

11-5-2003

Optimizing Cost and Data Entry for Assignment of Patients to Clinical Trials Using Analytical and Probabilistic Web-Based Agents

Bhavesh Dineshbhai Goswami
University of South Florida

Follow this and additional works at: <https://scholarcommons.usf.edu/etd>

 Part of the [American Studies Commons](#)

Scholar Commons Citation

Goswami, Bhavesh Dineshbhai, "Optimizing Cost and Data Entry for Assignment of Patients to Clinical Trials Using Analytical and Probabilistic Web-Based Agents" (2003). *Graduate Theses and Dissertations*.
<https://scholarcommons.usf.edu/etd/1378>

This Thesis is brought to you for free and open access by the Graduate School at Scholar Commons. It has been accepted for inclusion in Graduate Theses and Dissertations by an authorized administrator of Scholar Commons. For more information, please contact scholarcommons@usf.edu.

Optimizing Cost and Data Entry for Assignment of Patients to Clinical Trials Using
Analytical and Probabilistic Web-Based Agents

by

Bhavesh Dineshbhai Goswami

A thesis submitted in partial fulfillment
of the requirements for the degree of
Master of Science in Computer Science
Department of Computer Science & Engineering
College of Engineering
University of South Florida

Co-Major Professor: Lawrence O. Hall, Ph.D.
Co-Major Professor: Dmitry B. Goldgof, Ph.D.
Eugene Fink, Ph.D.

Date of Approval:
November 5th, 2003

Keywords: Expert System, Rule Based System, Breast Cancer.

© Copyright 2003, Bhavesh Goswami.

Table of Contents

List of Tables	iii
List of Figures	iv
ABSTRACT.....	v
Chapter 1 Introduction.....	1
1.1 Introduction	1
1.2 Significance	2
Chapter 2 Previous Work	3
2.1 Software Agents.....	3
2.2 History.....	4
2.3 Similar Systems	4
Chapter 3 System Design	6
3.1 System Overview	6
3.2 Knowledge Representation.....	8
3.2.1 Medical Tests	9
3.2.2 Questions	9
3.2.3 Eligibility Criteria	9
3.3 Algorithm	11
3.3.1 Test Reordering.....	13
3.3.1.1 Analytical Reordering Agent	13
3.3.1.2 Probabilistic Reordering Agent.....	15
3.4 Testing System.....	20
Chapter 4 User Interface.....	23
4.1 Interface Design.....	23
Chapter 5 Experiments and Results	30
5.1 Analytical Experiments.....	30
5.2 Cost Saving Experiments	33
5.3 Probabilistic Experiments.....	35
5.4 Eligibility Probability.....	39
Chapter 6 Conclusions and Future Work	44
6.1 Conclusions	44
6.2 Future Work.....	45

References.....46

List of Tables

Table 1: Results of Matching Retrospective 187 Patients to Clinical Trials	30
Table 2: Results of Matching Current 169 Patients to Clinical Trials	31
Table 3: Detailed Classification of New Matches Found.....	32
Table 4: Average Dollar Cost Savings for 187 Retrospective Patients.....	34
Table 5: Average Dollar Cost Savings for 169 Current Patients	34
Table 6: Individual Test Results for Ten-Fold Cross Validation.....	37
Table 7: Results of Ten-Fold Cross Validation Per Protocol	37
Table 8: Standard Deviation for Ten-Fold Cross Validation.....	38
Table 9: Results Excluding Patients who were Ineligible by Initial Questions.....	39
Table 10: Ten-Fold Cross-Validation for Heuristic that Uses Eligibility Probability	41
Table 11: Ten-Fold Cross-Validation for Reordering Agent that Uses the Bayes Method to Compute Eligibility Probability	42

List of Figures

Figure 1: System Architecture	7
Figure 2: Example of Tests and Questions	8
Figure 3: Eligibility Criteria and its Acceptance and Rejection Expressions	10
Figure 4: Acceptance Expression Graph	11
Figure 5: Rejection Expression Graph	12
Figure 6: Expansion of CNF to DNF	14
Figure 7: User Interface Design	23
Figure 8: Initial Patient Entry Page	24
Figure 9: Initial Questions Page	24
Figure 10: Clinical Trial Selection Page	25
Figure 11: Data Entry Page. The System Asks for More Information to Determine Patient's Eligibility.	26
Figure 12: Status of all Clinical Trials	27
Figure 13: Explanation of a Decision by Explanation Sub-System	28
Figure 14: System after Eligibility of all Clinical Trials is Decided	28
Figure 15: Status of all Available Tests for the Patient	29
Figure 16: Graph Showing Average Dollar Cost Savings for Retrospective Patients	34
Figure 17: Graph Showing Average Dollar Cost Savings for Current Patients	35

**Optimizing Cost and Data Entry for Assignment of Patients to Clinical Trials Using
Analytical and Probabilistic Web-Based Agents**

Bhavesh Dineshbhai Goswami

ABSTRACT

A clinical trial is defined as a study conducted on a group of patients to determine the effect of a treatment. Assignment of patients to clinical trials is a data and labor intensive task. Usually, medical personnel manually check the eligibility of a patient for a clinical trial based on the patient's medical history and current medical condition. According to studies, most clinical trials are under-enrolled which negatively affects their effectiveness. We have developed web-based agents that can test the eligibility of patients for many clinical trials at once. We have tested various heuristics for optimizing cost and data entry needed in assigning patients to clinical trials. Testing eligibility of a patient for many clinical trials is only feasible if it is cost and data entry efficient. Agents with different heuristics were then tested on data from current breast cancer patients at the Moffitt Cancer Center. Results with different heuristics are compared with each other and with that of the clinicians. It is shown that cost savings are possible in clinical trial

assignment. Also, less data entry is needed when probabilistic agents are used to reorder questions.

Chapter 1

Introduction

1.1 Introduction

A clinical trial is an experimental research study that evaluates a specific new treatment for a specific population of patients. The trial protocol is a rulebook which clearly identifies the criteria for a patient to be eligible for the trial. The criteria is based on the medical history and the present medical condition of the patient. Some criteria are general information such as the age and sex of the patient while others requires specific tests be done to determine if the patient matches them. Eligibility is usually checked manually by a clinician, nurse or trial coordinator. In the absence of enough information to determine the eligibility of a patient, clinicians order tests needed to obtain the required information. Each test has some cost associated with it and studies show that if tests are reordered correctly, cost savings are possible [1, 11, 14]. Although clinicians can potentially reduce costs by ordering inexpensive tests first, test reordering is a complex optimization problem which has many parameters like the actual cost of the test, number of eligibility criterion decided by it, number of protocols in which the test is needed and the probability of the test resulting in an eligibility decision for the patient.

1.2 Significance

Cancer is one of the major causes of death in United States resulting in 550,000 deaths annually. Cancer treatment is an active research area. Clinical trials are used to research new treatments. Accrual of patients to these clinical trials is of the utmost importance for the success of a trial. At any one time there are many active trials, and keeping track of them becomes a difficult task. Studies show that clinicians sometimes miss up to 50% to 60% of the eligible patients [11, 22, 33]. In such a scenario, having a web-based agent with a central database of clinical trials that can be accessed through the internet by any willing institute can drastically improve accrual of patients. It may also ease the process of trial sharing among various institutes, which is difficult otherwise.

Chapter 2

Previous Work

2.1 Software Agents

A software agent is defined as an autonomous software entity that can interact with its environment to accomplish certain objectives without any direct input or supervision from its user. Software agents are used in a variety of fields like VLSI design, search engines, e-commerce applications, business processes, medical research, etc. Earlier use of AI in medical domains was primarily through expert systems for diagnosis and treatment suggestions. The primary difference between expert system and agent is that the agent changes dynamically according to the changes in its environment and that many agents can work together to accomplish certain task(s). The analytical and probabilistic agents, which we developed, dynamically change their heuristic parameter values as new information is received. The analytical agent modifies the rule graphs and cost information when it acquires new information, while the probabilistic knowledge base is updated on acquisition of new information, which dynamically changes the heuristic parameter values of the probabilistic agent. Also the analytical and probabilistic agents can work together to reorder the questions and assign patients to clinical trials. Due to the above facts, we term them as agents.

2.2 History

Research in the use of Artificial Intelligence for medical diagnosis has been conducted since the early seventies. A variety of expert systems were developed for diagnosis and treatment of various diseases. The MYCIN system for diagnosis of bacterial diseases was developed by Shortliff and colleagues [13]. It evolved from a chemical expert system called DENDRAL [8, 24, 25]. It was a rule-based system with *if-then* rules using the backward chaining process. The system gathered information about the patient and provided treatment recommendations with explanations of how the conclusion was reached. MYCIN was found to be fairly accurate in the experiments conducted, which led to the development of various other medical expert systems. EMYCIN [7] was developed from MYCIN. It was a generic expert system shell where rules from any domain could be added to the knowledge base. A lung disease knowledge base was developed using EMYCIN and the resultant system was called PUFF [26]. An expert system called NEOMYCIN [9] was developed to train doctors by presenting them with practice cases, getting their diagnosis and correcting them when they made mistakes.

2.3 Similar Systems

Researchers have adopted many different approaches for systems to assign patients to clinical trials. A system called AIDS [28] was developed to assign patients to HIV clinical trial treatment protocols. It used Bayesian belief networks to manage uncertainty. It has a protocol driven mode, where each patient's probability of eligibility is checked automatically for a new clinical trial protocol using the data patient has. EON [30] was made with reusability as a high priority. Thus they made four reusable basic components of the system. The knowledge base could be generated for any protocol based therapy.

ONCODOC [4, 5, 6, 19, 20] was developed to “*enhance the accrual of patients in the best care plan.*” It had non-metastatic breast cancer protocols encoded as decision trees. The goal was to enhance the usability of the system by providing flexibility in the interpretation of the results by clinicians. Clinicians incrementally assign values to the decision parameters while navigating through the decision tree. Thus it is halfway between knowledge representation and a static written description and allows both formal and informal aspects of the protocols to be implemented. It means that a clinician can make an explicit decision if a test result is borderline or can navigate through the remaining decision tree to see the underlying logic and make his decision. An important drawback of using decision trees is that you cannot test a patient for multiple trials at once. One will have to test a patient separately for each trial, and thus the time consumed will increase drastically when the patient’s eligibility is checked for many clinical trials.

