

DECLARATION

We, Samir Sheriff and Satvik N bearing USN number 1RV09CS093 and 1RV09CS095 respectively, hereby declare that the dissertation entitled “**Decision Tree Classifier with GA based feature selection**” completed and written by us, has not been previously formed the basis for the award of any degree or diploma or certificate of any other University.

Bangalore

Samir Sheriff

USN:1RV09CS093

Satvik N

USN:1RV09CS095

R V COLLEGE OF ENGINEERING

(Autonomous Institute Affiliated to VTU)

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



CERTIFICATE

This is to certify that the dissertation entitled, “**Decision Tree Classifier with GA based feature selection**”, which is being submitted herewith for the award of B.E is the result of the work completed by **Samir Sheriff and Satvik N** under my supervision and guidance.

Signature of Guide

(Mrs. Shanta R)

Signature of Head of Department

(Dr. N K Srinath)

Name of Examiner

Signature of Examiner

1:

2:

ACKNOWLEDGEMENT

Any achievement, be it scholastic or otherwise does not depend solely on the individual efforts but on the guidance, encouragement and cooperation of intellectuals, elders and friends. A number of personalities, in their own capacities have helped me in carrying out this project work. We would like to take this opportunity to thank them all.

First and foremost we would like to thank Dr. B. S. Satyanarayana, Principal, R.V.C.E, Bengaluru, for his moral support towards completing my project work.

We deeply express my sincere gratitude to our guides Ms. Shanta Rangaswamy, Assistant Professor and Dr. Shobha G, Professor, Department of CSE, R.V.C.E, Bengaluru, for their able guidance, regular source of encouragement and assistance throughout this project.

We would like to thank Dr. N. K. Srinath, Head of Department, Computer Science Engineering, R.V.C.E, Bengaluru, for his valuable suggestions and expert advice.

We thank our Parents, and all the Faculty members of Department of Computer Science Engineering for their constant support and encouragement.

Last, but not the least, We would like to thank our peers and friends who provided me with valuable suggestions to improve our project.

Samir Sheriff

8th semester, CSE

USN:1RV09CS093

Satvik N

8th semester, CSE

USN:1RV09CS095

Vaishakh B N

8th semester, CSE

USN:1RV09CS114

ABSTRACT

A time series is a sequence of data points, measured typically at successive points in time spaced at uniform time intervals. Time series analysis comprises methods for analysing time series data in order to extract meaningful statistics and other characteristics of the data. In the context of statistics, the primary goal of time series analysis is forecasting. In the context of signal processing it is used for signal detection and estimation, while in the context of data mining, pattern recognition and machine learning time series analysis can be used for clustering, classification, query by content, anomaly detection as well as forecasting. This project is aimed making a time series data mining tool which can be used to accomplish the above goals.

@TODO Data sets names to be added

The project makes use of quite a few data sets in for time series analysis. these include Rainfall Data from Sita Nadi, electricity consumption data from XYZ city, ECG data and Sea Level Data. Anomaly detection, forecasting, Similarity detection is performed on these data sets using the algorithms available in literature.

@TODO Key findings and results come here

Contents

ACKNOWLEDGEMENT	i
ABSTRACT	ii
CONTENTS	ii
LIST OF FIGURES	vi
LIST OF TABLES	vii
1 INTRODUCTION	viii
1.1 DEFINITIONS AND USAGE	viii
1.2 LITERATURE SURVEY	viii
1.3 MOTIVATION	xv
1.4 PROBLEM STATEMENT	xv
1.5 OBJECTIVE	xv
1.6 SCOPE	xv
1.7 METHODOLOGY	xvi
1.8 ORGANIZATION OF REPORT	xvi
2 SOFTWARE REQUIREMENTS SPECIFICATION	xvii
2.1 PRODUCT PERSPECTIVE	xvii
2.2 PRODUCT FEATURES	xvii
2.3 CONSTRAINTS	xviii
2.4 ASSUMPTIONS AND DEPENDENCIES	xviii
2.5 SPECIFIC REQUIREMENTS	xviii
2.5.1 FUNCTIONAL REQUIREMENTS	xviii
2.5.2 SOFTWARE REQUIREMENTS	xix

2.5.3	HARDWARE REQUIREMENTS	xix
3	HIGH LEVEL DESIGN	xx
3.1	SYSTEM ARCHITECTURE	xx
3.2	DATA FLOW DIAGRAMS	xxi
3.2.1	DFD LEVEL 0	xxii
3.2.2	DFD LEVEL 1	xxiii
3.2.3	DFD LEVEL 2	xxiv
4	DETAILED DESIGN	xxvi
4.1	STRUCTURED CHART	xxvi
4.2	MODULES DESCRIPTION	xxviii
4.2.1	GUI MODULE	xxviii
4.2.2	SESSION KEY GENERATION MODULE	xxix
4.2.3	EMBEDDING ALGORITHM MODULE	xxx
4.2.4	DATA ENCRYPTION MODULE	xxxi
4.2.5	DATA DECRYPTION MODULE	xxxii
4.2.6	DATA EXTRACTION MODULE	xxxiii
5	IMPLEMENTATION	xxxv
5.1	PROGRAMMING LANGUAGE SELECTION	xxxv
5.2	PLATFORM	xxxvi
5.3	CODING STANDARDS	xxxvi
5.4	DIFFICULTIES ENCOUNTERED AND STRATEGIES USED TO TACKLE	xxxvi
6	Decision Tree Learning	xxxviii
6.1	Definition	xxxviii
6.2	The Basic Idea	xxxviii
6.3	Building the Decision Tree	xl
6.3.1	ID3 Algorithm	xl
6.3.2	Choosing the best attribute for a given node	xli
6.3.3	Entropy - a measure of homogeneity of the set of examples	xlii

6.3.4	Information gain measures the expected reduction in entropy . . .	xliii
6.3.5	Discretization	xliv
6.3.6	Limitation of Decision Tree Methods	xl v
7	Genetic Algorithms	xlvi
7.1	The Algorithm	xlvi
7.2	Genetic Operators	xl vii
7.2.1	Generation Zero	xl vii
7.2.2	Survival of the Fittest	xl vii
7.2.3	The Next Generation	xl viii
8	Decision Trees and Genetic Algorithms	li
8.1	Representation of Chromosomes	lii
8.2	Population	lii
8.3	Advantages	liii
9	Our State-of-the-art Object-Oriented System	liv
9.1	org.ck.sample	lv
9.2	org.ck.dt	lv
9.3	org.ck.ga	lvii
9.4	org.ck.gui	lvii
10	CONCLUSION AND FUTURE WORK	lx
10.1	Summary	lx
10.2	Limitations	lxi
10.3	Future enhancements	lxi
	BIBLIOGRAPHY	lxi
	APPENDICES	lxiii

List of Figures

3.1	System Architecture	xxi
3.2	Data Flow Diagram Level 0	xxii
3.3	Data Flow Diagram Level 0	xxiii
3.4	Data Flow Diagram Level 2: Embedding	xxiv
3.5	Data Flow Diagram Level 2: Embedding	xxv
4.1	Structure Chart	xxvii
4.2	GUI Module	xxix
4.3	Session Key Generator Module	xxix
4.4	Embedding Algorithm Module	xxx
4.5	Encryption Module	xxxi
4.6	Decryption Module	xxxii
4.7	Data Extraction Module	xxxiv
6.1	Sample Decision Tree	xxxix
6.2	ID3 Algorithm	xli
6.3	The entropy function relative to a binary classification, as the proportion of positive examples p_p varies between 0 and 1.	xliii
7.1	Random Genome	xlvi
7.2	Roulette Selection Pseudo-Code	xlvi
7.3	Crossover	xlix
7.4	Mutation	l
8.1	The data flow in DT/GA Hybrid Classifier Algorithm.	li

