

# A Knowledge-Based Framework for Information Extraction from Clinical Practice Guidelines

Corrado Loglisci, Michelangelo Ceci, and Donato Malerba

Dipartimento di Informatica, Università degli Studi di Bari  
via Orabona, 4 - 70126 Bari, Italy  
{loglisci, ceci, malerba}@di.uniba.it

**Abstract.** Clinical Practice Guidelines guide decision making in decision problems such as the diagnosis, prevention, etc. for specific clinical circumstances. They are usually available in the form of textual documents written in natural language whose interpretation, however, can make difficult their implementation. Additionally, the high number of available documents and the presence of information for different decision problems in the same document can further hinder their use. In this paper, we propose a framework to extract practices and indications considered to be important in a particular clinical circumstance for a specific decision problem from textual clinical guidelines. The framework operates in two consecutive phases: the first one aims at extracting pieces of information relevant for each decision problem from the documents, while the second one exploits pieces of information in order to generate a structured representation of the clinical practice guidelines for each decision problem. The application to the context of Metabolic Syndrome proves the effectiveness of the proposed framework.

**Keywords:** Information Extraction, Clinical Practices Guidelines, Medical Decision Making.

## 1 Introduction

Clinical practice guidelines (CPGs) are systematically developed statements to assist practitioner and patient decisions about appropriate health care for specific clinical circumstances [5]. They can support medical personnel in deciding the activities to follow for the diagnosis, treatment, prevention or management (decision problems, DP) of the specific health condition of a patient.

However, CPGs are usually available in the form of unstructured textual documents written in natural language whose interpretation often hinders their implementation at point of care [2]. Indeed, the typical lacking of structure in textual guidelines and the usual blending of information of several DP in a CPG can make the role of guidelines cumbersome and the work of practitioners unbearable.

One of the successful strategies for this issue is the integration of information technologies and research frameworks in health care environments. Works

reported in the literature follow two main research lines: 1) *model-centric*, which aims to model CPGs with directly computable formalisms (e.g., ontologies [10]), 2) *document-centric*, which aims to transform CPGs into a human interpretable structured format. In the latter large attention has been paid to approaches for manually marking-up the unstructured textual CPGs and transferring the literal content into a pre-defined representation format [9]. However, both research lines are likely to be resource-consuming because demand knowledge on both medical domains and modeling languages and require a lot of human intervention especially for the activity of marking-up the CPGs. A different approach is that of (semi)-automatically processing the content of unstructured textual CPGs and generate a structured representation of it based on pieces of information in natural language present in the text: CPGs are thus represented in a format more rigorous than that unstructured, which practitioners can more easily follow. The problem is often faced by adapting domain-independent frameworks of Information Extraction (IE) to the case of CPGs.

However, adapting IE techniques can be expensive and can suffer for low results accuracy because of the complexity of the CPGs. Instead, ad-hoc approaches to process CPGs are often designed for particular medical domains, specific sections of documents or specific DP (e.g., treatment [7]).

In this paper we present a two-phase framework which aims to extract information on practices and indications of any DP from all sections which compose the textual CPGs. It transforms CPGs in a human-interpretable and structured format, require little human intervention and is designed ad-hoc for CPGs. The framework operates in two consecutive phases. The first one localizes and recognizes relevant pieces of information from the sections composing CPGs through Text Processing techniques. The second one builds a structured representation of the CPG of a specific DP by filling a template structure with the previously extracted pieces of information: practices and indications will be represented as a composition of pieces of information. The paper is so organized. In the next section, after a brief overview of similar techniques, we point out the peculiarities of this work and introduce the proposed two-phase framework. A detailed description of the two phases is provided in the sections 3 and 4 respectively. In Section 5 we explore the application to the domain of Metabolic Syndrome and a quantitative evaluation is reported too. Finally conclusions close this work.

