

# School Dropout Screening through Artificial Neural Networks based Systems

Margarida Figueiredo, Lúcia Vicente, Henrique Vicente, and José Neves

**Abstract**—School dropout is one of the major concerns of our society. Indeed, it is a complex phenomenon, resulting in economic and social losses, either to the individual, family or the community to which the person belongs. Academic difficulty and failure, poor attendance, retention, disengagement from school together with family and socio-economic reasons can lead to such occurrence. In this work Logic Programming was used for knowledge representation and reasoning, letting the modeling of the universe of discourse in terms of defective data, information and knowledge. Artificial Neural Networks were used in order to evaluate potential situations of school dropout and the degree of confidence that one has on such a happening.

**Keywords**—Artificial Neuronal Networks, Knowledge Representation and Reasoning, Logic Programming, School Dropout.

## I. INTRODUCTION

EDUCATION is a powerful driver of development and one of the soundest instruments for reducing countries poverty. Although there has been great progress in the last decade, a large number of young people that finished their education did it without acquiring basic skills necessary for work and life. This is particularly detrimental when unemployment is high and labour markets are demanding more skills than ever before [1].

In the 21st century, in all developed countries, education of children and young people is recognized as one of the core values. However, school dropout numbers are much higher than would be desirable. According to data from the official statistics institute EU in 2012, Portugal had a dropout rate of 20.8%. Spain, with 24.9%, and Malta, with 22.6%, these were the only European countries which that year showed the highest values. Despite the decline between 2005 and 2011 (from 38.8% to 20.8%), this rate is still far from the national goal of a 10% dropout in secondary schools for 2020 [2].

This work is funded by National Funds through the FCT - Fundação para a Ciência e a Tecnologia (Portuguese Foundation for Science and Technology) within projects PEst-OE/EEI/UI0752/2014 and PEst-OE/QUI/UI0619/2012.

Margarida Figueiredo is with the Department of Chemistry & Évora Chemistry Centre, School of Science and Technology, University of Évora, Évora, Portugal (corresponding author to provide phone: +351-266745315; fax: +351-266745303; e-mail: mtf@uevora.pt).

Lúcia Vicente is with the Agrupamento de Escolas de Reguengos de Monsaraz, Reguengos de Monsaraz, Portugal (e-mail: lmrvicente@gmail.com).

Henrique Vicente is with the Department of Chemistry & Évora Chemistry Centre, School of Science and Technology, University of Évora, Évora, Portugal (e-mail: hvicente@uevora.pt).

José Neves is with the CCTC/Department of Informatics, University of Minho, Braga, Portugal (e-mail: jneves@di.uminho.pt).

The reasons why young people have poor skills when they finish their school careers are many and diverse. Solving this problem requires detailed knowledge of the causes that lead to this situation. School dropout is one of those causes. It is a complex phenomenon, resulting in economic and social losses, either to the individual, family or the community to which the person belongs. If school dropout is large in a country or in a developed region, the consequences will be mainly damaging in terms of economic competitiveness and social environmental degradation [3].

Young people drop out of school for multiple reasons, and such decision rarely is a spur of the moment. Indeed, in large number of cases, young people drop out of school following a long process of disengagement and academic struggle. The factors behind the school dropout phenomenon can be grouped in four main generic groups, namely student related factors, family related factors, school related factors and community related factors [4], as illustrated in Fig. 1.

With regard to student related factors, poor school performance is a strong indicator of dropping out of school. Indeed, low test scores, course failure, and grade retention have been found to be strongly associated with school dropout [4], [5]. The causes for this are multiple and complex. Sometimes students get involved with gangs, drugs, alcohol, get pregnant and commit crimes. Many have a poor school attitude and are frequently bored by school. They are disconnected to their families, school and life. They do not see the reasons why they need to go to school. They are not involved in school activities and have lack of self-esteem [4], [5].

