

DATA MINING IN HEALTHCARE: CURRENT APPLICATIONS AND ISSUES

By Ruben D. Canlas Jr., MSIT, MBA

5 August 2009

rcanlas@alumni.cmu.edu

Paper submitted to fulfill requirements for the
Master of Science in Information Technology

at the

Carnegie Mellon University

Australia

Data Mining in Healthcare: Current Applications and Issues

By Ruben D. Canlas Jr.

Abstract

The successful application of data mining in highly visible fields like e-business, marketing and retail have led to the popularity of its use in knowledge discovery in databases (KDD) in other industries and sectors. Among these sectors that are just discovering data mining are the fields of medicine and public health.

This research paper provides a survey of current techniques of KDD, using data mining tools for healthcare and public health. It also discusses critical issues and challenges associated with data mining and healthcare in general.

The research found a growing number of data mining applications, including analysis of health care centers for better health policy-making, detection of disease outbreaks and preventable hospital deaths, and detection of fraudulent insurance claims.

1. INTRODUCTION and RATIONALE

The successful application of data mining in highly visible fields like e-business, marketing and retail have led to its application in KDD in other industries and sectors. Among these sectors just discovering it healthcare.

This research paper intends to provide a survey of current techniques of knowledge discovery in databases using data mining techniques that are in use today in medical research and public health. We also discuss some critical issues and challenges associated with the application of data mining in the profession of health and the medical practise in general.

1.2 Objectives

The objectives of this paper are the following:

1. To enumerate current uses and highlight the importance of data mining in medicine and public health,
2. To find data mining techniques used in other fields that may also be applied in the health sector.
3. To identify issues and challenges in data mining as applied to the medical practise.

4. To outline some recommendations for discovering knowledge in electronic databases through data mining.

2. METHODOLOGY

Due to resource constraints and the nature of the paper itself, the main methodology used for this paper was through the survey of journals and publications in the fields of medicine, computer science and engineering. The research focused on more recent publications, with 2000 as the cut off year.

3. RESEARCH FINDINGS

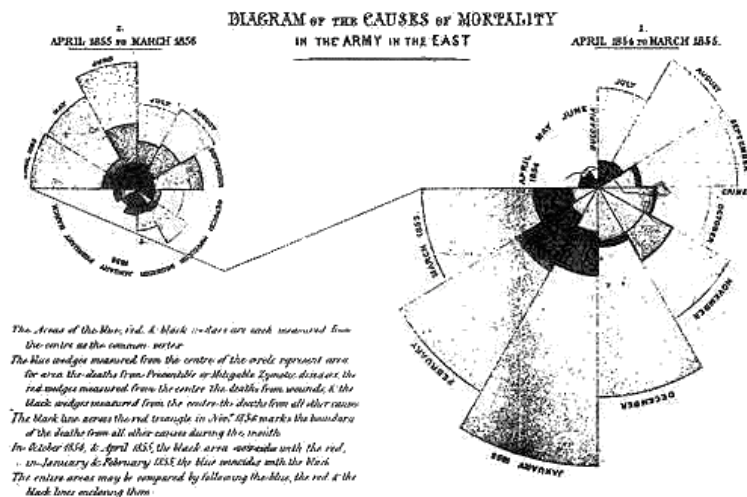
3.1 Data Mining in the Health Sector

The practise of using concrete data and evidence to support medical decisions (also known as evidence-based medicine or EBM) has existed for centuries. John Snow, considered to be the father of modern epidemiology, used maps with early forms of bar graphs in 1854 to discover the source of cholera and prove that it was transmitted through the water supply, below (Tufte 1997).



Snow counted the number of deaths and plotted the victim's addresses on the map as black bars. He discovered that most of the deaths clustered towards a specific water pump in London (center of the red circle in the map).

