

# Imputing Missing Values in Nuclear Safeguards Evaluation by a 2-Tuple Computational Model

Rosa M. Rodríguez<sup>1</sup>, Da Ruan<sup>2</sup>, Jun Liu<sup>3</sup>, Alberto Calzada<sup>1</sup>, and Luis Martínez<sup>1,\*</sup>

<sup>1</sup> University of Jaén, (Jaén-Spain)

{rmrodrig,martin,acalzada}@ujaen.es

<sup>2</sup> Belgium Nuclear Research Centre (SCK • CEN), (Mol-Belgium)

druan@esckcen.be

<sup>3</sup> University of Ulster, (Northern Ireland-UK)

j.liu@ulster.ac.uk

**Abstract.** Nuclear safeguards evaluation aims to verify that countries are not misusing nuclear programs for nuclear weapons purposes. Experts of the International Atomic Energy Agency (IAEA) evaluate many indicators by using diverse sources, which are vague and imprecise. The use of linguistic information has provided a better way to manage such uncertainties. However, missing values in the evaluation are often happened because there exist many indicators and the experts have not sufficient knowledge or expertise about them. Those missing values might bias the evaluation result. In this contribution, we provide an imputation process based on collaborative filtering dealing with the linguistic 2-tuple computation model and a trust measure to cope with such problems.

**Keywords:** missing values, nuclear safeguards, fuzzy sets, imputation, trust worthy.

## 1 Introduction

Nuclear safeguards are a set of activities accomplished by the International Atomic Energy Agency (IAEA) in order to verify that a State is living up to its international undertakings not to use nuclear programs for nuclear weapons purposes. The safeguards system is based on assessments of the correctness and completeness of the State's declarations to the IAEA concerning nuclear material and related nuclear activities [7]. As a part of the efforts to strengthen international safeguards, including its ability to provide credible assurance of the absence of undeclared nuclear material and activities, IAEA uses large amounts and different types of information about States' nuclear and related nuclear activities.

IAEA evaluates nuclear safeguards by using a hierarchical assessment system that assesses indicators about critical activities [6] in the nuclear fuel cycle and processes required for their activities. According to the existence of the processes is inferred a decision about the development of nuclear proposes for weapon purposes.

---

\* Corresponding author.

Experts and IAEA evaluate indicators on the basis of their analysis of the available information sources such as declarations of States, on-site inspections, IAEA non-safeguards databases [13,14]. This information is often uncertain and hard to manage by experts and it might happen that experts cannot provide either evaluations about some indicators or do so in an accurate way. In the latter case, the fuzzy linguistic approach [21] to deal with uncertain information provided good results [13]. However, when experts cannot provide all assessments for indicators, it is necessary to manage these missing values. To do this, there exist different ways in the literature, such as deletion, imputation or using as it is [8,17,18,19].

This contribution aims to present an imputation model based on collaborative filtering and the linguistic 2-tuple computational model to deal with missing values in nuclear safeguards. The imputed values are not real ones, therefore, we also introduce a trust measure to clarify the trustworthiness of the imputed values and of the final result.

This paper is structured as follows: In Section 2, we review some related works in nuclear safeguards problems. In Section 3, we briefly outline a linguistic background that will be used in our model. In Section 4, we propose both an imputation model and a trust measure. In Section 5, we show a numerical example of the proposed approach, and finally, we conclude the research in Section 6.

## 2 Related Works

Different approaches for the evaluation and synthesis of nuclear safeguards have been proposed. We have focused on those that deal with nuclear safeguards problems and have briefly reviewed some of them. In [6] the IAEA Physical Model provides a structure for organizing the safeguards relevant information, which is used by IAEA experts to evaluate in a better way the safeguards significance of information on some State's activities. In [13] an evaluation model for the treatment of nuclear safeguards used linguistic information based on a hierarchical analysis of States under activities in a multi-layer structure. This proposal is the basis of our model in this contribution. The hierarchical model based on the IAEA Physical Model is divided into several levels with lower complexity, from which a global assessment by using a multi-step linguistic aggregation process is obtained. A latest proposal to manage the nuclear safeguards problem was proposed in [14], where a framework for modeling, analysing and synthesising information on nuclear safeguards under different types of uncertainty was presented. Such a framework makes use of the multi-layer evaluation model presented in [13] and a new inference model based on a belief inference methodology (RIMER) to handle hybrid uncertain information in nuclear safeguards evolution process. Recently, the focus on nuclear safeguards has moved to the management of missing values [9]. Different proposals about dealing with missing values for various purposes have been published in the literature [3,16,18,20]. We have focused our proposal on this problem by using a collaborative filtering view.

