

Structured Statistical Models of Inductive Reasoning: An Application

Akash Pahariya & Vivek Agarwal

November 13, 2010

Abstract

Common concepts such as animals, cities et cetera confer to highly structured arrangements in their respective spaces. We as human beings continuously construct and improve these structures based on our daily activities and experiences. The construction of these structures can be explained by the probabilistic inference about the relations and similarities between different entities of a domain.

In this project we have tried to form an efficient organizational structure for some of the diseases based the extent to which they show similar symptoms. We have attempted to understand how a doctor decides on what disease an individual is suffering form based on the already existing mental structures. In this attempt we also propose a simple model which could be further used to create **Doctor-Bots** with artificial Intelligence.

We have used the model of structure learning specifically 'taxonomic model', described in the paper ('Structured Statistical Models of Inductive Reasoning', Kemp & Tenenbaum).

The code for structure learning is available on:

<http://www.psy.cmu.edu/~ckemp/code/structstatmodels.html>

The code utilized for this project in specific with all relevant changes and corrections is available on:

<http://home.iitk.ac.in/~agvivek/se367/project/code.zip>

1 Introduction

"Going beyond given information", there are two fundamental ways in which this can be achieved viz. **Deductively** or **Inductively**. In Deductive inferences, we draw conclusions that were not explicitly stated but can be derived from the given information. Whilst, in Inductive inferences we conclude in a systematic manner which is one of the most likely outcomes in the given context.

In Inductive Inference, the authors[1] view can be developed as a two step process; a combination of Statistical Inference & Structured Learning. Both of these models have previously been implemented individually, but their combination is an intuitive extension, since they tend to complement each other. Structured Learning is insensitive to noise, however brings out the underlying core efficiently. On the other hand Statistical Inference is sensitive to noise however fails to bring out the underlying essence.

The approach can thus be summarized as:

- (i) Generation of a prior structure which captures the relevant relationships between the categories in the domain;
- (ii) Identification of a stochastic process that captures the knowledge about how these properties tend to distribute themselves.

For Example: A natural representation for animals in the biological context(evolutionary) is a hierarchal tree structure where they are located at the leaves of the tree. Also the stochastic property is that of "Diffusion" as it ensures the symmetry of two way relationship between different animals is maintained. It also acknowledges the fact that the farther the animals lesser the similarities between them.

2 Kemp & Tenenbaum, 2009 (Appendix C): A Summary

A tree can be represented as a pair (ST, SL) where ST refers to the topology of the tree and SL is a vector specifying the lengths of its branches.

Our task is to find ST & SL in such a way so that we can maximize

$$p(ST, SL | D, T, F) = p(D | ST, SL, T) * p(ST, SL | F)$$

where,

- T** is the diffusion process,
- F** is the unknown structure,
- S** is the structures(tree) &
- D** is the data set

To achieve this we first identify the best topology ST and then search for the best set of branch lengths.

- (i) We employ a *greedy search* to identify the topology which includes first starting with the most basic topology in which every entity is connected through an edge to single node. Every time after witnessing some training data we try to

split one of the current nodes. And after each split we rearrange some of the nodes in order to improve the posterior $p(ST | D, F, T)$. This process continues till we get a best structure.

(ii) After the above step, a *gradient based search* is used to find the SL, which maximize $p(SL | ST, D, T)$

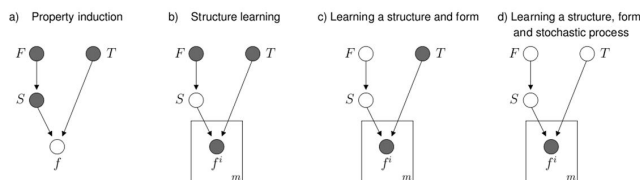


Figure One: Courtesy, Kemp & Tenenmaun, 2009

3 Project Description & Challenges

The motivation behind choosing this project has been to gain a deeper insight into the "Structured Statistical Model", which offers a comprehensive blend of the Statistical Approach along with the Structured Approach.

As a part of the project, we have attempted to understand & recreate the mental structure of how a doctor thinks, and based on all inputs(symptoms) conclude what the problem(disease) is. This has also opened a way to creating **Doctor-Bot** with **artificial intelligence** to diagnose & recommend corrective action.