A system to assign patients to clinical trials was designed by Papaconstantinou and colleagues [10] using Bayesian networks. They were successful in implementing three clinical trials but concluded that the complexity of the Bayesian network immensely increased when new trials were being added. Another drawback was the difficulty in getting the prior probability relationships between different nodes of the network. Also, propagation of probabilities with data evidence was very slow. Another system DIAVAL [8] was also a Bayesian expert system for echocardiography. It had the same problems discussed above. Smith and McNeely [2, 31] built an expert system to reorder tests for laboratory investigations of patients for general clinical problems. They used the ACQUIRE suite of expert system tools to design this system. Their results shows that the laboratory costs were reduced from mean \$232 to \$194 on tested patients.

Chapter 3

System Design

3.1 System Overview

The system is divided into three main parts (a) Knowledge entry system, (b) Agent and, (c) Testing system. Nikiforou [36] implemented the knowledge entry system. It is a web-based system that has user-friendly interface for encoding the clinical trials into a form that is understood by the agents. Fletcher, Kokku [33] and colleagues have built the initial agent for matching of patients to clinical trials. We modified the way cost was dealt with in the system and conducted systematic experiments to evaluate the cost-savings by the system. We also added the probabilistic agents to the system and conducted experiments to compare the results. Figure 1 represents the basic structure of the system.

As shown in Figure 1, the user interacts with the system through the web-based interface. A user can access old patient data, add new patients and enter data for existing patients. Each time new information is obtained about the patient, her eligibility is checked for the protocols selected. The patient is eligible, ineligible or more information is required to decide on eligibility for a particular protocol. If more information is needed to determine patient's eligibility, the probabilistic agent calculates the eligibility probability based on the current data of the patient and probabilistic knowledge base of the system.

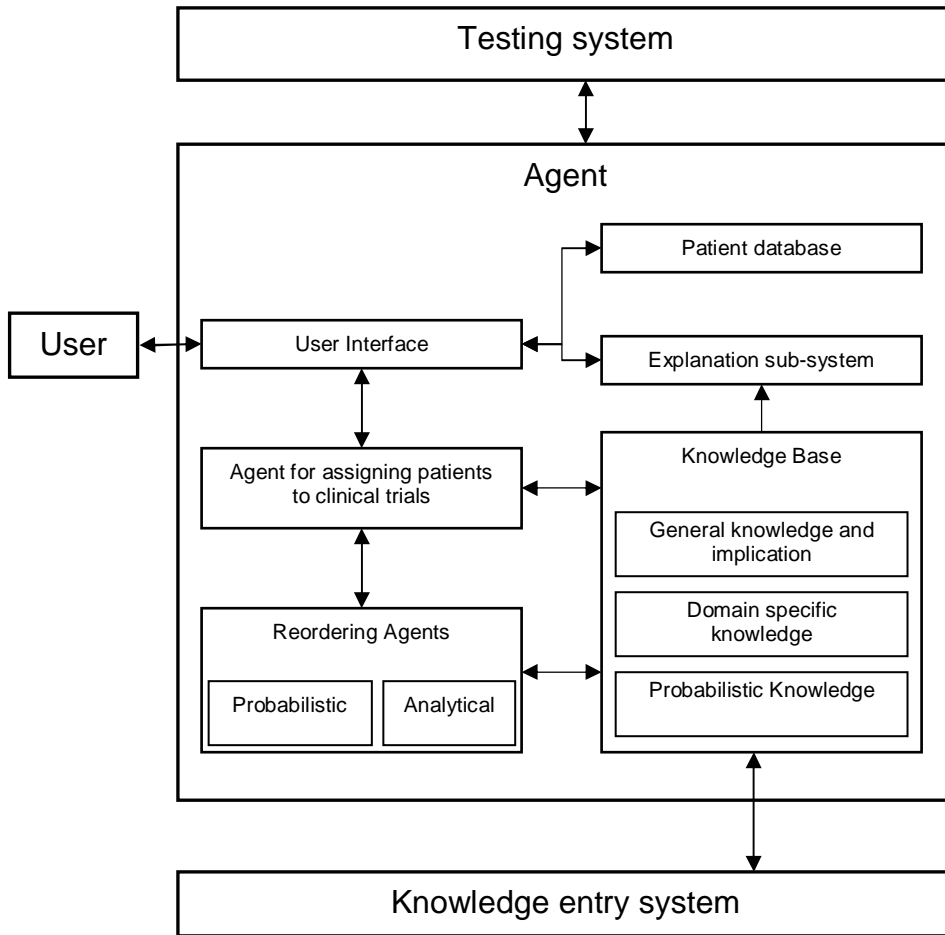


Figure 1: System Architecture

The system augments its probabilistic knowledge base by gathering information entered for the current patients and its effect on the eligibility of the patient. After checking eligibility, if the patient has protocols for which eligibility is not yet decided upon, the reordering subsystem reorders the relevant questions using analytical and probabilistic agents. At all times, the system has an explanation subsystem which can provide an explanation for system's decisions. More protocols can be added to the system using the knowledge entry system. The testing system is used to test new heuristics by

creating many new patients and testing the effectiveness of the system without human intervention.

3.2 Knowledge Representation

The system incorporates domain knowledge in the form of medical tests, questions and eligibility criteria. A test can be entered into the system with its name and cost. A list of questions that are answered by that test are then entered. For example, as shown in Figure 2, Mammogram test answers questions i) What is the cancer stage? ii) Does the patient have invasive cancer? Clinical trial protocols can then be entered into the system, which are essentially a set of rules composed of the entered questions. The eligibility of the patient for a particular protocol can then be tested using the agent. There is a web-based interface to accomplish all these tasks.

General Questions	Cost: \$0
What is patient's sex?	<input type="checkbox"/> Female <input type="checkbox"/> Male
What is patient's age?	<input type="text"/>
Mammogram	Cost: \$150
What is the cancer stage?	<input type="checkbox"/> I <input type="checkbox"/> II <input type="checkbox"/> III <input type="checkbox"/> IV
Does the patient have invasive cancer?	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Unknown
CT scan (Head)	Cost: \$850
Are there any symptoms of metastatic disease in brain?	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Unknown

Figure 2: Example of Tests and Questions

3.2.1 Medical Tests

Medical tests that are performed on patients to obtain information needed to determine their eligibility are entered into the system along with the cost incurred to perform that test. Each test, when performed, gives us information about the patient which is used to determine her eligibility. This information is stored in the form of questions.

3.2.2 Questions

A question can be any of the three forms: (a) Yes/No, (b) Multiple choice and (c) Numeric. A Yes/No question is the one where the answer to the question is yes, no or unknown. Similarly, multiple choice questions have a list of two or more answers to select the answer from. A numeric question has some numeric value as its answer. Figure 2 shows some tests and the questions associated with them.

3.2.3 Eligibility Criteria

The clinical trial eligibility criteria is essentially a logical expression of questions. Figure 3 shows an eligibility criteria and the corresponding logical expression. The system keeps asking for data until either the expression in Figure 3(a) is TRUE or the expression in Figure 3(b) is FALSE, which means that the patient is eligible or ineligible respectively. For example, let's assume that a patient's eligibility is checked for the given trial which has the eligibility criteria shown in Figure 3.

1. *Female of age 18 to 50 years.*
2. *Must be postmenopausal or using contraceptive.*
3. *Should not have metastatic disease in brain.*
4. *Cancer should not be invasive.*

(a) Acceptance Expression

sex = FEMALE AND
age ≥ 18 AND
age ≤ 50 AND
{ postmenopausal = YES OR use_contraceptive = YES } AND
invasive_cancer = NO AND
brain_metastatic = NO

(b) Rejection Expression

sex = MALE OR
age ≤ 18 OR
age ≥ 50 OR
{ postmenopausal = NO AND use_contraceptive = NO } OR
invasive_cancer = YES OR
brain_metastatic = YES

Figure 3: Eligibility Criteria and its Acceptance and Rejection Expressions

Assume the system already has the data that the patient is female and that she is postmenopausal. The remaining pertinent information to be obtained is (a) Age, (b) If there is metastatic disease in the head and (c) If cancer is invasive. We do not need to perform any tests to know a patient's age so that information is free to us. On the other hand, we need to perform a mammogram to check if the cancer is invasive and it costs \$150. Similarly, to check if there is metastatic disease in brain, we need to perform a CT scan of the head, which costs \$850. The cost optimization problem here is to see if it is

more beneficial to order a CT scan of the head or a mammogram. A mammogram looks like a better alternative at the first glance as it costs much less than CT scan. However there are other issues which are important too.

For example if most patients have invasive cancer, then it will be more efficient to go for a CT scan, even though it costs more, as the probability of a patient having invasive cancer is high and it requires us to do a CT scan anyways after a mammogram is performed. Also, when one tests for multiple protocols at once, many protocols might need information which can be obtained from a single test, so that test gets priority. Thus test ordering is a complex optimization task. The system keeps asking questions with continued reordering until it can either determine that the acceptance expression is TRUE or rejection expression is TRUE.

3.3 Algorithm

As described above, the system stores eligibility criteria of different trials in the form of rules which are described as logical expressions of different questions. While computing eligibility we store the rules in the form of nodes.

