	0	0	0	0
A <sub>2</sub>	0	0	1	1	0	0
A <sub>3</sub>	0	0	0	0	1	1

S	k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>
d <sub>1</sub>	0,9	0,6	0	0	0	0
d <sub>2</sub>	0,7	0	0,7	0	0	0
d <sub>3</sub>	0	0,4	0	0,6	0	0
d <sub>4</sub>	0	0	0	0,8	0,5	0
d <sub>5</sub>	0,3	0	0	0	0	0,4
d <sub>6</sub>	0	0,8	0	0	0	0,9

T <sub>1</sub>	k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>
A <sub>1</sub>	0,9	0,6	0,7	0	0	0
A <sub>2</sub>	0	0,4	0	0,8	0,5	0
A <sub>3</sub>	0,3	0,8	0	0	0	0,9

T <sub>2</sub>	k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>
A <sub>1</sub>	0	0	0	0	0	0
A <sub>2</sub>	0	0	0	0	0	0
A <sub>3</sub>	0	0	0	0	0	0

**LEVEL OF ADEQUACY (LEVEL OF FUZZY RELATION)**

A new concept of LEVEL OF ADEQUACY is introducing. This is a number into [0,1] which is set by the user and characterizes the relevancy between searched keyword (phrase) and proposed documents by the system.

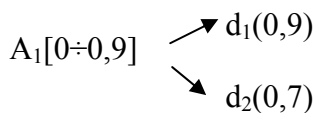
Let the searched word be **ks**. Let the searched level of adequacy to this word be  $\alpha_{ks}$ . Consequently the given level of adequacy is the level  $\alpha_{ks}$  of the fuzzy relation **S** of **DxK**, respectively  $\mu_S(d,k) \geq \alpha_{ks}$ .

The results under the conditions mentioned above can be presented as a hierarchical structure which elements have this type *Domain\_name[pessimistic÷optimistic\_expectation]*.

Because the min-max composition is always zero and to be clear and easy visible the given level of adequacy, the result may be presented this type: *Domain\_name[level\_of\_adequacy÷optimistic\_expectation]* or  $A_i[\alpha_{ks} \div \max \min \mu_{R \circ S}(A_i, ks)]$ .

*For example:*

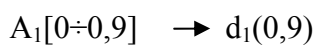
If the searched word is **ks=k<sub>1</sub>** and the level of adequacy is  $\alpha_{ks}=0,7$ , then the result from the example mentioned above will be only the documents **d<sub>1</sub>** and **d<sub>2</sub>** under these conditions and can represented as:



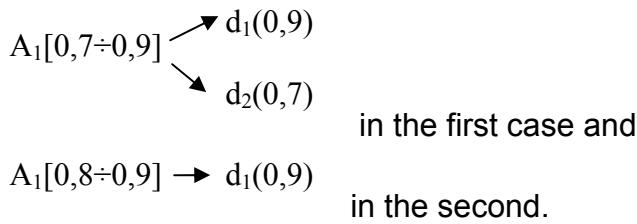
, where:

**A<sub>1</sub>** - this is the domain where the documents belong to and **[0÷0,9]** is the interval of which the results have to be.

If  $\alpha_{ks}=0,8$  then we have in return only one document - **d<sub>1</sub>**. The result is:



If the result is presented as  $A_i[\alpha_{ks} \div \max \min \mu_{R \circ S}(a_i, ks)]$  then we have:



**COMPLEX REQUEST**

A complex request is a request that has words such as “and”, “or”, “not” and are corresponding to fuzzy sets operations – “conjunction”, “disjunction” and “complement”. They can be summarized in:

AND	$\wedge$ (conjunction)	min
OR	$\vee$ (disjunction)	max
NOT	$\bar{k}_i$ (complement)	$1 - k_i$

By representing the searched word in this way we can find the fuzzy set **DS** which contains “all documents with the keyword **ks**” and where the level of adequacy is set to  $\alpha_{ks}$ .