To achieve this we have employed the taxonomic model; since this best captures the underlying relationship between diseases.

4 iDoctor

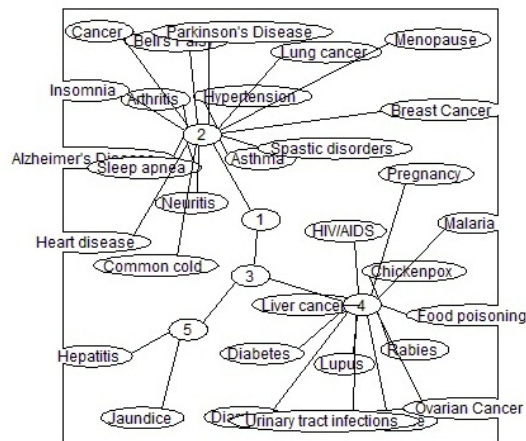
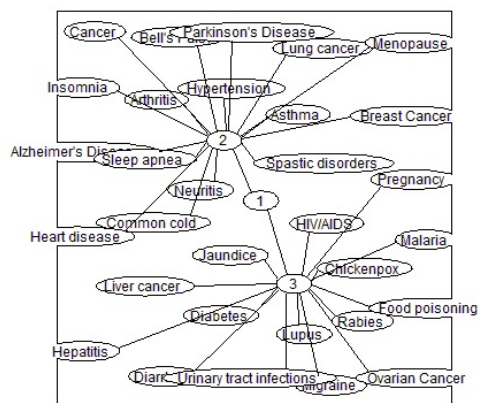
iDoctor works out of a basic set of 31 disease-81 symptom database. In its current format the user inputs his/her symptoms on a scale of 0 to 1. After Execution of the model a tree-leaf structure is obtained. By locating the disease in a particular branch; one can conclude the kind of disease one might be infected with.

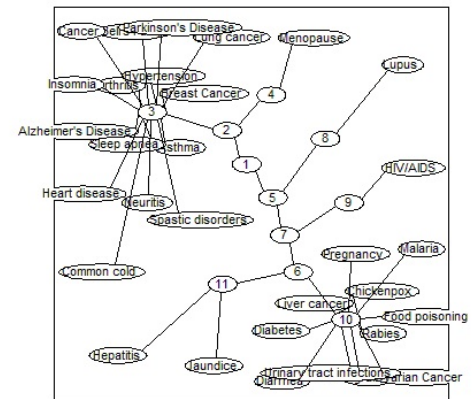
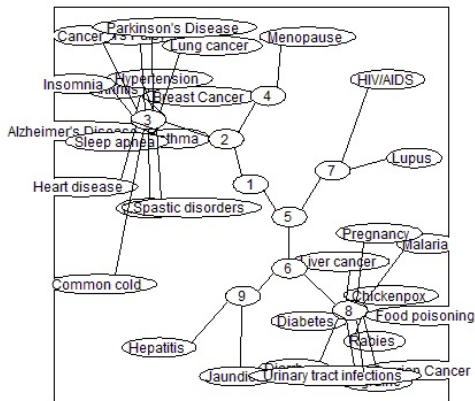
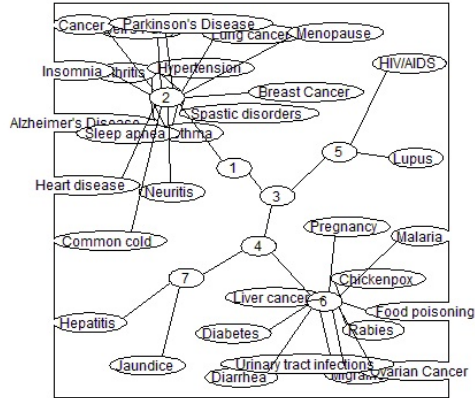
We have obtained the "Master Data Set" based on the inputs of three doctors. Each of whose data set is appended in **Appendix A**. All three of the data sets where then combined to create a "graded" feature set. The model then uses this as an input to try & understand how the diseases may be related & plots a graph.

