# Identification of Disease-Treatment Relationship in Bio-Science Text

## A PROJECT REPORT

Submitted by

## Kartheek Anumolu
Reg. No. 10MSE1027

in partial fulfillment for the award of the degree of

Master of Science

in

Software Engineering



# School of Computing Science and Engineering
VIT University
Vandalur - Kelambakkam Road, Chennai - 600 127

May - 2015

## School of Computing Science and Engineering

# DECLARATION

I hereby declare that the project entitled **Identification of Disease-Treatment Relationship in Bio-Science Text** submitted by me to the School of Computing Science and Engineering, VIT Chennai, 600 127 in partial fulfillment of the requirements of the award of the degree of **Master of Science** in **Software Engineering** is a bona-fide record of the work carried out by me under the supervision of **Prof. A Muralidhar**. I further declare that the work reported in this project, has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma of this institute or of any other institute or University.

Place: Chennai                                          Signature of Candidate
Date:                                                      (Kartheek Anumolu)

## School of Computing Science and Engineering

# CERTIFICATE

This is to certify that the report entitled **Identification of Disease-Treatment Relationship in Bio-Science Text** is prepared and submitted by **Kartheek Anumolu (Reg. No. 10MSE1027)** to VIT Chennai, in partial fulfullment of the requirement for the award of the degree of **Master of Science** in **Software Engineering** is a bona-fide record carried out under my guidance. The project fulfills the requirements as per the regulations of this University and in my opinion meets the necessary standards for submission. The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma and the same is certified.

**Guide/Supervisor**                                    **Program Chair**

Name:  Prof. A Muralidhar                    Name:  Dr. N.Maheswari
Date:                                                            Date:

**Examiner**                                              **Examiner**

Name:                                                          Name:
Date:                                                            Date:

(Seal of SCSE)

# Acknowledgement

I would like to extend my gratitude towards Prof. A Muralidhar (School of Computer Science and Engineering), it was under his guidance and supervision that I was able to complete my project with a satisfactory note. I would like to thank Prof. N. Maheswari (Program Manager, MS Software Engineering) and Dr.L.Jeganathan (Professor and Dean, School of Computer Science and Engineering) , VIT University, Chennai for their enormous support and advice which acted as a catalyst in successful completion of my project.

<div align="right">

Kartheek Anumolu
Reg.No.10MSE1027

</div>

# Abstract

People are very much cautious of their health and want to be well informed about the recent advancements in medical domain.There is an enormous increase in number of people searching disease related information over web. Health related information are very sensitive to trust and standardizing the data would increase the truthfulness of the results. There is a need to develop a tool that can assist people in finding the trusted information.One of the good ideas would be to integrate machine learning tools and Natural language processing techniques into medical domain to yield better results and user experience. My application can extract disease-treatment relationships from published Medline abstracts that are really sensible to trust. After extracting the data this application provides the users with an interface to select the disease and find its related treatment and also information about the treatment if the user has no idea regarding it

# Contents

# List of Figures

# Chapter 1

# Introduction

## 1.1   Background

Machine learning is one of the domain that is being widely used at present. Integrating Machine learning in medical domain can result in finding a good application that can be used by general users and medical related people. This can also be used in Medical domain to yield better results. Machine learning is a subfield of artificial intelligence domain that concerns about study of systems that can learn from data. Supervised and unsupervised learning algorithms are used in many of applications like classifying emails and predicting the severity of a disease etc. Some machine learning systems attempt to eliminate the need for human intuition in data analysis, while others adopt a collaborative approach between human and machine. Because of this factor developing machine learning systems has increased in recent years.

Medline has health related journals that contain all the disease information and recent developments in their treatments. Using Medline one can retrieve all the documents related to a disease. But this retrieval lack reliable and faster access to required information because this retrieval displays user all the related documents that can be many in number and sometimes may display documents that are not what the user requires. This may be quite frustrating to the user when the user need only required information alone. Performing relation extraction and using learning algorithms can give the user required information about the query.

