



# Interpretation of Best Medical Coding Practices by Case-Based Reasoning—A User Assistance Prototype for Data Collection for Cancer Registries

Michael Schnell<sup>1,2</sup>(✉), Sophie Couffignal<sup>1</sup>, Jean Lieber<sup>2</sup>, Stéphanie Saleh<sup>1</sup>,  
and Nicolas Jay<sup>2,3</sup>

<sup>1</sup> Department of Population Health, Luxembourg Institute of Health,  
1A-B, rue Thomas Edison, 1445 Strassen, Luxembourg

{michael.schnell,sophie.couffignal,stephanie.saleh}@lih.lu

<sup>2</sup> UL, CNRS, Inria, Loria, 54000 Nancy, France

jean.lieber@loria.fr

<sup>3</sup> Service d'évaluation et d'information médicales,

Centre Hospitalier Régional Universitaire de Nancy, Nancy, France

n.jay@chru-nancy.fr

**Abstract.** In the fight against cancer, cancer registries are an important tool. At the heart of these registries is the data collection and coding process. This process is ruled by complex international standards and numerous best practices, which can easily overwhelm (coding) operators. In this paper, a system assisting operators in the interpretation of best medical coding practices and a short evaluation are presented. By leveraging the arguments used by the coding experts to determine the best coding option, the proposed system answers coding questions from operators and provides a partial explanation for the proposed solution.

**Keywords:** Interpretation of best practices · Interpretive case-based reasoning · Coding standards · Cancer registries · User assistance · Decision support

## 1 Introduction

There are numerous cancer registries around the world collecting data about cancers diagnosed and/or treated in a given area. This data is used to monitor cancer (incidence rates, survival rates, etc.) and to evaluate cancer care (diagnosis, treatment, etc.). To produce comparable data, common definitions (e.g. terminologies like the International Classification of Diseases (ICD)) and coding practices [5] have to be followed. However, the broadness and complexity of these standards make the work of the medical staff in charge of coding (operators) more difficult.

The aim of this research is to address this complexity, by assisting both operators and coding experts in the interpretation of coding best practices.

As an illustrating example, let us consider the case denoted by `exmpl` of a particular woman. In 2016, multiple pulmonary opacities were discovered within her right lung. A CT scan indicated no mediastinal adenopathy.<sup>1</sup> A histological analysis of a sample identified the morphology<sup>2</sup> of the cancer as adenocarcinoma. The TTF1 marker test was positive. After further testing, another tumor is found in the ovaries. An operator might wonder which topography<sup>3</sup> should be coded (lung or ovaries?) and can request help. For the Luxembourg National Cancer Registry (NCR), operators ask their questions using an online ticketing system. With free text description provided by operators, coding experts provide a solution, i.e. an answer with their reasoning in the form of a motivated argument.

Section 2 describes an approach to assist the data collection process for cancer registries and how case-based reasoning (CBR [1]) is applied. In Sect. 3, a prototype and preliminary results are discussed. Section 4 presents a conclusion and points out what further efforts need to be undertaken in the future.

## 2 Case-Based Interpretation of Best Practices

This article summarizes the work presented in [9] and adds a description of the developed prototype and some preliminary results.

### 2.1 Preliminaries

RDFS<sup>4</sup> is a knowledge representation language of the semantic web. SPARQL (See footnote 4) is a query language for RDFS web.

A case (`srce, sol(srce)`) is composed of two parts: (1) `srce` is a problem given by a question (i.e. a subject) and a patient record, and (2) `sol(srce)` a solution for the problem `srce`.

The question indicates the subject (incidence date, topography, tumor nature, etc.). In the example, the question is about the topography.

The patient record represents the data from the hospital patient record (patient features, tumors, exams, treatments, etc.) needed to answer the question. The relevant data depends on the subject and is defined by coding experts. The patient record is represented by an RDFS graph [3] (see Fig. 1). Body parts and cancer morphologies use classes from the SNOMED Clinical Terms<sup>5</sup> ontology.