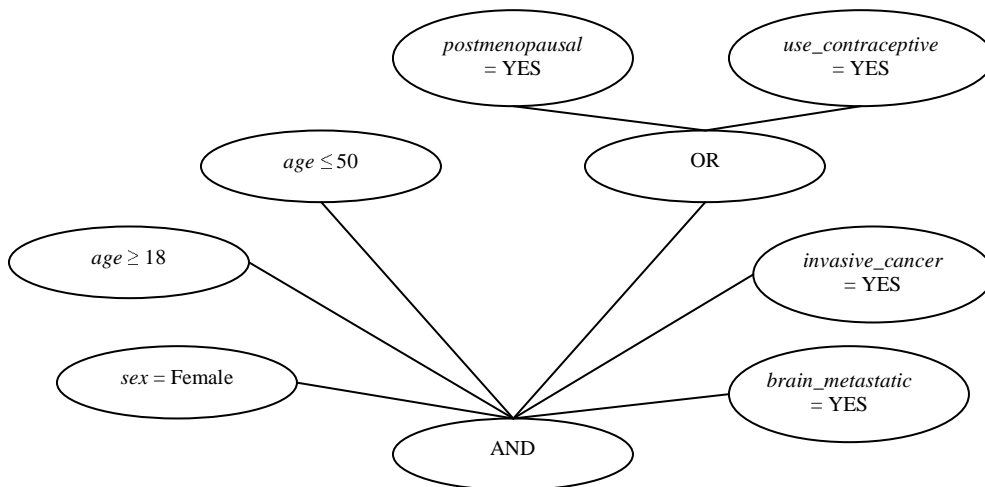


Figure 4: Acceptance Expression Graph

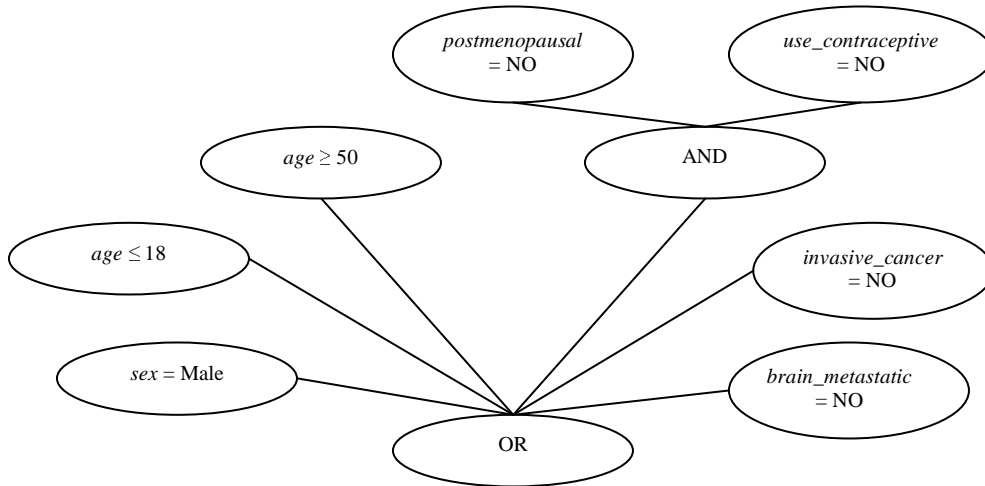


Figure 5: Rejection Expression Graph

Two graphs are created, one for the acceptance expression and the other for the rejection expression. The base node of the acceptance expression graph is the AND node which represents the eligibility of a patient. It represents the conjunction of conditions that are needed to be satisfied in order for a patient to be eligible for a protocol. Figure 4 shows the acceptance expression graph for eligibility criteria defined in Figure 3. Figure 5 is the rejection expression which represents the ineligibility criteria. The base node of this graph is an OR node. The OR node represents a disjunction of conditions that rule out a patient for a trial. The patient is eligible if the AND node of the acceptance expression is TRUE, ineligible if the OR node of the rejection expression is TRUE and the eligibility is not decided if neither of them is true. Only one of them can be TRUE at a time for a given patient. The system will keep asking questions until one of the base nodes is determined to be TRUE.

3.3.1 Test Reordering

Analytical and probabilistic agents are used to reorder tests for the optimization of cost and data entry. Analytical heuristics are based on test cost and the structure of the rejection and acceptance expressions. Probabilistic agents use probabilistic data accumulated over time by the system.

3.3.1.1 Analytical Reordering Agent

The analytical reordering agent uses three heuristic parameters to reorder tests. The first parameter is the most fundamental one, which is the cost of a test. The more expensive the test, the less likely it is for the system to obtain information associated with that test. Less expensive tests get priority over the more expensive ones. Although this heuristic looks straightforward, it has inherent drawbacks. For example imagine a scenario where the agent has to decide between ordering two tests, test A and test B, where test B costs twice as much as test A. Although the first thought would be to order test A ahead of test B, that may not always be the case. Imagine a scenario where test A's results are almost always favorable for the patients, and almost no patient gets ruled out of a protocol due to test A's results. In such a scenario it will be more efficient to order test B ahead of test A. When a patient is tested for multiple trials at once, a certain test might be needed to determine eligibility for more than one trial. That is the second heuristic parameter. The number of clinical trials that need a particular test be done, to determine the eligibility, is taken as a heuristic parameter. The probability of a test being ranked highest is linearly proportional to the number of trials that need it.

The third heuristic parameter is the number of clauses that include questions answered by that test in the acceptance expression. To calculate this, the acceptance

expression needs to be converted into Disjunctive Normal Form (DNF). Converting a graph into DNF can take a considerable amount of processing time if the graph has many OR nodes.

For example consider a simple expression $(a \vee b \vee c) \wedge (e \vee f \vee g) \wedge (h \vee i \vee j)$ with nine AND nodes and 1 OR node. After converting it into DNF the expression looks like that in Figure 6

$$\begin{aligned}
 &(a \wedge e \wedge h) \vee (a \wedge e \wedge i) \vee (a \wedge e \wedge j) \vee (a \wedge f \wedge h) \vee (a \wedge f \wedge i) \vee (a \wedge f \wedge j) \vee \\
 &(a \wedge g \wedge h) \vee (a \wedge g \wedge i) \vee (a \wedge g \wedge j) \vee (b \wedge e \wedge h) \vee (b \wedge e \wedge i) \vee (b \wedge e \wedge j) \vee \\
 &(b \wedge f \wedge h) \vee (b \wedge f \wedge i) \vee (b \wedge f \wedge j) \vee (b \wedge g \wedge h) \vee (b \wedge g \wedge i) \vee (b \wedge g \wedge j) \vee \\
 &(c \wedge e \wedge h) \vee (c \wedge e \wedge i) \vee (c \wedge e \wedge j) \vee (c \wedge f \wedge h) \vee (c \wedge f \wedge i) \vee (c \wedge f \wedge j) \vee \\
 &(c \wedge g \wedge h) \vee (c \wedge g \wedge i) \vee (c \wedge g \wedge j)
 \end{aligned}$$

Figure 6: Expansion of CNF to DNF

As shown in Figure 6, the expression now has 27 AND nodes and one OR node. For large expressions, this conversion takes considerable computing time and thus can be a drawback. After the conversion, for each test, the number of clauses that depend on a question belonging to that test is calculated. The tests that are needed in more clauses get priority.

After the system gets any new information, it regenerates the graphs and recalculates heuristics values. After calculating all three parameters of all tests, the agent does a linear combination of the parameter values and orders the tests using it.

3.3.1.2 Probabilistic Reordering Agent

Most people who have used a probabilistic approach towards expert systems or agents have used Bayesian Networks. There are many inherent problems with this approach and the prime ones are: (a) Bayesian Networks are very complex and it takes a lot of time to build new networks when new trials are added (b) It is difficult to get initial relational probabilities between nodes of the network (c) It takes a lot of computational time to modify the conditional probability between nodes of the network in the wake of new evidence [10, 15]. In spite of the above drawbacks, they also have some critical advantages over analytical rule-based systems. The most important advantage probabilistic systems have is their ability to estimate eligibility probability in the absence of some information. Rule based systems cannot estimate the probability even in the absence of a single piece of information. Consider a protocol which has 20 eligibility rules, and 19 of them are met, which seemingly makes it likely for a patient to be eligible. A rule based system will still not be able to predict anything about the probability of the patient's eligibility while a probabilistic system can estimate the eligibility of the patient using the prior probabilities that the system has. Probability estimation becomes important in certain cases, especially when you have many trials. The Moffitt Cancer Center at USF has about 15 active breast cancer trials at once. It is time-consuming and expensive to test patients for all trials. Thus rather than testing the patients for all trials, we can just test her eligibility for trials which show high initial eligibility probability. We can also use probabilistic knowledge accumulated over time by the system to reorder the tests.

To attempt to exploit the advantages of the probabilistic systems and to avoid its drawbacks, we used the probabilistic methods discussed below. The basic idea is to try to classify a patient as ineligible as soon as possible. The idea is that if a patient is ineligible, the information that is most likely to determine her ineligible should be obtained first. This would optimize the data-entry needed to decide upon a patient's eligibility. The system gathers this probabilistic knowledge over time. For each question in the system, it keeps a log of how many times a question is asked and how many times a patient is ruled out for a particular trial after that question is answered. This gives us the probability of that question, when asked, ruling out a patient for a particular trial. So the approach would be to ask the question which has the highest probability of ruling out a patient first. When multiple trials are tested at once and same question is in the acceptance expression of more than one trial, we add the separate probabilities of the question for each trial. Note that the result of the summation can be greater than one, thus it is not a probability anymore. Rather it is summation value which, when higher, suggests a higher likelihood of determining a patient ineligible.

These probabilities are also used to estimate the eligibility probability. Here we make an important assumption that all questions have independent probabilities. Although this assumption is not entirely true, it seems close enough. Most questions are either completely dependent on each other or are not at all dependent. For example if a patient has no positive lymph nodes, that means that the cancer stage is either 0 or I. Thus the two questions "Does the patient have positive lymph nodes?" and "What is the cancer stage?" are dependent on each other. We can take care of such situations by an implication sub-system. We can add implication rules in the system such as "If cancer

stage 0 or I then it implies that patient has no positive lymph nodes”. When system has information that the cancer stage is either 0 or I, the implication subsystem automatically generates information that no lymph nodes are positive and the system does not ask for that information. Thus, all such completely dependent questions are taken care of by the implication subsystem. There are very few questions which have a partial implication, like “If a patient is ER positive then there is an 80% chance of positive lymph nodes”. We ignore such conditional probabilities among questions and treat all questions as either completely dependent or independent.

When we assume all questions to be independent, the eligibility criteria for a protocol becomes a set of independent questions which must have favorable answers for a patient to be eligible. Also we have the probability of how many times a question, when answered, was responsible for a patient being ineligible for a protocol. Thus we have probability of each question being answered favorably for eligibility. A patient will be eligible for a protocol when each answer is favorable. Thus the eligibility probability of a protocol will be the product of the favorable probabilities of all questions. When a question is answered and the answer meets the eligibility rules of the clinical trial, we can exclude it from our product and thus we get the new eligibility probability. As more questions answered fit the eligibility criteria, the eligibility probability increases. When a question’s answer does not fit the eligibility profile, the patient becomes ineligible and the eligibility probability becomes zero. When the patient becomes eligible the eligibility probability is 1. The eligibility probability varies between 0 and 1 when the eligibility of a patient is undecided.

Although the above mentioned approach gave good results when the computed probabilities were used to reorder questions (as discussed in the “Experiments and Results” section), the probabilities computed were not normalized and thus were very small numbers that did not give any meaningful feedback to the user. To compute probabilities that provide meaningful feedback, the Naïve Bayes approach appeared to fit well. We can think of the patient enrolment procedure as a classification problem. The classification classes will be “Eligible” or “Ineligible”. The attributes will be the questions and the values for the question are “Favorable for eligibility” or “Unfavorable for eligibility” for each clinical trial. We have a probability for each question to be favorable and unfavorable for each clinical trial. Thus we have a classification problem where we have probabilities for the occurrence of each attribute value. To use Naïve Bayes we also needed probabilities of occurrence of each classification type, which is “Eligible” or “Ineligible” in our case. We modified the system so that for each clinical trial, the system recorded how many patients were tested for that clinical trial and how many patients were decided to be eligible and ineligible. Thus, we could now use Naïve Bayes to calculate the eligibility probability of a patient with partial information.

