

# Improving the Secondary Utilization of Clinical Data by Incorporating Context

Leonard D'Avolio

University of California, Los Angeles  
Dept. of Information Studies  
Dept. of Medical Imaging Informatics  
924 Westwood Blvd, Suite 420  
Los Angeles CA, 90024  
1(323)219-2335  
ldavolio@ucla.edu

Galya Rees

University of California, Los Angeles  
Dept. of Medical Imaging Informatics  
924 Westwood Blvd, Suite 420  
Los Angeles CA, 90024  
1(310)794-8956  
grees@mii.ucla.edu

Lousine Boyadzhyan, MD

University of California, Los Angeles  
Dept. of Radiology  
Dept. of Medical Imaging Informatics  
924 Westwood Blvd, Suite 420  
Los Angeles CA, 90024  
1(310)794-8977  
lboyadzhyan@mednet.ucla.edu

## ABSTRACT

There is great potential in the utilization of existing clinical data to assist in decision support, epidemiology, and information retrieval. As we transition from evaluating systems' abilities to accurately capture the information in the record, to the clinical application of results, we must incorporate the contextual influences that affect such efforts. A methodology is proposed to assist researchers in identifying strengths and weaknesses of clinical data for application to secondary purposes. The results of its application to three ongoing clinical research projects are discussed.

## 1. Secondary Utilization of Clinical Data

Increasingly, researchers in medical informatics are using natural language processing, information extraction, and data mining methods to utilize existing clinical data for secondary purposes. Researchers have employed such techniques in decision support efforts [1, 2], to encode medical reports [3-6], to provide structured data for data mining and clinical research [7-9], in order entry [10] and in information retrieval [11, 12].

Concern with the utility of medical information toward these efforts has primarily been concerned with formatting information in a manner amenable to computation. In information extraction efforts, a primary concern is the ability to extract medical concepts of interest with high rates of recall and precision. When information extraction efforts act as the foundation for stochastic data mining or machine learning approaches, the concern becomes how best to represent these extracted concepts as mathematical input. The issues raised by this computational perspective include the volume and heterogeneity of medical data, its poor mathematical characterization, and the lack of canonical form for medical concepts [13].

The factors introduced by this perspective center around achieving results that accurately reflect *what was contained in the record*. However, when the intent is to apply this extracted data in the context of patient care, our concern must expand beyond capturing what is contained in the record to include the utility of that information toward a specific clinical goal. This calls for a reconciling of the context surrounding the original capture of the information and the context surrounding its intended secondary use. Three general classes of contextual influences that affect the

application of clinical data for secondary purposes can be derived from existing literature.

**Clinical Context.** The influences that shape the inclusion and format of data recorded in the process of caring for patients are considered products of the clinical context. Several influences of note arise within this context. First, medical data is often considered 'soft' and of uncertain value [14], a far cry from the measurable set of 'facts' required for an objective analysis [15]. Second, medical data is mutually elaborative, and cannot be understood as isolated, atomic givens [16]. To analyze something as seemingly objective as a lab value independent of consideration of other relevant information is a flight from interpretation that is doomed to fail [17]. This is because medical data is interpreted and reinterpreted in light of surrounding data [18]. Medical data is also interpreted based on the credibility of its source [19] with omissions of data from a resident physician likely to be perceived differently than omissions from an attending [20]. Multiple uses of records proliferate the lack of uniform standards leading to a high degree of variation even in what values are considered worth capturing in records.

**Economic Context.** While there isn't proof that physicians intentionally fabricate medical data to appease insurers' cost control regulations, in one survey nearly 70% of physicians indicated that they would misrepresent a woman's condition in order to assure reimbursement for a mammography screening [21]. Coding of diagnoses and procedures (such as CPT codes) are used primarily for reimbursement, affecting their granularity and hence utility for research purposes. Further complicating their use, Hsia found that a statistically significant amount of improperly coded records led to higher reimbursement for hospitals [22]. As a result, the current medical environment has been labeled the 'misinformation era' with the economics 'influenc[ing] every record entry' [23].