## 2 Related Works and Contribution

As stated before, extracting relevant information from textual CPGs in order to support their application is an approach investigated in the literature either adapting domain-independent IE techniques or building ad-hoc IE systems. In both cases the extraction is guided by a rule set previously acquired (typically called *extraction patterns*). The mode of obtaining the rule set divides these techniques in two threads: *Automatic Learning*, when the rules are induced from example documents through a learning process [3], and *Knowledge Engineering*

(KE), when the rules are defined by exploiting knowledge on the documents [7,12]. Due to the difficulty to prepare a large training set and the high variability and complexity of the sections of textual CPGs, KE is often preferred to Automatic Learning [1]. A representative paper of the KE approach which is similar to ours has been proposed by Kaiser et al. [7]. In their work the goal was to extract information on clinical treatment processes and use it to support the manual modeling of the guidelines concerning that specific DP. The method outputs representations of the CPGs at different level of refinement by exploiting hand-made heuristics (extraction patterns) defined at phrase-level and discourse-level. However, although this approach allows to structure guidelines even at levels of single activities, the usage of extraction patterns tailored for the clinical treatment make it adequate only for that specific DP (in this case treatment).

Differently, our framework does not limit its application to a specific DP and extracts information from the several sections composing the CPGs. This allows us to take into account an important characteristic of the textual guidelines: a CPG contains practices and indications of any DP disseminated in the several sections of the document. For instance, the practices of any DP are summarized and reported together in a section while they are detailed and extended in another section.

The framework processes original CPGs with Text-Processing techniques [11] in order to extract relevant pieces of information in the several sections (first phase). Next, for each DP of a specific health condition, it builds a structured representation by arranging pieces of information in a template structure with a KE approach which exploits a production system (second phase). In this work a piece of information (afterward, pINF) has to be intended as a textual unit composed of (at least) a sentence. This textual unit expresses a well-defined activity to follow for a specific DP. A practice or an indication of a specific DP is then defined as a set of activities.

### 3 Extraction of Relevant Pieces of Information

The first phase is performed through natural language processing functionalities which analyze the textual CPGs at both section-level and sentence-level and output pINF and annotations. These annotations are represented in the form of attribute-value pairs and express information concerning the sections and pINF. Analysis proceeds in two steps: the first one (TS1) merely segments CPGs into sections and subsections, and obtains relational properties on them (i.e., ordering among (sub)sections, membership of a subsection to a section), while the second one (TS2) performs a linguistic analysis on the obtained (sub)sections aiming at capturing the semantics of pINF and representing it in the form of annotations.

More precisely, TS1 starts by transforming initial CPGs into semi-structured documents with tagged sections since guidelines are often released as marked-up and irregularly formatted documents (e.g., HTML documents). This is done

integrating the Xerces Parser tool<sup>1</sup>. Next, a segmentation on tagged sections is performed to split them into subsections and then to identify sets of pINF. Trivially, a section contains several subsections which, in their turn, contain sets of pINF. TS1 generates also instances of two kinds of annotations: *GROUP* and *VALUE* (see Table 1). *GROUP* annotations represent sections and subsections, and express relational properties on them. *VALUE* annotations rather represent information on sections, subsections and the contained pINF.

For instance, given the section *Scope* containing the subsection *Interventions and Practices considered* of the CPG concerning *Essential Hypertension* below reported obtained after the transformation in a semistructured document, TS1 generates the *GROUP* and *VALUE* annotations illustrated in Table 2.

```
<doc_id> <section_lbl>Scope</section_lbl>
  <field_lbl>Interventions and Practicesconsidered</field_lbl>...
  <strong> Treatment/Management </strong> <ol start="1" type="1">
  <li> Drug therapy <ul type="disc">
    <li> Diuretics (thiazide or loop) </li>
    <li> Beta blockers </li>
    <li> Angiotensin converting enzyme (ACE) inhibitors </li> ... </doc_id>
```

A further segmentation on the sets of pINF (namely, *CONTENT* of the *VALUE*) is performed w.r.t. delimiters of text in order to recognize single pINF and their relational properties. Here three annotations are produced: *LIST*, *ITEM*, *PARAGRAPH* (see Table 1). *LIST* describes the placement of a collection of pINF affiliated to a subsection, while *ITEM* and *PARAGRAPH* annotate pINF respectively contained in a *LIST* and not contained but included in the *CONTENT* attribute of the *VALUE*. In addition *LIST*, *ITEM*, *PARAGRAPH* express relational properties on pieces of information present in each subsection. By following the example above, from the attribute *CONTENT* of *VALUE* the *LIST* and *ITEM* annotations illustrated in Table 2 can be generated.