Let us assume that we have a clinical trial T with three questions Q_1 , Q_2 and Q_3 . Out of 100 patients that were tested on the clinical trial, we found 40 to be eligible and 60 to be ineligible. Question Q_1 was asked 90 times and it disqualified patients 10 times, Q_2 was asked 80 times and disqualified patients 5 times, and Q_3 was asked 70 times and disqualified patients 15 times. Thus $P(T_E)$, the probability of patient being eligible for

protocol T is $\frac{40}{100}$. $P(Q_1)$, the probability that questions Q_1 is answered favorably for

clinical trial T is $\frac{80}{90}$. Similarly $P(Q_2) = \frac{75}{80}$ and $P(Q_3) = \frac{55}{70}$.

So when we don't have any information, i.e. we don't have answers for any questions, then the probability of a given patient being eligible for clinical trial T is $\frac{40}{100}$. Now assume that we have answers to questions Q_1 and Q_2 and the answers are favorable for eligibility. Thus the new eligibility probability will be $P(T_E | Q_1, Q_2)$. According to Naïve Bayes:

$$P(T_E | Q_1, Q_2) = P(T_E) P(Q_1 | T_E) P(Q_2 | T_E)$$

$P(Q_1 | T_E)$ is the probability that question Q_1 was answered favorably for eligibility given the fact T_E , which means that patient is eligible for clinical trial T. Now, if a patient is eligible for a clinical trial, all the questions must be answered favorably for eligibility. Thus $P(Q_1 | T_E) = 1$. Similarly $P(Q_2 | T_E) = 1$. Thus $P(T_E | Q_1, Q_2) = P(T_E)$.

As shown, $P(T_E | Q_1, Q_2 \dots Q_n)$ will always be $P(T_E)$ regardless of the number of questions answered. Thus Naïve Bayes cannot be used in this case as the probabilities do not change dynamically when the acquisition of new information. For enhancing the estimation technique for the clinical trial probability, we then used Bayes rule. Using Bayes rule

$$P(T_E | Q_1, Q_2) = \frac{P(T_E) P(Q_1, Q_2 | T_E)}{P(Q_1, Q_2)}$$

As explained before $P(Q_1, Q_2 | T_E)$ will be the probability that questions Q_1, Q_2 are favorable for eligibility given that the patient is eligible for protocol T. If the patient is

eligible then the questions answered will always be favorable and thus $P(Q_1, Q_2 | T_E) = 1$.

Also we have assumed that all questions are independent and thus

$$P(T_E | Q_1, Q_2) = P(T_E | Q_1) P(T_E | Q_2)$$

Substituting all this values, we get the new equation as

$$P(T_E | Q_1, Q_2) = \frac{P(T_E)}{P(Q_1) P(Q_2)}$$

Substituting the values of $P(T_E)$, $P(Q_1)$ and $P(Q_2)$ we get the value 0.48 for the equation. Thus after answering two questions Q_1 and Q_2 , the eligibility probability of a patient for clinical trial T is 48%.

The generic equations for eligibility probability for a clinical trial when we have favorable answers for eligibility for n questions will be

$$P(T_E | Q_1, Q_2, \dots, Q_n) = \frac{P(T_E)}{P(Q_1) P(Q_2) \dots P(Q_n)}$$

3.4 Testing System

Every time a new reordering agent is implemented or the current agent is modified, its effectiveness needs to be determined. To obtain data that can be used to compare different agents' effectiveness, a certain number of patients need to be assigned to clinical trials using the clinical trial assignment system that incorporates the reordering agents in question. Performing the operation of repeatedly assigning numerous patients to clinical trials using the system is a time consuming and labor intensive task. The manual data entry required severely increases the time needed to perform such experiments and limits the number of patients that can be tested. To overcome such shortfalls, a testing system was developed.

The testing system takes a list of patient names, a list of clinical trials and path to the patient assignment system as its arguments. It then autonomously checks the eligibility of the list of patients for the list of clinical trials using the assignment system which is accessible at the path specified in the argument. The patients in the patient list already have had their eligibility manually determined for a set of clinical trials. For each such patient in the patient list, the testing system creates a new test-patient and attempts to determine the eligibility of the patient for the list of clinical trials provided using the clinical trial assignment system whose link is passed as a parameter. The clinical trial assignment system reorders questions and presents the question list with the topmost question being the most important, according to the reordering parameters. The testing system answers the topmost question of that list, by obtaining the answer from the patient's data that was manually added into the system. After the answer is obtained, the clinical trial assignment system regenerates the reordering parameter values and reorders the questions.

At the end of the execution, the testing system generates a summary file which contains information for each patient such as the number of questions asked, cost, eligibility status of the clinical trials, etc. Thus the testing system makes it fast, effective and easy to test different heuristics without human intervention. The testing system works autonomously, but it requires human intervention in certain cases. When the clinical trial assignment system presents the reordered questions list and the testing system does not have the answer to the topmost question on the list, it requires human intervention to find out the answer to the question. This generally happens because the original patient data do not contain answers for all questions. The clinical trial assignment system stops

presenting further questions for a trial when the patient is either eligible or ineligible for that trial. So a patient's data generally do not contain answers to all the questions. When a new heuristic is implemented, questions are reordered differently and often a question which is unanswered in the original patient data is presented. The testing system would not be able to find data for it and thus it seeks human intervention. In the summary file, all such patients are listed. A user can then look at medical database of the hospital or patient's chart to obtain the information and enter it into the system.

The code for the testing system is written in PERL, while the rest of the system is coded in C. The code can be found in the CD ROM presented with the thesis.

Chapter 4

User Interface

4.1 Interface Design

The system provides a user-friendly web-based interface for the user. The whole system is designed in such a way that there is a central system and clinical personnel can access it from anywhere via the internet. A user can enter new patients and check their eligibility for all the available trials. Figure 7 shows the control flow diagram of the user interface.

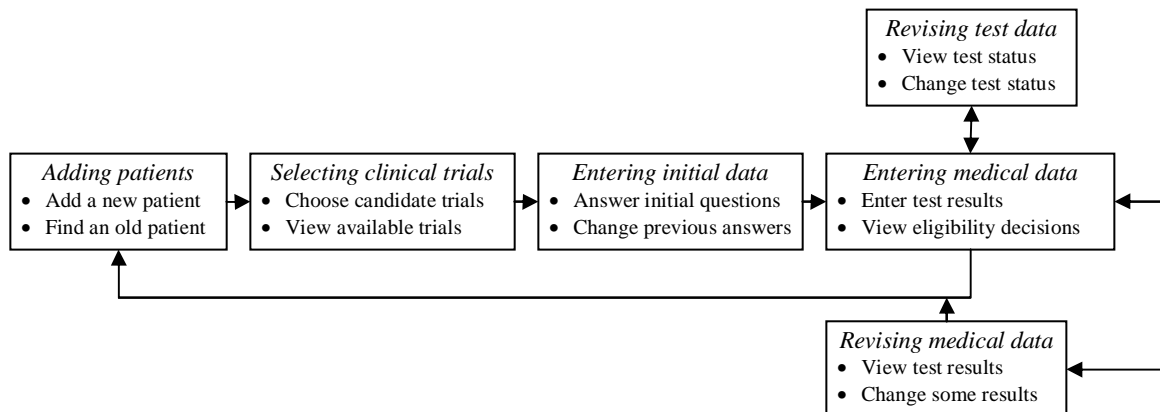


Figure 7: User Interface Design

The user starts by entering the name and identification number for the patient as shown in Figure 8. The combination of name and identification number should be unique for all the patients. The user also has search options to find patients whose data was

entered previously and may resume entering data for them. There is an option to optimize cost or data entry when testing. As shown in Figure 8, if the user selects “Cost Version”, cost is optimized and if the “No Cost Version” is selected, data entry is optimized.

Patient Name <input type="text" value="Christina Adams"/>	<input type="button" value="NEW PATIENT"/> Click to enter a new patient <input type="button" value="PATIENT SEARCH"/> Click if this patient was previously entered
Patient ID Number <input type="text" value="275113"/>	
Version <input type="radio"/> CostVersion <input checked="" type="radio"/> No Cost	
<input type="button" value="CLEAR"/>	

Figure 8: Initial Patient Entry Page

What is the stage of the breast cancer? (\$ 0.00) <input type="text" value="IIA"/>	What is the patients gender? (\$ 0.00) <input type="text" value="Female"/>
What is the patient's age in years? (\$ 0.00) <input type="text" value="45"/>	How many nodes are positive? (\$ 0.00) <input type="text" value="2"/>
What is the greatest diameter of the tumor in cm? (\$ 0.00) <input type="text" value="1.2"/>	What is the estrogen receptor status? (\$ 0.00) <input type="text" value="Positive"/>
What is the progesterone receptor status? (\$ 0.00) <input type="text" value="Positive"/>	Is the patient pregnant or nursing? (\$ 0.00) <input type="radio"/> Yes <input checked="" type="radio"/> No <input type="radio"/> Defer
Is the patient considered a candidate for adjuvant or first-line hormonal therapy? (\$ 0.00) <input type="radio"/> Yes <input checked="" type="radio"/> No <input type="radio"/> Defer	Has the patient had surgery for breast cancer? (\$ 0.00) <input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
Did the surgery include an axillary dissection? (\$ 0.00) <input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Has the patient been administered any therapy for cancer? (\$ 0.00) <input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
If the patient was administered therapy for cancer, what types were used? (\$ 0.00) <input type="text" value="Hormonal"/>	If the patient has had surgery for breast cancer, what was the most extensive type? (\$ 0.00) <input type="text" value="Breast-sparing_procedure_or_lumpectomy"/>
<input type="button" value="PROCESS"/> Click to submit your answers to the system	

Figure 9: Initial Questions Page

After the patient's name and identification are added, the system presents 14 initial questions as shown in Figure 9. These are general information questions like patient's age, sex, etc. and are generally used in most of the clinical trials. After the user has entered the initial information, the next screen is a list of available trials for which the patient's eligibility can be checked. At least one of the clinical trials must be selected. In Figure 10 five clinical trials of the available 17 were selected.

