

MEDICAL WORD RECOGNITION USING A COMPUTATIONAL SEMANTIC LEXICON

Robert Milewski

Center of Excellence for Document Analysis and Recognition
University at Buffalo, The State University of New York
UB Commons Suite 202
Amherst, NY 14228
milewski@cedar.buffalo.edu

Venu Govindaraju

govind@cedar.buffalo.edu

Abstract

Artificial Intelligence(AI) plays the following two crucial roles in medical form analysis: recognition, as an input, of the New York State (NYS) Prehospital Care Report(PCR), and data inferences as an output. The PCR provides medical, legal, and quality assurance (QA) data (approximately 2-3 years behind in storage and analysis) that needs to be efficiently centralized to aid health care. Automating NYS PCR analysis will facilitate a more efficient and useful description of a patient being admitted to a hospital emergency room (ER). ER environments can be highly stressful on the human body given the time constraints of bio-terrorism, trauma and/or disease. The recognition task will allow these ER health care professionals to evaluate all data and emergency techniques performed by paramedics and emergency medical technicians (EMT's). A computer screen, presenting diagrams, descriptions and inferences of a human body, representing the patient, will be updated with the corresponding handwritten PCR information. This information can then be transported to a central data bank where other hospitals can determine if there are possible outbreaks due to bio-terrorism, disease, hazardous materials incident or other non-obvious mass casualty incidents (MCI). Currently, it may take several days or even weeks, when it is clearly too late, to discover a massive atrocity. The recognition process will involve a method for reducing the size of the lexicon by integrating semantic knowledge with pattern recognition data.

1 Overview

Ambulance services are called to the scene of an accident where patients are rescued, evaluated, monitored, and transported to emergency room hospitals. Medical forms,

describing the full situation of the patient, contain machine and human printed fields. These New York State(NYS) Prehospital Care Reports(PCR) (Figure 2), used as both a medical and legal document, are used by emergency room physicians to evaluate the initial circumstances for the patients condition. Then the forms will be forwarded to the Western Regional Emergency Medical Services (WREMS) for quality assurance analysis. Unfortunately, these forms are approximately 3 years behind in being manually entered and evaluated into computers. The New York State Department of Health (DOH) is investigating several ways for automating these procedures: the use of expensive electronic boards and the automated analysis of existing forms. Ideally, the sharing of this patient medical data between emergency rooms (Figure 1) would be available; unfortunately, this technology is not being used in many states.

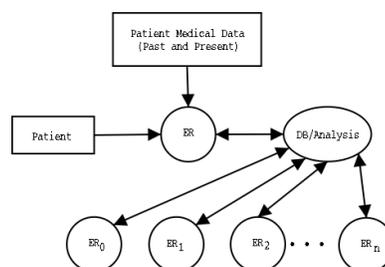


Figure 1. Input/Output Overview

The most challenging and crucial part of this development effort is the automated recognition of various handwriting in emergency conditions. There can be a mix of characters, digits, symbols, standard and non-standard abbreviations mixed with cross outs, cursive and print mixed, letters carrying over between lines, words being crushed to

fit on the form, messy handwriting and misspelled words. Compounded with over 50,000 possibilities of medical words and tens of thousands of arrangements of text input, the accurate shrinking of the lexicon, during character recognition, is essential.



Figure 2. An example NYS PCR

Both generic and customized systems are in competition with one another; it is generally believed that the most specific a priori knowledge will be required to accomplish this task; hence, a customized system. Intuitively, the need for semantic relationships between words, an understanding of the patient's situation, prior probabilistic analysis for similar situations, etc... will be needed to constrain the huge lexicon when executing form recognition algorithms. Generic algorithms involve the analysis of words independently with little or no semantic analysis between the words. The expectation of data in several specific places on different medical forms severely hinders the use of the generic algorithms.

2 Objective

The ultimate goal of this paper is to construct a hybrid semantic network and computational mind capable of taking related words from an NLP oracle machine and producing, with higher confidence than that of existing CEDAR recognition algorithms (e.g. WMR) [5] [11] [1] [7], the probabilistic semantic match of that word, to a known word, with known meaning in the field of emergency medicine.