## 1.2 Motivation

All the health management systems are evidence based systems and not reliable. But we need better, faster and more reliable access to information. In the medical domain one of the richest and most used source of information is Medline, a database of extensive life science published articles. So we are using Medline data repository because it contains more standardized information about disease and treatment relationships and mostly used by clinicians and scholars. Our application need to extract the disease treatment relation between the from Medline using Pattern Recognition algorithm. Thus the user gets the required information alone which saves his time and improves the quality of the result.
One of the main use of this application is that it can also be used by doctors and medical persons to analyze the treatment for a particular disorder and to proof check his view on a particular disease disorder. It can help users to let them know about the disease they are most likely to suffer from and know about the treatment relations which can be used as first aid or precautionary measures they can take when the doctor is not available immediately. This application can be embedded in a e-commerce website related to health domain. Still there is no such reliable application that can provide standardized solution with the treatment information that user wants.

## 1.3 Objective

The main objective of this project is to prepare useful data from the medical abstracts that can be useful to the people. This data is the disease treatment relationship. The initial data are taken from the abstracts that are taken from Medline.

## 1.4 Challenges

One of the main challenges is here in my work a lot of data should be handled. My methodology invloves a series of data processing phases where the data should be properly managed. Even if a small part of data goes missing then it will in turn reflect the entire data preprocessing phase. Moreover we are handling a lot of data so the time taken to process this data should also be taken into consideration and

thus we should use efficient technique for faster data processing.

## 1.5 Essence of the Approach

The approach that is being used in this work is my own idea which is mainly based on the concept of pattern matching and tagger classification. Here i have prepared my own set of data from the different resources and then used the above mentioned concepts to process the data obtained from the medical repository.

## 1.6 About MedLine

MEDLINE OR MEDLARS Online (Medical Literature Analysis and Retrieval System Online) is a bibliographic database of life sciences and biomedical information. It includes bibliographic information for articles from academic journals covering medicine, nursing,pharmacy, dentistry, veterinary medicine, and health care. MEDLINE also covers much of the literature in biology and biochemistry, as well as fields such as molecular evolution. It is compiled by the United States National Library of Medicine (NLM), MEDLINE is freely available on the Internet and searchable via PubMed and NLM's National Center for Biotechnology Information's Entrez system. MEDLARS (Medical Literature Analysis and Retrieval System) is a computerised biomedical bibliographic retrieval system. It was launched by the National Library of Medicine in 1964 and was the first large scale, computer based, retrospective search service available to the general public.
Since 1879, the National Library of Medicine had published Index Medicus, a monthly guide to medical articles in thousands of journals. The huge volume of bibliographic citations were manually compiled. By 1960 a detailed specification was prepared and by the spring of 1961 a request for proposals was sent out to many companies to develop the system. As a result a computer was developed to run MEDLARS was delivered to the NLM in March 1963. MEDLARS cost 3 million dollars to develop and at the time of its completion in 1964, no other publicly available, fully operational electronic storage and retrieval system of its magnitude existed. In late 1971, an online version called MEDLINE ("MEDLARS Online") became available as a way to do online searching of MEDLARS

from remote medical libraries. This early system covered 239 journals.

The database contains more than 21.6 million records from 5,639 selected publication covering Bio-medicine and health from 1950 to the present. Originally the database covered articles starting from 1965, but this has been enhanced, and records as far back as 1950/51 are now available within the main index. The database is freely accessible on the Internet via the PubMed interface and new citations are added Tuesday through Saturday. For citations added during 1995-2003: about 48 percent are for cited articles published in the U.S., about 88 percent are published in English, and about 76 percent have English abstracts written by authors of the articles. MEDLINE functions as an important resource for biomedical researchers and journal clubs from all over the world. MEDLINE facilitates evidence-based medicine. Most systematic review articles published presently build on extensive searches of MEDLINE to identify articles that might be useful in the review. MEDLINE influences researchers in their choice of journals in which to publish. More than 5,500 biomedical journals are indexed in MEDLINE. New journals are not included automatically or immediately. They are selected based on scientific scope and quality of a journal.