Available Protocols ([Find out about each protocol](#)):

- 10822 : Phase III Comparison of Adjuvant Chemotherapy with High-Dose Cyclophosphamide Plus Doxorubicin (AC) Versus Sequential Doxorubicin Followed by Cyclophosphamide (A->C) in High-Risk Breast Cancer Patients with 0-3 Positive Nodes.
- 10840 : Doxorubicin Dose Escalation, with or without Taxol, as Part of the Ca Adjuvant Chemotherapy Regimen for Node Positive Breast Cancer: A Phase III Intergroup Study(INT-0148, CALGB 9344, ECOG C9344, NCCTG 94-30-51, SWOG 9410).
- 11072 : A Phase II Study of Cyclophosphamide, Thiotepa, And Carboplatin (CTC) As A Preparative Regimen for Autologous Hematopoietic Stem Cell Transplant in Stage II OR III Breast Cancer (Excluding Inflammatory Breast).
- 11132 : Evaluation of Weight Gain in Breast Cancer Patients Receiving Adjuvant Radiation Therapy or Chemotherapy as a Function of Changes in Metabolic, Hormonal, Psychological and Lifestyle Factors.
- 11378 : A Comparison of Intensive Sequential Chemotherapy using Doxorubicin Plus Paclitaxel Plus Cyclophosphamide with High-Dose Chemotherapy and Autologous Hematopoietic Progenitor Cell Support for Primary Breast Cancer in Women with 4-9 Axillary Lymph Nodes (SWOG 9623).
- 11931 : NAFTA (The North American Fareston versus Tamoxifen Adjuvant Trial for Breast Cancer).
- 11971 : Phase II Neoadjuvant Trial of Sequential Doxorubicin and Docetaxel for the Treatment of Stage III Breast Cancer Measuring STAT Activation as a Predictor of Response To Therapy.
- 11992 : The Advanced Breast Biopsy Instrumentation (ABBI) Procedure for Lumpectomy Margin Determination.
- 12100 : A Randomized Trial of Axillary Node Dissection in Women With CLINICAL T1 - 2 N0 M0 Breast Cancer who have a Positive Sentinel Node.
- 12101 : A Prognostic Study of Sentinel Node and Bone Marrow Micrometastases in Women with Clinical T1 or T2 N0 M0 Breast Cancer.
- 12385 : Axillary Lymph Node Staging in Breast Cancer. Comparison of Sentinel Node Biopsy and Positron Emission Tomography (PET) Scanning.
- 12601 : The pharmacokinetics of weekly Docetaxel in older patients with metastatic breast cancer.
- 12643 : Breast Cancer Susceptibility Markers.
- 12757 : Phase III Trial of Doxorubicin Cyclophosphamide Followed by Taxol/ Trastuzumab as Adjuvant Treatment for Patients with HER2 Overexpressing Node Positive Breast Cancer.
- 12775 : Radioactive Seed Localization Breast Biopsy.Elimination of Post-Biopsy Specimen Mammogram.
- 12777 : Development and Validation of the Personal Changes Questionnaire for Persons with Cancer.
- 12885 : High Dose Cyclophosphamide, Thiotepa, and Carboplatin Followed by Hematopoeitic Blood Stem Cell Transplantation in Patients with High Risk, Locally Advanced Breast Cancer or Metastatic Breast Cancer.

Figure 10: Clinical Trial Selection Page

The system then checks the eligibility of the patient for the selected trials. If it needs more information to determine the eligibility then it makes a list of missing information questions and reorders the questions optimizing cost or data-entry as selected by the user. It then presents the top ten questions of the reordered question list on the next page as shown in Figure 11.

PROTOCOL-----	STATUS-----	QUESTIONS REMAINING-----	PERCENTAGE OF QUESTIONS ANSWERED
12100	More Information Needed	11	21 <input type="button" value="Why?"/>
How many months have passed since the diagnosis of the patient's cancer? (\$ 0.00)	<input type="text" value="6"/>	Does the patient have at least a sixth grade education? (\$ 0.00)	<input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
Is the patient able to understand, read, and write English? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Does the patient have a suspicious non-palpable breast lesion requiring breast biopsy for diagnosis? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
Is the patient is deemed by their treating physician to be at low risk for recurrence from prior malignancies? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Has there been any evidence of any prior malignancies for at least five (5) years? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
Did the patient undergo potentially curative therapy for all prior malignancies? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Does the patient have a medial quadrant lesion? (\$ 0.00)	<input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
Does the patient have life expectancy of 10 year or more? (\$ 0.00)	<input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Is the patient available for follow-up? (\$ 0.00)	<input checked="" type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer
<input type="button" value="PROCESS"/> Click to submit your answers to the system. <input type="button" value="REVIEW"/> Click to review and change your previous answers			
<input type="button" value="RETURN"/> Click to return to the entry page.			

Figure 11: Data Entry Page. The System Asks for More Information to Determine Patient's Eligibility.

PROTOCOL-----	STATUS-----	QUESTIONS REMAINING-----	PERCENTAGE OF QUESTIONS ANSWERED		
	12100	More Information Needed	11	21	Why?
	12100	More Information Needed	11	21	
How many months have passed since diagnosis of breast cancer? (\$ 0.00)	12101	More Information Needed	21	8	
	12601	Ineligible	30	3	sixth grade education? (\$ 0.00)
	12775	More Information Needed	1	50	
	12777	More Information Needed	3	25	
Is the patient able to understand, read, and write English? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Does the patient have a suspicious non-palpable breast lesion requiring breast biopsy for diagnosis? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer				
Is the patient is deemed by their treating physician to be at low risk for recurrence from prior malignancies? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Has there been any evidence of any prior malignancies for at least five (5) years? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer				
Did the patient undergo potentially curative therapy for all prior malignancies? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Does the patient have a medial quadrant lesion? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer				
Does the patient have life expectancy of 10 year or more? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer	Is the patient available for follow-up? (\$ 0.00) <input type="radio"/> Yes <input type="radio"/> No <input type="radio"/> Defer				
<input type="button" value="PROCESS"/> Click to submit your answers to the system. <input type="button" value="REVIEW"/> Click to review and change your previous answers					
<input type="button" value="RETURN"/> Click to return to the entry page.					

Figure 12: Status of all Clinical Trials

At any time the user can check the eligibility status of the clinical trials by clicking the drop down box as shown in Figure 12. The system keeps on presenting questions until the eligibility of all the clinical trials for the given patient is known. The user can stop entering data at any point. Data entry for the patient can be reinitiated at any later time by using the search option shown in Figure 8.

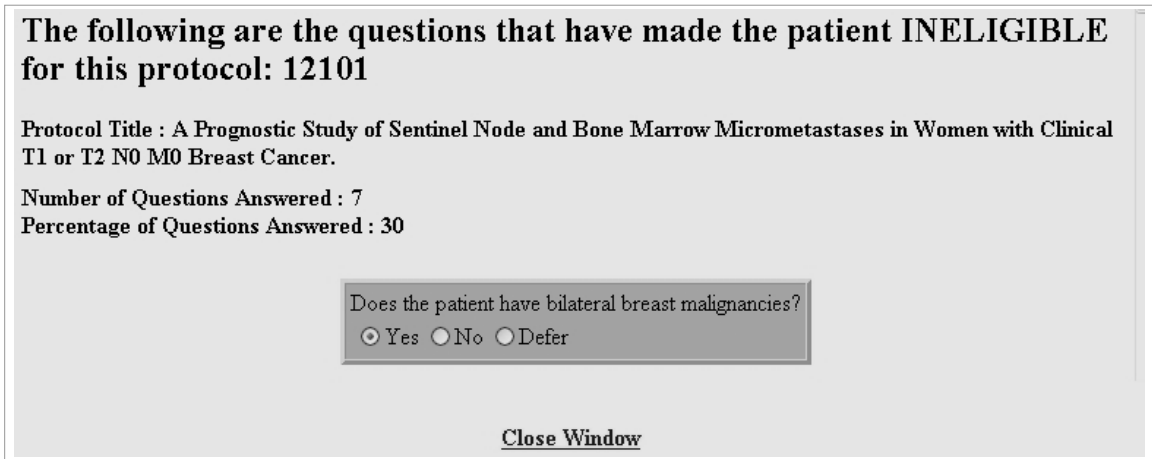


Figure 13: Explanation of a Decision by Explanation Sub-System

At any time you can click on why button and the explanation sub-system will justify the system’s decisions as shown in Figure 13. Figure 14 appears when the patient’s eligibility for all the clinical trials is decided and there is no more information needed.

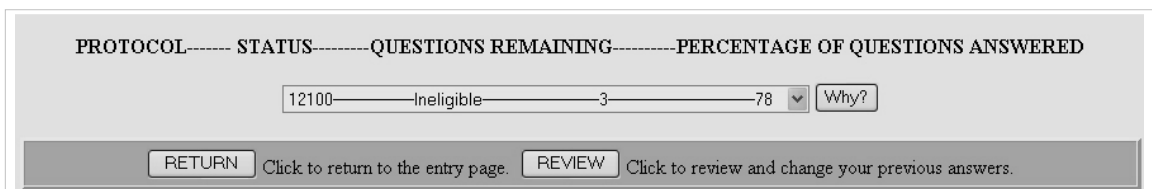


Figure 14: System after Eligibility of all Clinical Trials is Decided

Name: Standard Questions Cost: \$0 <input type="radio"/> Done Before <input checked="" type="radio"/> Done <input type="radio"/> Not Done	Name: CBC with Diff Cost: \$129 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: CBC Cost: \$95 <input checked="" type="radio"/> Done Before <input type="radio"/> Done <input type="radio"/> Not Done	Name: Standard Questions Cost: \$0 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: BMP Cost: \$65 <input checked="" type="radio"/> Done Before <input type="radio"/> Done <input type="radio"/> Not Done	Name: CMP Cost: \$286 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: CT Scan (Head) Cost: \$866 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done	Name: CT Scan (Pelvis) Cost: \$1800 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: CT Scan (Abdomen) Cost: \$2223 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done	Name: CT Scan (Chest) Cost: \$2215 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: Diffusing Lung Capacity Test (DLCO) Cost: \$115 <input type="radio"/> Done Before <input checked="" type="radio"/> Done <input type="radio"/> Not Done	Name: Pulmonary Function Tests DLCO Cost: \$764 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: Pulmonary Function Test FVC Cost: \$764 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done	Name: T-MUGA Cost: \$979 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done
Name: HIV Antibody Test Cost: \$50 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done	Name: T-CBC with Diff Cost: \$129 <input type="radio"/> Done Before <input type="radio"/> Done <input checked="" type="radio"/> Not Done

Click to submit your answers to the system.