8.2	Sample Representation	lii
10.1	Application Window - Welcome Screen	lxiv
10.2	Decision Tree Constructor Window.	lxiv
10.3	Decision Tree Classifier Window.	lxv
10.4	Decision Tree Construction with GA based Feature Selector.	lxv
10.5	Decision Tree Construction with manual feature selection.	lxvi
10.6	The project source code on github public repository.	lxvi

List of Tables

Chapter 1

INTRODUCTION

A time series is a set of observations X_t , each one being recorded at a specific time t . Discrete-time time series is one in which the set T of times at which observations are made is a discrete set. Continuous-time time series are obtained when observations are recorded continuously over some time interval, e.g., when T_0 belongs $[0,1]$. Examples of time series are the daily closing value of the ECG readings and the annual flow volume of the Nile River at Aswan. Time series are very frequently plotted via line charts. Time series analysis comprises methods for analysing time series data in order to extract meaningful statistics and other characteristics of the data. Time series forecasting is the use of a model to predict future values based on previously observed values.

1.1 DEFINITIONS AND USAGE

Definition of Time Series: An ordered sequence of values of a variable at equally spaced time intervals. TSDM : Time Series Data Mining tool. @TODO ADD DETAILS OF THIS SECTION

1.2 LITERATURE SURVEY

A time series is a collection of observations made sequentially through time. At each time point one or more measurements may be monitored corresponding to one or more

attributes under consideration. The resulting time series is called univariate or multivariate respectively. In many cases the term sequence is used in order to refer to a time series, although some authors refer to this term only when the corresponding values are non-numerical. Throughout this paper the terms sequence and time series are being used interchangeably. The most common tasks of time series data mining methods are: indexing, clustering, classification, novelty detection, motif discovery and rule discovery. In most of the cases, forecasting is based on the outcomes of the other tasks. A brief description of each task is given below.

Indexing: Find the most similar time series in a database to a given query time series.

Clustering: Find groups of time series in a database such that, time series of the same group are similar to each other whereas time series from different groups are dissimilar to each other.

Classification: Assign a given time series to a predefined group in a way that is more similar to other time series of the same group than it is to time series from other groups.

Novelty detection: Find all sections of a time series that contain a different behavior than the expected with respect to some base model.

Motif discovery: Detect previously unknown repeated patterns in a time series database.

Rule discovery: Infer rules from one or more time series describing the most possible behaviour that they might present at a specific time point (or interval).

The temporal aspect of data arises some special issues to be considered and/or imposes some restrictions in the corresponding applications. First, it is necessary to define a similarity measure between two time series and this issue is very important in TSDM since it involves a degree of subjectivity that might affect the final result. A lot of research has focused on defining different similarity measures in order to improve the performance of the corresponding methods. Second, it is necessary to apply a representation scheme on the time series data. Since the amount of data may range from a few megabytes to terabytes, an appropriate representation of the time series is necessary in order to manipulate and analyze it efficiently. The desirable properties that this approach should hold are: (a) the completeness of feature extraction (b) the reduction of the dimensionality curse [1]. More specifically, the method of extraction features should guarantee that there would be no

pattern missed, the number of patterns falsely identified as interesting will be minimized and the dimensionality reduction will be substantial. In many cases also, the objective is to take advantage of the specific characteristics of a representation that make specific methods applicable (i.e. inducing rules, Markov models). Consequently, the majority of the researchers are focused on defining novel similarity measures and representation schemes in order to improve indexing performance. Clustering and classification of time series rely heavily on the similarity measure and the representation scheme selected, thus, there are very few papers proposing a novel algorithm [2]. A recent survey on clustering time series is provided by Liao [3]. Novelty detection is a very important task in many areas. Several alternative terms for novelty have been used, such as, anomaly, interestingness, surprising, faults to name a few. Moreover, many problems of finding periodic patterns can be considered as similar problems. The important point here is to provide a clear and concise definition of the corresponding notion. For instance, Keogh et al. [14] describe a pattern as surprising if the frequency with which it appears, differs greatly from that expected given previous experience. The authors present a novel algorithm, called Tarzan, and provide useful pointers to relevant literature. Recently, Aref et al. [4] focus on discovering partial periodic patterns in one or more databases. They present algorithms for incremental mining (how to maintain discovered patterns over time as the database is being expanded). Motif discovery has only recently attracted the interest of the data mining community [9]. Motifs are defined to be previously unknown, frequently occurring patterns in a time series. These patterns may be of particular importance to other data mining tasks, such as, rule discovery and novelty detection. The recent work of Tanaka et al. [7] proposes a new method for identifying motifs from multi-dimensional time series. They apply Principal Component Analysis to reduce dimensionality and perform a symbolic representation. Then, the motif discovery procedure starts by calculating a description length of a pattern based on the Minimum Description Length principle.

- Indexing Indexing approaches are mostly influenced by the pioneer work of Agrawal et al. [1], generalized by Faloutsos et al. [12]. The emerged framework from these papers, referred as GEMINI, can be summarized in the following steps [11]: extract k essential features from the time series map into a point in k -dimension feature

space organize points with off-the-shelf spatial access method. discard false alarms

The first and second step suggests the application of a representation scheme in order to reduce the dimensionality. However, this mapping should guarantee that it would return all the qualifying objects. This implies that the similarity measure in the k-dimension feature space should lower bound the corresponding similarity measure in the original space [8]. The third step is an opened selection, however most of the times R-tree structures are used. Other indexing structures may be vp-trees [7] [9], hB-trees and grid-files. The fourth step is a consequence of the fact that this approach can not guarantee that there will not be returned unqualified objects, thus these false alarms should be discarded in a post processing phase. Recently, Vlachos et al. [8] presented an external memory indexing method for discovering similar multidimensional time series under time warping conditions. The main contribution of this work is the ability to support various distance measures without the need to reconstruct the index. Two approaches with respect to distance measures are taken under consideration, namely, the Longest Common Subsequence (LCS) and the Dynamic Time Warping (DTW). Their indexing technique works by splitting a set of multiple time series in multidimensional Minimum Bounding Rectangles (MBR) and storing them in an R-tree. For a given query, a Minimum Bounding Envelope (MBE) is constructed, that covers all the possible matching areas of the query under time warping conditions. This MBE is decomposed into MBRs and then probed in the R-tree index.