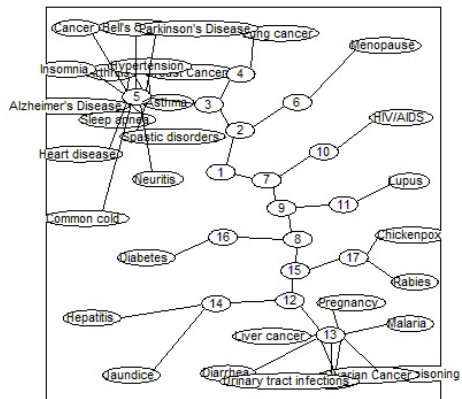
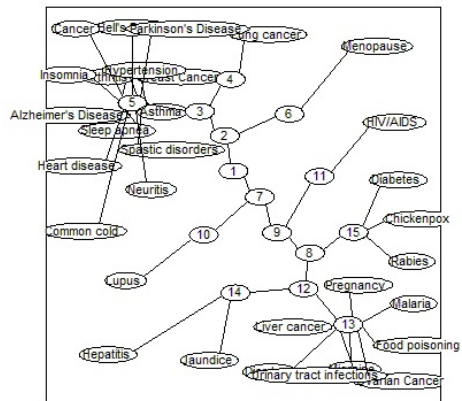
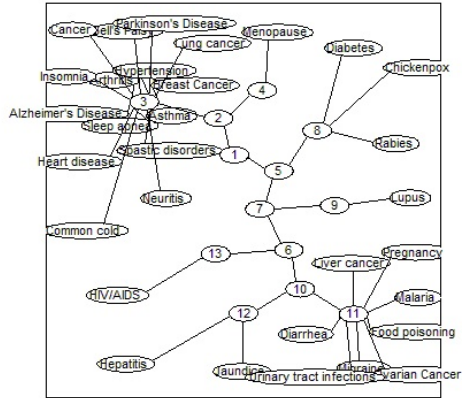
To diagnose a particular disease we request an individual to input symptoms on a scale of 0 to 1 into the dataset. The model is run again, the location of the Disease to be found is located on a branch. It is then concluded that the disease is most likely to be like the members of its branch. A close examination with the original tree helps concluding the disease.

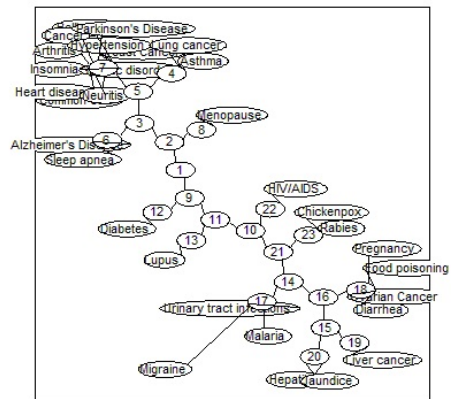
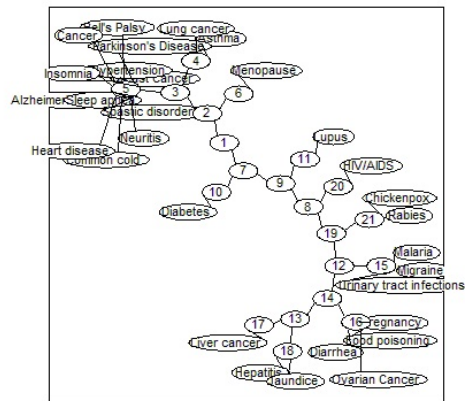
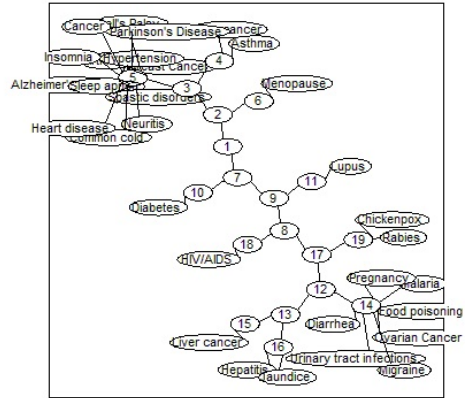
5 Results