<sup>1</sup> An adenopathy is an enlargement of lymph nodes, likely due to cancer.

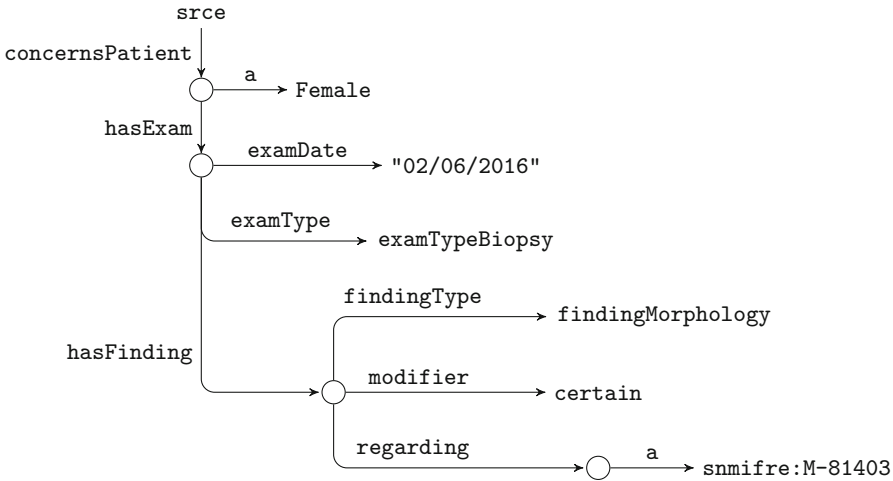
<sup>2</sup> The morphology describes the type and behavior of the cells that compose the tumor.

<sup>3</sup> The topography is the location where the tumor originated.

<sup>4</sup> <https://www.w3.org/TR/rdf-schema/> and <https://www.w3.org/TR/sparql11-query>.

<sup>5</sup> <https://bioportal.bioontology.org/ontologies/SNOMEDCT>.

The solution contains the answer to the question and the most important arguments in favor of (**pros**) and against (**cons**) this answer. In the example, the answer is to consider the topography to be the ovaries. The presence of multiple pulmonary opacities is an argument in favor, as they are indicative of a lung metastasis and thus the tumor is unlikely to have originated in the lungs.



**Fig. 1.** Short patient record in RDFS. This graph represents a woman with a single biopsy (exam), identifying the tumor as adenocarcinoma (which is coded as M-8140/3). The circles represent blank nodes.

The arguments have two uses. They help explain the answer to operators and serve as a reminder for coding experts. They are also used in the proposed approach during the retrieval step. Three types of arguments will be considered: strong pros, weak pros and weak cons. The difference between a strong and a weak argument comes from their reliability for a given conclusion. A strong argument is considered to be a sufficient justification for an answer, unlike a weak argument which is more of an indication or clue. It can be noted that there are no strong cons in the source cases. Indeed, such an argument would be an absolute argument against the given answer. Formally, an argument is a function that associates a Boolean to a case and is stored as a SPARQL ASK query. The following shows an argument **arg**, followed by an explanation:

```

arg(case) = ASK {
  case concernsPatient ?patient .
  ?patient hasExam ?exam_morpho .
  ?exam_morpho hasFinding ?finding .
  ?finding findingType findingTypeFindMorphology .
  ?finding modifier certain .
  ?finding regarding [ a smifr:M-81403 ] .
  ?patient hasExam ?exam_ttf .
  ?exam_ttf hasFinding ?finding .
  ?finding findingType findingTypeFindTTF1Marker .
  ?finding present yes .
}

```

arg says that a TTF1 positive adenocarcinoma is in favor of a primitive lung cancer. The argument checks that the morphology of the tumor is of type adenocarcinoma and that the tumor is positive for the TTF1 marker. This argument applies for the example described in the introduction, i.e. `arg(exmpl) = TRUE`.

### 2.2 Global Architecture

The proposed approach uses a 4-R cycle (retrieve, reuse, revise, retain) adapted from [1] and four knowledge containers [8] (case base, domain knowledge, retrieval knowledge, adaptation knowledge), as shown in Fig. 2.