**Legal Context.** One possible use of medical records is as legal evidence in malpractice suits. This has been blamed for the use of "hedged" language [24], depersonalization, use of passive voice, and treating the technology as an agent [25]. As a result, attempts to determine, for example, the potential causes of adverse outcomes must consider the fact that evidence is unlikely appear in a straightforward manner.

## 2. Methodology

The contextual categories outlined above and their enumerated considerations were used to create a methodology for discovering the utility of clinical data toward a specific clinical goal. First, a list of the data items believed to be of interest was created (e.g. lab value, tumor size, surgical step). Next, a random sample of 15 records was drawn and the format and consistency of the items of interest noted. Finally, for each data item, a Secondary Utilization Matrix (SUM) analysis was conducted. Inspired by business' SWOT analysis, the Y-axis of the grid features the categories; *strengths*, *weaknesses*, *strategy*. Across the X-axis are the enumerated considerations of the three categories of contextual influence (clinical, economical, legal). At each intersection in the matrix, potential effects of contextual influences are listed. Critical to the completion of a SUM analysis by a non-clinical researcher is a review of the results with a clinician responsible for creating the reports. The results facilitate an understanding of potential obstacles and opportunities of data from both a computational and clinical perspective. SUM analyses were conducted on three ongoing clinical research efforts; surgical prostate cancer removal, brain tumor treatment, and treatment of female pelvic prolapse. Abbreviated findings are discussed below.

## 3. Results

**Prostate Cancer Removal.** This project is designed to discover surgical processes that correlate to specific surgical outcomes using operative notes. Experience of dictating surgeons was discovered to be responsible for significant variation in data included in reports. A follow up analysis of 42 reports showed residents averaged twice as many words in their reports as attending surgeons (1594.54  $\pm$ 129.64, 95% C.I. versus 817.37  $\pm$ 26.67, 95% C.I.). Inclusion of certain semi-structured fields (e.g. estimated blood loss) varied by surgeon. The legal context implied that one may have to "read between the lines" to discover correlations to complications. Pattern recognition techniques are currently in use to identify "red flag" observations.

**Brain Tumor Treatment.** The goal of this project is to create a comprehensive representation of the progression of a brain tumor disease with the help of a naive Bayesian network. For this purpose, evidence was extracted from radiology reports. In the clinical context category, the lack of standards surrounding the description of tumor size (a key factor in the proposed model) was problematic with only 74% addressing the size of the lesion in a sample of 40. In addition, there was great diversity in the representation of the size of the tumor in radiology reports with 15% of the cases providing verbal description ('large'/'small').

**Female Pelvic Prolapse Treatment.** The selecting of appropriate surgical candidates and interventions is an issue of primary concern. In attempting to create a statistical model of ideal versus non-ideal surgical candidates, clinical context prevented key values from appearing consistently. This materialized in radiologists' notes in the form of inconsistencies in the inclusion of measurements and the gradings of prolapse included in the MRI report. A review of the standards employed in the clinical context showed physician personal preference in macro (dictation template) style contributing to included or excluded values of clinical interest.