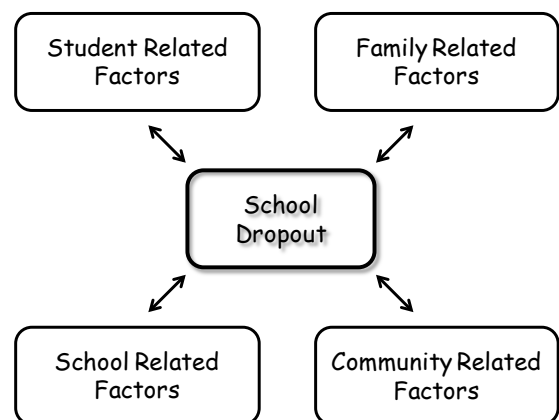


Fig. 1 Causes behind of the school dropout phenomenon

Family related factors affecting heavily the school dropout situations. Students coming from families with low socio-economic backgrounds, where there are many other children, older children often have to go to work or have to stay at home to take care of younger siblings. Frequently, in these families, parents have also a history of school dropout [4], [5].

The policies followed by the schools are important in the involvement of students in the educational project. High number of students per class and overloaded schedules can contribute to disinterest students. Problematic students need adequate guidance counselling and alternative curricula properly tailored. The characteristics of the community where the student lives and where the school is located can also contribute to school dropout. Many students live in places where education is not valued, where there is low demand for qualifications by employers and where drugs, gangs and violence abound [4], [5].

Solving problems related to dropout requires a proactive strategy able to take into account all these factors. Thus, the development of models to evaluate the risk of dropout may be a way to solve or minimize the problem. This work introduces a computational system to evaluate potential situations of school dropout centred on logic programming to knowledge representation and reasoning, complemented with a computational framework based on Artificial Neural Networks.

## II. QUALITY-OF-INFORMATION VERSUS DEGREE OF CONFIDENCE

Due to the growing need to offer user support in decision making processes some studies have been presented [6], [7], related to the qualitative models and qualitative reasoning in Database Theory and in Artificial Intelligence research. With respect to the problem of knowledge representation and reasoning in Logic Programming (LP), a measure of the *Quality-of-Information (QoI)* of such programs has been object of some work with promising results [8], [9]. The *QoI* with respect to the extension of a predicate  $i$  will be given by a truth-value in the interval  $[0,1]$ , i.e., if the information is *known (positive)* or *false (negative)* the *QoI* for the extension of  $predicate_i$  is 1. For situations where the information is unknown, the *QoI* is given by:

$$QoI_i = \lim_{N \rightarrow \infty} \frac{1}{N} = 0 \quad (N \gg 0) \quad (1)$$

where  $N$  denotes the cardinality of the set of terms or clauses of the extension of  $predicate_i$  that stand for the incompleteness under consideration. For situations where the extension of  $predicate_i$  is unknown but can be taken from a set of values, the *QoI* is given by:

$$QoI_i = 1/Card \quad (2)$$

where  $Card$  denotes the cardinality of the *abducibles* set for  $i$ , if the *abducibles* set is disjoint. If the *abducibles* set is not disjoint, the *QoI* is given by:

$$QoI_i = \frac{1}{C_{Card}^{Card} + \dots + C_{Card}^{Card}} \quad (3)$$

where  $C_{Card}^{Card}$  is a card-combination subset, with  $Card$  elements. The next element of the model to be considered is the relative importance that a predicate assigns to each of its attributes under observation, i.e.,  $w_i^k$ , which stands for the relevance of attribute  $k$  in the extension of  $predicate_i$ . It is also assumed that the weights of all the attribute predicates are normalized, i.e.:

$$\sum_{1 \leq k \leq n} w_i^k = 1, \forall_i \quad (4)$$

where  $\forall$  denotes the universal quantifier. It is now possible to define a predicate's scoring function  $V_i(x)$  so that, for a value  $x = (x_1, \dots, x_n)$ , defined in terms of the attributes of  $predicate_i$ , one may have:

$$V_i(x) = \sum_{1 \leq k \leq n} w_i^k \times QoI_i(x)/n \quad (5)$$

It is now possible to engender all the possible scenarios of the universe of discourse, according to the information given in the logic programs that endorse the information depicted in Fig. 3, i.e., in terms of the extensions of the predicates *General Characterization*, *Student Related Factors*, *Family Related Factors*, *School Related Factors* *Community Related Factors* and *School Dropout*.