# Chapter 2

# Overview / Literature Review

The goal of machine learning is to build computer systems that can adapt and learn from their experience. Every machine learning algorithm has both a computational aspect how to compute the answer and a statistical aspect (how to ensure that future predictions are accurate). The ultimate aim of my study is to provide base of information technology model called as Health care Information Model (HIM) or frame work that helps the society to identify and disseminates health care information. There are various approaches (Rosario et al., 2004) (T.F. Michael Raj and S. Prasanna) of classifying the Medline abstracts http:structuredabstracts.nlm.nih.gov.

## 2.1   Existing Methods

### 2.1.1   Bag of words

The bag of words representation is commonly used for text classification tasks. It is a representation in which features are chosen among the words that are present in the training data. Because we deal with a feature appears more than once in a sentence, this means that it is important and the frequency value representation will capture thisthe features value will be greater than that of other features. We keep only the words that appeared at least three times in the training collection, contain at least one alphanumeric character, are not part of an English list of stop words [15] and are longer than three characters. Words that have length of two or one character are not considered as features because of two other reasons possible

incorrect tokenization and problems with very short acronyms in the medical domain that could be highly ambiguous (could be an acronym or an abbreviation of a common word).

### 2.1.2   Support Vector Machines(SVM)

SVMs are a new learning method introduced by Vapnik (1995). A Support Vector Machine (SVM) is a concept in statistics and computer science for a set of related supervised learning methods that analyze data and recognize patterns, used for classification and regression analysis. The standard SVM takes a set of input data and predicts, for each given input, which of two possible classes forms the input, making the SVM a non-probabilistic binary linear classifier. They are well founded in terms of computational learning theory and very open to theoretical understanding and analysis. Support vector machines are based on the Structural Risk Minimization principle from computational learning theory. The idea of structural risk minimization is to find a hypothesis h for which we can guarantee the lowest true error. The true error of his the probability that h will make an error on an unseen and randomly selected test example. An upper bound can be used to connect the true error of a hypothesis h with the error of h on the training set and the complexity of H (measured by VCDimension), the hypothesis space containing h (Vapnik, 1995). Support vector machines find the hypothesis h which (approximately) minimizes this bound on the true error by effectively and efficiently controlling the VC-Dimension of H.

### 2.1.3   Decision based models(decision trees)

Decision trees are powerful and popular tools for classification and prediction. Decision trees represent rules, which can be understood by humans and used in knowledge system such as database. Key requirement of this model is:

## 2.1.4 Attribute-value description

Object or case must be expressible in terms of a fixed collection of properties or attributes (e.g., hot, mild, cold).
Predefined classes:
The target function has discrete output values (bollean or multiclass)
Sufficient data:
Enough training cases should be provided to learn the model Decision trees classify instances or examples by starting at the root of the tree and moving through it until a leaf node. It is used in hand crafted models and Is suitable for short texts.
Decision tree is a classifier in the form of a tree structure:
Then a decision node which specifies a test on a single attribute and a Leaf node which indicates the value of the target attribute. There is an Arc/edge which splits an attribute and a path which is a disjunction of test to make the final decision

## 2.1.5 Data Chunking

This involves chunking the data step by step. The process involved here is

1)**Removal of Stop words:**
In general the sentences consists of words like (who, they, them) etc so here we try to remove all the stop words. A list of stop words are present thus we remove them

2)**Removal of Prepositions and Conjunctions:**
Now after the removal of stop words then we try to remove the prepositions that are present in the data. Thus this step involves the removal of all prepositions

3)**Removal of other words:**
Now after the step 2 the preprocessed data consists of proper nouns, adjectives and bio-medical words and other remaining data like dates numbers etc. Thus we are struck here in this step and it takes more number of processing methods to slowly extract the biomedical words. Thus this is a very long time taking process and also the efficiency is also very poor.

## 2.2   My Methodology

**Pattern Matching and Taggers**:

In my methodology i have taken the medical related words and then pattern matched them with the abstracts that are obtained from the MedLine. This pattern matching process is done in different steps. Thus the data is trained with a different set of medical related words and thus when ever a word is matched then its particular tag is inserted. Similarly for different data taggers are inserted. Then after the taggers are inserted then based on the taggers found then a particular taggers are inserted at the end of the line. Thus after steps of preprocessing we get the processed data.