## 4. References

- [1] Fiszman, M. and Haug, P. Using medical language processing to support real-time evaluation of guidelines. in Proceedings of the American Medical Informatics Association Symposium. 2000.
- [2] Friedman, C., Knirsch, C., et al. Automating a severity score guideline for community-acquired pneumonia employing medical language processing of discharge summaries. in Proceedings of the American Medical Informatics Association Symposium. 1999.
- [3] Gunderson, M., Haug, P., et al., development and evaluation of a computerized admission diagnoses encoding system. *Computers in Biomedical Research*, 1996. 29: p. 351-372.
- [4] Blanquet, A. and Zweigenbaum, P. A lexical method for assisted extraction and coding of ICD-10 diagnosis from free text patient discharge summaries. in Proceedings of the American Medical Informatics Association Symposium. 1999.
- [5] Lussier, Y., Shagina, L., and Friedman, C. Automating SNOMED Coding using medical language understanding: A feasibility study. in Proceedings of the American Medical Association Symposium. 2001.
- [6] Heinze, D., Morsch, M., and Holbrook, M. Mining free-text medical records. in Proceedings of the American Medical Informatics Association Annual Symposium. 2001.
- [7] Doddi, S., Marathe, A., et al., Discovery of association rules in medical data. *Medical Informatics and the Internet in Medicine*, 2001. 26(1): p. 25-33.
- [8] Hripcsak, G., Friedman, C., and Et, Unlocking Clinical Data from Narrative Reports: A Study of Natural Language Processing. *Annals of Internal Medicine*. 122(9): p. 681.
- [9] Wilcox, A. and Hripcsak, G. Medical test representations for inductive learning. in Proceedings of the American Medical Association Symposium. 2000.
- [10] Lovis, C., Chapko, M., et al., Evaluation of a command-line parser-based order entry pathway for the Department of Veterans Affairs electronic patient record. *Journal of American Medical Informatics Association*, 2001. 8(5): p. 486-498.
- [11] Hersh, W., Mailhot, M., et al., Selective automated indexing of findings and diagnoses in radiology reports. *Journal of Biomedical Informatics*, 2001. 34(4): p. 262-273.
- [12] Chu, S. and Cesnik, B., Knowledge representation and retrieval using conceptual graphs and free text document self-organization techniques. *International Journal of Medical Informatics*, 2001. 62(2-3): p. 121-133.
- [13] Cios, K. and Moore, G., Uniqueness of medical data mining. *Artificial Intelligence in Medicine*, 2002. 26(1-2): p. 1-24.
- [14] Musen, M., The strained quality of medical data. *Methods of Information in Medicine*, 1989. 28: p. 123-125.
- [15] Greenhalgh, T., Narrative based medicine: Narrative based medicine in an evidence based world. *British Medical Journal*, 1999. 318: p. 323-325.
- [16] Berg, M. and Goorman, E., The contextual nature of medical information. *International Journal of Medical Informatics*, 1999. 56(1): p. 51-60.

- [17] Leder, D., Clinical interpretation: the hermeneutics of medicine. *Theoretical Medicine*, 1990. 11: p. 9-24.
- [18] Whalen, J. Accounting for 'standard' task performance in the execution of 9-1-1 operations. in *Annual Meeting of the Americal Sociological Association*. 1993. Miami, FL.
- [19] Cicourel, A., The integration of distributed knowledge in collaborative medical diagnosis, in *Intellectual teamwork: Social and intellectual foundations of cooperative work*, J. Galegher, R. Kraut, and C. Egidio, Editors. 1990, Lawrence-Erlbaum Associates: Hillsdal, NJ. p. 221-242.
- [20] Bosk, C., *Forgive and remember: Managing medical failure*. 2nd ed. 2003, Chicago: University of Chicago Press.
- [21] Novack, D., Detering, B., et al., Physicians' attitudes toward using deception to resolve difficult ethical problems. *Journal of the American Medical Association*, 1989. 261: p. 2980-2985.
- [22] Hsia, D., Krusgat, W., et al., Accuracy of diagnostic coding for Medicare patients under the prospective-payment system. *New England Journal of Medicine*, 1988. 318: p. 352-355.
- [23] Burnum, J., The misinformation era: The fall of the medical record. *Annals of Internal Medicine*, 1989. 110: p. 482-484.
- [24] Prince, E.F., Frader, J., and Bosk, C. On hedging in physician-physician discourse. in *Linguistics and the professions: Proceedings of the second annual Delaware symposium on language studies*. 1982: Ablex.
- [25] Anspach, R., Notes on the sociology of medical discourse. *Journal of Health and Social Behavior*, 1988. 28: p. 357-375.