## 3 Linguistic Approach Background

Nuclear safeguards deals with a huge amount of information that usually involves uncertainties, which are mainly related to human cognitive processes. The use of linguistic

information has provided good results for managing such types of information in nuclear processes [13,15]. We shall use the linguistic modeling for nuclear safeguards evaluation by extending the work in [13,14]. The Fuzzy Linguistic Approach [21] represents the information with linguistic descriptors and their semantics. The former can be chosen by supplying a term set distributed on a scale with an order [2]. For example, a set of seven terms,  $S$ , could be:

$$S = \{s_0 : \text{nothing}(n), s_1 : \text{very low}(vl), s_2 : \text{low}(l), s_3 : \text{medium}(m), s_4 : \text{high}(h), s_5 : \text{very high}(vh), s_6 : \text{perfect}(p)\}$$

Usually, in these cases, it is required that in the linguistic term set there exist the operators: (i) Negation:  $\text{Neg}(s_i) = s_j$  such that  $j = g - i$  ( $g + 1$  is the cardinality), (ii) Maximization:  $\max(s_i, s_j) = s_i$  if  $s_i \geq s_j$ , (iii) Minimization:  $\min(s_i, s_j) = s_i$  if  $s_i \leq s_j$ .

The semantics of the terms is represented by fuzzy numbers defined in the interval  $[0, 1]$ , described by membership functions. A way to characterize a fuzzy number is to use a representation based on parameters of its membership function [2].

The use of linguistic information implies processes of computing with words (CW). There exist different computational models to accomplish them. We use the 2-tuple computational model presented in [5] to improve the precision of such processes of CW.

The 2-tuple model represents the linguistic information by means of a pair of values, called 2-tuple,  $(s_i, \alpha)$ , where  $s$  is a linguistic term and  $\alpha$  is a numerical value representing the symbolic translation.

**Definition 1.** *The symbolic translation is a numerical value assessed in  $[-0.5, 0.5]$  that supports the “difference of information” between a counting of information  $\beta$  assessed in the interval of granularity  $[0, g]$  of the term set  $S$  and the closest value in  $\{0, \dots, g\}$  which indicates the index of the closest linguistic term in  $S$ .*

This linguistic model defines a set of functions to carry out transformations between numerical values and 2-tuple [5].

$$\begin{aligned} \Delta : [0, g] &\longrightarrow S \times [-0.5, 0.5] \\ \Delta(\beta) &= (s_i, \alpha), \text{ with } \begin{cases} i = \text{round}(\beta), \\ \alpha = \beta - i, \end{cases} \end{aligned} \quad (1)$$

where *round* is the usual round operation,  $s_i$  has the closest index label to  $\beta$ , and  $\alpha$  is the value of the symbolic translation.

We note that  $\Delta$  is bijective [5] and  $\Delta^{-1} : S \times [-0.5, 0.5] \longrightarrow [0, g]$  is defined by  $\Delta^{-1}(s_i, \alpha) = i + \alpha$ . In this way, the 2-tuple of  $S$  is identified with the numerical values in the interval  $[0, g]$ .

Besides, together this representation model, a computational model based on the functions aforementioned was also introduced in [5].

## 4 A Model to Impute Missing Values in Nuclear Safeguards

So far, the main interest in nuclear safeguards has been focused on the development of evaluation processes whose general structure is shown in Fig. 1. The sub-factors are

aggregated for levels and then, these results are aggregated again to obtain a global assessment. However, recently the focus on nuclear safeguards evaluation has moved to the treatment of missing values [9], because it has been noted that such a treatment is the key for obtaining reliable results. The treatment of missing values can be considered in different ways: deletion, imputation and using as it is [8,17,19]. We are interested in an imputation process, in the literature can be found different imputation processes for different proposes [3,16,18,20]. Here, we propose an imputation process for imputing missing values in nuclear safeguards based on collaborative filtering and a trust measure to compute the reliability of the imputed values. Hence, the general structure of the nuclear safeguards evaluation process is extended according to the Fig. 2.