It is now feasible to rewrite the extensions of the predicates referred to above, in terms of a set of possible scenarios according to productions of the type:

$$predicate_i((x_1, \dots, x_n)) :: QoI \quad (6)$$

and evaluate the *Degree of Confidence (DoC)* given by  $DoC = V_i(x_1, \dots, x_n)/n$ , which denotes one's confidence in a particular term of the extension of  $predicate_i$ . To be more general, let us suppose that the Universe of Discourse is described by the extension of the predicates:

$$a_1(\dots), a_2(\dots), \dots, a_n(\dots) \text{ where } (n \geq 0) \quad (7)$$

Therefore, for a given *scenario*, one may have (where  $\perp$  denotes an argument value of the type unknown; the values of the others arguments stand for themselves):

$$\left\{ \begin{array}{l} \neg a_1(x_1, y_1, z_1) \leftarrow not a_1(x_1, y_1, z_1) \\ a_1([3, 5], \perp, 1.2) :: 0.85 \\ \underline{[0, 15][5, 35][0, 7.2]} \\ \text{attribute's domains for } x_1, y_1, z_1 \\ \neg a_2(x_2, y_2, z_2) \leftarrow not a_2(x_2, y_2, z_2) \\ a_2([13, 22], [8, 10], \perp) :: 0.7 \\ \underline{[10, 40] [6, 14] [2, 5]} \\ \text{attribute's domains for } x_2, y_2, z_2 \\ \vdots \end{array} \right.$$

↓ 1st interaction: transition to continuous intervals

$$\left\{ \begin{array}{l} \neg a_1(x_1, y_1, z_1) \leftarrow \text{not } a_1(x_1, y_1, z_1) \\ a_1([3, 5], [5, 35], [1.2, 1.2]) :: 0.85 \\ \quad \underline{[0, 15] [5, 35] [0, 7.2]} \\ \text{attribute's domains for } x_1, y_1, z_1 \\ \\ \neg a_2(x_2, y_2, z_2) \leftarrow \text{not } a_2(x_2, y_2, z_2) \\ a_2([13, 22], [8, 10], [2, 5]) :: 0.7 \\ \quad \underline{[10, 40] [6, 14] [2, 5]} \\ \text{attribute's domains for } x_2, y_2, z_2 \\ \\ \vdots \end{array} \right.$$

↓ 2nd interaction: normalization  $\frac{Y - Y_{min}}{Y_{max} - Y_{min}}$

$$\left\{ \begin{array}{l} \neg a_1(x_1, y_1, z_1) \leftarrow \text{not } a_1(x_1, y_1, z_1) \\ a_1\left(\left(\frac{3-0}{15-0}, \frac{5-0}{15-0}\right), \left(\frac{5-5}{35-5}, \frac{35-5}{35-5}\right), \left(\frac{1.2-0}{7.2-0}, \frac{1.2-0}{7.2-0}\right)\right) \equiv \\ a_1([0.2, 0.33], [0, 1], [0.17, 0.17]) :: 0.85 \\ \quad \underline{[0, 1] [0, 1] [0, 1]} \\ \text{attribute's domains for } x_1, y_1, z_1 \\ \\ \neg a_2(x_2, y_2, z_2) \leftarrow \text{not } a_2(x_2, y_2, z_2) \\ a_2\left(\left(\frac{13-10}{40-10}, \frac{22-10}{40-10}\right), \left(\frac{8-6}{14-6}, \frac{10-6}{14-6}\right), \left(\frac{2-2}{5-2}, \frac{5-2}{5-2}\right)\right) \equiv \\ a_2([0.1, 0.4], [0.25, 0.5], [0, 1]) :: 0.7 \\ \quad \underline{[0, 1] [0, 1] [0, 1]} \\ \text{attribute's domains for } x_2, y_2, z_2 \\ \\ \vdots \end{array} \right.$$

The *Degree of Confidence (DoC)* is evaluated using the equation  $DoC = \sqrt{1 - \Delta l^2}$ , as it is illustrated in Fig. 2. Here  $\Delta l$  stands for the length of the arguments' intervals, once normalized.