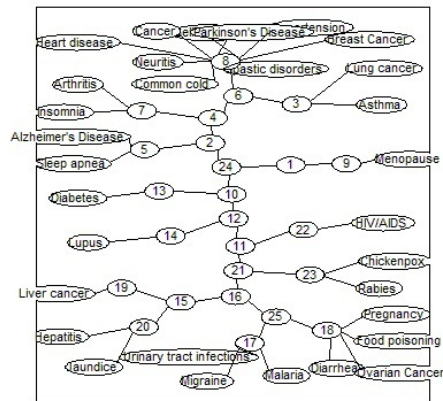
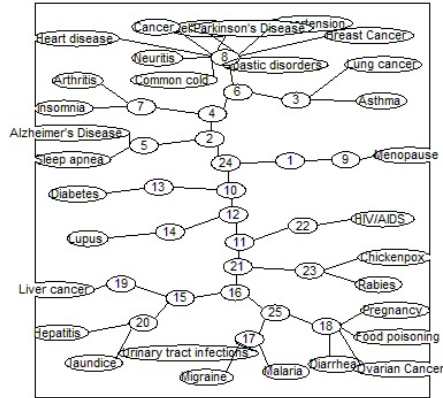
Initial Dataset







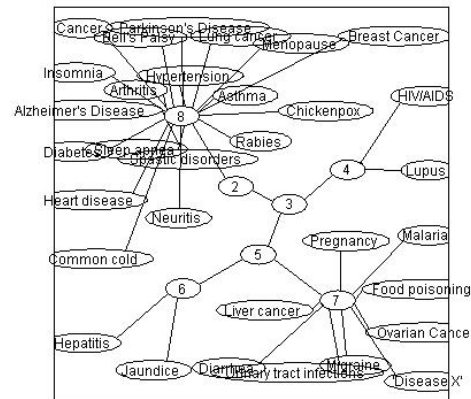
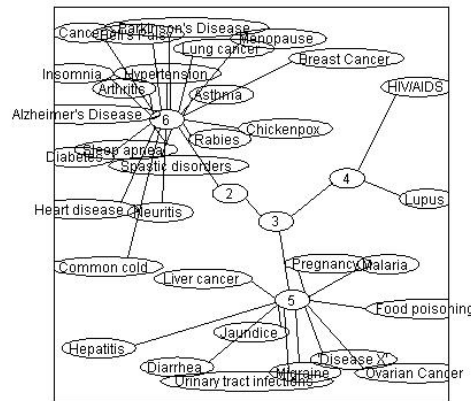
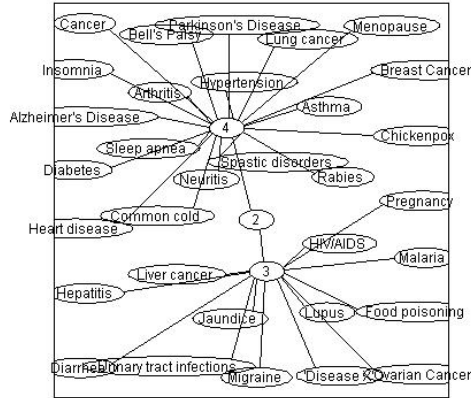


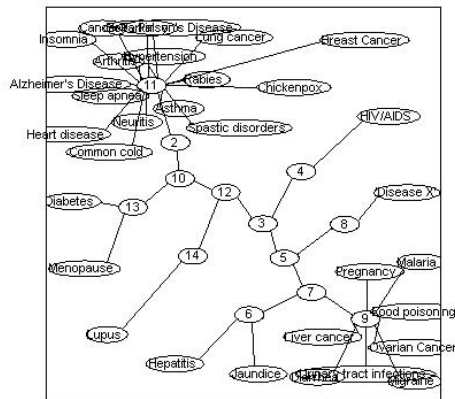
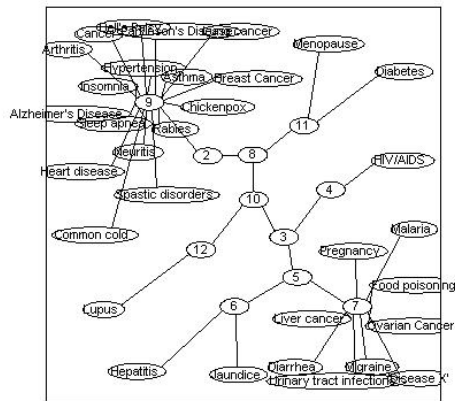
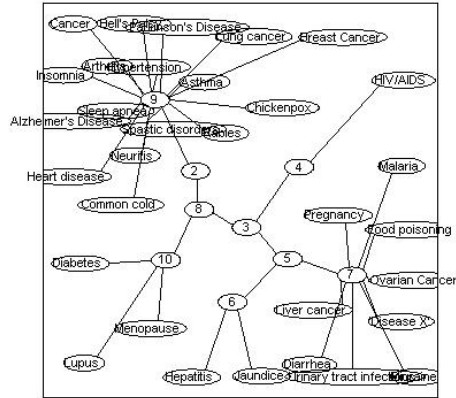