Now after obtaining the processed data now using a java program i can extract the required data based on the line taggers and word taggers. Then the data obtained from the taggers is stored into a local database. Thus we get the required related disease treatment relationship.
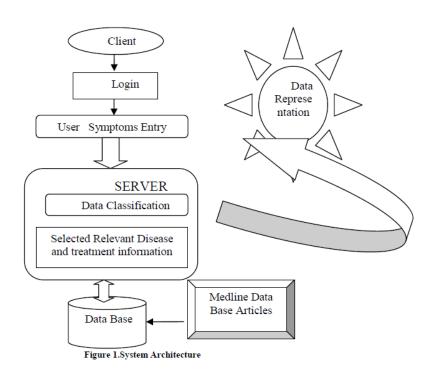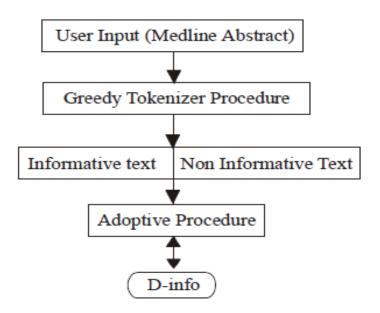
# Chapter 3

# System Design

## 3.1 System Architecture



Figure 3.1: System architecture

## 3.2   Process Flowchart



Figure 3.2: Process Flowchart
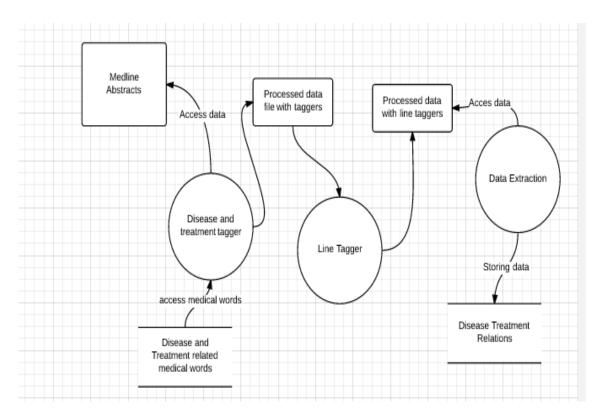
## 3.3   Data Flow Diagram



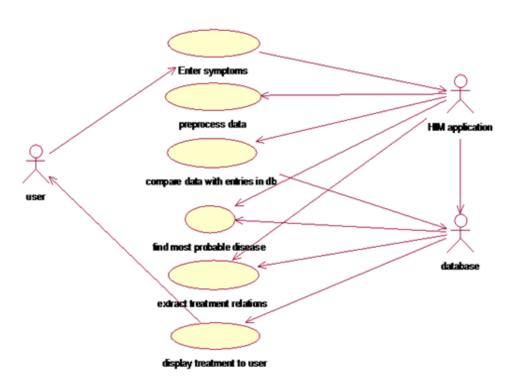Figure 3.3: Data Flow diagram

## 3.4 Usecase Diagram
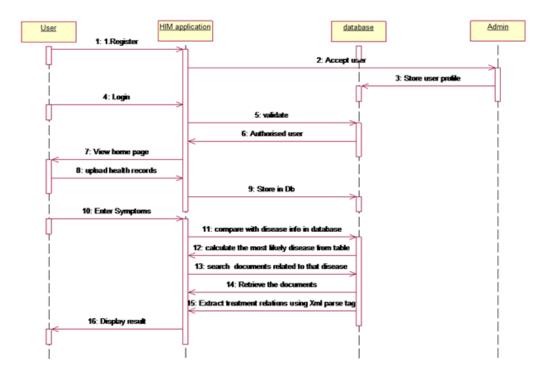


Figure 3.4: Usecase Diagram

## 3.5   Sequence Diagram



Figure 3.5: Sequence Diagram

## 3.6   Requirements/Pre-requisites

### 3.6.1   Software:

**Languages:** JAVA, php, SQL

**Server:** XXAMP

**Database**: phpmyadmin

**Tools Used:** Eclipse, Editplus

### 3.6.2   Hardware:

**Processor:** Intel

**Hard Disk:** More than 100GB

**RAM:** Greater than 512MB

# Chapter 4

# Implementation of System/ Methodology

MedLine abstracts are in the form of paragraphs. These paragraphs consists information about a particular kind of disease or about the occurrence or something about a disease. But these paragraphs also consists of treatment to a particular disease. In some cases only diseases are mentioned but there may or may not be a treatment to that particular disease. So our aim is to identify the disease and its related treatment.