- Time series representation There have been several time series representations proposed in the literature, mainly on the purpose of reducing the intrinsically high dimensionality of time series. We will refer to some of the most commonly used representations. Discrete Fourier Transform (DFT) [1] was one of the first representation schemes proposed within data mining context. DFT transforms a time series from the time domain into the frequency domain whereas a similar representation scheme, Discrete Wavelet Transform (DWT) [8], transforms it into the time/frequency or space/frequency domain. Singular Value Decomposition (SVD)

[5] performs a global transformation by rotating the axes of the entire dataset such that the first axis explains the maximum variance, the second axis explains the maximum of the remaining variance and is orthogonal to the first axis etc. Piecewise Aggregate Approximation (PAA) [3] divides a time series into segments of equal length and records the mean of the corresponding values of each one. Adaptive Piecewise Constant Approximation (APCA) [10] is similar to PAA but allows segments of different lengths. Piecewise Linear Approximation (PLA) approximates a time series by a sequence of straight lines. Recently, more representation schemes have been proposed in order to reduce dimensionality. The first class of these schemes consists of symbolic representations. Lin et al. [11] propose a Symbolic Aggregate Approximation (SAX) method, which uses as a first step the PAA representation and then discretizes the transformed time series by using the properties of the normal probability distribution. Bagnal[5] assess the effects of clipping original data on the clustering of time series. Each point of a series is mapped to 1 when it is above the population mean and to 0 when it is below. This representation is called clipping and has many advantages especially when the original series is long enough. It achieves adequate accuracy in clustering, it efficiently handles outliers and it provides the ability to employ algorithms developed for discrete or categorical data. Megalooikonomou et al. [30] introduce a novel dimensionality reduction technique, called Piecewise Vector Quantized Approximation (PVQA). This technique is based on vector quantization that partitions each series into segments of equal length and uses vector quantization to represent each segment by the closest codeword from a codebook. The original time series is transformed to a lower dimensionality series of symbols. This approach requires a training phase in order to construct the codebook, a data-encoding scheme and a distance measure. Cole et al. [9] provide a work that addresses the task of discovering correlated windows of time series (synchronously or with lags) over streaming data. They concentrate in the case where the time series are uncooperative, meaning that there does not exist a fundamental degree of regularity that would allow an efficient implementation of DFT transformations. The proposed method involves a combination of

several techniques sketches (random projections), convolution, structured random vectors, grid structures, and combinatorial design in order to achieve high performance. Gionis and Mannila [7] introduce a different approach, which is mainly motivated from research on human genome sequences. However, this approach is more general and involves multivariate time series. The notion behind their approach is that, the high variability that some time series very often exhibit, may be explained by the existence of several different sources that affect different segments of this series. More specifically, the task is to find a proper way to segment a time series into k segments, each of which comes from one of h different sources ($k \ll h$). This task is analogous to clustering the points of a time series in h clusters with the additional constraint that a cluster may change at most $k-1$ times. Gionis and Mannila provide three algorithms for solving this problem and they test them on synthetic and genome data. Finally, Vlachos et al. [3] propose to represent a time series by applying discrete Fourier transformations and retain the k best Fourier coefficients instead of the first few ones. Although this paper is motivated by mining knowledge from the query logs of the MSN search engine, the proposed methods may be applied for time series data mining in general.

- **Similarity Measures** The definition of novel similarity measures has been one of the most researched areas in the TSDM field. Generally, they are strongly related to the representation scheme applied to the original data. However, there are some similarity measures that appear frequently in the literature. Most of the researchers choices are based on the family of L_p norms, that include the Euclidean distance. Yi and Faloutsos [3] presented a novel and fast indexing scheme when the distance function is any of the arbitrary L_p norms ($p = 1, 2, \dots$). Another similarity measure that attracted a lot of attention, Dynamic Time Warping (DTW), comes from the speech recognition field [6]. The main advantage of this measure is that it allows acceleration-deceleration of a series along the time dimension (nonlinear alignments are possible), however it is computationally expensive. Markov models have been constructed and experimented. Another family of distance measures,

Longest Common Subsequence Measures (LCS), often used in speech recognition and text pattern matching. As an example of this approach, we refer to the work of Agrawal et al. [2] who define two sequences as similar when they have enough, non-overlapping, time-ordered pairs of subsequences that are similar. Li et al. [6] propose an algorithm for fast and efficient recognition of motions in multi-attribute continuous motion sequences. The main contribution of this paper is the definition of a similarity measure based on the analysis of Singular Value Decomposition (SVD) properties of similar multi-attribute motions. The proposed measure deals with noise and takes into account the different rates and durations of each motion. The authors also propose a five-phase algorithm for handling segmentation and recognition in real-time. Sakurai et al. [5] propose the Fast search method for dynamic Time Warping (DTW) that satisfies the following criteria: (a) it is fast (b) it produces no false dismissals (c) it does not pose any restriction on the series length (d) it supports for any, as well as for no restriction on warping scope. Their approach is based on a new lower bounding distance measure. They represent the sequence with approximate segments, not necessary of equal length, and operate on them. Three segments, the lower bound, the upper bound, and the time interval, correspond to each one of these approximate segments. In order to fulfill all of the above criteria, they provide algorithms for dynamic programming and searching adjusted to the properties of this representation. Fu et al. [14] propose a new technique to query time series that incorporates global scaling and time warping. The argument is that most real world problems require the ability to handle both types of distortion simultaneously. The approach is to scale the sequence by a bounded scaling factor and also to find nearest neighbor or evaluate range query by applying time warping. The authors provide definitions and proofs of the necessary lower bounds. Furthermore, there is the expected contribution to defining similarity measures by papers that propose novel representation schemes, since these two tasks are interrelated to each other.

1.3 MOTIVATION

mr. Manju. @TODO Add video link here..

1.4 PROBLEM STATEMENT

Make a paper and publish water data everywhere. :/ @TODO add problem statement.

1.5 OBJECTIVE

The ability to model and perform decision modelling and analysis is an essential feature of many real-world applications ranging from emergency medical treatment in intensive care units to military command and control systems. Existing formalisms and methods of inference have not been effective in real-time applications where trade-offs between decision quality and computational tractability are essential. The objective of this project is to fill the void that exists and help in proper analysis of time varying data.

1.6 SCOPE

The scope of a time series data mining tool is two fold. The first is to obtain an understanding of the underlying forces and structure that produced the observed data. The second is to fit a model and proceed to forecasting, monitoring or even feedback and feed forward control. The time series data mining tool can be used in the following fields.

- **Economic Forecasting**
- **Sales Forecasting**
- **Rainfall Analysis**
- **Stock Market Analysis**
- **Yield Projections**

- **Process and Quality Control**
- **Census Analysis**

1.7 METHODOLOGY

Time series analysis of data requires the user to be able to view the different algorithms and the result obtained from each algorithm along with the graphs which help the user understand the time varying nature of the data. Hence, the representation of data becomes very important. Having understood this requirement in the early phase of the project, we adopted a methodology that will accomplish the objectives in a neat and intuitive way. A GUI was developed in the form of Java Server Pages and the back end was coded in Java which helped us exploit the object oriented paradigm in design of algorithms.

1.8 ORGANIZATION OF REPORT

@TODO after all chapters are complete.

Chapter 2

SOFTWARE REQUIREMENTS SPECIFICATION

Software Requirement Specification (SRS) is an important part of software development process. It includes a set of use cases that describe all the interactions of the users with the software. Requirements analysis is critical to the success of a project.

2.1 PRODUCT PERSPECTIVE

Time Series Data Mining tool is a unique product that makes use of different algorithms to predict, view similarities, and points out the anomalies in different time varying data sets. It is built in a pluggable fashion where the only requirement at the users end is the browser and a working internet connection.

2.2 PRODUCT FEATURES

The time series data mining tool has many features that distinguishes it from the others available already in the open world. It provides accurate results using the similarity finding, anomaly finding algorithms. The back propagation neural network helps us predicting the future values. On the front end, the user has options to choose the algorithm

of her choice. Also, the charts which depict the output are carefully plotted using the google charts API which has been made available by Google Inc. Also, Java beans along with servlets and java server pages and best practices of coding have been followed.

2.3 CONSTRAINTS

During the development of this product, constraints were encountered. Some specific constraints under which the time series data mining tool has are :

- @TODO Add Constraints
- @TODO Add Constraints

2.4 ASSUMPTIONS AND DEPENDENCIES

- It is assumed that the user of this tool has basic understanding of time series data mining.
- Also, the user must have a decent knowledge of the interpretation of line graphs.

2.5 SPECIFIC REQUIREMENTS

This section shows the functional requirements that are to be satisfied by the system. All the requirements exposed here are essential to run this tool successfully.