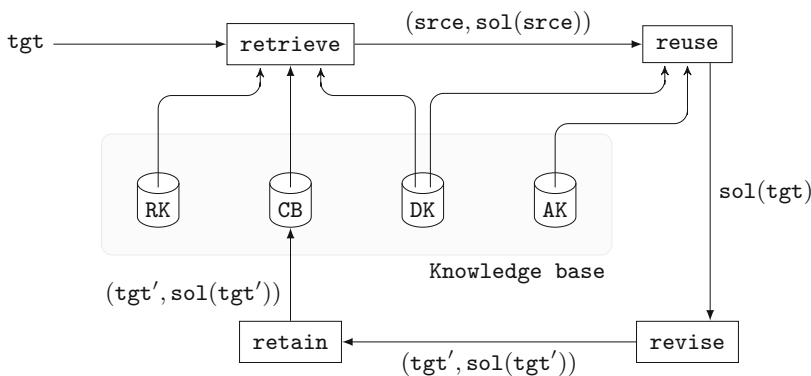


Fig. 2. Adapted 4-R cycle and knowledge containers for the proposed approach.

### 2.3 Retrieve

The proposed approach relies on arguments to find similar cases. Indeed, similar answers should be based on similar reasoning and thus the same arguments

should apply. Our method checks the applicability of arguments from source cases on the target problem **tgt** and uses this to determine the preferred source case to solve **tgt**. This preference relation is denoted by the preorder  $\preceq_{\mathbf{tgt}}$ . The comparison between two source cases **i** and **j** relies on three criteria,  $\mathcal{C}_s$  for strong arguments,  $\mathcal{C}_w$  for weak arguments and  $\mathcal{C}_{\text{dist}}$  for patient records.

An argument **arg** is applicable for a case **c** if the preconditions of the argument are met in the patient record of **c**. For the argument **arg** described in the preliminaries, **arg** applies for a case if the patient record contains at least two exams, one identifying the morphology as adenocarcinoma and another exam reporting a TTF1 positive tumor. Formally an argument **arg** is applicable for a case **c** if  $\mathbf{arg}(\mathbf{c}) = \text{TRUE}$ .

For the criterion  $\mathcal{C}_s$ , the source case with more applicable strong arguments is preferred. Formally,  $\mathcal{C}_s$  is met if  $\Delta_{i,j}^s > 0$ , where  $\Delta_{i,j}^s$  is defined as

$$\Delta_{i,j}^s = \mathcal{N}^{\text{sp}}(\text{srce}_i, \mathbf{tgt}) - \mathcal{N}^{\text{sp}}(\text{srce}_j, \mathbf{tgt})$$

and  $\mathcal{N}^{\text{args}}(\text{srce}, \mathbf{tgt})$  denotes the number of arguments of type **args** of a the source case **srce** which are applicable for a case **tgt**, i.e.

$$\mathcal{N}^{\text{args}}(\text{srce}, \mathbf{tgt}) = |\{\mathbf{a} \in \text{args}(\text{srce}) \mid \mathbf{a}(\mathbf{tgt}) = \text{TRUE}\}|$$

and **args**  $\in \{\text{sp}, \text{wp}, \text{wc}\}$  is an argument type.  $\text{sp}(\text{srce})$  is the set of strong pros,  $\text{wp}(\text{srce})$  the set of weak pros and  $\text{wc}(\text{srce})$  the set of weak cons of **srce**.

