Using Text Generation to Access Clinical Data in a Variety of Contexts

Olof TORGERSSON and Göran FALKMAN*

Department of Computing Science, Chalmers University of Technology, and Göteborg University, SE-412 96 Göteborg, Sweden *Department of Computer Science, University of Skövde PO Box 408, SE-541 28, Skövde, Sweden

Abstract. MedView is a joint project with participants from oral medicine and computer science. The aim of the project is to build a large database from patient examinations and produce computerised tools to access data in various ways. One way to access data is to read case descriptions generated from stored cases. We give a description of how documents are generated from data and how these are used in a variety of contexts to supply useful information.

1. Introduction

The introduction of formalised electronic storage of medical information opens up for new ways of accessing collective knowledge among clinicians students and patients. Within the MedView project [1] formalised protocols are used during all patient visits to build a large collection of cases in the area of oral medicine. Tools have been developed for use in the examination room, to enter data electronically and view cases, but also for subsequently visualising and analysing the entire database [5].

A common way of accessing medical data is to read case descriptions. With a formalised case base, such as the one available in MedView, documents can be *generated* from data tailored for various purposes. Among the uses being developed or tested so far within MedView are medical histories viewed by the clinician in the examination room, discharge summaries, texts for personalised patient information and texts for use in an online educational system for oral medicine.

The aim of this paper is to present the methods used for generating medical texts in MedView and to show some of the applications that have been built using it. We will also discuss the MedView approach in the context of Natural Language Generation (NLG).

2. The MedView System

MedView is primarily aimed at increasing the speed by which we may obtain new and valuable information within the field of oral medicine. A formalisation of clinical procedures and visualisation of information provide a possibility for recognising new trends and patterns otherwise hidden in large amounts of non-transparent clinical records.

When elaborating the MedView system, great care was taken to determine what clinical information could be defined as useful and constitute the foundation in the database. The result from these considerations was standardised protocols for input of clinical information, developed in close collaboration between participants from oral medicine and computer science. Case history and all clinical data are entered by use of predefined parameters from the mentioned protocols. Through this process a solid base for subsequent analysis and intelligible reasoning of results is obtained. The formalised protocols have a logic interpretation [6], which make them suitable for automated reasoning in a computerised system. At the same time, they are simple enough to have an obvious intuitive reading needing no further explanation.



Figure 1: General description of MedView.

Today, the system has been in use for a couple of years at several clinics. The database contains more than 2000 examinations and some 2500 images, which in the area of oral medicine is a significant contribution. The examination rooms are equipped with a PC on a custom-built table and a digital video camera used for taking images of the oral mucosa. The collected data are stored on a server. Seminars are held regularly within the network of users, discussing cases from the database. An overview of the system is shown in Fig. 1.

3. Text Generation in MedView

Natural Language Generation is the activity of generating text from some kind of sources. In principle, there are two approaches to generation, the *deep* and *shallow* approach. A deep system builds on a deep understanding of linguistics whereas shallow systems use simpler methods to generate text. The advantage of deep text-generators is that they are more domain independent and thus can be applied to various areas with relative ease. Shallow systems are typically specialised for a particular task and need not be more complicated than the task demands.

Typically, a NLG system is divided into three phases [4], (i) Content Determination: what the text should contain, (ii) Sentence Planning: planning the structure at sentence level, (iii) Surface Realisation: realising the desired structure into text. Other approaches are used as well. An overview of applications of NLG in health-care can be found in [3].

3.1. Generating Case Descriptions

462

The main focus during the development of the text-generation system used in MedView has been to create a very flexible system where users can experiment with different texts without having any linguistic expertise. Therefore, a shallow approach to NLG was chosen. Close to a simple mail-merge system, it can be classified as a slot-and-filler, or canned-text with knowledge base references system [2].