Once the pINF have been localized and annotated, TS2 can be performed. It resorts to linguistic analysis techniques [11] first to split pINF containing several sentences into single sentences, then to capture their semantics. Linguistic analysis exploits hand-coded controlled dictionaries and grammars and includes tokenization, sentence splitting, stemming, part-of-speech tagging and named-entity recognition to be executed in sequence. Three further annotations are produced: *SENTENCE*, *WORD*, *DOMAIN* which describe respectively i) placement of the sentence contained in the pINF, ii) lexical and morphological features of the words contained in the sentences and iii) domain-specific information on the *SENTENCE* annotations (see Table 1). For instance, from the *ITEM* annotations of the previous example the *SENTENCE* and *WORD* annotations in Table 2 can be generated. Annotations so produced will be exploited by Template Filling method to arrange pINF in a template form, namely the final structured representation of the CPGs.

---

<sup>1</sup> <http://xerces.apache.org/xerces-j/>

**Table 1.** Representation of the annotations generated by the first phase

<i>GROUP: [ID, TYPE, CONTENT, PARENT, POSITION]</i>
where ID is a unique identifier of the annotation, TYPE is the tag of the corresponding (sub)section, CONTENT is the tag value, PARENT is the identifier of father GROUP for the subsections, POSITION is a progressive index which univocally represents the section in the CPG (in the case of sections) or the subsection in the section (in the case of subsections).
<i>VALUE: [ID, CONTENT, GROUP]</i>
where CONTENT is the tag value, GROUP is the identifier of the group of affiliation of the VALUE.
<i>LIST: [ID, TYPE, GROUP, POSITION]</i>
where TYPE denotes whether the list is ordered or not, GROUP indicates of the group of affiliation of LIST, POSITION is a progressive index representing the position of the LIST in its GROUP.
<i>ITEM: [ID, CONTENT, LIST, POSITION]</i>
where CONTENT is the textual content of the represented element, LIST indicates the LIST which contains the ITEM, POSITION is a progressive index representing the ITEM in its LIST.
<i>PARAGRAPH: [ID, CONTENT, GROUP, POSITION]</i>
where CONTENT is the textual content of the represented element, GROUP indicates the GROUP which contains the PARAGRAPH, POSITION is the position of PARAGRAPH in its GROUP.
<i>SENTENCE: [ID, CONTENT, PARAGRAPH, ITEM, GROUP, POSITION]</i>
where CONTENT is the textual content, PARAGRAPH, ITEM, GROUP denote the object of the current annotation: only one of them can be instantiated. PARAGRAPH is instantiated when the sentence refers to a PARAGRAPH, ITEM is instantiated when the sentence refers to a ITEM and GROUP is instantiated when the sentence refers to a GROUP and none of the previous. POSITION is a progressive index of the SENTENCE.
<i>WORD: [ID, CONTENT, SENTENCE, CATEGORY, POS, PROPERTIES, POSITION]</i>
where CONTENT is the textual content, SENTENCE denotes the SENTENCE of the current annotation, CATEGORY represents a generalization of the concept expressed by CONTENT (e.g. organ is the category of liver) according to controlled vocabularies, POS denotes the part-of-speech tag of CONTENT, PROPERTIES expresses linguistic/orthographic properties of CONTENT (e.g. lowercase), POSITION is a progressive index of WORD in the corresponding SENTENCE.
<i>DOMAIN: [ID, CONTENT, SENTENCE, CATEGORY]</i>
where CONTENT is the textual content of the corresponding SENTENCE, SENTENCE denotes the SENTENCE of the current annotation, CATEGORY represents a generalization of the concept expressed by CONTENT (e.g. intended users is the category of Allied Health Personnel) according to controlled vocabularies.