For the criterion for weak arguments  $\mathcal{C}_w$ , a combination of pros and cons is used. Intuitively, if more weak pros and less weak cons are applicable, the source case is preferred. Formally,  $\mathcal{C}_w$  is met if  $\Delta_{i,j}^w > 0$ , where  $\Delta_{i,j}^w$  is defined as

$$\begin{aligned} \Delta_{i,j}^w = & \lambda_p * (\mathcal{N}^{\text{wp}}(\text{srce}_i, \mathbf{tgt}) - \mathcal{N}^{\text{wp}}(\text{srce}_j, \mathbf{tgt})) \\ & - \lambda_c * (\mathcal{N}^{\text{wc}}(\text{srce}_i, \mathbf{tgt}) - \mathcal{N}^{\text{wc}}(\text{srce}_j, \mathbf{tgt})) \end{aligned}$$

where  $\lambda_p$  and  $\lambda_c$  are two nonnegative coefficients that are currently fixed to  $\lambda_p = 3$  and  $\lambda_c = 2$ . When more data are available, these parameters values will be reevaluated.

For the criterion  $\mathcal{C}_{\text{dist}}$ , a graph edit distance between patient record RDFS graphs is used [4]. Formally,  $\mathcal{C}_{\text{dist}}$  is met if  $\Delta_{i,j}^{\text{dist}} \geq 0$ , where  $\Delta_{i,j}^{\text{dist}}$  is defined as

$$\Delta_{i,j}^{\text{dist}} = \text{dist}(\text{srce}_j, \mathbf{tgt}) - \text{dist}(\text{srce}_i, \mathbf{tgt})$$

The three criteria are considered lexicographically, first  $\mathcal{C}_s$ , then  $\mathcal{C}_w$  and finally  $\mathcal{C}_{\text{dist}}$  (see [9]).  $\text{srce}_i$  is preferred over  $\text{srce}_j$ , i.e.  $\text{srce}_i \preceq_{\mathbf{tgt}} \text{srce}_j$ , ifq

$$\Delta_{i,j}^s > 0 \text{ or } (\Delta_{i,j}^s = 0 \text{ and } (\Delta_{i,j}^w > 0 \text{ or } (\Delta_{i,j}^w = 0 \text{ and } \Delta_{i,j}^{\text{dist}} \geq 0)))$$

## 2.4 Reuse

Once an appropriate source case has been found, the solution associated to the source case is copied:  $\text{sol}(\mathbf{tgt}) := \text{sol}(\text{srce})$ . The arguments that do not apply to the target problem, if any, are removed.

## 2.5 Revise and Retain

The newly formed case  $(tgt, sol(tgt))$  can be reviewed by a coding expert, to modify the answer, the arguments and/or the patient record. A coding expert may choose to remove unnecessary information from the patient record, removing unwanted specificity. Thus,  $(tgt, sol(tgt))$  is substituted by  $(tgt', sol(tgt'))$ , where  $tgt'$  is more general than  $tgt$ .  $(tgt', sol(tgt'))$  is a generalized case that has a larger coverage than  $(tgt, sol(tgt))$  [6].

## 3 Prototype and Preliminary Results

The prototype designed for the NCR serves as a ticketing system, where operators ask coding questions and experts provide answers. It assists operators in structuring questions, making it easier for the NCR and coding experts to find similar questions later. For topography questions, it will also provide a tentative answer. This answer is calculated using the approach described in [9]. All the answers are reviewed by experts. The prototype presents itself as a single page

Questions > Question 653

Sujets : Topographie  
 Type de cancer : Cancer gynéco, Cancer du poumon

Posé par Pierre Sevérick (Hopital du Luxembourg, Registre National du Cancer)  
 Posé le 28/07/2017 (11:24)

Description	Réponse
<p>Femme (age non-renseigné)</p> <p>Imagerie (5/1/2016)</p> <ul style="list-style-type: none"> <li>• Image en lâcher de ballons</li> </ul> <p>CT scan (9/1/2016)</p> <ul style="list-style-type: none"> <li>• Pas adénopathies trouvées : ganglion lymphatique médiastinal</li> </ul> <p>Biopsie (13/2/2016)</p> <ul style="list-style-type: none"> <li>• Marqueur TTF1 : positif</li> <li>• Morphologie : adénocarcinome, sai</li> </ul> <p>PET scan (1/3/2016)</p> <ul style="list-style-type: none"> <li>• Lésion tumorale : ovaire, sai</li> </ul>	<p>C80.9 : origine primitive inconnue</p> <p>Arguments</p> <ul style="list-style-type: none"> <li>• <span style="color: green;">✔</span> Aucun argument fort en faveur.</li> <li>• <span style="color: green;">✔</span> On observe une image en lâcher de ballons, résultat typique d'une métastase pulmonaire.</li> <li>• <span style="color: red;">✘</span> Le marqueur TTF1 est positif (souvent pour les tumeurs primitives du poumon, ce marqueur est positif).</li> </ul>