Florence Nightingale invented polar-area diagrams in 1855 (below) to show that many army deaths could be traced to unsanitary clinical practises and were therefore preventable. She used the diagrams to convince policy-makers to implement reforms that eventually reduced the number of deaths (Audain 2007). (Diagram from Nightingale 1858.)



Snow and Nightingale were able to personally collect, sift through and analyze the mortality data during their times because the volume of information was manageable. Today, the size of the population, the amount of electronic data gathered, along with globalization and the speed of disease outbreaks make it almost impossible to accomplish what the pioneers did.

This is where data mining becomes useful to healthcare. It has been slowly but increasingly applied to tackle various problems of knowledge discovery in the health sector.

Data mining and its application to medicine and public health is a relatively young field of study. In 2003, Wilson et al began to scan cases where KDD and data mining techniques were applied in health databases. They found confusion in the field regarding what constituted data mining. “Some authors refer to data mining as the process of acquiring information, whereas others refer to data mining as utilization of statistical techniques within the knowledge discovery process.” (Wilson et al. 2003)

Because of misconceptions still going on in the medical community about what data mining comprises, let us first define what we mean by it. The generally accepted definition of data mining today is the set of procedures and techniques for discovering and describing patterns and trends in data (Witten and Frank 2005). We shall use this definition throughout the paper.

3.2 The Importance and Uses of Data Mining in Medicine and Public Health

Despite the differences and clashes in approaches, the health sector has more need for data mining today. There are several arguments that could be advanced to support the use of data mining in the health sector, covering not just concerns of public health but also the private health sector (which, in fact, as can be shown later, are also stakeholders in public health).

Data overload. There is a wealth of knowledge to be gained from computerized health records. Yet the overwhelming bulk of data stored in these databases makes it extremely difficult, if not impossible, for humans to sift through it and discover knowledge (Cheng, et al 2006).

In fact, some experts believe that medical breakthroughs have slowed down, attributing this to the prohibitive scale and complexity of present-day medical information. Computers and data mining are best-suited for this purpose. (Shillabeer and Roddick 2007).

Evidence-based medicine and prevention of hospital errors. When medical institutions apply data mining on their existing data, they can discover new, useful and potentially life-saving knowledge that otherwise would have remained inert in their databases. For instance, an ongoing study on hospitals and safety found that about 87% of hospital deaths in the United States could have been prevented, had hospital staff (including doctors) been more careful in avoiding errors (HealthGrades Hospitals Study 2007). By mining hospital records, such safety issues could be flagged and addressed by hospital management and government regulators.

Policy-making in public health. Lavrac et al. (2007) combined GIS and data mining using among others, Weka with J48 (free, open source, Java-based data mining tools), to analyze similarities between community health centers in Slovenia. Using data mining, they were able to discover patterns among health centers that led to policy recommendations to their Institute of Public Health. They concluded that “data mining and decision support methods, including novel visualization methods, can lead to better performance in decision-making.”

The preceding factors remind us of an incident in the Philippines at the Rizal Medical Center in Pasig City in October 2006. Failing to implement strict sanitation and sterilization measures the hospital contributed to the death of several new-born babies due to neonatal sepsis (bacterial infection). No one really knew what was going on until the deaths became more frequent. Upon examining hospital records, the Department of Health (DOH) found that 12 out of 28 babies born on October 4, for example, died of sepsis (Tandoc 2006). With an integrated database and the application of data mining the DOH could detect such unusual events and curtail them before they worsen.

More value for money and cost savings. Data mining allows organizations and institutions to get more out of existing data at minimal extra cost. KDD and data mining have been applied to discover fraud in credit cards and insurance claims (Kou et al. 2004). By extension, these techniques could also be used to detect anomalous patterns in health insurance claims, particularly those operated by PhilHealth, the national healthcare insurance system for the Philippines.

Early detection and/or prevention of diseases. Cheng, et al cited the use of classification algorithms to help in the early detection of heart disease, a major public health concern all over the world. Cao et al (2008) described the use of data mining as a tool to aid in monitoring trends in the clinical trials of cancer vaccines. By using data mining and visualization, medical experts could find patterns and anomalies better than just looking at a set of tabulated data.