**Table 2.** Annotation examples generated by the first phase

<i>GROUP: [ID1, Section, Scope, null, 1]</i>
<i>GROUP: [ID2, Field, Interventions and Practices Considered, ID1, 1]</i>
<i>VALUE: [ID3, &lt;strong&gt; Treatment/Management &lt;/strong&gt; &lt;ol start="1" type="1"&gt; &lt;li&gt; Drug therapy &lt;ul type="disc"&gt; &lt;li&gt; Diuretics (thiazide or loop) &lt;/li&gt; &lt;li&gt; Beta blockers &lt;/li&gt; &lt;li&gt; Angiotensin converting enzyme (ACE) inhibitors &lt;/li&gt;, ID2]</i>
<i>LIST: [ID4, ordered, ID2, 1]</i>
<i>ITEM: [ID5, Drug therapy, ID4, 1]</i>
<i>LIST: [ID6, unordered, ID2, 1]</i>
<i>ITEM: [ID7, Angiotensin converting enzyme (ACE) inhibitors, ID6, 3]</i>
<i>SENTENCE: [ID13, Drug therapy, null, ID5, null, 1]</i>
<i>SENTENCE: [ID9, Angiotensin converting enzyme (ACE) inhibitors, null, ID7, null, 1]</i>
<i>WORD: [ID8, Drug therapy, ID13, therapeutic procedure, complex POS, null, 1]</i>
<i>WORD: [ID10, Angiotensin converting enzyme, ID9, pharmacological substance, complex POS, null, 1]</i>
<i>WORD: [ID11, ACE, ID9, pharmacological substance acronym, nn, upperCase, 2]</i>
<i>WORD: [ID12, inhibitor, ID9, word, nns, lowerCase, 3]</i>

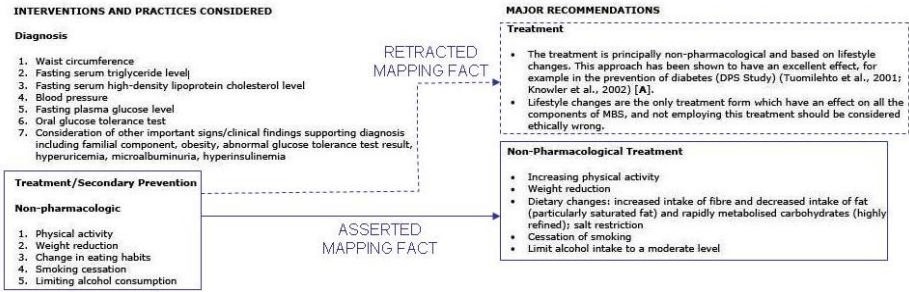
## 4 Template Filling for Structuring CPGs

A final CPG has a pre-defined template structure composed of five main slots named as *Disease/ Health Conditions*, *Target Population*, *Benefits*, *Harms*, *Practices*. Each of them, in its turn, can contain nested slots: for instance, the *Practices* slot is composed of as many inner slots as the practices of the corresponding CPG are, while each inner slot contains, in its turn, as many basic slots as the activities of the corresponding practice are. It is thus expected that two different CPGs have different numbers of inner and basic slots albeit they present the same number of main slots.

The filling method constructs the final CPG through a bottom-up strategy which fills basic slots and merges them to return inner slots up to obtain the main slots: initially it acquires information on the pINF (annotations) for a particular DP of a specific health condition, then it infers the content of the slots *Disease/ Health Conditions*, *Target Population*, *Benefits*, *Harms* (which express information required by standard conventions or expected by practitioners). Subsequently, it derives the content of the *Practices* slot by first identifying the sections of the original CPG which potentially contain the single practices and, then, by inferring, for each practice, the content of its slots, namely the activities described in summarized and detailed form.

The filling procedure is based on a production system with forward chaining inferential mechanism [6], which is composed of an inferential engine, facts-list and knowledge base. Facts-list contains two kind of assertions: the previously generated annotations represented in a suitable language and facts deduced during the inferential process which determine the pINF which fill the slots. Deducted facts, in their turn, are divided into *mapping* and *matching* facts. Mapping facts identify the sections which can contain pINF for the slots, while matching ones indicate pINF which can fill the slots. Each mapping or matching fact presents a numerical score based on the annotations of the associated pINF. During the inferential process different matching facts can compete for the same slot and the final assignment is decided on the basis of the higher score value. Knowledge base consists of *if-then* rules whose conditional part is described in terms of facts, while the consequential part specifies particular operations to do on the facts-list. Rules are divided into standard conventions rules (SCr), control rules (Cr), filling rules (Fr). SCr define characteristics that final CPG must have and information expected by practitioners. For instance, one of the SCr rules states that in the final CPGs the summarized description of each practice has to be followed by the detailed description. Cr regulate the activation of Fr by retracting facts which can fire Fr when other facts with higher score have been previously inferred. Finally, Fr can infer matching facts (namely, facts associated to the pINF which can fill the slots) and mapping facts (namely, facts on the sections of the original CPG which can contain potential pINF for the slots).