**Résolue en réutilisant la question 609**

Sujets : Topographie  
 Type de cancer : Cancer gynéco, Cancer du poumon

Posé par Pierre Sevérick (Hopital du Luxembourg, Registre National du Cancer)  
 Posé le 08/03/2016 (13:40)

Description	Réponse
<p>Avis RCP</p> <ul style="list-style-type: none"> <li>• Lésion tumorale : origine primitive inconnue</li> </ul> <p>CT scan</p> <ul style="list-style-type: none"> <li>• Aucune lésion tumorale</li> </ul> <p>PET scan</p> <ul style="list-style-type: none"> <li>• Lésion tumorale : poumon droit, sai, poumon gauche, sai</li> </ul> <p>Biopsie</p> <ul style="list-style-type: none"> <li>• Marqueur TTF1 : positif</li> <li>• Morphologie : adénocarcinome, sai</li> </ul> <p>Imagerie</p> <ul style="list-style-type: none"> <li>• Image en lâcher de ballons</li> </ul>	<p>C80.9 : origine primitive inconnue</p> <p>Arguments</p> <ul style="list-style-type: none"> <li>• <span style="color: green;">✔</span> Aucun argument fort en faveur.</li> <li>• <span style="color: green;">✔</span> En RCP, on a conclu pour une origine primitive inconnue.</li> <li>• <span style="color: green;">✔</span> On observe une image en lâcher de ballons, résultat typique d'une métastase pulmonaire.</li> <li>• <span style="color: red;">✘</span> Le marqueur TTF1 est positif (souvent pour les tumeurs primitives du poumon, ce marqueur est positif).</li> <li>• <span style="color: red;">✘</span> Aucune autre tumeur synchrones n'a été découverte.</li> </ul>

**Fig. 3.** Example of a solved case. The top displays the new question asked and the provided solution. The bottom displays the source case used to solve the new question.

application built using Angular<sup>6</sup> with a backing REST API built with Go (See footnote 6) and the Gin framework.<sup>7</sup> The data is stored in a triple store Apache Jena and exposed as a SPARQL endpoint using Apache Fuseki.<sup>8</sup> Figures 3 and 4 show screenshots of the prototype.

**Fig. 4.** Form used to describe coding questions and patient records. The French labels for body parts and morphologies are taken from the SNMIFRE (a French translation of SNOMED, <http://bioportal.lirmm.fr/ontologies/SNMIFRE>).

The prototype was tested internally, to perform a first assessment of its usability and utility. Some old cases concerning the topography were formalized and coded, with some domain knowledge. For the arguments, great care was given during modeling in order to make them more broadly applicable. Then new questions were presented to the system, and the proposed solution compared with the expected ones. While the prototype answered every question, not all of them were correct. The main reasons for the difference were the small amount of cases (15 originally, however the case base will be enriched by routine usage) and the simple reuse method used at this stage. Indeed, as the arguments have been formalized to be more general, some of the provided answers might be slightly incorrect (e.g. answering upper lung lobe instead of lower lung lobe). Despite this, as the prototype displays the reused source case, an operator should be able to make the necessary adaptation to the provided solution. For the questions concerning other subjects, the prototype relies entirely on the coding experts to provide answers.

<sup>6</sup> <https://angular.io>.

<sup>7</sup> <https://golang.org>, <https://github.com/gin-gonic/gin>.

<sup>8</sup> <https://jena.apache.org/> and <https://jena.apache.org/documentation/fuseki2/>.