Below, one has the expected representation of the universe of discourse, where all the predicates' arguments are nominal. They speak for one's confidence that the unknown values of the arguments fit into the correspondent intervals referred to above.

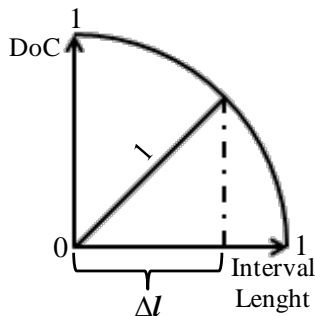


Fig. 2 Degree of Confidence evaluation

$$\left\{ \begin{array}{l} \neg a_1(x_1, y_1, z_1) \leftarrow \text{not } a_1(x_1, y_1, z_1) \\ a_1(0, 0.75, 0.8) :: 0.9 \\ \quad [0, 1] [0, 1] [0, 1] \\ \\ \neg a_2(x_2, y_2, z_2) \leftarrow \text{not } a_2(x_2, y_2, z_2) \\ a_2(0.65, 0.7, 0) :: 0.85 \\ \quad [0, 1] [0, 1] [0, 1] \\ \\ \vdots \end{array} \right.$$

### III. KNOWLEDGE REPRESENTATION AND REASONING

Many approaches for knowledge representation and reasoning have been proposed using the *Logic Programming (LP)* paradigm, namely in the area of Model Theory [10]–[12], and Proof Theory [13], [14]. We follow the proof theoretical approach and an extension to the *LP* language, to knowledge representation and reasoning. An *Extended Logic Program (ELP)* is a finite set of clauses in the form:

$$p \leftarrow p_1, \dots, p_n, \text{not } q_1, \dots, \text{not } q_m \quad (8)$$

$$?(p_1, \dots, p_n, \text{not } q_1, \dots, \text{not } q_m) \quad (n, m \geq 0) \quad (9)$$

where  $?$  is a domain atom denoting falsity, the  $p_i$ ,  $q_j$ , and  $p$  are classical ground literals, i.e., either positive atoms or atoms preceded by the classical negation sign  $\neg$  [14]. Under this representation formalism, every program is associated with a set of *abducibles* [10] [12], given here in the form of exceptions to the extensions of the predicates that make the program. Once again, LP emerged as an attractive formalism for knowledge representation and reasoning tasks, introducing an efficient search mechanism for problem solving. Therefore, and in order to exemplify the applicability of our model, we will look at the relational database model, since it provides a basic framework that fits into our expectations [15], and is understood as the genesis of the LP approach to knowledge representation and reasoning.

Consider, for instance, the scenario where a relational database is given in terms of the extensions of the relations or predicates depicted in Fig. 3, which stands for a situation where one has to manage information about school dropout. Under this scenario some incomplete data is also available. For instance, in *General Characterization* database the *Household incomes* for case 3 is unknown, while for example 1 ranges in the interval [2800,3200].

The values of the *Student*, *Family*, *School* and *Community Issues* in *School Dropout* database are the sum of the respective values, ranging between [0,7], [0,11], [0,6] and [0,5] respectively.

Now, we may consider the extensions of the relations given in Fig. 3 to populate the extension of the *dropout* predicate, given in the form:

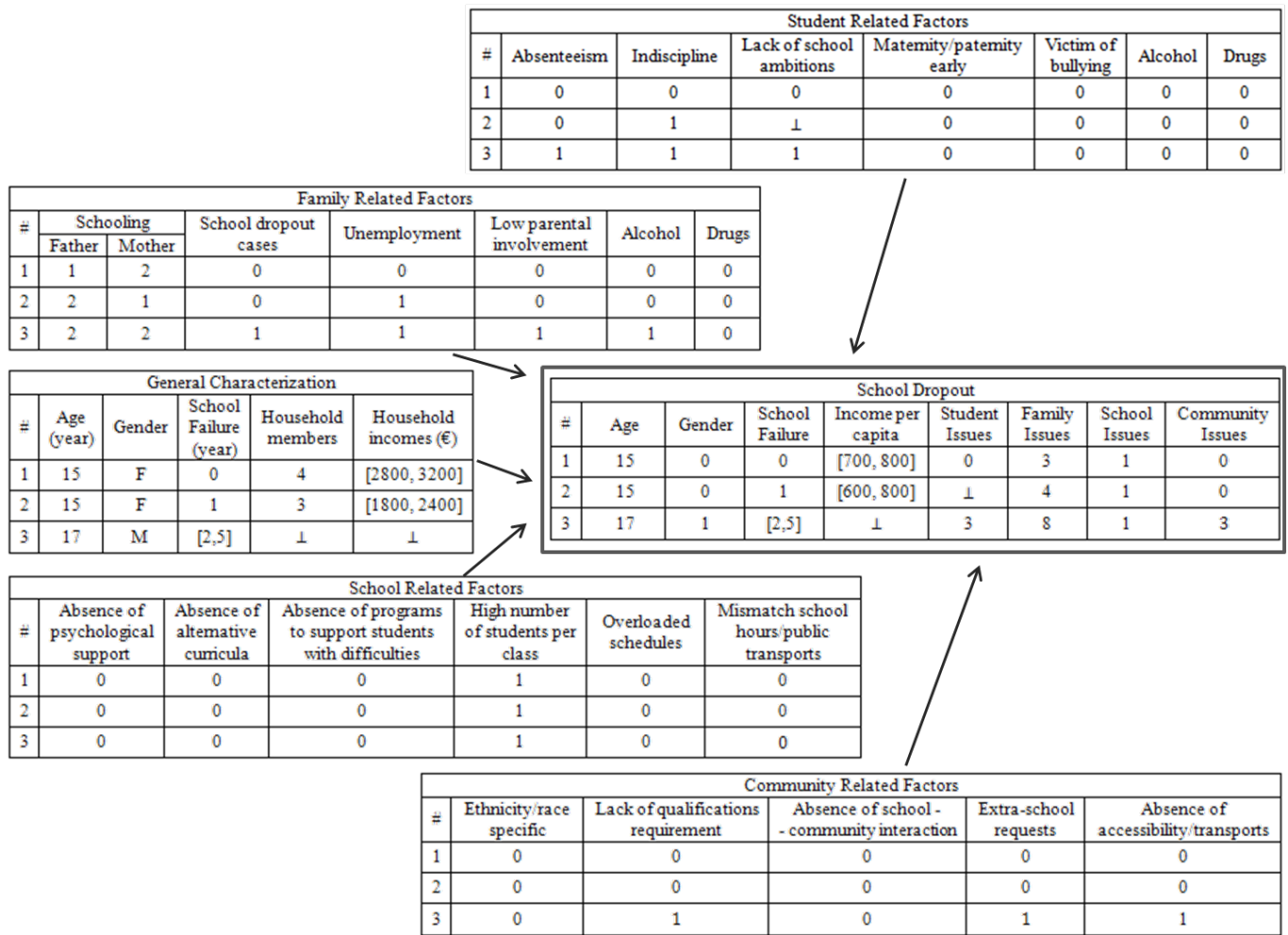


Fig. 3 An extension of the relational database model. In column *Gender* of *School Dropout* database, 0 (zero) and 1 (one) stand, respectively, for *female* and *male*. In column *Schooling* of *Family Related Factors* database, 0 (zero), 1 (one), 2 (two) and 3 (three) denote, respectively, *university course*, *secondary education*, *primary education* and *illiterate*. In the remaining columns of *Student*, *Family*, *School* and *Community Related Factors* 0 (zero) denotes *absence/no* and 1 (one) denotes *presence/yes*