A concrete example is reported in Table 3 and illustrates the inference of facts for the original CPG reported in Figure 1. The subsections *Interventions and Practices considered* (IPc) and *Major Recommendations* (MR) contain respectively the summarized and detailed description of the practices for



**Fig. 1.** Example of activation of rules in the knowledge base: the rule retracts a mapping fact previously inferred and asserts another one having higher numerical score

any DP. Suppose that IPc and MR have been annotated with GROUP1 and GROUP2 respectively. The rule in Table 3 derives that the practices for the non-pharmacological treatment DP have to be sought in the subsections *Non-pharmacologic* of IPc (annotated with GROUP1\_1) and *Non-pharmacological Treatment* of MR (GROUP2\_2): more precisely, it retracts the mapping fact (dotted line in Figure 1) on *Non-pharmacologic* of IPc and *Treatment* of MR (GROUP2\_1) and infers the new mapping fact (solid line in Figure 1) on the basis of the occurrence of the same DP in the SENTENCE annotations (i.e., Treatment DP) and the occurrence of the shared terms in WORD annotations (e.g., Non-pharmacologic).

## 5 Application to the Metabolic Syndrome CPGs

In order to prove the viability of the proposed framework we applied it to the context of CPGs concerning Metabolic Syndrome. The phase of extraction of relevant pINF exploits the facilities of GATE (General Architecture for Text Engineering) system [4], while the production system used for template filling is developed as an expert system in CLIPS language. Moreover, the vocabularies<sup>3</sup> used for linguistic analysis are built considering the dictionaries available

**Table 3.** Example of Cr used for regulating the activation of Fr on a mapping fact

<p><b>IF</b> word_set(SENTENCE(GROUP1_1))∩word_set(SENTENCE(GROUP2_2))≠ ∅ and DP(GROUP1_1)==DP(GROUP2_1) and affiliationToSummary(GROUP1_1) and affiliationToDetail(GROUP2_2) and any_mapping_exists(GROUP1_1)</p> <p><b>THEN</b> retract (any_mapping(GROUP1_1)) and assert(mapping(GROUP1_1,GROUP2_2))</p> <p><b>FACTS-LIST:</b> mapping_fact(GROUP1_1,GROUP2_1), [SENTENCE1, Non-pharmacologic, null, null, GROUP1_1, GROUP1, 1], [SENTENCE3, Non-pharmacological Treatment, null, null, GROUP2_2, GROUP2, 1], [SENTENCE2, Treatment, null, null, GROUP2_1, GROUP2, 1], [WORD1, Non-pharmacologic, SENTENCE1, GROUP1_1, word, adjective, upperInitial, 1], [WORD2, Non-pharmacologic, SENTENCE3, word, adjective, upperInitial, 1], [WORD3, Treatment, SENTENCE3, Treatment DP, noun, upperInitial, 2], [WORD4, Treatment, SENTENCE2, Treatment DP, noun, upperInitial, 1]</p>
---

**Table 4.** Experimental results: percentage values of precision and recall

Guideline Title	Decision Problem	#avs	#tfs	#cfs	recall	precision
Osteoporosis in gastrointestinal disease	diagnosis	17	18	17	94	100
	treatment	27	29	22	81	75.8
	management	27	29	22	81	75.8