Figure 15: Status of all Available Tests for the Patient

At any time you can also check the status of the tests know for the current trials. Figure 15 shows the test status of all the tests for the given patient. There are three statuses for the tests: a) Done before, b) Done, and, c) Not done. “Done before” means that the test was already done before the patient was checked for enrollment in a clinical trial. This gives the correct cost incurred for determining eligibility for the patient as the system will not count the cost of such tests in total cost for the patient. “Done” tests are the ones done to obtain more information to decide a patient’s eligibility for clinical trials. Similarly “Not done” are the tests that are available in the system but are still not done for the patient. At any point the user can change the status of any test.

Chapter 5

Experiments and Results

5.1 Analytical Experiments

Data of patients treated at the Moffitt Cancer Center was used to compare the system's results with that of clinicians. Results from the application of various heuristics are also compared. The first set of experiments were done on retrospective data from 187 patient charts. Table 1 shows the comparison of automated assignment by the system and manual assignment by clinicians at Moffitt Cancer Center. These experiments were started by Kokku [33]. The results are also published in IEEE SMC 2003 conference [14].

Table 1: Results of Matching Retrospective 187 Patients to Clinical Trials

Clinical Trial	Same Matches	New Matches	Missing Data
10822	10	5	0
10840	0	19	3
11072	48	26	19
11378	4	19	3
11992	5	6	0
12100	8	20	13
12101	20	30	0

Table 1 shows the number of patients that were eligible for a particular trial. "Same Matches" column indicates the patients who were found eligible by the system and were also found eligible by the clinicians. The "New Matches" column indicates the patients

who were found eligible by the system but were potentially missed by the clinicians. The “Missing Data” column indicates the patients that did not have enough information to determine their eligibility. All the patients are not tested for all the clinical trials and thus some of them do not have some tests performed that are essential to determine eligibility for a specific clinical trial. Also the data was retrospective, thus the patients were not undergoing treatment at the Moffitt Cancer Center.

Some questions, which were generally consent-based questions, were answered “Yes” while performing the experiments as we had no actual way to determine their answers. Below are a few such questions:

1. Is the patient willing to sign the consent form?
2. Does the patient have at least sixth grade education?
3. Is the patient willing to use contraceptives?

The next set of experiments were done on 169 patients who were getting treatment at Moffitt Cancer Center at the time of experiments. The results are shown in Table 2.

Table 2: Results of Matching Current 169 Patients to Clinical Trials

Clinical Trial	Number of Tested Patients	Same Matches	Missing Info	New Matches
11132	7	4	1	1
11931	169	2	5	26
11971	159	4	0	0
12100	162	0	0	5
12101	166	11	6	52
12385	42	0	0	19
12601	162	0	3	1
12643	63	18	2	34
12757	58	1	4	3
12775	133	23	6	17

Table 3: Detailed Classification of New Matches Found

Clinical Trial	New Matches	New Matches with no enrollment	New Matches with no conflict	New Matches with conflict	% of "Assignable" New Matches over number of checked patients	% of "Actual" Enrollment over number of checked patients
11132	1	0	1	0	14.29%	57.14%
11931	26	15	11	0	15.38%	1.78%
11971	0	0	0	0	0.00%	2.52%
12100	5	4	1	0	3.09%	0.00%
12101	52	33	6	13	23.49%	6.63%
12385	19	9	2	8	26.19%	0.00%
12601	1	1	0	0	0.62%	0.00%
12643	34	24	10	0	53.97%	30.16%
12757	3	2	1	0	5.17%	1.72%
12775	17	11	3	3	10.53%	18.80%

Table 2 is similar to Table 1 which was explained previously. Table 3 categorizes the new matches found by the system. The third column of Table 3 specifies the number of new matches in which the patient was not enrolled in any clinical trial. Thus this number represents the patients who were found eligible for some clinical trial by the system, but were never enrolled in any clinical trial by the clinicians. The fourth column indicates the number of matches in which the patient was assigned to a clinical trial but there was still a clinical trial to which she could be assigned. Many clinical trials are compatible and thus a patient can be in two clinical trials at once. Such clinical trials are called compatible clinical trials. Thus the numbers in fourth column indicates the patients who were assigned to a clinical trial but there was an available compatible clinical trial to which they were not assigned. The fifth column indicates the number of matches found in which the patient was already enrolled in an incompatible trial. Thus this match has no significance as the patient could not have been assigned to the matching trial anyway. The sixth column indicates the percentage of assignable matches found over the number

of patients checked. The number of assignable matches found, is the summation of the third and fourth columns. This indicates the matches in which the patient could have been safely enrolled in a new clinical trial to which the patient was not previously assigned. The last column is the percentage of actual enrollment in that clinical trial. Note that in the majority of protocols, the percentage of new assignable matches found is greater than the actual enrollment.

The results clearly show that the system was able to detect many matches which were potentially missed by clinicians. Given that fact that many clinical trials fail due to under-recruiting [11, 22, 33], the system can potentially play a crucial role in the success of a clinical trial.

5.2 Cost Saving Experiments

In case of lack of information to determine a patient's eligibility for a clinical trial, the clinicians order specific tests to be performed on the patient to obtain relevant information. Research suggests that test costs can be saved if we optimize the test ordering [1, 11, 14]. The system's reordering agent reorders the questions presented by the system to the user. This, in turn, results in reordering of the tests as tests needs to be performed to obtain answers to the questions. Reordering heuristics were designed to optimize the cost as explained in the system design section. Tables 4 and 5 show the cost saving information for the retrospective 183 patients and current 169 patients respectively.

Table 4: Average Dollar Cost Savings for 187 Retrospective Patients

Clinical Trial	Average Dollar Cost	
	Without Test Reordering	With Test Reordering
10822	\$70	\$11
10840	\$0	\$0
11072	\$209	\$60
11378	\$35	\$19
11992	\$0	\$0
12100	\$0	\$0
12101	\$0	\$0

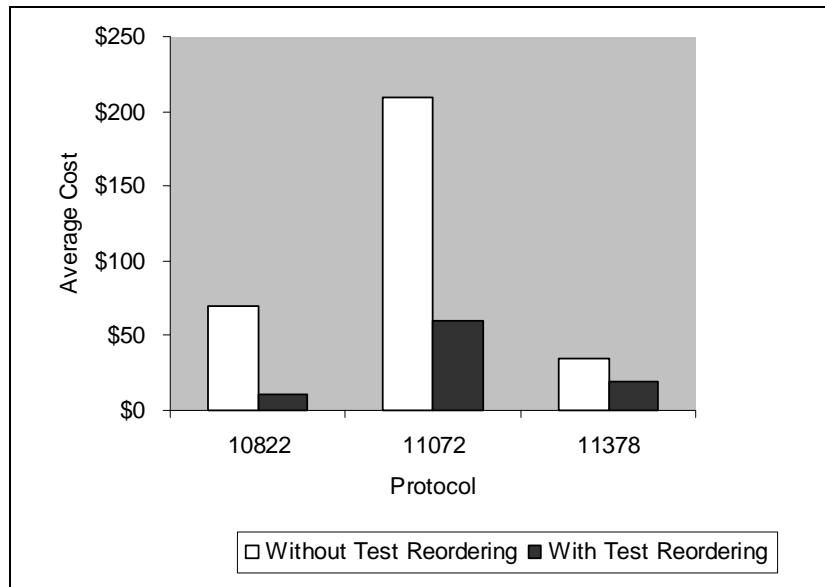


Figure 16: Graph Showing Average Dollar Cost Savings for Retrospective Patients

Table 5: Average Dollar Cost Savings for 169 Current Patients

Clinical Trial	Average Dollar Cost	
	Without Test Reordering	With Test Reordering
11132	\$0	\$0
11931	\$0	\$0
11971	\$192	\$192
12100	\$0	\$0
12101	\$0	\$0
12385	\$0	\$0
12601	\$36	\$3
12643	\$0	\$0
12757	\$107	\$107
12775	\$0	\$0

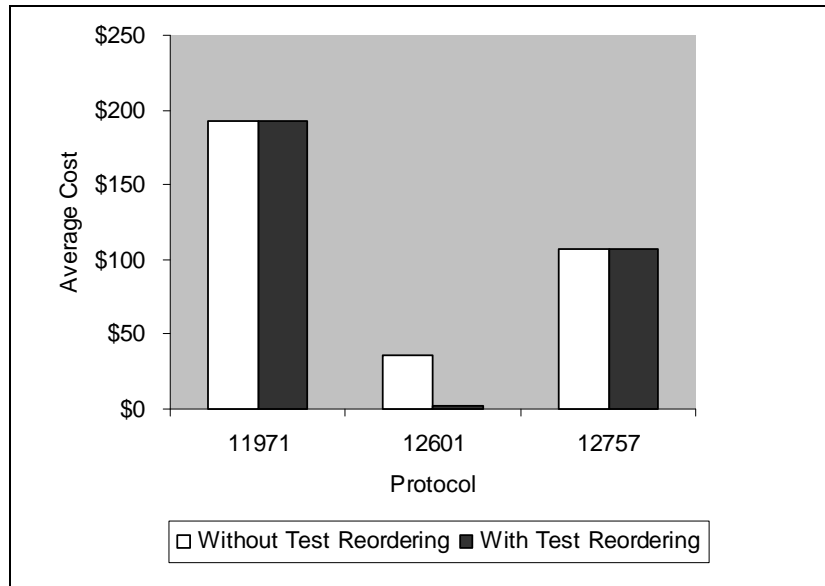


Figure 17: Graph Showing Average Dollar Cost Savings for Current Patients

5.3 Probabilistic Experiments

It was noted that medical test results were highly probabilistic in nature. Most test results can be estimated with reasonable accuracy before the test is performed. Given enough data this knowledge can be used to reorder tests and to optimize data entry needed to determine the eligibility of a patient. Also it was noted that over time, clinicians typically get a good estimate of what information is likely to be most crucial in determining the eligibility of a patient and prefer to get that information first. We developed a probabilistic knowledge base on the same principles. Whenever a question is answered by a user, the system records it. If answering the question made a patient ineligible for a clinical trial, this is recorded too. Thus we get a basic ratio of the number of times a question was asked, to the number of times it made a patient ineligible for a particular protocol. Treating this as a probability we reorder the questions, with the questions having a higher probability of ruling a patient out of a clinical trial being asked first. The

idea is that if a patient is ineligible then the question which is most likely to determine she is ineligible, from past experience, is asked first and this will reduce data entry. This method also closely mimics the method adopted by clinicians and thus can increase acceptability of the system by medical personnel.