The text-generator takes as input a *template* describing the texts to generate. This template consists of a number of files providing (i) an RTF, HTML or LaTeX template text, (ii) a file that classifies the attributes of the database into a number of categories, (iii) a file that classifies the attributes of the database into a number of categories, (iii) a file that defines the text-fragments to use as slot-fillers for attribute values. The template text contains a number of *sections*, where each section is made up of a number of *fragments* with slots to be to be filled in depending on the values for attributes in an examination record. In the template text the different fragments are separated either by a full-stop (at sentence level) or by some special character. The contents of a template is administered by the end-users themselves, thereby allowing for simple customisation. A template does not contain any procedural information but can be read in a purely declarative manner. The procedural behaviour is defined by a few simple rules. The format of the resulting document is the same as the format used in the template text.

To generate a case description for a particular task the user first selects the appropriate template. The system then performs the following steps:

- Content planning. Depending on the user's choice it is decided what template to use and which sections of the template should be included in the text.
- Sentence planning. Depending on which attributes have values, it is decided which fragments should be included. Fragments for which values are missing are omitted.
- Surface Realisation. Depending on the values for attributes in the database particular text-fragments are selected and used to fill slots in selected fragments

Example. An examination record in MedView forms a tree structure with top level nodes representing the different main tasks from which information is gathered. However, for the purpose of this paper it is sufficient to view it as a set of equations where a number of attributes are defined by certain values. For instance a very small part of a record could be:

Occup = Lärare Ref-in = Tandläkare Ref-cause = Slemhinneförändring. Civ-stat = Gift. Born = Sverige. Checkup = Ja.

where the actual values are given in Swedish. A section of a template (where full-stop is used to separate fragments) used to create a summary in English for this patient could be

§DISEASE HISTORY§

<age> year old <sex> <occup> who is referred by <ref-in> because of a <ref-cause>. The patient is <civ-stat> and comes originally from <born>. <dis-now>. <checkup>.

Now, if value-text maps include the following:

| Occup: | Ref-in: |
|---------------------------------|---|
| Lärare = teacher. | Tandläkare = "a general dentist". |
| Civ-stat: | Ref-cause: |
| Gift $=$ married. | Slemhinneförändring = "mucosal lesion". |
| Born: | Checkup: |
| Sverige = Sweden. | Ja = "Attends medical check-ups regularly". |
| Dis-now. | |
| Nej = | |
| itself = "Current diseases: ?". | |
| | |

the summary becomes:

DISEASE HISTORY

58 year old female teacher, who is referred by a general dentist because of a mucosal lesion. The patient is married and comes originally from Sweden. Attends medical check-ups regularly.

The slot <dis-now> is omitted since it has no value. If the equations, Disnow=Hypertension, and Dis-now=Glaucoma, had been present the sentence "Current diseases: Hypertension and Glaucoma" would have been added to the text. The special marker --- means that for this value the system should behave as no value at all was defined. It is also possible to use intervals when defining substitutions for numerical values.

The result of generation is determined over time through experimentation with properties like the size of text fragments used as replacements for attributes, what values to include, the order and so on. Our overall experience is that due to the rather static structure of case descriptions, elaborating with different templates and canned-text chunks of different sizes works well enough in most cases.

4. Application Areas

MedView makes a clear distinction between the activities of entering information and viewing collected data. To gain access to data a *viewer* application is used. A number of these that use text generation to display clinical data are briefly mentioned here.

MedSummary. The view of the database presented by MedSummary is that of a summary of one or more examination records together with any associated images. MedSummary is used in the examination room to access the available clinical data for the patient being examined. The text can then printed and used for things like providing a detailed medical history if the patient is sent to another clinician, eliminating the need for dictation. Through the years the application has been used to generate literally thousands of documents.

MedViewer. MedViewer is an interactive search and analysis tool used for exploring the database. It makes use of several templates for providing information at different levels. Typically, when an interesting case is found the user first demands a quick overview text and then, if more information is needed, a full description together with associated images as illustrated in Fig. 2.



Figure 2: Text Generation in MedViewer. By clicking a point in the graph a full description is shown.