$dropout : Age, Gender, SchoolFailure, Income\ per\ capita,$   
 $Student\ Issues, Family\ Issues, School\ Issues,$   
 $Community\ Issues \rightarrow \{0,1\}$

where 0 (zero) and 1 (one) denote, respectively, the truth-values *false* and *true*. It is now possible to give the extension of the predicate *dropout*, in the form:

```

{
  ¬dropout (Age, G, SF, Inc, Stud, Fam, Sch, Com)
  ← not dropout (Age, G, SF, Inc, Stud, Fam, Sch, Com)
  dropout (15, 0, 0, [700,800], 0, 3, 1, 0) :: 1
    [6,22][0,1][0,8][200,3000][0,7][0,11][0,6][0,5]
  dropout (15, 0, 1, [600,800], 1, 4, 1, 0) :: 1
    [6,22][0,1][0,8][200,3000][0,7][0,11][0,6][0,5]
}
    
```

In this program, the first clause denotes the closure of predicate *dropout*. The next clauses correspond to two terms taken from the extension of the *dropout* relation. It is now possible to have the arguments of the predicates extensions normalized to the interval [0, 1], in order to compute one's confidence that the nominal values of the arguments under considerations fit into the intervals depicted previously. One may have:

```

{
  ¬dropout (Age, G, SF, Inc, Stud, Fam, Sch, Com)
  ← not dropout (Age, G, SF, Inc, Stud, Fam, Sch, Com)
  dropout ([0.56,0.56], [0,0], [0,0], [0.18,0.21],
    [0,0], [0.27,0.27], [0.17,0.17], [0,0]) :: 1
    [0,1] [0,1] [0,1] [0,1]
    [0,1] [0,1] [0,1] [0,1]
}
    
```

```

dropout([0.56,0.56],[0,0],[0.12,0.12],[0.14,0.21],
        [0,1],[0.36,0.36],[0.17,0.17],[0,0]) :: 1
        [0,1] [0,1] [0,1] [0,1]
        [0,1] [0,1] [0,1] [0,1]
}
    
```

The logic program referred to above, is now presented in the form:

```

{
  ¬dropoutDoC(Age, G, SF, Inc, Stud, Fam, Sch, Com)
← not dropoutDoC(Age, G, SF, Inc, Stud, Fam, Sch, Com)

  dropoutDoC(1,1,1,0.9995,1,1,1,1) :: 1
  dropoutDoC(1,1,1,0.9975,0,1,1,1) :: 1
}
    
```

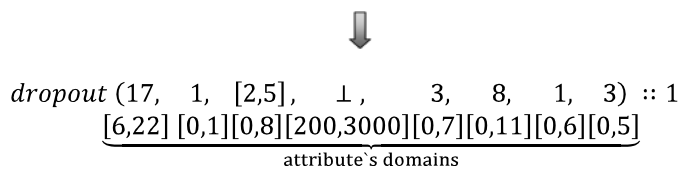
where its terms make the training and test sets of the Artificial Neural Network given below (Fig. 4).

#### IV. ARTIFICIAL NEURAL NETWORKS

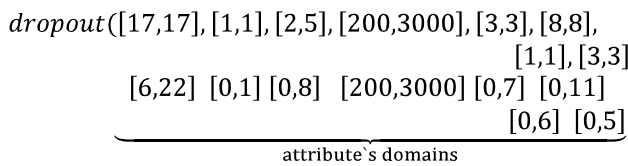
In [16]–[18] it is shown how Artificial Neural Networks (ANNs) could be successfully used to model data and capture complex relationships between inputs and outputs. ANNs simulate the structure of the human brain being populated by multiple layers of neurons. As an example, let us consider the third case presented in Fig. 3, where one may have a situation that may lead to school dropout, which is given in the form:

```

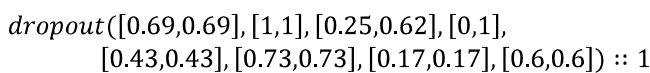
{
  dropout attributes (Age, G, SF, Inc, Stud, Fam, Sch, Com)
}
    
```



↓ 1st interaction: transition to continuous intervals



↓ 2nd interaction: normalization  $\frac{Y - Y_{min}}{Y_{max} - Y_{min}}$



↓

```

dropoutDoC(1,1,0.93,0,1,1,1,1) :: 1
}
    
```

In Fig. 4 it is shown how the normalized values of the interval boundaries and their *DoC* values work as inputs to the ANN. The output translates the chance of a situation of school dropout to occur, and *DoC* the confidence that one has on such a happening. In addition, it also contributes to build a database of study cases that may be used to train and test the ANNs.