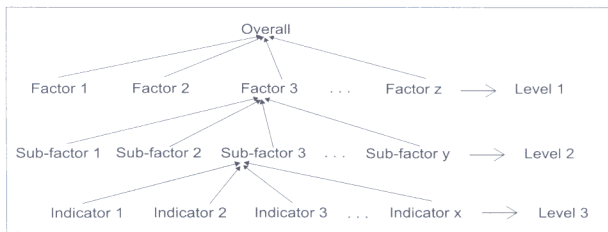


Fig. 1. Structure of the overall evaluation

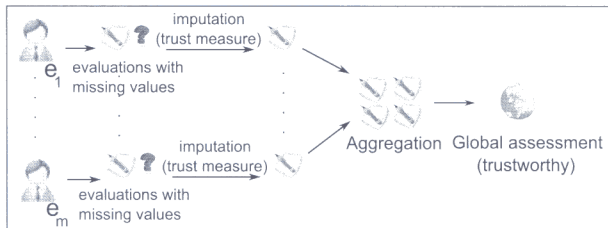


Fig. 2. General steps for nuclear safeguards evaluation with missing values

In the following section we present the CF process to estimate the imputed values and define a trust measure to compute the trustworthiness of such imputed values in order to know how reliable will be the final result.

### 4.1 Imputation Process

The imputation process is based on a k-NN scheme and an estimation similar to the process used by a collaborative recommender system [1,4] (see Fig. 3).

The main idea is to group the indicators according to their similarities by using a similitude measure on expert assessments. To do this, a new proposal for nuclear safeguards which utilises a collaborative filtering technique based on the k-NN algorithm (K nearest neighbours) is applied. To obtain those similarities is used the cosine distance [11] (see Eq. 1). In order to impute a plausible value for the missing one of an

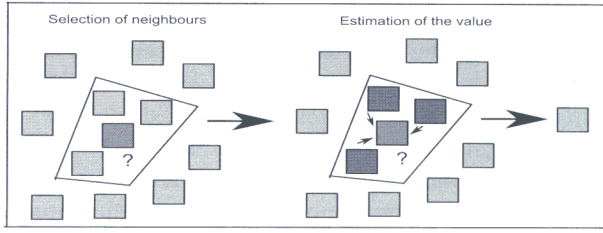


Fig. 3. Scheme of a Collaborative Recommender System

indicator is used the weighted mean (see Eq. 2). This imputation is expressed by means of a linguistic 2-tuple [5] to improve the precision of the linguistic imputed value.

$$w(i, j) = \cos(\bar{v}_i \bar{v}_j) = \frac{\bar{v}_i \bar{v}_j}{\|\bar{v}_i\|^2 \|\bar{v}_j\|^2} \quad (1) \quad v(e, i) = \frac{\sum_{j=1}^{j=k} w_{i,j} v_{e,j}}{\sum_{j=1}^{j=k} w_{i,j}} \quad (2)$$

In Eq. 1 and Eq. 2  $i$  and  $j$  are indicators and  $e$  the expert who has not provided his/her assessment.

The imputed values are not real expert's values, but rather an approximation to them. So, there exist some imprecision. Different precision metrics for collaborative filtering can be found in the literature such as MAE (Mean Absolute Error) [12], ROC [10] and so on. These metrics compute an error between the imputed values and real values. But that is not enough in our problem, because we need to know how trustworthy are those imputed values rather than its possible error.

## 4.2 Trustworthy of the Imputed Values

Therefore, in order to know the trustworthy of the imputed values computed by the previous imputation process, we define the following trust measure:

**Definition 2.** Assume that an expert provides his/her vector of linguistic assessments for the  $m$  indicators in nuclear safeguards evaluation,  $X = \{x_1, \dots, x_m\}$ ,  $x_i \in S = \{s_0, \dots, s_g\}$  and, there exists a set of missing values  $\bar{X} = \{\bar{x}_1, \dots, \bar{x}_n\} \subseteq X$ , that has been imputed in the imputation process. The trustworthy of an imputed value,  $T(\bar{x}_j)$ , is defined as:

$$T(\bar{x}_j) = (1 - \bar{S}_j)V + \bar{S}_j \frac{k}{K}, \quad T(\bar{x}_j) \in [0, 1] \quad (2)$$

where  $\bar{S}_j = \frac{\sum_{l=1}^{l=k} w(j,l)}{k}$  is the arithmetic mean of the similarities among the indicator,  $j$ , and the  $k$  nearest indicators. And  $V = \frac{g - sd(x_j)}{g}$ , being  $sd$  the standard deviation of the assessments used to compute the imputed value and  $g + 1$  the granularity of  $S$ . Eventually  $k$  indicates the real number of neighbours involved in the computation of  $\bar{x}_j$  from the initial  $K$  computed by the  $k$ -NN algorithm.