For example:

If the request is  $ks = k_1 \vee k_2 \wedge \bar{k}_3$  then we can symbolize it as  $\min[\max(k_1, k_2), 1 - k_3]$ . Consequently  $\mu_D(ks) = \min[\max(\mu_D(k_1), \mu_D(k_2)), 1 - \mu_D(k_3)]$  and the fuzzy set **DS** is:

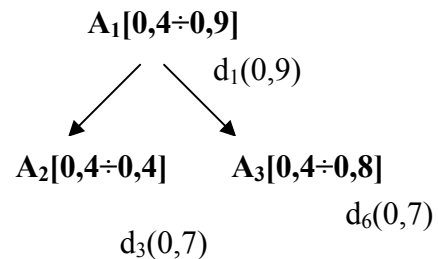
DS=

<b>d<sub>1</sub></b>	<b>d<sub>2</sub></b>	<b>d<sub>3</sub></b>	<b>d<sub>4</sub></b>	<b>d<sub>5</sub></b>	<b>d<sub>6</sub></b>
0,9	0,3	0,4	0	0,3	0,8

If we set a level  $\alpha_{ks} = 0,4$  then we will find out that the documents **d<sub>1</sub>**, **d<sub>3</sub>** и **d<sub>6</sub>** cover this condition. When we find the max-min and min-max composition of **R** и **DS** we have:

And the result will be:

T <sub>1</sub>	<b>ks</b>	T <sub>2</sub>	<b>ks</b>
<b>A<sub>1</sub></b>	0,9	<b>A<sub>1</sub></b>	0
<b>A<sub>2</sub></b>	0,4	<b>A<sub>2</sub></b>	0
<b>A<sub>3</sub></b>	0,8	<b>A<sub>3</sub></b>	0

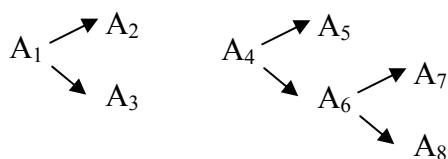


**GENERALE CASE IN REAL SYSTEM**

The general case is possible as one keyword describes documents of more than one domain and their sub domains, and / or as one document has been classified in more than one domain / sub domain. Consequently when is searched for such a keyword, the result will be two or more different hierarchies grouping the result thematically in domains and their sub domains.

For example:

We have the following hierarchies:



The documents where they consist of are described with binary relation  $R=UxD$ , where  $U$  is universe of discourse containing all domains and sub domains we need and  $D$  is a set of all documents classified in some of sets  $A_1, A_2 \dots A_8$  (in this case).

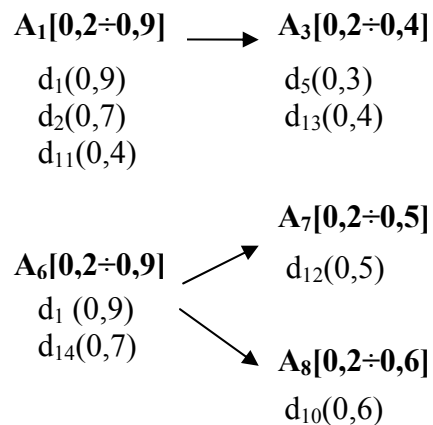
The relation  $R$  is presented as:

R	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d <sub>4</sub>	d <sub>5</sub>	d <sub>6</sub>	d <sub>7</sub>	d <sub>8</sub>	d <sub>9</sub>	d <sub>10</sub>	d <sub>11</sub>	d <sub>12</sub>	d <sub>13</sub>	d <sub>14</sub>
A <sub>1</sub>	1	1	0	0	0	0	0	0	0	0	1	0	0	0
A <sub>2</sub>	0	0	0	1	0	0	0	0	0	0	0	0	0	0
A <sub>3</sub>	0	0	0	0	1	0	0	1	0	0	0	0	1	0
A <sub>4</sub>	0	0	0	0	0	1	1	0	0	0	0	0	0	0
A <sub>5</sub>	0	0	0	1	0	0	0	0	1	0	0	0	0	0
A <sub>6</sub>	1	0	0	0	0	0	0	0	0	0	0	0	0	1
A <sub>7</sub>	0	0	1	0	0	0	0	0	0	0	0	1	0	0
A <sub>8</sub>	0	0	0	0	0	1	0	0	0	1	0	0	0	0