The system begins by assigning a 50% probability to all the questions. The probability is modified as the system sees new evidence. To test the effectiveness of this technique we did a ten-fold cross validation. We randomly selected 90% of the patients and their data was used to generate a probabilistic knowledge base for the system. The remaining 10% of the patients were tested using the system which used this probabilistic knowledge base to reorder the questions. This process was repeated ten times with each 10% of the testing patients being unique. Six clinical trials were used in the experiments. These six clinical trials were selected out of about 15 possible clinical trials, as these trials were in “open” status for a long enough duration during the experiments to have an adequate number of patients being tested for those clinical trials. In the probabilistic experiments it was necessary to train the system on an adequate number of patients before it can be used to reorder questions for other patients. We selected 90 patients at random from our list of patients and used their data in experiments. As mentioned above, a ten-fold cross validation was carried out, so that the system was trained on 81 patients and the remaining 9 patients were tested using the system.

Table 6 shows the results of each fold of the cross validation. Table 7 shows results of the cross validation per clinical trial.

Table 6: Individual Test Results for Ten-Fold Cross Validation

Ten-fold cross validation				
Test Number	Average number of questions			Difference %
	Probabilistic System	Analytical System	Difference	
1	16.67	20.75	4.08	19.68
2	15.17	17.00	1.83	10.78
3	15.83	17.58	1.75	9.95
4	15.75	18.25	2.50	13.70
5	13.83	16.67	2.83	17.00
6	15.58	17.75	2.17	12.21
7	15.83	18.25	2.42	13.24
8	15.50	16.83	1.33	7.92
9	16.50	18.50	2.00	10.81
10	15.83	19.17	3.33	17.39
Total	15.65	18.08	2.43	13.42

Table 7: Results of Ten-Fold Cross Validation per Protocol

Ten-fold cross validation				
Protocol	Average number of questions			Difference %
	Probabilistic System	Analytical System	Difference	
11931	15.35	18.90	3.55	18.78
12100	13.85	13.95	0.10	0.72
12101	21.65	24.75	3.10	12.53
12521	14.75	19.05	4.30	22.57
12601	13.90	15.70	1.80	11.46
12777	14.40	16.10	1.70	10.56
Average	15.65	18.08	2.43	13.42

As seen in Table 7 the probabilistic system reduces data entry by 13.42% on average as compared to the analytical system. The average number of questions asked by the system is reduced by 2.43. One important observation is that the probabilistic system **always** asks fewer questions than the analytical system for all the patients tested. Using the t-test, the probabilistic system is statistically significantly better at the 99.99% confidence interval.

Table 8: Standard Deviation for Ten-Fold Cross Validation

Standard deviation		
Protocol	New System	No Cost Version
11931	7.44	6.11
12100	1.39	1.61
12101	7.44	6.11
12521	5.39	6.54
12601	3.67	5.55
12777	5.69	7.26

Also note that as shown in the system’s interface design diagrams, the system asks 14 basic questions after one enters a patient into the system. These questions are general questions which are used in almost all trials, like patient’s age, sex, etc. If the patient’s eligibility is still not determined for all the trials for which her eligibility is tested, the system asks for further information. Thus technically you answer a minimum of 14 questions for a patient as it is very unusual to not to have information about any of those questions. For such patients who get ruled out of a clinical trial after the questions on the initial page, the number of questions asked in both the systems remains the same as there is no reordering of questions. Table 9 shows the average number of questions when we exclude such patients from our experiment. The probabilistic system asks 3.69 fewer questions on an average than analytical system which is 18.53% less than analytical system. It saves up to 33% of the data entry for 12601, which is the maximum.

Table 9: Results Excluding Patients who were Ineligible by Initial Questions

Ten-fold cross validation				
Protocol	Average number of questions			Difference %
	Probabilistic System	Analytical System	Difference	
11931	16.71	21.79	5.07	23.28
12100	14.75	15.25	0.50	3.28
12101	22.05	25.32	3.26	12.89
12521	14.75	19.05	4.30	22.57
12601	14.60	21.80	7.20	33.03
12777	14.42	16.21	1.79	11.04
Average	16.21	19.90	3.69	18.53

5.4 Eligibility Probability

A major disadvantage of using an analytical approach to finding eligibility is that even in the absence of a single piece of evidence, we cannot predict anything about the eligibility of a patient. On the other hand, with a probabilistic approach, we can make some predictions about the eligibility. Papaconstantinou [10] showed in his experiment that by using Bayesian Networks, we can determine a patient’s eligibility in the absence of some evidence. In our rule based system, we tried to use the independent probabilities that we had for each question, to generate a prediction about a patient’s eligibility probability in the absence of evidence. As explained before in the system design section, we assume that all questions are independent, so the eligibility probability for a protocol is the product of the individual probabilities that a question does not rule out a patient, of all the questions of that protocol. As we enter more information in the system, if the patient is ruled ineligible by the system, her eligibility is 0, and if she is eligible then her eligibility probability is 1. If we still need more information, we recalculate the eligibility probability as the product of eligibility probability of all the unanswered questions.

One of the uses of the eligibility probability can be to try to quickly assign a patient to a clinical trial. This can be achieved by generating an initial eligibility probability for all the available trials for that patient. We then check the eligibility of the patient for the trial which has the highest eligibility probability. After every piece of information, the probabilities are regenerated and the system asks for more information about the trial with the highest eligibility probability until the patient is found eligible for a trial or is determined ineligible for all the trials. The system stops seeking further information after the patient is found eligible for a clinical trial as the purpose of these experiments is to find a single matching protocol with least the number of questions being answered. To test the effectiveness of this approach we did a ten-fold cross validation on the available patients using this approach. We used six clinical trials. We compare the results with the analytical system in which we stop answering questions when it finds a clinical trial for which the patient is eligible. Results are shown in Table 10.

Table 10: Ten-Fold Cross-Validation for Heuristic that Uses Eligibility Probability

Ten-fold cross validation				
Test Number	Average number of questions			Difference %
	Probabilistic System	Analytical System	Difference	
1	22.33	28.67	6.33	22.09
2	30.67	34.33	3.67	10.68
3	31.33	24.33	-7.00	-28.77
4	30.00	33.00	3.00	9.09
5	19.67	25.00	5.33	21.33
6	20.67	31.67	11.00	34.74
7	28.33	33.00	4.67	14.14
8	24.00	36.67	12.67	34.55
9	21.67	22.67	1.00	4.41
10	17.67	24.33	6.67	27.40
Average	24.63	29.37	4.73	16.12

As shown in Table 10, using the eligibility probability in reordering heuristic saves 16.12% of data entry. For the third test in ten-fold cross validation, the probabilistic system asks more questions than the analytical system. This is because the protocols for which the patient is more likely to be eligible are tested first in the probabilistic system and the patients were ineligible for that protocol. Thus the number of questions asked increases as compared to the analytical system. Using the t-test, the probabilistic system is statistically significantly better at the 95% confidence interval.

Although the method gave good results, it had an inherent drawback that the probabilities were not normalized and thus the resultant probabilities were a very low number that were inappropriate for providing feedback to users. As discussed in the “System Design” section, we used Bayes probabilities to compute probabilities that were more suitable to give appropriate feedback. The same set of 90 patients that were used in the previous experiments were used to conduct 10-fold cross validation experiments. Figure 11 shows the results.

Table 11: Ten-Fold Cross-Validation for Reordering Agent that Uses the Bayes Method to Compute Eligibility Probability

Ten-fold cross validation				
Test Number	Average number of questions			Difference %
	Probabilistic System	Analytical System	Difference	
1	20.67	28.67	8.00	27.91
2	29.00	34.33	5.33	15.53
3	31.67	24.33	-7.33	-30.14
4	26.33	33.00	6.67	20.20
5	22.33	25.00	2.67	10.67
6	18.67	31.67	13.00	41.05
7	25.67	33.00	7.33	22.22
8	22.67	36.67	14.00	38.18
9	19.33	22.67	3.33	14.71
10	17.33	24.33	7.00	28.77
Average	23.37	29.37	6.00	20.43

The results show that Bayes method of computing eligibility probability and using it in reordering the questions reduces the data entry needed by 20.43 % on an average. Also the probabilities generated were more appropriate for feedback as the probabilities computed were real probabilities unlike the previous method where they were not normalized and thus were very small numbers. Using the t-test, the probabilistic system is statistically significantly better at the 95% confidence interval.

As shown by the results, the reordering agent using Bayes method for computing eligibility probability reduces data entry more than the reordering agent using the previous method. Although the Bayes method is known to provide with more accurate probability, the previous method was tried as it had an important feature that the eligibility probability was dependent on the number of questions in the clinical trial. As the individual probabilities of each question was very small number, the more questions in the eligibility probability, the lesser the eligibility probability gets as it was the product of individual probability of each question. This method is biased towards the clinical

trials with lesser number of questions, which can help optimize data entry. Out of the six clinical trial used in the experiments, clinical trial number 12777 has four questions and clinical trial 12521 have 34 questions. Is such a scenario it might be optimal to check the eligibility of clinical trial 12777 before 12521 even if 12521 has higher rate or enrollment. If 12521 have better rate or enrollment then 12777 then the reordering agent using Bayes method for calculating eligibility probability will estimate a higher eligibility probability for 12521 and ask for more information about that clinical trial. Whereas, the reordering agent that estimates eligibility probability by product of individual questions will most certainly estimate a higher eligibility probability for 12777. The experiment results suggest that when the reordering agent uses Bayes method to estimate eligibility probability, it results in a more data entry efficient agent.

Chapter 6

Conclusions and Future Work

6.1 Conclusions

Recruiting patients to clinical trials is very time and labor intensive work. Many clinical trials fail due to under recruitment. The system presented here has found a substantial number of matches potentially missed by clinicians. Thus the system can play a critical part in the success of a clinical trial. Also the web-based interface of the system makes it possible to have a central system which can be accessed by clinical personnel from any medical institute around the country. All large research centers have clinical trials of their own and it is very hard to exchange the trial information between them as the same trials can be interpreted in a different way by different clinicians. Having an electronic version of the clinical trials encoded using our knowledge entry system can make sharing of clinical trials between different hospitals very convenient and effective. This in turn can increase accrual for clinical trials as the pool of potential participants increases.