How do we identify that a particular word from a paragraph is a disease or some other bio-medical term. The methodology used by me is Pattern Matching and Taggers. Thus through a repeated process of data training the processed data is obtained. Thus the modules present in this are:

# 4.1   Module 1

Here I have taken all the disease related words from the medical dictionary and put them into a file. There were nearly 5000 disease related words that are present

```
        ---+----1----+----2----+----3----+----4----+----5----+----6----+----7----+----8----+----9----+----0----+----1----+----2-
4206  Tobacco, Chewing
4207  Toddlers, Sleep
4208  Toe, Broken
4209  Toenail Fungus
4210  Toenails, Ingrown
4211  Toilet Substitutes for Incontinence
4212  Tomography, Computerized Axial
4213  Tongue Cancer
4214  Tongue Problems
4215  Tonic Contractions
4216  Tonic Seizure
4217  Tonic Spasms
4218  Tonic-Clonic Seizure
4219  Tonometry
4220  Tonsillectomy
4221  Tonsils
4222  Tonsils and Adenoids
4223  Tooth Damage
4224  Tooth Decay
4225  Tooth Pain
4226  Tooth, Infected
4227  Toothache Overview
4228  Toothpastes
4229  Torn ACL
4230  Torn Meniscus
4231  Tornadoes
4232  Torsion Dystonia
4233  Torsion, Testicle
4234  Torticollis
4235  Total Abdominal Hysterectomy
4236  Total Hip Replacement
4237  Total Knee Replacement
4238  Totipotent Stem Cells
4239  Tounge Thrusting
4240  Tourette Syndrome
4241  Toxemia
```
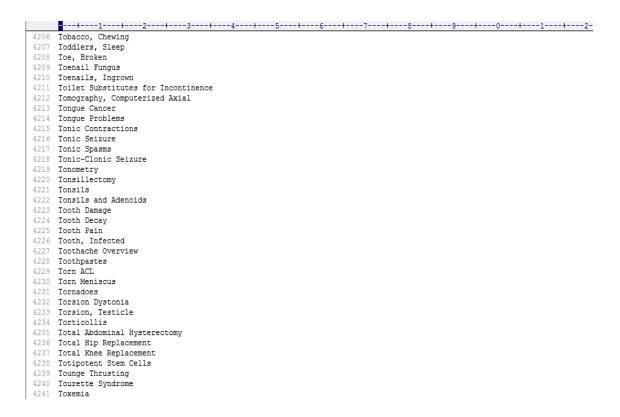
Figure 4.1: Disease related medical words

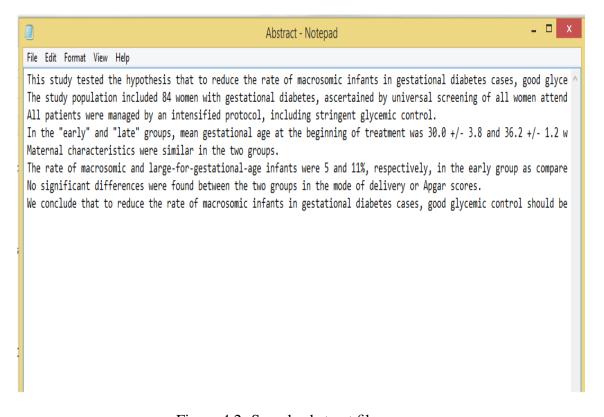Now take an abstract file. This is a sample abstract file