The definition of  $T(\bar{x}_j)$  is based on the different cases of study whose results show that the more assessments are used to compute the imputed value the more trustworthy is.



Similarly, the more homogeneous are the assessments the more reliable is the imputed value. Thus, the more  $T(\overline{x_l})$  the more reliable is the imputed value. This measure will be used in the process presented in Fig. 3 in order to obtain the trustworthiness of the global assessment.

## 5 A Case Study

Here we present a case study that shows the results obtained by the proposed imputation model in a reduced dataset of four experts and 22 indicators in a nuclear safeguards problem (see Table 1). For the imputation model we have fixed the following parameters:  $K=15$ , the similarity measure is calculated by utilizing the cosine distance (Eq. 1) and for the imputation algorithm the weighted mean (Eq. 2). Therefore, the model imputes values and the results obtained are shown in Table 2. Additionally, in Table 2 is shown the trustworthiness of each imputed value as well.

**Table 1.** Experts evaluations

ind.	1	2	3	4	5	6	7	8	9	10	11
e1	h	p	vh	m	?	h	m	p	m	l	?
e2	l	vh	?	m	l	m	l	vl	l	l	m
e3	h	h	p	vh	m	vh	vh	h	?	m	vh
e4	p	p	p	?	p	m	vh	vh	h	h	m
ind.	12	13	14	15	16	17	18	19	20	21	22
e1	vl	p	vh	?	l	h	m	m	p	h	m
e2	m	p	m	l	l	m	m	?	m	m	vl
e3	l	p	p	vl	l	p	m	vh	p	vh	?
e4	?	vh	h	l	m	?	l	vh	h	vh	l

Once we have obtained the imputed values and their trustworthiness, we keep applying the safeguards process shown in Fig. 2 based on [13], in which the experts assessments are synthesized in a multi-step aggregation process to obtain a global value and its trustworthiness is also obtained by aggregating the trustworthiness of each assessment. To obtain the global assessment, first it is computed a collective assessment for these indicators and then, a global assessment is obtained by aggregating such collective assessments. In this case, we have obtained a linguistic 2-tuple for the nuclear safeguards result (**high**, **0.23**) with a trustworthiness **0.86**.

If we use the safeguards model presented in Fig. 1 without the imputation of missing values the result obtained is (**medium**, **0.33**). We can then observe the relevance of the missing values. Therefore, the treatment of such values must be a cornerstone in nuclear safeguards evaluation.

Table 2. Predictions and trust measures

ind.	e1	e2	e3	e4	
	pred.	trust pred.	trust pred.	trust pred.	trust
3		(m,-.154) .866			
4				(h,.313) .866	
5	(h,-.46) .863				
9			(vh,-.361) .931		
11	(h,-.132) .93				
12				(h,.242) .796	
15	(m,.493) .798				
17				(vh,-.459) .866	
19		(m,-.16) .866			
22			(vh,-.4) .858		

6 Conclusions

Nuclear safeguards evaluation is a complex problem where the experts use different sources of information to evaluate related indicators. This evaluation is usually inaccurate due to the uncertainty of the sources of information and the huge amount of information to manage. Such uncertainties and inaccuracies make that experts sometimes cannot provide assessments for all indicators appearing missing values. In this contribution, we have presented a linguistic nuclear safeguards evaluation model that manages these missing values by means of an imputation process based on a collaborative filtering algorithm. Additionally, we have provided a trust measurement to measure the goodness of the imputed values.

Acknowledgements

This work is partially supported by the Research Project TIN-2009-08286, P08-TIC-3548 and FEDER funds.