#### V. CONCLUSIONS AND FUTURE WORK

The anticipation of situations that lead to school dropout is a hard and complex task, which needs to consider many different conditions with intricate relations among them. These characteristics put this problem into the area of problems that may be tackled by AI based methodologies and techniques to problem solving. Despite that, little to no work has been done in that direction. This work presents the founding of a computational framework that uses powerful knowledge representation and reasoning techniques to set the structure of the information and the associate inference mechanisms. This representation is above everything else, very versatile and capable of covering every possible instance by considering incomplete, contradictory, and even unknown data. The main contribution of this work is to be understood in terms of the evaluation of the *DoC*, and the possibility to address the issue of incomplete information. Indeed, the new paradigm of knowledge representation and reasoning enables the use of the normalized values of the interval boundaries and their *DoC* values, as inputs to the ANN. The output translates the chance of school dropout and the degree of confidence that one has on such a happening. Future work may recommend that the same problem must be approached using other computational frameworks like Case Based Reasoning or Particle Swarm, just to name a few.

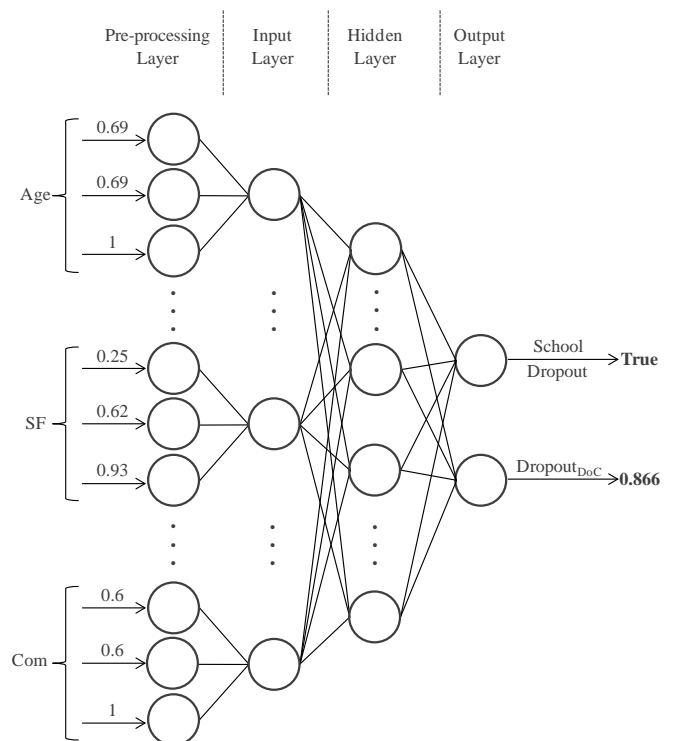


Fig. 4 A possible Artificial Neural Network topology

## REFERENCES

- [1] World Bank Group Education Strategy 2020. (2014, May 31). Learning for All. Investing in People's Knowledge and Skills to Promote Development [Online]. Available: [http://siteresources.worldbank.org/EDUCATION/Resources/ESSU/Education\\_Strategy\\_4\\_12\\_2011.pdf](http://siteresources.worldbank.org/EDUCATION/Resources/ESSU/Education_Strategy_4_12_2011.pdf)
- [2] European Commission. (2014, May 28). Progress in tackling early school leaving and raising higher education attainment - but males are increasingly left behind [Online]. Available: [http://europa.eu/rapid/press-release\\_IP-13-324\\_en.htm](http://europa.eu/rapid/press-release_IP-13-324_en.htm)
- [3] T. Andrei, D. Teodorescu and B. Oancea, "Characteristics and causes of school dropout in the countries of the European Union", *Procedia - Social and Behavioral Sciences*, vol. 28, pp. 328-332, 2011.
- [4] C. F. V. Castro "School Dropout - Factors and Strategies and Combat" Master dissertation, Department of Education and Psychology of University of Trás-os-Montes e Alto Douro, Vila Real, 2010.
- [5] J. H. Tyler and M. Lofstrom, "Finishing High School: Alternative Pathways and Dropout Recovery", *The Future of Children*, vol. 19, pp. 77-103, Spring 2009.
- [6] J. Halpern, *Reasoning about uncertainty*. Massachusetts: MIT Press, 2005.
- [7] B. Kovalerchuck and G. Resconi, "Agent-based uncertainty logic network", in *Proceedings of the IEEE International Conference on Fuzzy Systems – FUZZ-IEEE 2010*, Barcelona, Spain, 2010, pp. 596-603.