in Unified Medical Language System (UMLS) specific for Metabolic Syndrome provided by domain experts. The set of initial unstructured guidelines was retrieved by submitting the query “Metabolic Syndrome” to National Guideline Clearinghouse (NGC)<sup>2</sup> search engine: from the returned CPGs (more than 100) we selected 39 documents which did not present text in tabular structures and which concern at least one of the following decision problems: Diagnosis, Treatment, Management, Prevention, Risk Assessment, Evaluation. A subset of 23 documents was thus used to develop the knowledge base of the production system whose rules were hand-coded by exploiting knowledge of the practitioners on which sections are present within a CPG, which sections are useful for them and standard conventions of CPGs. Knowledge base comprises a total set of 40 rules<sup>3</sup> so partitioned: SCr (5), Cr (17), Fr (18). The similar organization of these documents permits us to analyze a small set of them and apply the derived rule set also to others. The remaining 16 documents were processed by the framework which returned 43 templates so distributed: Diagnosis (15), Management (11), Treatment (11), Prevention (3), Screening (1), Evaluation(2). Final CPGs (namely, filled templates) were evaluated according to *Precision*, *Recall*, *macroaveraged Precision* ( $\pi^M$ ) and *macroaveraged Recall* ( $\rho^M$ ) [8]. Recall estimates the number of filled slots w.r.t. the total number of available slots ( $\#avs$ ), while Precision estimates the number of correctly filled slots ( $\#cfs$ ) against all filled slots ( $\#tfs$ ). The values of  $\pi^M$  and  $\rho^M$  amount respectively to 79.39% and 80.79% w.r.t the total number of the final CPGs. However, for limitations of space, we report only the most significant CPG (see Table 4).

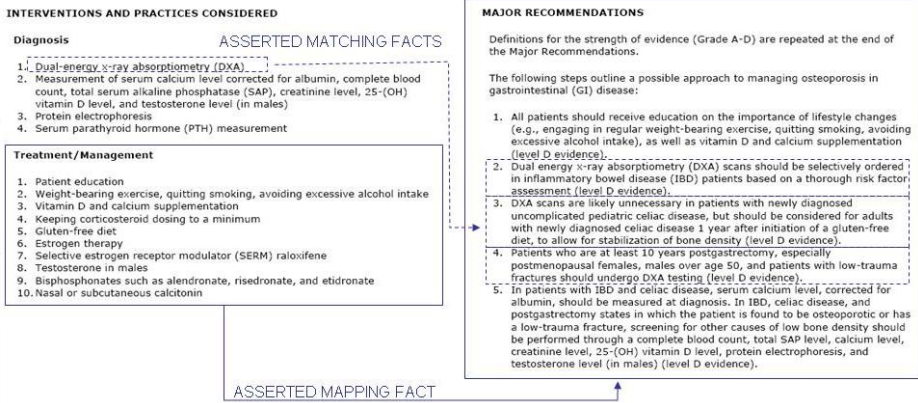
By analyzing the results of Table 4 it emerges that templates of different DP can present different recall and precision values even when extracted from the same initial guideline (e.g., diagnosis and treatment of *Osteoporosis in gastrointestinal disease*). This can be attributed to the irregular organization of the detailed description of the practices, namely MR section (see Figure 1). For instance, in the case of *Osteoporosis in gastrointestinal disease* (see Figure 2), the detailed description of the Treatment practices is blended in the MR section with the practices of other DPs. In this case, after that the annotations for the sections IPc and MR have been instantiated, the procedure of template filling first searches for the summarized practices of the Treatment DP in IPc, then infers that the Treatment practices have in the MR section a dedicated and subsection within which they are detailed, while that subsection actually does not exist. Hence, it derives that the detailed practices have to be sought in the entire

<sup>2</sup> <http://www.guidelines.gov/>

<sup>3</sup> Downloadable at

<http://www.di.uniba.it/~malerba/software/FELIPE/CPGresources/>





**Fig. 2.** Summarized and detailed descriptions of Treatment and Diagnosis practices for *Osteoporosis in gastrointestinal disease*

section MR and infers the mapping fact between the GROUP of Treatment of IPc and the GROUP of MR (solid line in Figure 2) which anyway contains also the practices of other DP. These may represent “noisy” practices for the generation process of the Treatment template and may negatively influence the final precision and recall (75.8%, 81% w.r.t. the averages 79.39% and 80.79%).