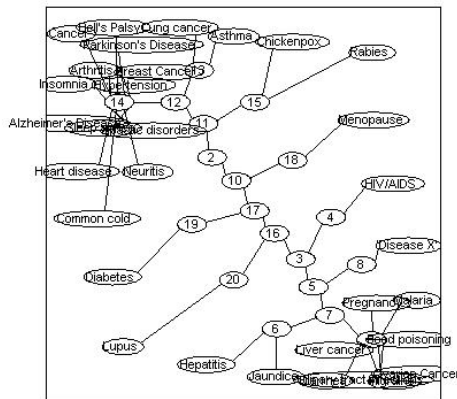
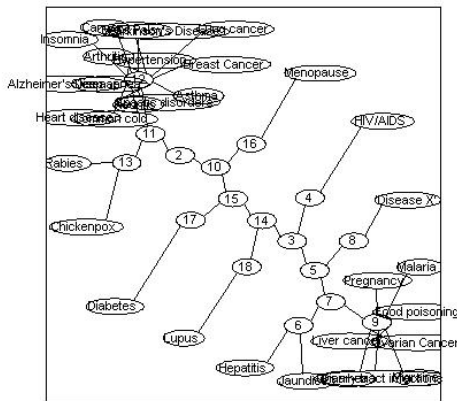
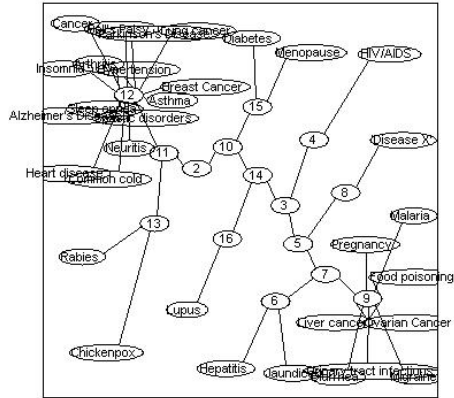
For Detection

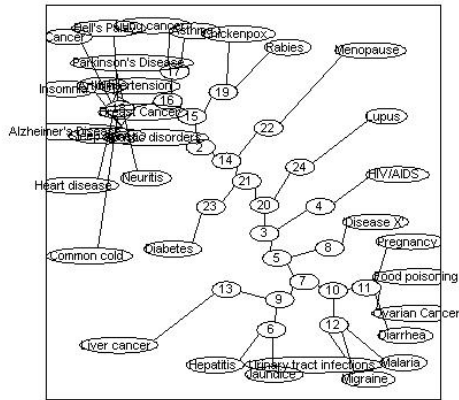
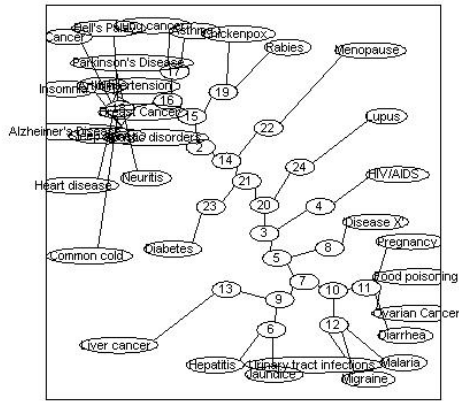
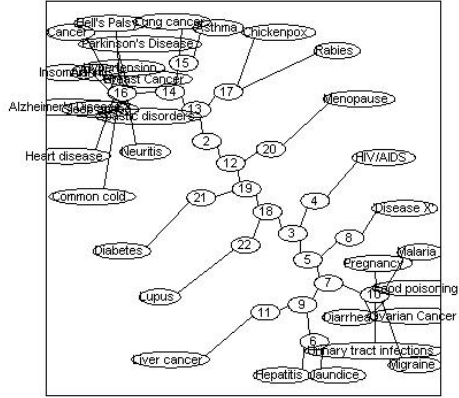
Disease X: had the following Symptoms:

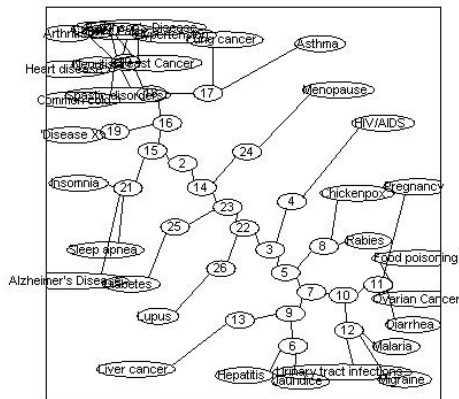
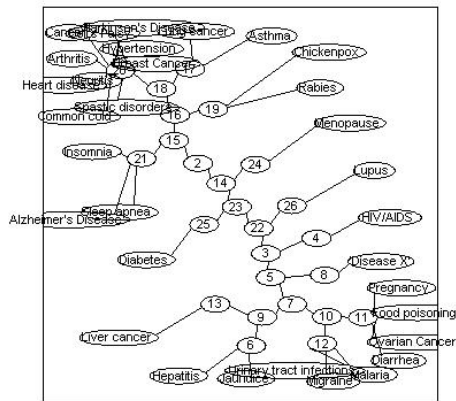
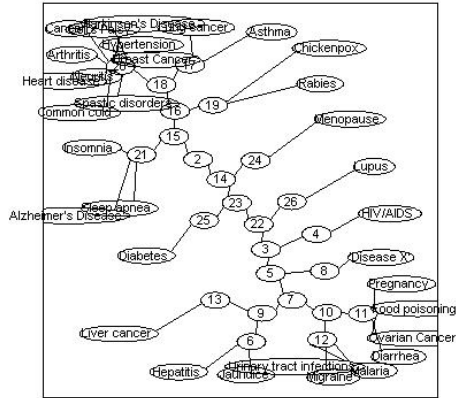
Abdominal Symptoms (Bloating, Pain, Cramps, Discomfort), Anemia, Dark Colored, Urine, Diarrhea, Dizziness, Facial Symptoms (Expression Changes, Flushing, Muscle Spasm, Pain, Swelling, Tingling), Knee Symptoms, Loose stool, Muscle Syptoms (Stiffness, Weakness, Spasms), Palpitations, Sexual problems, Shortness of breath, Skin Symptoms, Swelling, Vomiting, Wheezing, Yellow Skin

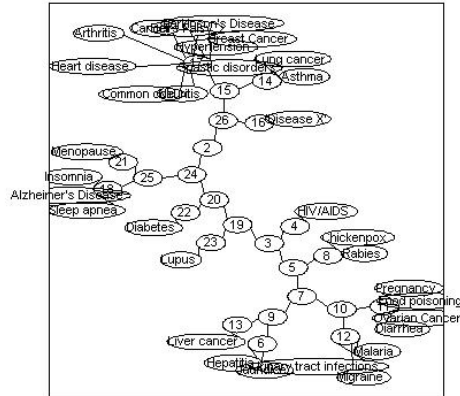












6 Inference

Based on the results the following conclusions can be drawn:

(i) Most of the Diseases tend to pair up well. Some of them pair up based on their origin like Hepatitis, Jaundice. There are also pairs which are disconnected Pregnancy, Food Poisoning, Ovarian Cancer, Diarrhoea. Of these Pregnancy & Ovarian Cancer and Food Poisoning & Diarrhoea are connected in pairs but not as a whole. This can be attributed to the limitations of the dataset & can be bettered with more noise free dataset.