Early detection and management of pandemic diseases and public health policy formulation. Health experts have also begun to look at how to apply data mining for early detection and management of pandemics. Kellogg et al. (2006) outlined techniques combining spatial modeling, simulation and spatial data mining to find interesting characteristics of disease outbreak. The analysis that resulted from data mining in the simulated environment could then be used towards more informed policy-making to detect and manage disease outbreaks.

DOH orders probe after Rizal hospital tragedy

Sanitation regulations stressed

By Edson C. Tandoc Jr.

HEALTH SECRETARY FRANCISCO DUQUE III has stressed the need for government hospitals to observe strictly sanitation and sterilization procedures.

Duque issued the reminder following the death of seven babies born at the Rizal Medical Center in Pasig City apparently of infections.

The health secretary created a five-person investigating committee to determine how the babies got infected. He added that the incident should "serve as a lesson to other government hospitals."

Parents of the children who died of neonatal sepsis or bacterial infection, days after their birth in the hospital on Oct. 4, have blamed the medical center for the deaths.

But health officials said it was also possible some of the babies, who had already been taken home, could have gotten the infection either outside the hospital, or from their mothers who could have passed on the infections.

The team of four government health officials and a private doctor has five days to complete the investigation. "We want to assure the public the investigation will be transparent and independent," Duque told the INQUIRER.

Wong et al. (2005) introduced WSARE, an algorithm to detect outbreaks in their early stages. WSARE, which is short for “What’s Strange About Recent Events” is based on association rules and Bayesian networks. Applying WSARE on simulation models have been claimed to result to relatively accurate predictions of simulated disease outbreaks. Of course, these sorts of claims always come with warnings to take precaution when applying these models in real life.

Non-invasive diagnosis and decision support. Some diagnostic and laboratory procedures are invasive, costly and painful to patients. An example of this is conducting a biopsy in women to detect cervical cancer. Thangavel et al (2006) used the K-means clustering algorithm to analyze cervical cancer patients and found that clustering found better predictive results than existing medical opinion. They found a set of interesting attributes that could be used by doctors as additional support on whether or not to recommend a biopsy for a patient suspected of having the cervical cancer.

Gorunescu (2009) described how computer-aided diagnosis (CAD) and endoscopic ultrasonographic elastography (EUSE) were enhanced by data mining to create a new non-invasive cancer detection. In the traditional approach, doctors look at the ultrasound movie and decide on whether a patient is to be subjected to a biopsy.

The physician’s judgment is primarily subjective, depending mostly on the her interpretation of the ultrasound video (see sample video screenshot, next page). Gorunescu approached this problem in a different way, using data mining. He did not study patient demographics. Instead his team focused on the ultrasound movies. They first trained a classification algorithm using a multi-layer perceptron (MLP) on known cases of malignant and benign tumors.

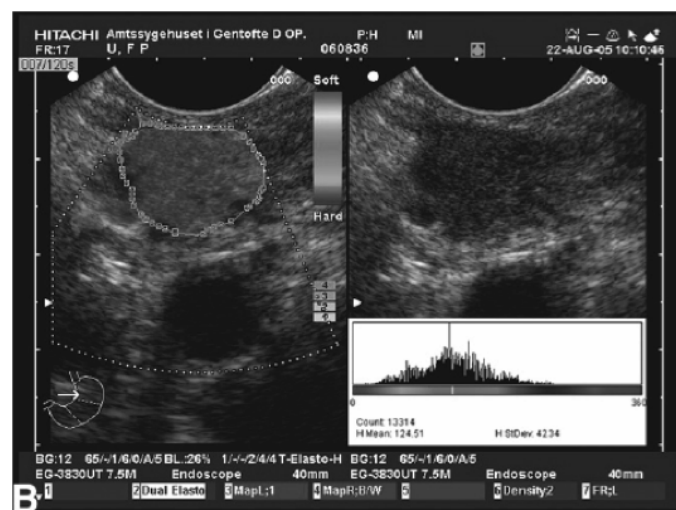


Fig. 1 EUSE sample movie frame with corresponding histogram

The model analyzed the pixels and their RGB content to find sufficient patterns to distinguish between malignant and benign tumors. Then the team applied the resulting model to other cases. They found that their model resulted to high accuracy in diagnosis with only a small standard deviation.