References

1. Adomavicius, G., Tuzhilin, A.: Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Transactions on Knowledge and Data Engineering* 17(6), 734–749 (2005)
2. Bonissone, P.P., Decker, K.S.: Selecting Uncertainty Calculi and Granularity: An Experiment in Trading-Off Precision and Complexity. In: Kanal, L.H., Lemmer, J.F. (eds.) *Uncertainty in Artificial Intelligence*. North-Holland, Amsterdam (1986)
3. Dubois, D., Prade, H.: Incomplete conjunctive information. *Computers and Mathematics with Applications* 15(10), 797–810 (1988)
4. Herlocker, J.L., Konstan, J.A., Riedl, J.: An Empirical Analysis of Design Choices in Neighborhood-based Collaborative Filtering Algorithms. In: *Information Retrieval*, pp. 287–310. Kluwer Academic, Dordrecht (2002)

5. Herrera, F., Martínez, L.: A 2-tuple fuzzy linguistic representation model for computing with words. *IEEE Transactions on Fuzzy Systems* 8(6), 746–752 (2000)
6. Physical model. Int. Atomic Energy Agency, IAEA, Vienna, Rep. STR-314 (1999)
7. Nuclear Security and Safeguards. In: IAEA Bulletin, Annual Report. IAEA, vol. 43 (2001)
8. Jiang, N., Gruenwald, L.: Estimating Missing Data in Data Streams. In: Kotagiri, R., Radha Krishna, P., Mohania, M., Nantajeewarawat, E. (eds.) *DASFAA 2007*. LNCS, vol. 4443, pp. 981–987. Springer, Heidelberg (2007)
9. Kabak, Ö., Ruan, D.: A cumulative belief-degree approach for nuclear safeguards evaluation. In: *Proc. of IEEE Conference on Systems, Man and Cybernetics*, San Antonio, TX-USA (2009)
10. Landgrebe, T.C.W., Duin, R.P.W.: Efficient multiclass ROC approximation by decomposition via confusion matrix perturbation analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 30(5), 810–822 (2008)
11. Lee, T.Q., Park, Y., Pack, Y.T.: A Similarity Measure for Collaborative Filtering with Implicit Feedback. In: Huang, D.-S., Heutte, L., Loog, M. (eds.) *ICIC 2007*. LNCS (LNAI), vol. 4682, pp. 385–397. Springer, Heidelberg (2007)
12. Lin, J.H., Sellke, T.M., Coyle, E.J.: Adaptive stack filtering under the mean absolute error criterion. In: *Advances in Communications and Signal Processing*, pp. 263–276 (1989)
13. Liu, J., Ruan, D., Carchon, R.: Synthesis and evaluation analysis of the indicator information in nuclear safeguards applications by computing with words. *International Journal of Applied Mathematics and Computer Science* 12(3), 449–462 (2002)
14. Liu, J., Ruan, D., Wang, H., Martínez, L.: Improving nuclear safeguards evaluation through enhanced belief rule-based inference methodology. *Int. J. Nuclear Knowledge Management* 3(3), 312–339 (2009)
15. Martínez, L.: Sensory evaluation based on linguistic decision analysis. *International Journal of Approximated Reasoning* 44(2), 148–164 (2007)
16. Nowicki, R.: On combining neuro-fuzzy architectures with the rough set theory to solve classification problems with incomplete data. *IEEE Transactions on Knowledge and Data Engineering* 20(9), 1239–1253 (1989)
17. Olman, L.B., Yahia, S.B.: Yet another approach for completing missing values. In: Yahia, S.B., Nguifo, E.M., Belohlavek, R. (eds.) *CLA 2006*. LNCS (LNAI), vol. 4923, pp. 155–169. Springer, Heidelberg (2008)
18. Pawlak, M.: Kernel classification rules from missing data. *IEEE Transactions on Information Theory* 39(3), 979–988 (1993)
19. Siddique, J., Belin, R.: Using and approximate bayesian bootstrap to multiply impute nonignorable. *Computational Statistics and Data Analysis* 53(2), 405–415 (2008)
20. Slowinski, R., Stefanowski, J.: Rough classification in incomplete information-systems. *Mathematical and Computer Modelling* 12(10–11), 1347–1357 (1989)
21. Zadeh, L.A.: The concept of a linguistic variable and its applications to approximate reasoning. *Information Sciences, Part I, II, III* 8, 8, 9, 199–249, 301–357, 43–80 (1975)