2.1 Strategy

Figure 3 illustrates the lexicon reduction flow by using the data on the PCR form. Starting at the top left of

the figure, the Ambulance and/or Hospital personnel place the PCR form in a scanner. The computer reads in the form and segments the PCR into several blocks corresponding directly to the form titles (e.g. Presenting Problem...Dispatch, PMH, and Patient Identification). There needs to be a method for minimizing the high quantity of possible words, used by a medic, to describe a patient's condition. A lexicon database will contain a list of English and medical words which are weighted according to the popularity of that word over time (i.e. as more and more PCR's get analyzed, the popular words become weighted to assist in probable recognition results). The subjective and patient complaint area contain the smallest amount of text on the PCR. The machine printed checkboxes, in conjunction with the lexicon database, will be used to determine general analysis path(s) for the patient status. This a priori data will be used for further recognition in the larger handwriting regions. The objective and comments region contain a lot of varying abbreviations, symbols and numbers in conjunction with regular handwriting. Therefore, we can use the general paths to narrow in on specific problems. This data is then sent to a data compiler, which collects similar cases for further minimization. The data compiler will also mine the patients PMH, the PMH of other family members, and other similar patients (from older PCR's). At this stage we have minimized the patient problems to a finite set of possibilities. In order to determine which of these hypothesis truly classify the patients condition the best, we chose the Naive Bayesian Classifier. The classifier will take as an input, the reduced data and previous inferences from a database and produce scores for possible medical conditions. The patient analysis will be stored in the inference database for future inference matching. The data is then sent back to the data compilation module which checks to see if the data routes can be further minimized. If there exist greater or equal to one hypothesis that cannot be further divided, the system will conclude multiple problems. The hypothesis set will be converted to data that a human can comprehend. Specifically, a graphical user interface (GUI) visualizer will color code the hypothesis based on confidence values produced by the system.

We have built a primitive semantic network as the infrastructure to a more computationally expensive lexicon network, called the Java Constrained Object Inference Net (JCOIN). The following basic steps will be required for the analysis infrastructure:

- 1) The meanings of words, abbreviations, symbols, etc... need to be stored
- 2) The frequency with which these words occur on these specific forms in particular patterns
- 3) Feature vector analysis for the words which train this hybrid network

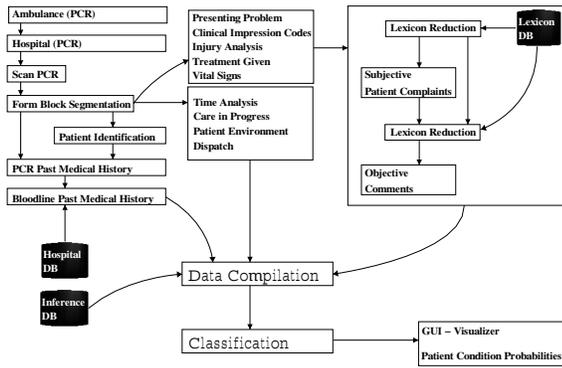


Figure 3. Recognition Flow

- 4) The general meanings as a generic (i.e. super category) to the words in a network
- 5) The probability that network words semantically relate to each other

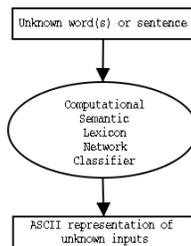


Figure 4. Network classifier input/output

Observe Figure 4; given the additions of this data to the network, it will then be possible to build a search algorithm which takes a given sentence or sequence of unknown words, in digital form, and outputs the ASCII meaning in the context of the medical document. The context will not be constrained by ordinary semantic analysis, but also include a machine learning strategy for adapting the current word meanings, features, and probabilistic analysis into the network. This will make the hybrid network the lexicon, and will illustrate the need for dynamic intelligent lexicons.

2.2 Semantic Lexicon Example

This system expects certain data elements at certain locations on the PCR; for example, the patient's chief complaint, subjective, and objective hand-written areas on the form, contain detailed medical and environmental analysis. These areas can contain words from various dictionaries: English, Medical, and Pharmaceutical are the

top 3; there are also symbols and numbers. Fortunately, the PCR contains checkboxes and boxes containing digits; which have a high recognition probability. This example will show one possible combination of PCR snippets to assist with the handwritten recognition pieces. Please note Figure 2 as the example PCR.

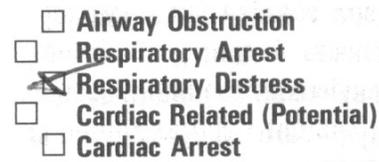


Figure 5. PCR Presenting Problem Snippet

The presenting problems section (Figure 5) is relatively straight forward; this snippet illustrates 5 checkboxes, however, the section contains about 35 checkboxes (not shown), 2 of which may contain further handwritten information. In this case, the “respiratory distress” checkbox is checked, indicating the system requirement to verify the patients problem (i.e. the system needs to find PCR data which places the patients condition in the respiratory category).