Figure 4.2: Sample abstract file

Now I have taken each word and compared it each word in a single abstract and if the word is present then I have put a tag to it namely <DIS >. For example if the word in from our list was diabetes and if it was found in the abstract, then it would be made like this <DIS >diabetes </DIS >.
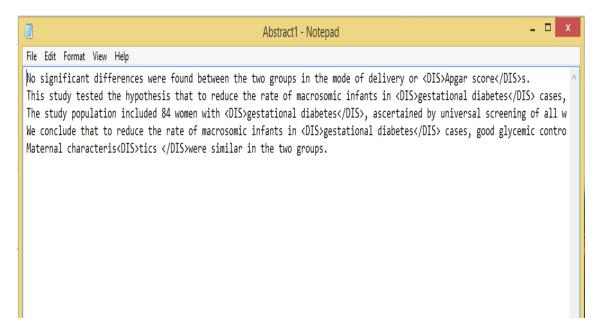


Figure 4.3: After Disease tagger

## 4.2   Module 2

So similarly for the treatment also I have prepared a list of words that are taken from the medical dictionary and similarly the word is identified in the abstracts and then the tag is inserted as <TREAT >.
List of Treatments

```
    ----+----1----+----2----+----3----+----4----+----5----+----6----+----7----+----8----
132  Corneal transplant
133  Corticosteroid injections
134  Corticosteroids
135  Cortisone
136  Coumadin
137  Counseling
138  Counselling
139  COX-2 inhibitors
140  Crutches
141  Cryosurgery
142  Curettage
143  Cyclophosphamide
144  Cyclosporine
145  Dairy-avoidance diet
146  Decongestants
147  Dermabrasion
148  Diabetes pills
149  Diabetic blood sugar control
150  Diazepam
151  Diet changes
152  Dieting
153  Digitalis
154  Digoxin
155  Diuretics
156  Dopamine agonists
157  Doxycycline
158  DTP vaccination
159  Ear surgery
160  Education
161  Elastic stockings
162  Electrical stimulation
163  Electroconvulsive therapy
```

Figure 4.4: Treatment List

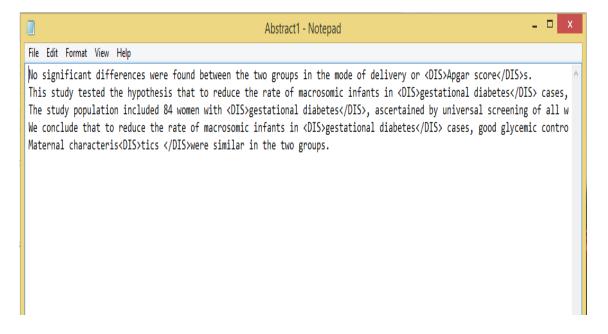Input for the second module



Figure 4.5: Input Abstract

Data after processing the treatment

ween the two groups in the mode of delivery or <DIS>Apgar score</DIS>s.

reduce the rate of macrosomic infants in <DIS>gestational diabetes</DIS> cases, <TREAT>good glycemi

th <DIS>gestational diabetes</DIS>, ascertained by universal screening of all women attending the a

crosomic infants in <DIS>gestational diabetes</DIS> cases, <TREAT>good glycemic control</TREAT> sho

similar in the two groups.

Figure 4.6: After Treatment tagger

## 4.3  Module 3

We now take the consized data that is only the lines that consists of the tags. Now how to we know whether treatment is for a particular disease. Now we take the processed data that consists of tags and then read each line.

Now if a line contains only <TREAT >then at the end of the line we give it as TREATONLY

Physico-chemical studies on the stability of <TREAT> penicillin salts </TREAT> .||TREATONLY

Amino acid metabolism and its relation to the biosynthesis of <TREAT> penicillin </TREAT> .||

The contribution of <TREAT> thoracic surgery </TREAT> to our discipline .||TREATONLY

The scope of <TREAT> sphincterotomy </TREAT> in <TREAT> biliary and pancreatic surgery </TRE

<TREAT> Sarcoidosis </TREAT> involving the veriform appendix .||TREATONLY

Figure 4.7: Lines with only Treatment

Similarly if the line contains only $<$DIS $>$then we give it as $DIS_ONLY$

Dextran of low molecular weight in <DIS> peripheral arterial insufficiency </DIS> ||DISONLY