Adverse drug events (ADEs). Some drugs and chemicals that have been approved as non-harmful to humans are later discovered to have harmful effects after long-term public use. Wilson et al. (2003) revealed that the US Food and Drug Administration uses data mining to discover knowledge about drug side effects in their database. This algorithm called MGPS or Multi-item Gamma Poisson Shrinker was able to successfully find 67% of ADEs five years before they were detected using traditional ways.

We have seen how data mining applications could be used in early detection of diseases, prevention of deaths, the improvement of diagnoses and even detecting fraudulent health claims. However, there are caveats to the use of data mining in healthcare.

4. ISSUES and CHALLENGES

Applying data mining in the medical field is a very challenging undertaking due to the idiosyncracies of the medical profession. Shillabeer and Roddick's work (2007) cite several inherent conflicts between the traditional methodologies of data mining approaches and medicine.

In medical research, data mining starts with a hypothesis and then the results are adjusted to fit the hypothesis. This diverges from standard data mining practise, which simply starts with the data set without an apparent hypothesis.

Also, whereas traditional data mining is concerned about patterns and trends in data sets, data mining in medicine is more interested in the minority that do not conform to the patterns and trends. What heightens this difference in approach is the fact that most standard data mining is concerned mostly with describing but not explaining the patterns and trends. In contrast, medicine needs those explanations because a slight difference could change the balance between life or death.

For example, anthrax and influenza share the same symptoms of respiratory problems. Lowering the threshold signal in a data mining experiment may either raise an anthrax alarm when there is only a flu outbreak. The converse is even more fatal: a perceived flu outbreak turns out to be an anthrax epidemic (Wong et al 2005). It is no coincidence that we found that, in most of the data mining papers on disease and treatment, the conclusions were almost-always vague and cautious. Many would report encouraging results but recommend

further study. This failure to be conclusive indicates the current lack of credibility of data mining in these particular niches of healthcare.

The confusion about the definition of data mining also complicates the issue. For example, we found a couple of papers with the keywords “data mining” in their titles but turned out to be the simple use of graphs. Shillabeer (2009) said that this misunderstanding is prevalent in the relatively young existence of data mining in healthcare.

Even if data mining results are credible, convincing the health practitioners to change their habits based on evidence may be a bigger problem. Ayres (2008) reports a couple of cases where hospital doctors refused to change hospital policy even when confronted with evidence. In one case, it was found that doctors coming out of autopsy without washing hands and led to a high probability of deaths in the patients they treated after the autopsy. Presented with this evidence, doctors still refused to change their habits until only much later.

Shillabeer (2009) also reported most doctors (at least in Australia) prefer to listen to a respected opinion leader in the medical profession, rather than to the result of data mining. Shillabeer’s observation can be validated by us, since we have worked with doctors in a medical school in our capacity as an organizational management consultant.

Privacy of records and ethical use of patient information is also one big obstacle for data mining in healthcare. For data mining to be more accurate, it needs a sizeable amount of real records. Healthcare records are private information and yet, using these private records may help stop deadly diseases.

5. CONCLUSION and RECOMMENDATIONS

This survey of data mining applications in medicine and public health provided only an overview of current practises and challenges. Health care organizations and agencies could look into these applications to find ideas on how to extract knowledge from their own database systems.