2.5.1 FUNCTIONAL REQUIREMENTS

The functionality requirements for a system describe the functionality or the services that the system is expected to provide. This depends on the type of software system being developed. The requirements that are needed for this project are :

- The data sets should be normalized so that the algorithms can be applied effectively.

- A good representation of the results should be made available to the users through proper representation media like graphs.
- @TODO ADD SOMETHING HERE.

2.5.2 SOFTWARE REQUIREMENTS

DEVELOPERS MACHINE

- Front End: Java Enabled Browser
- Back End: Java
- Operating System: Windows 7
- JDK 7.0
- Apache Tomcat Server version 7.0
- Eclipse IDE for J2EE Developers
- Active Internet Connection
- JQuery UI and Ajax Libraries

END USERS MACHINE

Java Enabled Browser

Active Internet Connection

2.5.3 HARDWARE REQUIREMENTS

- Processor: Intel Pentium 4 or higher version
- RAM: 512MB or more
- Hard disk: 5 GB

SOFTWARE INTERFACES

The Java Runtime Environment (JRE) is required to run the software.

Chapter 3

HIGH LEVEL DESIGN

The software development usually follows Software Development Life Cycle (SDLC). The second stage of SDLC is the design phase. The design stage involves two substages namely High level design and Detailed level design.

High level design gives an overview of how the system works and top level components comprising the system.

3.1 SYSTEM ARCHITECTURE

This section provides an overview of the functionality and the working of the time series data mining tool. The overall functionality of the application is divided into different modules in an efficient way. The system architecture is shown in Figure 3.1

3.2 DATA FLOW DIAGRAMS

A DFD is a figure which shows the flow of data between the different processes and how the data is modified in each of the process. It is very important tool in software

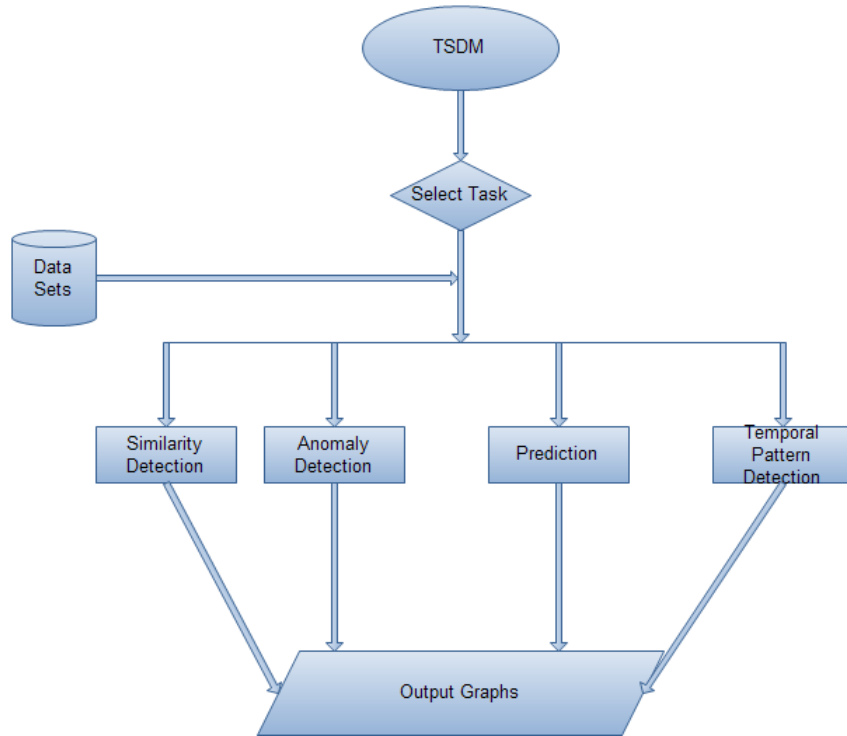


Figure 3.1: System Architecture

engineering that is used for studying the high level design.

There are many levels of DFDs. Level 0 gives the general description and level 1 gives the detailed description. Going higher in the level numbers greater description of the processes will be given.

3.2.1 DFD LEVEL 0

The level 0 DFD is shown in Fig. 3.2 below which gives the general operation of the steganographic system. There are three major components. Two external entities called sender and receiver and one major system called the steganographic system.

- Sender : The sender is the one responsible to send the data to the receiver. The data to be sent is hidden using the Steganography system. The data is hidden in an image.
- Receiver : The receiver receives the data that is sent to him in the embedded format. The receiver then extracts the data from the embedded format to get the actual data. The receiver makes use of the Steganography system to extract the data.
- Steganography System : This the software which is used to embed the data into the image and also to extract the data from the image.

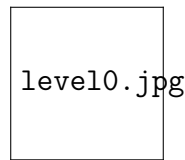


Figure 3.2: Data Flow Diagram Level 0

3.2.2 DFD LEVEL 1

There are two major processes in level 1 DFD as shown in Figure 3.3. The processes involved in here are :

- Data Embedding: This process represents the actual embedding the data. The inputs given to this process are the data files, cover image. Input image is where the data is to be embedded and the key must also be shared between the sender and the receiver.
- Data Extraction: This process is used to extract the data back from the stego image. This is the reverse process of embedding. Here the input is the stego image and the key.

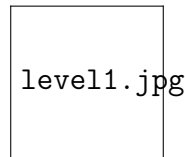


Figure 3.3: Data Flow Diagram Level 0

3.2.3 DFD LEVEL 2

Embedding Phase

The major processes in Level 2 DFD of Embed process on the sender side are shown in Figure 3.4.

- Data Holder: Data holder contains all the data bits, which upon invocation will give the required number of bits of data need to be embedded.
- Pixel Retriever: This retrieves the individual pixels from the cover image.
- Processor: Takes the key as the input, encrypts the data and embeds it into the pixels received from the retriever.
- A final image is formed as a result of this process.

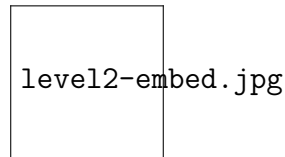


Figure 3.4: Data Flow Diagram Level 2: Embedding

Extraction Phase

The major processes in the Level 2 DFD of the Extraction module on the receiver side is shown in Figure 3.5.

- Data Retriever: The data retriever module extracts the data file from the morphed image which it receives. The key should be present for proper retrieval.
- Pixel Retriever: Retrieves the pixels from the morphed image and provides it to the Data Retriever.

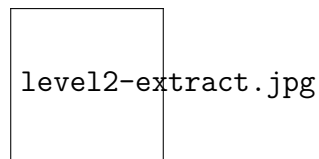


Figure 3.5: Data Flow Diagram Level 2: Embedding

Chapter 4

DETAILED DESIGN

4.1 STRUCTURED CHART

Structure charts are used to specify the high level design or architecture of a computer program. As a design tool, they help the programmer in dividing and conquering a large software problem, i.e. recursively breaking a problem down into parts that are small enough to be understood by a human brain. The process is called top-down design or functional decomposition.

Programmers use a structure chart to build a program in a manner similar to how an architect uses a blueprint to build a house. In the design stage, the chart is drawn and used as a method for the client and various software designers to communicate. During the actual building of the program, the chart is continuously referred to as master plan. Often, it is modified as programmers learn new details about the program. After a program is completed, the structured chart is used to fix bugs and to make changes.

The entire program starts with user entering his choice of input as to either embed or extract. Based on this, the GUI module responds appropriately with success or failure. The digital media chosen by user is taken as input for embedding phase. The digital media is encrypted using AES algorithm before embedding to stego image. This encrypted data is embedded onto image using embedding algorithm to get the stego image containing

secret data. A key is generated using Diffie Hellman Key exchange which acts as the input for the embedding algorithm.