Clinico-statistical considerations on <DIS> vescicular mole </DIS> ||DISONLY

On the significance of histologic findings in <DIS> neoplasms </DIS> of the trophoblast ( case contribu

Diagnostic evaluation of the patient with <DIS> high blood pressure </DIS> . ||DISONLY

Epidemiology of <DIS> hypertension </DIS> . ||DISONLY

<DIS> Renal vascular hypertension </DIS> ; diagnosis and treatment . ||DISONLY

The use of radioactive isotopes in the diagnosis of <DIS> hypertension </DIS> . ||DISONLY

The angiogram in the study of <DIS> hypertension </DIS> . ||DISONLY

Figure 4.8: Lines with only Disease

Whereas if the line consists of both <DIS >and <TREAT >then at the end of the line we give it as TREATMENT

for <DIS> recurrent spontaneous abortion </DIS> . || TREAT_FOR_DIS

or treatment of <DIS> chronic stable angina pectoris </DIS> . || TREAT_FOR_DIS

s </TREAT> for treatment of <DIS> osteoarthritis </DIS> of the knee . || TREAT_FOR_DIS

Figure 4.9: Disease and its Treatment

After preparing the dataset with tag elements and the classification we now extract the treatment to a particular disease by taking the line with ——TREATMENT and then extract the data between the tag parts and store them into a database. Later we use this database to display to the user a treatment to a particular disease After displaying the treatment we also provide the user with an idea or description of the treatment.

# Chapter 5

# Results and Discussions

## 5.1   Technical Results

The step by step process involved here is pattern matching and the taggers part. First the data from the medical abstracts is taken and then compared with the diseases list and then taggers part is done. Similarly it is also compared with the treatment list. Then from that another line tagger part is done. Then a database is obtained from the data with in the tagger part and thus this data is stored in the database. Now i have linked the database to the interface part which is useful to the people where the people can select the disease and then get its related treatment.

Figure 5.1: Disease Tagger

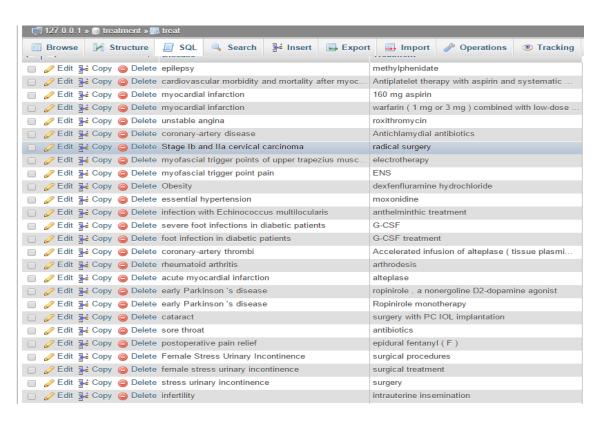Figure 5.2: Writing into the database

Figure 5.3: Database

## 5.2   Performance Analysis



```
                    }
                    // writer.println(text);
            }
            br.close();
        }
        System.out.println(count);
        for (i = 1; i <= count; i++) {
            System.out.println(start[i]);
            System.out.println(end[i]);


        }
        writer.close();
        long y=System.currentTimeMillis();
        System.out.println("The starting time is");
        System.out.println(x);
```

```
Problems  @ Javadoc  Declaration  Console ⊠  Debug  Search
<terminated> DiseaseMatch [Java Application] C:\Program Files (x86)\Java\jre7\bin\javaw.exe (10-May-2015 12:29:24 AM)
82
102
44
64
61
81
20
25
The starting time is
1431197964806
The ending time is
1431197965259
The time taken to compare a single abstract      453
```

Figure 5.4: Abstract Performance

Figure 5.5: Data insertion

# Chapter 6

# Conclusion and Future Work

This project mainly emphasizes on the data preprocessing of the data that is taken from the MedLine and this is what the main objective of my work. This project also provides us with a database of the disease treatment relationships and thus we use this database in the website which is designed as an interface for the people. Here the people can select the particular disease and thus get its related treatment. The methodology used in this work is simple and easily understood by everyone and also the this methodology helped in obtaining the desired results. In this project i have also included the description or information about a particular treatment just to provide the users with some basic knowledge of what the treatment involves.