(ii) The method allows to detect a disease, and potentially can also be used to identify new diseases. In the above case, the disease is an outlier and can be only traced to closest neighbours for some hints. Which suggests that this is a new disease.

7 Future Work

In our modelling we have only been able to utilize the structured approach to the fullest, the statistical approach that should have been used to arrive at the working database has been used in a very limited fashion. With more **organized** inputs one can use the *"Bayesian Inference Engine"* to build the *"master database"*. This will help us in exploring the potential of the model to the fullest. In the current version of *iDoctor* a very limited & rudimentary data set has been employed. Detail & correctness of this data set determines the accuracy of the model. One of the many & essential improvements would be to reorganize the symptoms in such a fashion that first a major branch of disease can be arrived at and then further detailed analysis can be carried out.

This will not only save computation time, but also help in looking at greater details. Very certainly a lot can also be done with the interface.

8 Conclusion

The "Structured Statistical Approach" has myriad applications & immense potential to unlock many underlying phenomena. In our humble attempt we have been able to model this approach in part by implementation of the structured approach in complete & the statistical approach in part we have gained a better understanding of the "Structured Statistical Approach".

9 Acknowledgement

(i) Dr. Amitabh Mukherjee
Department of Computer Science & Engineering,
Indian Institute of Technology, Kanpur

(ii) Dr. Divya Balasundaram

(iii) References:

(1) "Structured Statistical Models of Inductive Reasoning", Kemp & Tenenbaum (2009)

10.2 Appendix B

List of Diseases:

Alzheimer's Disease
Arthritis
Asthma
Bell's Palsy
Breast Cancer
Cancer
Chickenpox
Common cold
Diabetes
Diarrhea
Food poisoning
Heart disease
Hepatitis
HIV/AIDS
Hypertension
Insomnia
Jaundice
Liver cancer
Lung cancer
Lupus
Malaria
Menopause
Migraine
Neuritis
Ovarian Cancer
Parkinson's Disease
Pregnancy
Rabies
Sleep apnea
Spastic disorders
Urinary tract infections

10.3 Appendix C

List of Symptoms: Amenorrhea (Absence of periods)
Abdominal Symptoms (Bloating, Pain, Cramps, Discomfort)
Aches
Anemia
Anxiety
Bacterial Infections
Behavioral Changes
Bladder Problems
Blisters

Bowel Symptoms
Breast Symptoms (Lump, Shape Change, Thickening, Skin Ridges)
Breathing difficulty
Cardiac Symptoms
Chest Ache
Chills
Confusion
Cough
Dark Colored Urine
Dehydration
Depression
Diarrhea
Difficult Swallowing
Difficulty Speaking
Digestive Symptoms
Dizziness
Drooling
Ear Symptoms
Emotional Symptoms
Eye Symptoms (Conjunctivitis)
Facial Symptoms (Expression Changes, Flushing, Muscle Spasm, Pain, Swelling, Tingling)
Failure to thrive
Fainting
Fatigue
Fever
Fluid Retention
Flu-like symptoms
Foamy Urine
Foot Symptoms (Numbness, Tingling)
Forgetfulness
Hair Loss
Headache
Hoarseness
Indigestion
Irritability
Itching
Joint Symptoms (Pain, redness, Stiffness, Swelling, Tenderness)
Kidney Symptoms (Failure, Infection, Inflammation)
Knee Symptoms
Lethargy
Light Sensitivity
Light-headedness
Liver failure
Loose stool
Loss of appetite

Malaise
Memory Problems
Mental Symptoms
Menstrual Symptoms (Changes, Absence, Bleeding)
Mood Symptoms (Change, Swing)
Mouth Symptoms (Blisters, Infections, Sores, White Spots)
Muscle Symptoms (Stiffness, Weakness, Spasms)
Nausea
Muscle Symptoms (Stiffness, Weakness, Spasms)
Neurologic Symptoms (Problems, Disorders)
Pale feces
Palpitations
Sexual problems
Shortness of breath
Skin Symptoms
Sleep Problems
Sore throat
Sound sensitivity
Swelling
Vague symptoms
Vomiting
Weakness
Weight Gain
Weight Loss
Wheezing
Yellow eyes
Yellow Skin