At the receiver end, the Diffie Hellman key is generated once again. From the stego images, secret data is extracted. It will be in encrypted form. The original input data is extracted by using decompression module of AES algorithm. This gives the hidden file's original content.

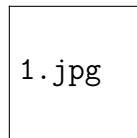


Figure 4.1: Structure Chart

4.2 MODULES DESCRIPTION

This section describes the main modules that are used in developing the project. This will help us in understanding the working of individual components.

4.2.1 GUI MODULE

Definition:

This module is the core module which takes the users input to decide upon which operation to be performed. Based upon user choice, the necessary inputs will be taken in this module. If any errors are generated during any operation, then suitable report is displayed to the user.

Resources:

The input files are bitmap images and are chosen to embed data. A session key is generated using Diffie Hellman protocol. Also, the data file to be hidden is taken.

Functionality:

This is the main flowchart which shows all the functions provided by the software. In the beginning, the user is given two options, according to which the user either embeds extracts or exits from the application. This is shown in the decision symbol. According to the decision taken, the operations are performed. If the user chooses to embed, then the functions to select cover image, the secret data, and key generation are performed. Before embedding the secret data encryption is done. After performing these operations, the stego images are sent to the receiver using any wireless medium or LAN. At the receivers end, the receiver implements the functions to extract stego image values, extract secret digits data. After the implementation of these functions, if the user wants to exit from the software, that particular option is also provided to the user. This is shown in Fig 4.2.

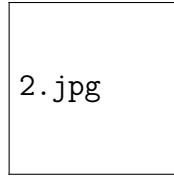


Figure 4.2: GUI Module

4.2.2 SESSION KEY GENERATION MODULE

Definition

This module generates a session key used for embedding using the principles of Diffie Hellman protocol.

Input Resources

Takes two global parameters Q and α .

Output Resources

It returns a session key k .

Functionality

Using the parameters Q and α , calculations are performed on these and finally a key is generated which is used for embedding the data.

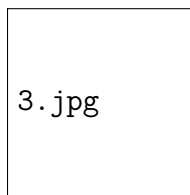


Figure 4.3: Session Key Generator Module

4.2.3 EMBEDDING ALGORITHM MODULE

Input Resources

This module takes the cover image, input file and the key as input.

Output Resources

The stego image is returned as the output of this process.

Functionality

First the size of the file is converted to base 8 and each value is embedded in first 32 bytes using Session Key. Next, embed 3 bits using Session Key bit order and choose positions in cover image to replace with the Data bits. Continue this until the end of file.

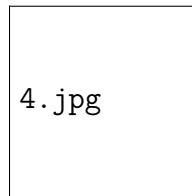


Figure 4.4: Embedding Algorithm Module

4.2.4 DATA ENCRYPTION MODULE

Input Resources

This module takes the file to be embedded as input.

Output Resources

The module returns the encrypted form of the file as output.

Functionality

The module encrypts the file using Advanced Encryption Standard (AES) using 128 bit key.

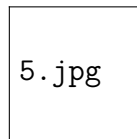


Figure 4.5: Encryption Module

4.2.5 DATA DECRYPTION MODULE

Input Resources

This module takes the encrypted file to be embedded as input.

Output Resources

The module returns the actual file by decrypting it.

Functionality

The module decrypts the the encrypted form of the file using Advanced Encryption Standard (AES) decryption by taking a 128 bit key.

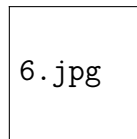


Figure 4.6: Decryption Module

4.2.6 DATA EXTRACTION MODULE

Input Resources

This module takes the stego image as the input.

Output Resources

The module returns the original data file as output.

Functionality

The module reads 3 pixels and obtains the lower half of the byte. Obtain the secret digits according to the current bit order of the Session Key. The final step is to decrypt the extracted data to get original data.

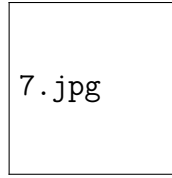


Figure 4.7: Data Extraction Module

Chapter 5

IMPLEMENTATION

The implementation phase of any project development is the most important phase and yields the final solution which solves the problem at hand. The implementation phase involves the actual materialization of the ideas, which are expressed in a suitable programming language. The factors concerning the programming language selection and platform chosen are described in the following sections.

5.1 PROGRAMMING LANGUAGE SELECTION

The programming language chosen must reflect the necessities of the project to be completely expressed in terms of the analysis and the design documents. Therefore before choosing the language, features to be included in the project are decided. The time series data mining project needs the following features in a language to be implemented. Some of the features required are stated as follows:

- J2EE provides us with servlets and JSP which help in dynamically constructing web pages.
- J2EE provides us with Java Beans which help in proper data manipulation.
- JSP and servlets make use of Java backend in a very optimal manner. They have special tags which help us exploit these features.

- Java's core classes are designed from scratch to meet the requirements of an object oriented system.

With these necessities in mind, J2EE is selected as the optimal programming language to implement the project.

5.2 PLATFORM

The TSDM tool was built and designed on Windows Operating system family. They were specifically tested on Windows 7 with Google Chrome and Mozilla Firefox browsers. Because the product is browser based, any user with the browsers mentioned above will be able to run the tool. The product is hence platform independent in the true sense.

5.3 CODE CONVENTIONS

The code standards for the Java programming Language document contains the standard conventions that follows. It includes file names, file organizations, indentation, comments, declarations, naming conventions and programming practices. Code conventions improve the readability of the software.

5.3.1 Naming Conventions

@TODO Add shitty things here

5.3.2 File Organization

@TODO Add shitty things here

5.3.3 Class Declarations

@TODO Add shitty things here

5.3.4 Comments

@TODO Add shitty things here

5.4 DIFFICULTIES ENCOUNTERED AND STRATEGIES USED TO TACKLE

There were a number of challenges that were faced while implementing the Time Series Data Mining tool. Some challenges were challenging and ended up in helping us think innovatively and come up with efficient solutions. Some major problems that were encountered have been stated in brief along with their solutions.

Problem 1

In the initial stages of the project charts4j libraries were used to plot graphs. There were some internal problems with the URL rendering.

Solution

This Problem was solved later by making use of Google's Charts API and java script.

Problem 2

@TODO Add Difficulties Encountered here.. more bullshitting to be done here :P ...

Initially the 10 digit key was being entered by the user could not be transmitted to the other end.

Solution

Diffie Hellman key exchange module was later developed which helped in key distribution.

Problem 3

Initially it was difficult to view both the morphed and the original image at the same time.

Solution

This problem was later solved by developing a module in SWT to compare the two images side by side.

Chapter 6

Decision Tree Learning

6.1 Definition

Decision tree is the learning of decision tree from class labeled training tuples. A decision tree is a flow chart like structure, where each internal (non-leaf) node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf (or terminal) node holds a class label. The topmost node in tree is the root node.