For example, DOH could coordinate with government-operated hospitals, PhilHealth and the National Statistics Office to collate and analyze public health indicators. They could apply data mining techniques to find trends in disease outbreaks or deaths (eg, infant mortality), per region and per hospital.

DOH could uncover hidden patterns in deaths or disease that could lead to better health policies like better vaccination planning, identification of disease vectors like malaria, prevention

of hospital errors (the case of Rizal Medical Center comes to mind), and the inexplicable but sporadic outbreaks of flu and cholera in certain areas that are supposed to have eradicated these diseases. PhilHealth could also apply data mining to find and stop anomalous insurance claims.

Before embarking on data mining, however, an organization must formulate clear policies on the privacy and security of patient records. It must enforce this policy with its partner-stakeholders and its branches and agencies.

Public health concerns like rapid pandemic outbreaks, the need to detect the onset of disease in a non-invasive, painless way, and the need to be more responsive to its customers -- all these add up to an increasing need for health organizations to integrate data and apply data mining to analyze these data sets.

Bibliography

- Audain, C. (2007). Florence Nightingale. On-line: <http://www.scottlan.edu/lriddle/women/nitegale.htm>. Accessed 30 July 2009.
- Ayres, I (2008). *Super Crunchers*. New York: Bantam Books.
- Bailey-Kellog, C. Ramakrishnan, N. and Marathe, M. Spatial Data Mining to Support Pandemic Preparedness. *SIGKDD Explorations* (8) 1, 80-82.
- Cao, X., Maloney, K.B. and Brusica, V. (2008). Data mining of cancer vaccine trials: a bird's-eye view. *Immunome Research*, 4:7. DOI:10.1186/1745-7580-4-7
- Cheng, T.H., Wei, C.P., Tseng, V.S. (2006) Feature Selection for Medical Data Mining: Comparisons of Expert Judgment and Automatic Approaches. *Proceedings of the 19th IEEE Symposium on Computer-Based Medical Systems (CBMS'06)*.
- Health Grades, Inc. (2007). *The Fourth Annual HealthGrades Patient Safety in American Hospitals Study*.
- Kou, Y., Lu, C.-T., Sirwongwattana, S., and Huang, Y.-P. (2004). Survey of fraud detection techniques. In *Networking, Sensing and Control, 2004 IEEE International Conference on Networking, Sensing and Control*. (2) 749-754.
- Nightingale, F (1858). *Notes on Matters Affecting the Health, Efficiency and Hospital Administration of the British Army*.
- Shillabeer, A (29 July 2009). *Lecture on Data Mining in the Health Care Industry*. Carnegie Mellon University Australia.
- Shillabeer, A. and Roddick, J (2007). Establishing a Lineage for Medical Knowledge Discovery. *ACM International Conference Proceeding Series*. (311) 70, 29-37.
- Tandoc, E.S (14 October 2006). DOH order probe after Rizal hospital tragedy -- Sanitation regulations stressed. *Philippine Daily Inquirer*, p. A19.
- Thangavel, K., Jaganathan, P.P. and Easmi, P.O. Data Mining Approach to Cervical Cancer Patients Analysis Using Clustering Technique. *Asian Journal of Information Technology* (5) 4, 413-417.
- Tufte, E. (1997). *Visual Explanations. Images and Quantities, Evidence and Narrative*. Connecticut: Graphics Press.
- Wong, W.K., Moore, A., Cooper, G. and Wagner, M (2005). What's Strange About Recent Events (WSARE): An Algorithm for the Early Detection of Disease Outbreaks. *Journal of Machine Learning Research*. 6, 1961-1998.
- Wilson A., Thabane L., Holbrook A (2003). "Application of data mining techniques in pharmacovigilance". *British Journal of Clinical Pharmacology*. (57) 2, 127-134.
- Witten, I. H. and Frank, E. (2005). *Data mining : practical machine learning tools and techniques*. Morgan Kaufmann series in data management systems. Morgan Kaufman.