Coming to the future work, this is a medical based project and there is always scope for more improvement. In this project i have included only information for a finite number of treatments which can be more increased. Moreover as i have processed the disease treatment relationship we can also develop a new kind of system which takes the users symptoms and then finds out the occurrence of a disease and then its corresponding treatment. Moreover this work of mine includes extracting treatments but these may be a treatment or a prevention thus the next challenge would be to identify what it is whether a treatment or a prevention or a side effect.

# Appendices

## Sample Code

```java
import java.io.*;
import java.util.*;
import java.util.regex.Matcher;
import java.util.regex.Pattern;
public class DiseaseMatch
public static void main(String[] args) throws Exception
String text = "";
int[] start = new int[20];
int[] end = new int[20];
int i = 1;
int count = 0;
PrintWriter writer = new PrintWriter(("C:/Users/Kartheek/Desktop/review-1/pendrive/Abstract2.txt")
BufferedReader brr = new BufferedReader(new FileReader "C:/Users/Kartheek/Desktop/review-
1/pendrive/DisList.txt"));
String dis = "";
long x=System.currentTimeMillis();
System.out.println(System.currentTimeMillis());
while ((dis = brr.readLine()) != null)
System.out.println(dis);
BufferedReader br = new BufferedReader(new FileReader("C:/Users/Kartheek/Desktop/review-
1/pendrive/Abstract.txt"));
while ((text = br.readLine()) != null)
String patternString = dis;
Pattern pattern = Pattern.compile(patternString,Pattern.CASEINSENSITIVE);
Matcher matcher = pattern.matcher(text);
```

```
while (matcher.find())
count++;
System.out.println("found: " + count + " : "+ matcher.start() + " - " + matcher.end());
start[i] = matcher.start();
end[i] = matcher.end();
String sub1 = text.substring(matcher.start(), matcher.end());
System.out.println(sub1);
String newLine = text.substring(0, start[i]) + "¡DIS¿"+ text.substring(start[i], end[i])
+ "¡/DIS¿"+ text.substring(end[i], text.length());
String newLine2 = text.substring(end[i], text.length());
System.out.println(newLine);
writer.println(newLine);
i++;
System.out.println(count);
for (i = 1; i ¡= count; i++)
System.out.println(start[i]);
System.out.println(end[i]);
writer.close();
long y=System.currentTimeMillis();
System.out.println(System.currentTimeMillis());
System.out.println(y-x);
```

Appendices are provided to give supplementary information, which is not included in the main text may serve as a separate part contributing to main theme.

# Bibliography

[1] Oana Frunza, Diana Inkpen and Thomas Tran
*A Machine Learning Approach for Identifying Disease-Treatment Relations in Short Texts" vol. 23, 2011.*

[2] B. Rosario and M.A. Hearst
*Semantic Relations in Bioscience Text, Proc. 42nd Ann. Meeting on Assoc. for Computational Linguistics, vol. 430, 2004*

[3] Mark Craven
*Learning to Extract Relations from Medline, Carnegie Mellon University*

[4] R. Bunescu and R. Mooney
*A Shortest Path Dependency Kernel for Relation Extraction, Proc. Conf. Human Language Technology and Empirical Methods in EMNLP), pp. 724-731, 2005*

[5] T.Sakthimurugan , S.Poonkuzhali
*An Effective Retrieval of Medical Records using Data Mining Techniques*

[6] A.M. Cohen and W.R. Hersh, and R.T. Bhupatiraju
*Feature Generation, Feature Selection, Classifiers, and Conceptual Drift for Biomedical Document Triage, Proc. 13th Text Retrieval Conf. 2004*

[7] http://www.biomedcentral.com/1471-2105/5/146

[8] http://www.nlm.nih.gov/bsd/licensee/2013stats/baselinemedfilecount.html