The system also effectively reorders the tests and reduces the cost incurred in determining eligibility. We developed the probabilistic agent which accumulates probabilistic knowledge over time. Data entry was successfully optimized by as much as 30% using the probabilistic reordering agent. The testing system makes it possible to

efficiently check the performance of different reordering heuristics. The system also mimics human behavior which may increase its acceptability to medical personnel.

6.2 Future Work

The system was designed to match patients to clinical trials of any kind, but currently only breast cancer clinical trials are being implemented. The system currently reorders to optimize cost and data entry but it also has provisions to incorporate pain measures for each test and reorder to minimize pain for the patients. These options still need to be explored.

The probabilistic reordering agent use independent probabilities of each question and estimate eligibility probability using Bayes rule . Other probabilistic methods need to be explored to try to improve the reordering and save more data entry and costs. In the experiments conducted, the probabilistic agents were trained on 81 patients. The effectiveness of the agents needs to be investigated when different number of patients are used in training.

The system currently is a web-based CGI application. Thus every time the system is accessed, it has to regenerate all the graphs and load parts of knowledge into memory. Currently we are implementing a port-listening version of the system, which keeps all the graphs and needed knowledge in memory until needed. This will eventually improve the responsiveness of the system even when the patient is being tested for a large number of clinical trials at once.

The reordering heuristics either use the probabilistic or the analytical approach. Heuristics that use both approaches together still need to be implemented and tested.

References

- [1] A.H. Wu, “*Reducing the inappropriate utilization of clinical laboratory tests*”, *Conn Med.*, 61, pp. 15-21, 1997.

- [2] Beverly J. Smith and Michael D. McNeely. *The influence of an expert system for test ordering and interpretation on laboratory investigations*. *Clinical Chemistry*, 45(8):1168–1175, 1999.

- [3] Brigitte S´eroussi, Jacques Bouaud, and Eric-Charles Antoine. *Enhancing clinical practice guideline compliance by involving physicians in the decision process*. In Werner Horn, Yuval Shahar, Greger Lindberg, Steen Andreassen, and Jeremy C. Wyatt, editors, *Artificial Intelligence in Medicine*, pages 76–85. Springer-Verlag, Berlin, Germany, 1999.

- [4] Brigitte S´eroussi, Jacques Bouaud, and Eric-Charles Antoine. *Users’ evaluation of ONCODOC, a breast cancer therapeutic guideline delivered at the point of care*. *Journal of the American Medical Informatics Association*, 6(5):384–389, 1999.

- [5] Brigitte S´eroussi, Jacques Bouaud, and Eric-Charles Antoine. *ONCODOC: A successful experiment of computer-supported guideline development and implementation in the treatment of breast cancer*. *Artificial Intelligence in Medicine*, 22(1):43–64, 2001.

- [6] Brigitte S´eroussi, Jacques Bouaud, Eric-Charles Antoine, Laurent Zelek, and Marc Spielmann. *Using ONCODOC as a computer-based eligibility screening system to improve accrual onto breast cancer clinical trials*. In Werner Horn, Yuval Shahar, G. Lindberg, Steen Andreassen, and J. Wyatt, editors, *Artificial Intelligence in Medicine*, pages 413–430. Springer-Verlag, Berlin, Germany, 2001.

- [7] Bruce G. Buchanan and Edward H. Shortliffe. *Rule Based Expert Systems: The mycin Experiments of the Stanford Heuristic Programming Project*. Addison-Wesley, Reading, MS, 1984.

- [8] Bruce G. Buchanan, Georgia L. Sutherland, and Edward A. Feigenbaum. *Heuristic dendral: A program for generating explanatory hypotheses in organic chemistry*. In

Bernard Meltzer and Donald Michie, editors, *Machine Intelligence*, volume 4, pages 209-254. Edinburgh University Press, Edinburgh, United Kingdom, 1969.

- [9] Clancey, W.J. *Intelligent tutoring systems, a tutorial survey*. Universite de L'Etat, Belgium, 1986.

- [10] Constantinos Papaconstantinou, Georgios Theocharous, and Sridhar Mahadevan. *An expert system for assigning patients into clinical trials based on Bayesian networks*. *Journal of Medical Systems*, 22(3):189–202, 1998.

- [11] Cyrus Kotwall, Leo J. Mahoney, Robert E. Myers, and Linda Decoste. *Reasons for non-entry in randomized clinical trials for breast cancer: A single institutional study*. *Journal of Surgical Oncology*, 50:125-129, 1992.

- [12] D. Bareford and A. Hayling. *Inappropriate use of laboratory services: Long term combined approach to modify request patterns*. *British Medical Journal*, 301(6764):1305–1307, 1990.

- [13] Edward H. Shortliffe. *MYCIN: A Rule-Based Computer Program for Advising Physicians Regarding Antimicrobial Therapy Selection*. PhD thesis, Computer Science Department, Stanford University, 1974.

- [14] Eugene Fink, Lawrence O. Hall, Dmitry B. Goldgof, Bhavesh D. Goswami, Matthew Boonstra, and Jeffrey P. Krischer. *Experiments on the automated selection of patients for clinical trials*. In *Proceedings of the IEEE International Conference on Systems, Man, and Cybernetics*, 2003.

- [15] Francisco J. Diez, Jos´e Mira, E. Iturralde, and S. Zubillaga. *DIIVAL, a Bayesian expert system for echocardiography*. *Artificial Intelligence in Medicine*, 10(1):59–73, 1997.

- [16] Franco Perraro, Paolo Rossi, Carlo Liva, Adolfo Bulfoni, G. Ganzini, and Adriano Giustinelli. *Inappropriate emergency test ordering in a general hospital: Preliminary reports*. *Quality Assurance Health Care*, 4:77–81, 1992.

- [17] Henrik Eriksson. *Specification and generation of custom-tailored knowledge-acquisition tools*. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*, volume 1, pages 510–518, 1993.

- [18] Ian R. Morrison, B. A. Schaefer, and Beverly J. Smith. *Knowledge acquisition: The ACQUIRE approach*. In Proceedings of the First Semi-Annual Conference in Policy Making and Knowledge Systems, 1991.
- [19] Jacques Bouaud, Brigitte S'erotoussi, Eric-Charles Antoine, Mary Gozy, David Khayat, and Jean-Francois Boisvieux. *Hyper textual navigation operationalizing generic clinical practice guidelines for patient-specific therapeutic decisions*. Journal of the American Medical Informatics Association, 5(suppl.):488–492, 1998.
- [20] Jacques Bouaud, Brigitte S'erotoussi, Eric-Charles Antoine, Laurent Zelek, and Marc Spielmann. *Reusing ONCODOC, a guideline-based decision support system, across institutions: A successful experiment in sharing medical knowledge*. In Proceedings of the American Medical Informatics Association Annual Symposium, volume 7, 2000.
- [21] Jihie Kim and Yolanda Gil. *Acquiring problem-solving knowledge from end users: Putting interdependency models to the test*. In Proceedings of the Seventeenth National Conference on Artificial Intelligence, pages 223–229, 2000.
- [22] Jim Blythe, Jihie Kim, Surya Ramachandran, and Yolanda Gil. *An integrated environment for knowledge acquisition*. In Proceedings of the International Conference on Intelligent User Interfaces, pages 13–20, 2001.
- [23] John H. Gennari and Madhu Reddy. *Participatory design and an eligibility screening tool*. In Proceedings of the American Medical Informatics Association Annual Fall Symposium, pages 290-294, 2000.
- [24] Joshua Lederberg. *dendral-64: A system for computer construction, enumeration and notation of organic molecules as tree structures and cyclic graphs. Part II*. Technical Report N66-14074, nasa Scientific and Technical Aerospace Reports, 1965.
- [25] Joshua Lederberg. *How dendral was conceived and born*. In Proceeding of the ACM Symposium on the History of Medical Informatics. National Library of Medicine, 1987.
- [26] J. S. Aikins, et al., PUFF: *An Expert System for Interpretation of Pulmonary Function Data*. Computers and Biomedical Research 16 (1983) 199--208.

- [27] Kenneth G. Keppel, Jeffrey N. Percy, and Diane K. Wagener. *Trends in racial and ethnic specific rates for the health status indicators: United States, 1990–98*. Healthy People 2000, Statistical Notes, 23, 2002.
- [28] Lucila Ohno-Machado, Eduardo Parra, Suzanne B. Henry, Samson W. Tu, and Mark A. Musen. *AIDS: A decision-support tool for decreasing physicians' uncertainty regarding patient eligibility for HIV treatment protocols*. In Proceedings of the Seventeenth Annual Symposium on Computer Applications in Medical Care, pages 429–433, 1993.
- [29] Mark A. Musen. *Automated generation of model-based knowledge acquisition tools*. Morgan Kaufmann, San Mateo, CA, 1989.
- [30] Mark A. Musen, Samson W. Tu, Amar K. Das, and Yuval Shahar. *EON: A component based approach to automation of protocol-directed therapy*. Journal of the American Medical Informatics Association, 3(6):367–388, 1996.
- [31] Michael D. McNeely and Beverly J. Smith. *An interactive expert system for the ordering and interpretation of laboratory tests to enhance diagnosis and control utilization*. Canadian Medical Informatics, 2(3):16–19, 1995.
- [32] M. Korver and A. R. Janssens. *Development and validation of HEPAR, an expert system for the diagnosis of disorders of the liver and biliary tract*. Medical Informatics, 16(3):259–270, 1993.
- [33] Princeton K. Kokku, Lawrence O. Hall, Dmitry B. Goldgof, Eugene Fink, and Jeffrey P. Krischer. *A cost-effective agent for clinical trial assignment*. In Proceedings of the IEEE International Conference on Systems, Man, and Cybernetics, 2002.
- [34] Sanjuncta Bhanja, Lynn M. Fletcher, Lawrence O. Hall, Dimtry B. Goldgof, and Jeffrey P. Krischer. *A qualitative expert system for clinical trial assignment*. In Proceedings of the Eleventh International Florida Artificial Intelligence Research Society Conference, pages 84–88, 1998.
- [35] Salim Yusuf, Peter Held, K. K. Teo, and Elizabeth R. Toretzky. *Selection of patients for randomized controlled trials: Implications of wide or narrow eligibility criteria*. Statistics in Medicine, 9:73–86, 1990.

[36] Savvas Nikiforou. Selection of clinical trials: *Knowledge representaton and acquisition*. Master's thesis, Department of Computer Science and Engineering, University of South Florida, 2002.