TIME	RESP	PULSE	B.P.
04:00	Rate: <input type="checkbox"/> Regular <input type="checkbox"/> Shallow <input checked="" type="checkbox"/> Labored	Rate: 102 <input type="checkbox"/> Regular <input type="checkbox"/> Irregular	102 68
04:10	Rate: <input checked="" type="checkbox"/> Regular <input type="checkbox"/> Shallow <input type="checkbox"/> Labored	Rate: 104 <input type="checkbox"/> Regular <input type="checkbox"/> Irregular	104 68

Figure 6. PCR Vital Signs Snippet

Figure 6 shows two vital sign measurements within 10 minutes: the first shows labored breathing and the second shows a regular breathing rate. This is probably impart to the 12 LPM of Oxygen given to the patient via a non-rebreather mask (these snippets are not shown). This increases the confidence of a respiratory issue.

Figure 7. PCR Past Medical History Snippet

The past medical history(PMH) section, illustrated in Figure 7, can be quite helpful in drawing correlations between a past condition and the current condition. This patient has a history of hypertension, diabetes and asthma; indicating a possible recurring condition of any or all of the three critical illnesses. This observation won't neglect other possibilities; certainly a patient with asthma can still break their leg. However, an individual with severe trauma and a respiratory PMH, has a greater chance of becoming unstable. This semantic reasoning, can assist a combinatorial based lexicon, in which two or more lexicons can be used to interpret the situation and assist with word selection constraints. Using the information from the PMH, we can use a pharmaceutical lexicon to reduce the selections of possible medications. After determining that albuterol is the medication, and with the prior knowledge that this is an asthma medication, we can note that asthma, with this patient, is the most serious.

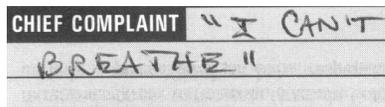


Figure 8. PCR Chief Complaint Snippet

After using the easier recognition portions of the forms as inputs to our system, it will then be possible to adjust the lexicon size as well as the scores for individual words. This will be accomplished by a hybrid engine using both machine learning (e.g. Bayesian laws) and knowledge representation (KR) techniques. The first handwriting snippet, is the chief complaint (Figure 8) and gives us the patients statement as to the chief problem. Using our prior analysis, we can expect words and phrases along the lines of: shortness of breath, difficulty breathing, restricted breathing, etc... Combining both the results of image processing and our hybrid semantic network, we can now begin to produce a more accurate ASCII representation of this text; in this case, the patient complains of not being able to breath.

Further breaking down other portions of handwriting text, we may be searching for words such as lungs, trachea, airway, mouth, breathing, etc... and given the knowledge that this patient has asthma, another possibility is the word inhaler. In Figure 10, the word inhalers is plural; this illustrates the importance in integrating word prefixes and suffixes into lexicon selection. In Figures 11 and 12, we see our first example of combining a medical symbol along with the respiratory related information. The word trachea in the phrase "negative tracheal shift", would have a higher score since we know that a medic must perform a physical exam,

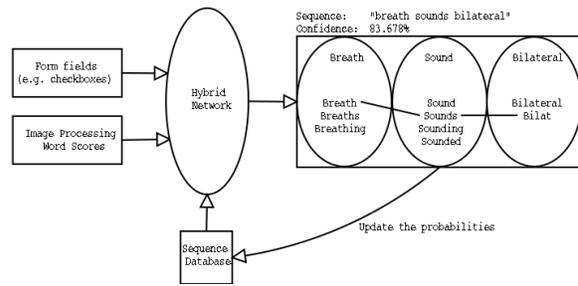


Figure 9. Hybrid network snippet analysis

and given the respiratory data, the trachea will be one of the words to be scanned. Similarly with the "decrease of bilateral breath sounds", all words are standard and frequently used by emergency personnel.

The hybrid network will make decisions based on (see

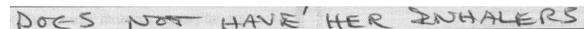


Figure 10. PCR Subjective Snippet: Does not have her inhalers

Figure 9):

- 1) The probability the sequence occurs relative to other PCR inputs.
- 2) A score measuring the relevance of the sequence based on medical knowledge.