To the best of our knowledge, few other research attempts to use arguments in the context of the retrieval process. The closest method found is a work by McSherry [7]. The proposed approach creates explanations afterwards, using the closest source case to provide the conclusion and the closest source case with the opposite conclusion to compute which attributes favor the conclusion and which attributes do not. Unlike our approach, each argument is linked to a single attribute. Thus they cannot show how the combination of attributes might influence a given outcome.

## 4 Conclusion

Recently there has been a growing interest for case-based reasoning applications in health sciences [2]. In this paper, an approach to assist operators in the interpretation of best medical coding practices has been proposed. This approach is based on discussions with operators and coding experts on actual coding problems. A dozen tricky problems were discussed in detail, among a hundred simpler problems. The coding questions asked by the operators are compared to previous questions and solved by reusing the pros and cons of previously given solutions. The results discussed are only preliminary and a more thorough evaluation, including the operators and coding experts, is planned.

At the moment the reasoning process is only partial. Arguments are only a part of a more complex reasoning process. The formalization of this process and the eventual integration of the coding standards remains an interesting avenue for future work.

After the prototype has been validated and improved by routine usage, a second version will be designed that is less domain-dependent. The objective is to build a generic system for argumentative case-based reasoning using semantic web standards.

**Acknowledgments.** The authors wish to thank the anonymous reviewers of the Joint Workshop on Artificial Intelligence in Health for their remarks which have helped in improving the quality of the paper. The first author would also like to thank the Fondation Cancer for their financial support.

## References

1. Aamodt, A., Plaza, E.: Case-based reasoning: foundational issues, methodological variations, and system approaches. *AI Commun.* **7**, 39–59 (1994)
2. Bichindaritz, I., Marling, C., Montani, S.: Case-based reasoning in the health sciences. In: *Workshop Proceedings of ICCBR (2015)*
3. Brickley, D., Guha, R.V.: RDF Schema 1.1, W3C recommendation, last consultation: March 2017 (2014). <https://www.w3.org/TR/rdf-schema/>
4. Bunke, H., Messmer, B.T.: Similarity measures for structured representations. In: Wess, S., Althoff, K.-D., Richter, M.M. (eds.) *EWCBR 1993*. LNCS, vol. 837, pp. 106–118. Springer, Heidelberg (1994). [https://doi.org/10.1007/3-540-58330-0\\_80](https://doi.org/10.1007/3-540-58330-0_80)



5. Tyczynski, J.E., Démaret, E., Parkin, D.M., European Network of Cancer Registries: Standards and Guidelines for Cancer Registration in Europe: the ENCR Recommendations. International Agency for Research on Cancer, Lyon (2003)
6. Maximini, K., Maximini, R., Bergmann, R.: An investigation of generalized cases. In: Ashley, K.D., Bridge, D.G. (eds.) ICCBR 2003. LNCS (LNAI), vol. 2689, pp. 261–275. Springer, Heidelberg (2003). [https://doi.org/10.1007/3-540-45006-8\\_22](https://doi.org/10.1007/3-540-45006-8_22)
7. McSherry, D.: Explaining the pros and cons of conclusions in CBR. In: Funk, P., González Calero, P.A. (eds.) ECCBR 2004. LNCS (LNAI), vol. 3155, pp. 317–330. Springer, Heidelberg (2004). [https://doi.org/10.1007/978-3-540-28631-8\\_24](https://doi.org/10.1007/978-3-540-28631-8_24)
8. Richter, M.M., Weber, R.O.: Case-Based Reasoning: A Textbook. Springer, Berlin (2013). <https://doi.org/10.1007/978-3-642-40167-1>
9. Schnell, M., Couffignal, S., Lieber, J., Saleh, S., Jay, N.: Case-based interpretation of best medical coding practices—application to data collection for cancer registries. In: Aha, D.W., Lieber, J. (eds.) ICCBR 2017. LNCS (LNAI), vol. 10339, pp. 345–359. Springer, Cham (2017). [https://doi.org/10.1007/978-3-319-61030-6\\_24](https://doi.org/10.1007/978-3-319-61030-6_24)