Another factor which can affect Precision and Recall lies in the domain-specific vocabularies used for the linguistic analysis: exploiting well curated vocabularies can facilitate the instantiation of annotations which better capture the semantics and domain knowledge of pINF (e.g., WORD, DOMAIN annotations in Table 1). For instance, in the case of the Diagnosis template in Figure 2, integrating domain acronyms (i.e., DXA) and hierarchies among background concepts (i.e., x-ray absorptiometry is a Diagnostic procedure) into vocabularies allows to generate WORD and DOMAIN annotations, for the subsection Diagnosis in IPc and for the section MR, which i) better express information of the Diagnosis practices, ii) mitigate the effect of the irregular organization of the MR section and iii) improve the final accuracy of Diagnosis template. A concrete illustration of that is represented by the inference of the matching fact (dotted line) of the practice 1 of the Diagnosis in Figure 2: indeed, its summarized description is extracted by the list item 1 in IPc while that detailed is localized in the items 2,3,4 of the MR section.

## 6 Conclusions

In this paper we have proposed a computational solution to support the usage and interpretation of CPGs. It based on the extraction of information deemed to be important and the representation of it in a structured mode which practitioners can more easily follow. The approach investigates two particular issues of CPGs: the presence of information (i.e., practices and indications) concerning

different decision problems and the typical lacking of structure of the textual guidelines. The purposeful aspect of this work is the automatic identification of practices relevant for a given problem in the several sections of CPGs and the presentation of these in the structured form where each practice is both synthetically and fully described. The application to the Metabolic Syndrome scenario shows the adaptability of the approach also to real contexts. As future work, we intend to integrate an Automatic Learning technique into the filling procedure to generate more accurate structured representation of the CPGs.

## Acknowledgments

This work is partial fulfillment of objective of ATENEO-2009 project “Modelli e metodi computazionali per la scoperta di conoscenza in dati biomedici”.

## References

1. Appelt, D.E.: Introduction to information extraction. *AI Communications* 12, 161–172 (1999)
2. Cabana, M.D., Rand, C.S., Powe, N.R., Wu, A.W., Wilson, M.H., Abboud, P.A., Rubin, H.R.: Why don't physicians follow clinical practice guidelines? A framework for improvement. *JAMA* 282(15), 1458–1465 (1999)
3. Califf, M.E., Mooney, R.J.: Relational Learning of Pattern-Match Rules for Information Extraction. In: *Proc. of AAAI/IAAI*, pp. 328–334 (1999)
4. Cunningham, H., Maynard, D., Bontcheva, K., Tablan, V.: Gate: A framework and graphical development environment for robust nlp tools and application. In: *40th Anniversary Meeting of the Association for Computational Linguistics* (2002)
5. Field, M.J., Lohr, K.H. (eds.): *Clinical Practice Guidelines: Directions for a New Program*. National Academy Press, Institute of Medicine (1990)
6. Ignizio, J.P.: *Introduction to Expert Systems: The Development and Implementation of Rule-Based Expert Systems*. McGraw-Hill, Inc., New York (1991)
7. Kaiser, K., Miksch, S.: Modeling Treatment Processes Using Information Extraction. *Advanced Computational Intelligence Paradigms in Healthcare* (1), 189–224 (2007)
8. Sebastiani, F.: Machine learning in automated text categorization. *ACM Comput. Surv.* 34(1), 1–47 (2002)
9. Svatek, V., Ruzicka, M.: Step-by-step formalisation of medical guideline content. *Int. Journal of Medical Informatics* 70(2-3), 329–335 (2003)
10. Wang, D., Peleg, M., Tu, S.W., Boxwala, A.A., Ogunyemi, O., Zeng, Q.T., Greenes, R.A., Patel, V.L., Shortliffe, E.H.: Design and implementation of the GLIF3 guideline execution engine. *Journal of Biomedical Informatics* 37(5), 305–318 (2004)
11. Weiss, S., Indurkha, N., Zhang, T., Damerou, F.: *Text Mining: Predictive Methods for Analyzing Unstructured Information*. Springer, Heidelberg (2004)
12. Yangarber, R., Grishman, R.: NYU: description of the Proteus/ PET system as used for MUC-7 ST. In: *Proc. of the 7th MUC*. Morgan Kaufmann, San Francisco (1998)