The fuzzy set  $K$  is a set of all keywords describing the documents. To keep it simple, we will work only with one keyword -  $\kappa_1$  and its membership function to  $D$ , which is shown in the binary relation  $S$ . Now we can find the binary relations  $T_1$  and  $T_2$  that are presenting the optimistic and pessimistic expectations of how  $\kappa_1$  is describing the following sets  $A_1, A_2 \dots A_8$ .

S	$\kappa_1$	T <sub>1</sub>	ks	T <sub>2</sub>	ks
d <sub>1</sub>	0,9	A <sub>1</sub>	0,9	A <sub>1</sub>	0
d <sub>2</sub>	0,7	A <sub>2</sub>	0	A <sub>2</sub>	0
d <sub>3</sub>	0	A <sub>3</sub>	0,4	A <sub>3</sub>	0
d <sub>4</sub>	0	A <sub>4</sub>	0	A <sub>4</sub>	0
d <sub>5</sub>	0,3	A <sub>5</sub>	0	A <sub>5</sub>	0
d <sub>6</sub>	0	A <sub>6</sub>	0,9	A <sub>6</sub>	0
d <sub>7</sub>	0	A <sub>7</sub>	0,5	A <sub>7</sub>	0
d <sub>8</sub>	0	A <sub>8</sub>	0,6	A <sub>8</sub>	0
d <sub>9</sub>	0				
d <sub>10</sub>	0,6				
d <sub>11</sub>	0,4				
d <sub>12</sub>	0,5				
d <sub>13</sub>	0,4				
d <sub>14</sub>	0,7				

Let  $ks = \kappa_1$  and  $\alpha_{ks} = 0,2$ . Then the result will be:



This result in a real system would look like a hierarchy of hyperlinks to adjacent list of documents, for example:

- Domain1 [0,2÷0,9]
- Sub-domain3 [0,2÷0,5]
- Domain4
- Sub-domain6 [0,2÷0,9]
- Sub-domain7 [0,2÷0,5]
- Sub-domain8 [0,2÷0,6]

In this case “*Domain4*” is visualized because it is giving the information about subsequent sub domains.

If we follow the link “*Sub-domain3 [0,2÷0,5]*” for example, the result will be:

*Sub-domain3:*

*Document5 (0,3)*

*keyword# (0,3); keyword# (#,#); keyword# (#,#);...*

*Document13 (0,4)*

*keyword# (0,4); keyword# (#,#); keyword# (#,#);...*

Where “*Document13 (0,4)*” is hyperlink pointing the real document.

## **CONCLUSIONS**

Knowledge representation using the fuzzy sets theory gives the opportunity to define in a more adequate and clear way the characteristics of an object and classes where it consists of. The presented model gives the possibility to present the result that is classified and helps the user to find the needed domains where he is interested of, grouping the results and in this case filtering out the one’s that are not in the required domain. It is given the possibility to set the level of adequacy of the given results.

## **REFERENCES**

- [1] Бърнев П., П. Станчев. Размити множества. Народна просвета, София, 1987. 110с.
- [2] Гаврилова Т. А., В. Ф. Хорошевский. Базы знаний интеллектуальных систем. Питер, Санкт-Петербург, 2000. 384с.
- [3] Дюбуа Д., А. Прад. Теория возможностей – приложения к представлению знаний в информатике. Радио и связь, Москва, 1990. 288с.
- [4] Леоненков А.. Нечеткое моделирование с среде METLAB и fuzzyTECH. БХВ-Петербург, Санкт-Петербург, 2003. 736с.

## **ABOUT THE AUTHORS**

Eng. Hristina Moneva, PhD student, Department of Computer Systems and Technology, University of Veliko Turnovo, Phone: + 359 88 9513288, E-mail: [xmoneva@ieee.org](mailto:xmoneva@ieee.org)