- [8] P. Lucas, "Quality checking of medical guidelines through logical abduction", in *Proceedings of AI-2003 (Research and Developments in Intelligent Systems XX)*, F. Coenen, A. Preece and A. Mackintosh, Eds. London: Springer, 2003, pp. 309-321.
- [9] J. Machado, A. Abelha, P. Novais, J. Neves and J. Neves, "Quality of service in healthcare units", *International Journal of Computer Aided Engineering and Technology*, vol. 2, pp. 436-449, 2010.
- [10] A. Kakas, R. Kowalski and F. Toni) "The role of abduction in logic programming", in *Handbook of Logic in Artificial Intelligence and Logic Programming*, vol. 5, D. Gabbay, C. Hogger and I. Robinson, Eds., Oxford: Oxford University Press, 1998, pp. 235-324.
- [11] M. Gelfond and V. Lifschitz, "The stable model semantics for logic programming", in *Logic Programming – Proceedings of the Fifth International Conference and Symposium*, R. Kowalski and K. Bowen, Eds. Cambridge: MIT Press, 1988, pp. 1070-1080.
- [12] L. Pereira and H. Anh, "Evolution prospection", in *New Advances in Intelligent Decision Technologies – Results of the First KES International Symposium IDT 2009*, K. Nakamatsu, G. Phillips-Wren, L. Jain and R. Howlett Eds. Studies in Computational Intelligence, vol. 199, Berlin: Springer, 2009, pp. 51-64.
- [13] J. Neves, J. Machado, C. Analide, A. Abelha and L. Brito, "The halt condition in genetic programming", in *Progress in Artificial Intelligence - Lecture Notes in Computer Science*, vol 4874, J. Neves, M. F. Santos and J. Machado Eds. Heidelberg: Springer, 2007, pp. 160-169.
- [14] J. Neves, "A logic interpreter to handle time and negation in logic data bases", in *Proceedings of the 1984 annual conference of the ACM on the fifth generation challenge*, R. L. Muller and J. J. Pottmyer Eds. New York: Association for Computing Machinery, 1984, pp. 50-54.
- [15] Y. Liu and M. Sun, "Fuzzy optimization BP neural network model for pavement performance assessment", in *Proceedings of the 2007 IEEE International Conference on Grey Systems and Intelligent Services*, Nanjing, China, 2007, pp. 18-20.
- [16] H. Vicente, S. Dias, A. Fernandes, A. Abelha, J. Machado, and J. Neves, "Prediction of the Quality of Public Water Supply using Artificial Neural Networks", *Journal of Water Supply: Research and Technology – AQUA*, vol. 61, pp. 446-459, 2012.
- [17] C. Salvador, M. R. Martins, H. Vicente, J. Neves, J. M. Arteiro and A. T. Caldeira, "Modelling Molecular and Inorganic Data of Amanita ponderosa Mushrooms using Artificial Neural Networks". *Agroforestry Systems*, vol. 87, pp. 295-302, 2013.
- [18] A. T. Caldeira, J. M. Arteiro, J. C. Roseiro, J. Neves, and H. Vicente, "An Artificial Intelligence Approach to Bacillus amyloliquefaciens CCM1 1051 Cultures: Application to the Production of Antifungal Compounds", *Bioresource Technology*, vol. 102, pp. 1496-1502, 2011.