Web Access. Since the real treasure of MedView is the database being built we have developed some basic web applications to let clinicians world-wide access data using their native language. Essentially these create the same kind of documents as MedSummary but in HTML format instead. A demo version in Swedish and English is accessible [7].

PPI. The MedView Personalised Patient Information System lets patients log in using the www to get information about their medical records, diseases and upcoming medical events. To present information about various diseases rather large chunks of canned text can be used. Due to lack of resources we have not yet performed any real testing of the system.

mEduWeb. Currently under development, mEduWeb is intended to make use of the MedView database to provide a large number of cases and exercises for students of oral medicine in a web based setting. When creating a sample case the course administrator uses text generation to create a case description from data. If needed the text can be edited before it is made available to the students. The students can also search and view the entire database having cases presented through generated documents.

5. Discussion

Of course, the basis for the possibility to easily present case descriptions using natural language generation is the formalisation of clinical data. The presence of a well-organised data source facilitates the development of templates for use in a template based text generation system. As stated earlier, ease of use by non-experts has been deemed more important than producing optimal text quality or using linguistically motivated methods. Several choices made in the development of text generation in MedView are essentially orthogonal to the approaches suggested for NLG, as discussed below.

The corpus-based approach [4] advocates that the first step in the construction of a NLG system is to build a corpus of example texts. This corpus is then analysed for linguistic and information content. Our approach has instead been to build a system where the users, through experiments with given tools, can decide the texts themselves.

It is often argued that a major advantage of sophisticated NLG over template systems is that deeper systems are more flexible and easier to maintain. In the development of MedView we have selected to use a simple template approach to achieve flexibility. Of course, this is related to the fact that it is necessary that the *end-users* themselves can modify what the generated texts should be. The templates used are simple enough to be modified by end-users. To expect that they would be able to easily control the workings of a sophisticated NLG system is not realistic.

Whether the higher quality produced by systems building on linguistic knowledge is needed depends on the application domain. The structure of medical record text is typically very static and uses a rather formal language. Furthermore, there is no need to produce text with great variation. On the contrary, too much variation might be disturbing since clinicians reading the texts expect them to follow certain patterns. This indicates that for the MedView application domain a template-based approach is sufficient. We have noticed that for text intended for use in the examination room there is a clear trend among users towards less sophisticated language in favour of a more telegrammed style.

While a deep NLG system does not appear to be needed, using a *hybrid system* would be quite useful. A hybrid system combines templates with deeper NLG techniques. One obvious technique being a candidate for inclusion in the future version is aggregation used to combine related phrases and sentences together in a linguistically correct manner.

Finally, we note that *multi-modality* is of increasing importance in document generation. We need to be able to include diagrams, tables, and other graphics into patient summaries. Creating fully multi-modal documents is an interesting challenge that will further increase the value of generated documents.

References

- [1] Y. Ali, G. Falkman, L. Hallnäs, M. Jontell, N. Nazari, and O. Torgersson. Medview: Design and adoption of an interactive system for oral medicine. In A. Hasman, et al, *Medical Infobahn for Europe: Proceedings of MIE2000 and GMDS2000*. IOS Press, 2000.
- [2] S. Busemann and H. Horacek. A flexible shallow approach to text generation. In E. Hovy, editor, Proceedings of the Ninth International Natural Language Generation Workshop (INLG'98), 1998.
- [3] Cawsey. B. Webber and R. Jones. Natural Language Generation in Health Care, Journal of the Medical Informatics Society, 4, 1997.
- [4] R. Dale and E. Reiter. Building Applied Natural Language Generation Systems, Journal of Natural Language Engineering, 3:55--87, 1997.
- [5] G. Falkman. Information visualization in clinical odontology: multidimensional analysis and interactive data exploration. *Artificial Intelligence in Medicine*, 22(2):133--158, 2001.
- [6] L. Hallnäs, Partial Inductive Definitions. TCS 87(1):115-142, 1991.
- [7] MedView WebDemos, http://www.cs.chalmers.se/~oloft/MedView/webdemos.htm.