After getting these two scores, the system will make a final probabilistic decision that the given image snippet has the proposed ASCII translation.

Clearly, there are many possible combinations of words and phrases, between multiple dictionaries, on medical forms. Constructing a general recognition system will not have as high the performance as with the addition of semantic analysis. If we are expecting words to show up, given a scenario, it makes sense to use those words in lexicon analysis.



Figure 11. PCR Objective Snippet: Negative Tracheal Shift

2.3 Prior and Competing Work

CEDAR is famous for the design, analysis, and implementation of document analysis algorithms; particularly in the areas of printed forms, word/character/digit recognizers and lexicon pruning [1] [10] [7]. Page segmentation al-

↓ BREATH SOUNDS' BILAT

Figure 12. PCR Objective Snippet: Decreased Breath Sounds Bilateral

gorithms have several techniques which can be applied to form recognition and the breakdown of sub-blocks within forms [8]. Machine learning techniques, such as Bayesian probabilities [9] have been demonstrated in several document analysis applications. Analysis on degraded images, that need recovery, play a role in the feature vector analysis of such forms as well [3] [4]. These PCR forms, which contain the standard characters and numbers, also contain several medical symbols and word combinations, making this analysis similar to that of foreign character recognition, such as Chinese and Japanese [6].

Products, such as (Notfall-Organisations- und Arbeitshilfe) NOAH, SafetyPAD, TCPR (Transportable computer based patient record), PenComputerSolutions Inc., and SOAPware all use high tech electronic devices, such as notepad computers, for data transfer between an ambulance and a hospital [14] [16] [18] [15] [17]. Each product shares its own unique characteristics: SafetyPADmobile runs on a "Pentium-class pen-based computer" allowing medics to remain highly versatile in emergency settings. SafetyPADmobile can also operate on a notebook computer, allowing the medic to electronically register and transfer patient data to the hospital through a wireless data link. Similarly, NOAH which was recently implemented in Germany, uses a wireless data transfer through a data communication network called "Modacom". Analogous to NOAH, PenComputerSolutions Inc. has hand held devices to transfer patient data from the emergency setting to the hospital, through a wireless channel. TCPR uses an extensive object oriented design to assess and prioritize the patient's condition. SOAPware uses an existing database to electronically store patient information and past history. Even though there are differences in the products, we do see the recurrence of using a wireless means of data transfer, and also the use of high tech equipment such as computers and hand held devices.

These methods possessing their unique characteristics theoretically satisfy the needs of transferring patient data prior to patient arrival, but they are not apposite to ambulatory settings, where medics are in a spate to scribe data concerning the patient's condition. In the same scope, electronic boards which are extremely expensive, as illustrated by the products sold by PenComputerSolutions, cannot be accommodated into all ambulatory organizations, where financial resources are limited, and the need for a cost-effective approach is being pushed. Costs could approximate to millions when boards, computers, other elec-

tronic devices and teachers are purchased to supply medical organizations. Other concerns also connected with electronic data transfer include mal-functioning of the product, or damage to the product in case of an accident. There is little assurance of constructing a damage proof product with advanced technology. Research involving optimized patient analysis using existing PCR forms does not compromise financial resources or dependability.

The Canadian Institute of Health Sciences is funding a new Beta Test Pilot, to restructure the PCR form, and also to improve efficiency of data transfer between the ambulatory setting and the hospital setting. Research is being conducted at the University of Rochester, where the focus is centralized on hospitals receiving patient data prior to patient arrival. The office of prehospital care at the university is constructing a detailed analysis of patient care in the pre-hospital setting by using the PCR [13].

2.4 Future Work

The first version of the system will be specific to the NYS PCR forms. To allow for new PCR forms, the image processing needs to be adapted for different organizations. Therefore, a mechanism for allowing new PCR forms needs to be incorporated into the recognition phase. The data compilation step would only need to be modified if a new form contained unique data (i.e. the form contained data that did not exist in any other PCR the system understands). Beyond the handwritten recognition component, other systems could have the ability to have voice recognition in the near future. Electronic boards are also part of the future of emergency medicine. The ability to transfer data through wireless means, and to use high tech computers to record patient data all show promise. Incorporating other technologies, such as the "Virtual Emergency Room" [2], could further assist ER training. Quality Assurance could also be optimized by flagging PCR's that might not have followed protocol; there exist many flowcharts in BLS and ALS that act as guidelines and law [12].