There are many specific decision-tree algorithms. Notable ones include:

- **ID3** (Iterative Dichotomiser 3)
- **C4.5** algorithm, successor of ID3
- **CART** (Classification And Regression Tree)
- **CHi-squared Automatic Interaction Detector** (CHAID). Performs multi-level splits when computing classification trees.
- **MARS**: extends decision trees to better handle numerical data

6.2 The Basic Idea

Decision tree is a classifier in the form of a tree structure (as shown in Fig. 3.1, where each node is either:

1. A **leaf node** - indicates the value of the target attribute (class) of examples (*In Fig. 3.1, the nodes containing values $K=x$, $K=y$*)
2. A **decision node** - specifies some test to be carried out on a single attribute-value, with one branch and sub-tree for each possible outcome of the test. *In Fig. 3.1, the nodes containing attributes A , B and C*

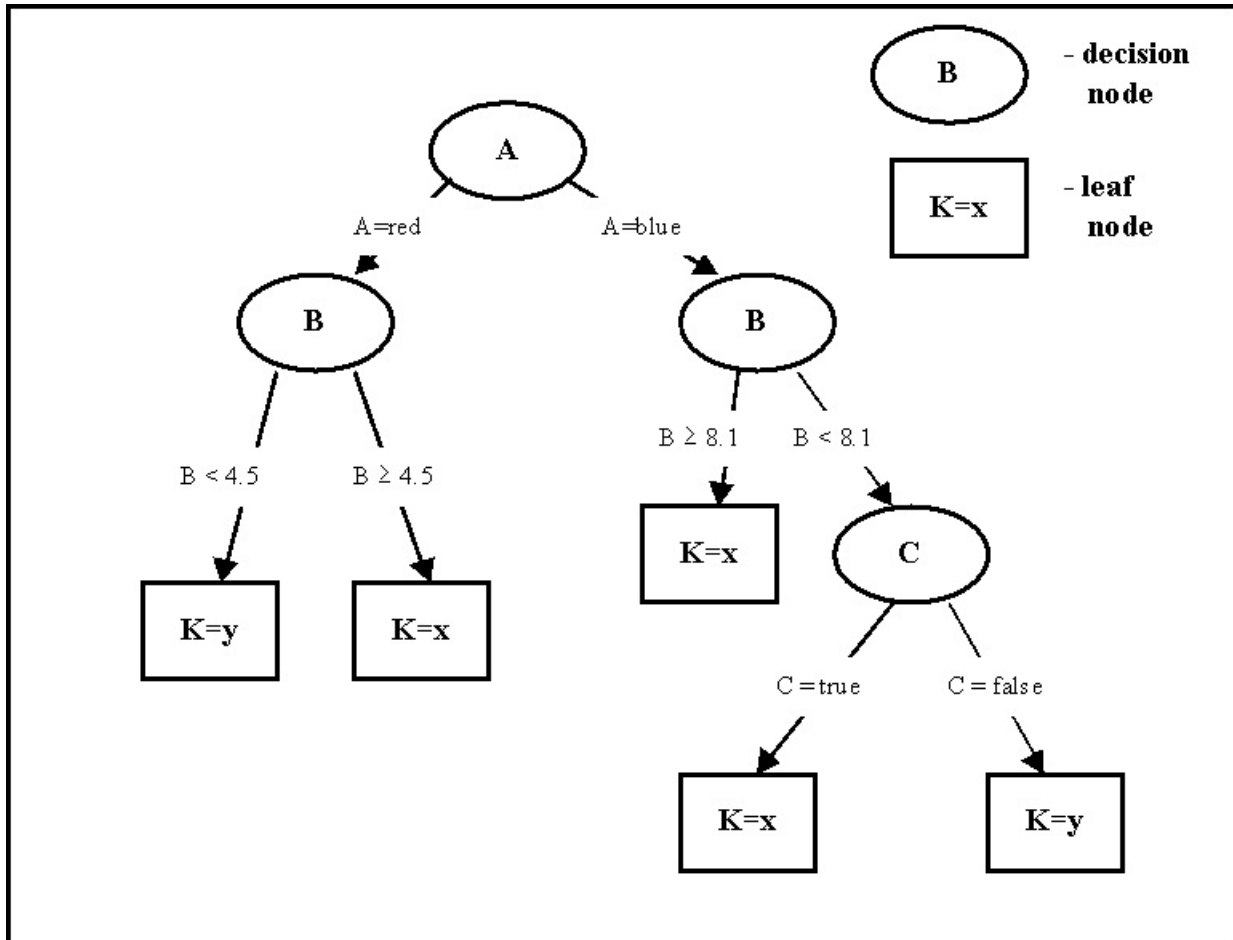


Figure 6.1: Sample Decision Tree

A decision tree can be used to classify an example by starting at the root of the tree and moving through it until a leaf node, which provides the classification of the instance.

Decision tree induction is a typical inductive approach to learn knowledge on classification. The key requirements to do mining with decision trees are:

- **Attribute-value description:** object or case must be expressible in terms of a fixed collection of properties or attributes. This means that we need to discretize continuous attributes, or this must have been provided in the algorithm. (*A, B and C, in Fig. 3.1*)
- **Predefined classes (target attribute values):** The categories to which examples are to be assigned must have been established beforehand (supervised data) (*Classes X and Y in Fig. 3.1*).
- **Discrete classes:** A case does or does not belong to a particular class, and there must be more cases than classes.
- **Sufficient data:** Usually hundreds or even thousands of training cases.

6.3 Building the Decision Tree

Most algorithms that have been developed for learning decision trees are variations on a core algorithm that employs a top-down, greedy search through the space of possible decision trees. Decision tree programs construct a decision tree T from a set of training cases.

6.3.1 ID3 Algorithm

J. Ross Quinlan originally developed ID3 at the University of Sydney. He first presented ID3 in 1975 in a book, *Machine Learning*, vol. 1, no. 1. ID3 is based on the Concept Learning System (CLS) algorithm. ID3 searches through the attributes of the training instances and extracts the attribute that best separates the given examples. If the attribute perfectly classifies the training sets then ID3 stops; otherwise it recursively operates on the m (where m = number of possible values of an attribute) partitioned subsets to get their "best" attribute. The algorithm uses a greedy search, that is, it picks the best attribute and never looks back to reconsider earlier choices. Note that ID3 may misclassify data.

```

function ID3
Input:   (R: a set of non-target attributes,
         C: the target attribute,
         S: a training set) returns a decision tree;
begin
  If S is empty, return a single node with
    value Failure;
  If S consists of records all with the same
    value for the target attribute,
    return a single leaf node with that value;
  If R is empty, then return a single node
    with the value of the most frequent of the
    values of the target attribute that are
    found in records of S; [in that case
    there may be errors, examples
    that will be improperly classified];
  Let A be the attribute with largest
    Gain(A,S) among attributes in R;
  Let {aj | j=1,2, ..., m} be the values of
    attribute A;
  Let {Sj | j=1,2, ..., m} be the subsets of
    S consisting respectively of records
    with value aj for A;
  Return a tree with root labeled A and arcs
    labeled a1, a2, ..., am going respectively
    to the trees (ID3(R-{A}, C, S1), ID3(R-{A}, C, S2),
    ....., ID3(R-{A}, C, Sm);
  Recursively apply ID3 to subsets {Sj | j=1,2, ..., m}
    until they are empty
end

```

Figure 6.2: ID3 Algorithm

6.3.2 Choosing the best attribute for a given node

The estimation criterion in the decision tree algorithm is the selection of an attribute to test at each decision node in the tree. The goal is to select the attribute that is most useful for classifying examples. A good quantitative measure of the worth of an attribute is a statistical property called information gain that measures how well a given attribute separates the training examples according to their target classification. This measure is used to select among the candidate attributes at each step while growing the tree.

6.3.3 Entropy - a measure of homogeneity of the set of examples

In order to define information gain precisely, we need to define a measure commonly used in information theory, called entropy, that characterizes the (im)purity of an arbitrary collection of examples. Given a set S , containing only positive and negative examples of some target concept (a 2 class problem), the entropy of set S relative to this simple, binary classification is defined as:

$$\text{Entropy}(S) = -p_p \log_2 p_p - p_n \log_2 p_n$$

where p_p is the proportion of positive examples in S and p_n is the proportion of negative examples in S . In all calculations involving entropy we define $0 \log 0$ to be 0.

To illustrate, suppose S is a collection of 25 examples, including 15 positive and 10 negative examples [15+, 10-]. Then the entropy of S relative to this classification is

$$\text{Entropy}(S) = -(15/25)\log_2(15/25) - (10/25)\log_2(10/25) = 0.970$$

Notice that the entropy is 0 if all members of S belong to the same class. For example, if all members are positive ($p_p = 1$), then p_n is 0, and:

$$\text{Entropy}(S) = -1\log_2(1) - 0\log_2 0 = -10 - 0\log_2 0 = 0.$$

Note the entropy is 1 (at its maximum!) when the collection contains an equal number of positive and negative examples. If the collection contains unequal numbers of positive and negative examples, the entropy is between 0 and 1. Figure 3.3 shows the form of the entropy function relative to a binary classification, as p_+ varies between 0 and 1.

One interpretation of entropy from information theory is that it specifies the minimum number of bits of information needed to encode the classification of an arbitrary member of S (i.e., a member of S drawn at random with uniform probability). For example, if p_p is 1, the receiver knows the drawn example will be positive, so no message need be sent, and the entropy is 0. On the other hand, if p_p is 0.5, one bit is required to indicate whether the drawn example is positive or negative. If p_p is 0.8, then a collection of messages can be encoded using on average less than 1 bit per message by assigning shorter codes to collections of positive examples and longer codes to less likely negative examples.

Thus far we have discussed entropy in the special case where the target classification is binary. If the target attribute takes on c different values, then the entropy of S relative

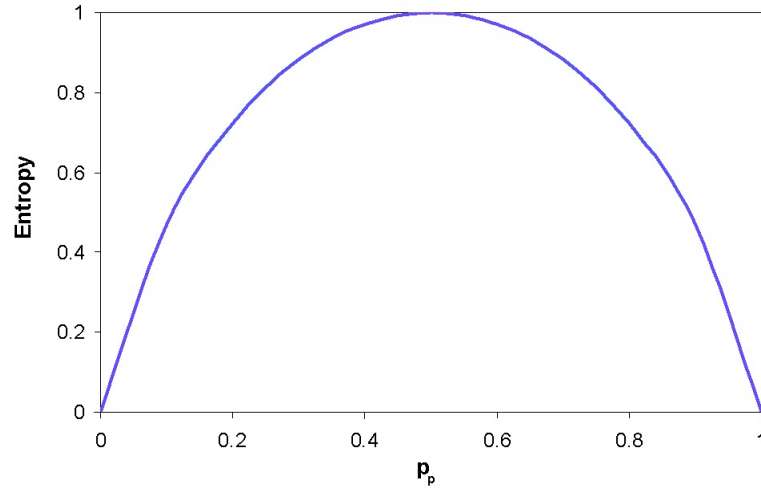


Figure 6.3: The entropy function relative to a binary classification, as the proportion of positive examples p_p varies between 0 and 1.

to this c-wise classification is defined as:

$$Entropy(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

where p_i is the proportion of S belonging to class i . Note the logarithm is still base 2 because entropy is a measure of the expected encoding length measured in bits. Note also that if the target attribute can take on c possible values, the maximum possible entropy is $\log_2 c$.

6.3.4 Information gain measures the expected reduction in entropy

Given entropy as a measure of the impurity in a collection of training examples, we can now define a measure of the effectiveness of an attribute in classifying the training data. The measure we will use, called information gain, is simply the expected reduction in entropy caused by partitioning the examples according to this attribute. More precisely, the information gain, $\text{Gain}(S, A)$ of an attribute A , relative to a collection of examples S , is defined as:

$$\text{Gain}(S, A) = Entropy(S) - \sum_{v \in \text{Value}(A)} \frac{S_v}{S} Entropy(S_v)$$

where $\text{Values}(A)$ is the set of all possible values for attribute A , and S_v is the subset of S for which attribute A has value v (i.e., $S_v = \{s \in S \mid A(s) = v\}$). Note the first term in the equation for Gain is just the entropy of the original collection S and the second term is the expected value of the entropy after S is partitioned using attribute A . The expected entropy described by this second term is simply the sum of the entropies of each subset S_v , weighted by the fraction of examples $|S_v|/|S|$ that belong to S_v . Gain (S, A) is therefore the expected reduction in entropy caused by knowing the value of attribute A . Put another way, Gain(S, A) is the information provided about the target attribute value, given the value of some other attribute A . The value of Gain(S, A) is the number of bits saved when encoding the target value of an arbitrary member of S , by knowing the value of attribute A .

The process of selecting a new attribute and partitioning the training examples is now repeated for each non-terminal descendant node, this time using only the training examples associated with that node. Attributes that have been incorporated higher in the tree are excluded, so that any given attribute can appear at most once along any path through the tree. This process continues for each new leaf node until either of two conditions is met:

1. Every attribute has already been included along this path through the tree
2. The training examples associated with this leaf node all have the same target attribute value (i.e., their entropy is zero).

6.3.5 Discretization

Decision Tree Learning requires a discrete feature space. To handle continuous feature spaces, the process of discretization has to be carried out on the feature space to obtain a discrete feature space, which can act as an input to the ID3 algorithm. One of the most famous algorithms for discretization is, perhaps, equal width interval binning. It involves sorting the observed values of a continuous feature and dividing the range of observed values for a variable into k equally sized bins, where k is a parameter supplied by the user. If a variable x is observed to have values bounded by x_{min} and x_{max} , then

this method computes the bin width: $\delta = \frac{x_{max} - x_{min}}{k}$ and constructs bin boundaries, or thresholds, at $x_{min} + i\delta$ where $i = 1, \dots, k-1$. This method is applied to each continuous feature independently.

6.3.6 Limitation of Decision Tree Methods

The weaknesses of decision tree methods

- Decision trees are less appropriate for estimation tasks where the goal is to predict the value of a continuous attribute.
- Decision trees are prone to errors in classification problems with many class and relatively small number of training examples.
- Decision tree can be computationally expensive to train. The process of growing a decision tree is computationally expensive. At each node, each candidate splitting field must be sorted before its best split can be found. In some algorithms, combinations of fields are used and a search must be made for optimal combining weights. Pruning algorithms can also be expensive since many candidate sub-trees must be formed and compared.
- Decision trees do not treat well non-rectangular regions. Most decision-tree algorithms only examine a single field at a time. This leads to rectangular classification boxes that may not correspond well with the actual distribution of records in the decision space.

Chapter 7

Genetic Algorithms

Nature seems to have an uncanny knack for problem-solving. Life began as a handful of simple, single-celled organisms barely equipped to survive the harsh environment of planet Earth. However, in the short span of a few billion years, nature has adapted and evolved them into beings complex enough to ponder their own origins. While this is indeed amazing, the truly incredible part is that it all happened according to a simple plan—allow individuals with favorable traits to survive and reproduce, and let die all the rest. This, in short, is the basis for a genetic algorithm.

7.1 The Algorithm

- Create an initial population of random genomes.
- Loop through the genetic algorithm, which produces a new generation every iteration.
 - Assess the fitness of each genome, stopping if a solution is found.
 - Evolve the next generation through natural selection and reproduction.
 - * Select two random genomes based on fitness.
 - * Cross the genomes or leave them unchanged.

- * Mutate genes if necessary.
- Delete the old generation and set the new generation to the current population.
- When a solution is found or a generation limit is exceeded, the loop breaks and the genetic algorithm is complete.