3 Contributions

There are several core objectives that benefit both health care and computer science research:

- 1) More efficient and accurate hospital patient care
- 2) Emergency room visual aids to assist human reasoning
- 3) WREMS quality assurance
- 4) Department of Health Counter-BioTerrorism Analysis
- 5) DOH efficiency of data entry with decreased cost
- 6) Computer science algorithm research
- 7) Optional digital device migration for health care

3.1 Medical Contributions

There are many factors that influence human decisions: medical ethics, stress, insurance, natural human bias to name a few. This system is intended to assist, but not replace, medical personnel with decisions. There exist many possibilities; signs and symptoms can be very similar. Given an infinite amount of future patients and an infinite realm of possibilities, it is reasonable to expect human error. Unfortunately, in the medical profession, human error may lead to a misdiagnosis, for example. This system is intended to minimize these errors for the patient and assist with human analysis to ease the load and minimize suffering.

3.2 Computer Science Contributions

Many existing pattern recognition systems do not integrate different computer models into one. This hybrid system combines the areas of image processing, machine learning, knowledge representation and data mining to enhance machine performance in the medical arena. The flexibility of this particular model allows other modules to be designed and imported.

References

- [1] D. Niyogi, S. Srihari, V. Govindaraju Handbook of Character Recognition and Document Image Analysis, "Analysis of Printed Forms." *World Scientific Publishing; Singapore. C1997. p. 485-502.*
- [2] Daniel Sorid, et al. *IEEE Spectrum*, "The Virtual Surgeon." <http://www.spectrum.ieee.org>; July 2000 Volume 37 Number 7.
- [3] Kanungo, Tapas, et al. "A Statistical, Nonparametric Methodology for Document Degradation Model Validation". *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 22, Number 11, November 2000; p1209-1223.*
- [4] Kato, Yoshiharu and Yasuhara, Makoto. "Recovery of Drawing Order from Single-Stroke Handwriting Images". *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 22, Number 9, September 2000; p938-949.*
- [5] Kim, Gyeonghwan, Govindaraju, Venu. "A Lexicon Driven Approach to Handwritten Word Recognition for Real-Time Applications". *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 19, Number 4, April 1997; p366-379.*
- [6] Liu, Cheng-Lin and Nakagawa, Masaki. "Precise Candidate Selection for Large Character Set Recognition by Confidence Evaluation." *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 22, Number 6, June 2000; p636-642.*
- [7] Madhvanath, S., Krpasundar, V. and Govindaraju, V. "Syntactic methodology of pruning large lexicons in cursive script recognition." *Published by Elsevier Science C2001. The Journal of the Pattern Recognition Society; Pattern Recognition 34(2001); p37-46.*
- [8] Mao, Song and Kanungo, Tapas. "Empirical Performance Evaluation Methodology and Its Application to Page Segmentation Algorithms." *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 23, Number 3, March 2001; p242-256.*
- [9] Mitchell, Tom M. Machine Learning. *McGraw Hill Higher Education. C1997.*
- [10] Park, Jaehwa, Govindaraju, Venu, Srihari, Sargur N. "OCR in a Hierarchical Feature Space." *IEEE Transactions on Pattern Analysis and Machine Intelligence; Volume 22, Number 4, April 2000; p400-407.*
- [11] Park, Jaehwa, Govindaraju, Venu. "Use of Adaptive segmentation in handwritten phrase recognition." *Published by Elsevier Science C2001. The Journal of the Pattern Recognition Society; Pattern Recognition 35(2002); p245-252.*
- [12] Young, Gregory, Hens, Virginia Adult & Pediatric Advanced Life Support Protocols Second Edition; *Published by Western Regional Emergency Medical Services (WREMS) C1995.*
- [13] "Beta testing, Canadian Institute of Health Sciences." C2000 <http://www.jppc.org/ambcare/beta/index.htm>
- [14] "NOAH (Notfall-Organisations- und Arbeitshilfe)". C2000 <http://www.ncbi.nlm.nih.gov>
- [15] "PenComputerSolutionsInc.". C2000 <http://www.pencomputersolutions.com/>
- [16] "SafetyPAD". C2000 <http://www.safetypad.com/>
- [17] "SOAPware". C2000 <http://www.docs.com/>
- [18] "TCPR (Transportable Computer-based Patient Record)". C2000 <http://akamai.tamc.amedd.army.mil/PacMedTek/Presentations/Thursday/vatinelle/tsld006.htm>