7.2 Genetic Operators

The basic genetic algorithm attempts to evolve traits that are optimal for a given problem. It has a wide variety of common uses, notably for balancing weights in neural networks.

7.2.1 Generation Zero

The first step in the genetic algorithm is to create an initial population, generation zero, that contains a set of randomized strings of genes. Each string of genes, illustratively called a genome or chromosome, represents a series of traits that may or may not be useful for the problem at hand. These “genes” are usually represented by either binary digits or real numbers.

Random Genome									
Bits (Genes)	0110	1100	1111	1011	0100	1010	0111	0101	1110
Values (Traits)	6	12	15	11	4	10	7	5	14

Figure 7.1: Random Genome

7.2.2 Survival of the Fittest

Every genome in the population must now be assigned a fitness score according to how well it solves the problem at hand. The process and approach to measuring a genomes fitness will be different for every problem. Determining the fitness measure is the most important and often most difficult part of developing a genetic algorithm.

7.2.3 The Next Generation

Once the fitness for every genome is determined, its time to start building the next generation of genomes based on probability and fitness. This is the main part of the genetic algorithm, where the strong survive and the weak perish. It usually consists of these three parts:

Selection

Two genomes are selected randomly from the current population (reselection allowed), with fitter genomes having a higher chance of selection. The selected genomes, which should have a relatively high fitness score, are guaranteed to pass some of their traits to the next generation. This means that the average fitness of each successive generation will tend to increase.

The best way to program the selection function is through a method creatively named roulette selection. First, a random number between zero and the sum of the populations fitness is generated. Imagine this value as a ball landing somewhere on a pie graph of the populations fitness. Then, each genomes fitness, or slice of the pie graph, is added one by one to a running total. If the ball ends up in that genomes slice, it is selected.

```
RouletteSelection()
{
    float ball = rand_float_between(0.0, total_fitness);
    float slice = 0.0;

    for each genome in population
    {
        slice += genome.fitness;

        if ball < slice
            return genome;
    }
}
```

Figure 7.2: Roulette Selection Pseudo-Code

Crossover

The two genomes now have a good chance of crossing over with one another, meaning that they will each donate a portion of their genes to form two offspring that become part of the next generation. If they do not cross over, they simply go on to the next generation unchanged. The crossover rate determines how often the genomes will cross over, and should be in the vicinity of 65-85

A crossover operation on the binary genomes in our example would begin by choosing a random position at which to cross them. The first part of the fathers genes and the second part of the mother's genes combine to form the first child, with a similar effect for the second child. The following shows a crossover operation with the crossover point at 12.

Before Crossing

Father 011110010011 001011011000111011010000

Mother 010100111110 010101111101000100010010

After Crossing

Child 1 011110010011 010101111101000100010010

Child 2 010100111110 001011011000111011010000

Figure 7.3: Crossover

Mutation

Just before the genomes are placed into the next generation, they have a slight chance of mutating. A mutation is simply a small, random change to one of the genes. With binary genes, mutation means flipping the bit from 1 to 0 or 0 to 1. With real number genes, a small, random perturbation is added to the gene.

The mutation rate determines the chances for each gene to undergo mutation, meaning that every individual gene should get a chance to mutate. The mutation rate should be roughly 1-5 percent for binary and 5-20 percent for real numbers.

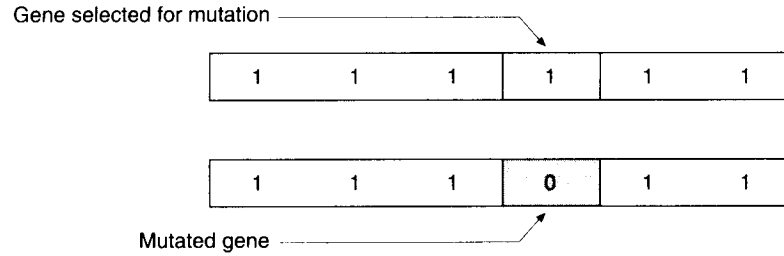


Figure 7.4: Mutation

The purpose of mutation in GAs is preserving and introducing diversity. Mutation should allow the algorithm to avoid local minima by preventing the population of chromosomes from becoming too similar to each other, thus slowing or even stopping evolution. This reasoning also explains the fact that most GA systems avoid only taking the fittest of the population in generating the next but rather a random (or semi-random) selection with a weighting toward those that are fitter.

Chapter 8

Decision Trees and Genetic Algorithms

In GA based DT Classifier, the search component is a GA and the evaluation component is a decision tree. A detailed description of this algorithm is shown in Figure 5.1.

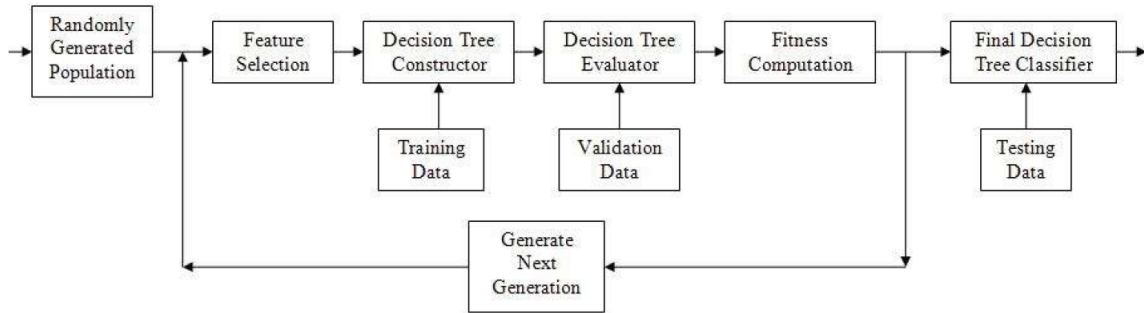


Figure 8.1: The data flow in DT/GA Hybrid Classifier Algorithm.

The basic idea of our hybrid system is to use GAs to efficiently explore the space of all possible subsets of a given feature set in order to find feature subsets which are of low order and high discriminatory power. In order to achieve this goal, fitness evaluation has to involve direct measures of size and classification performance, rather than measures such as the ranking methods such as information gain, etc.

An initial set of features is provided together with a training set of the measured feature vectors extracted from raw data corresponding to examples of concepts for which the decision tree is to be induced. The genetic algorithm (GA) is used to explore the

space of all subsets of the given feature set where preference is given to those features sets which achieve better classification performance using smaller dimensionality feature sets. Each of the selected feature subsets is evaluated (its fitness measured) by testing the decision tree produced by the ID3 algorithm. The above process is iterated along evolutionary lines and the best feature subset found is then recommended to be used in the actual design of the pattern classification system.

8.1 Representation of Chromosomes

Every individual of the population has N genes, each of which represents a feature of the input data and can be assigned to 1 or 0. 1 means the represented feature is used during constructing decision trees; 0 means it is not used. As a result, each individual in the population represents a choice of available features.

PH	Cu	HCo3	Co3	F	K	Ca	Cl	Na	Mg	No3	S04	Th	Si	EC
1	0	0	0	1	1	1	0	0	1	1	0	0	0	0

Figure 8.2: Sample Representation

In this figure, only the following features are being considered: PH, F, K, Ca, Mg, and No3, yielding the chromosome: 100011100110000

8.2 Population

The initial population is randomly generated. For each individual in the current population, a decision tree is built using the ID3 algorithm. This resulting decision tree is then tested over validation data sets, which generate classification error rates. In our application, we calculate fitness values using the weighted average of training and test classification errors. The lower the classification error rate, the better the fitness of the individual. In other words, the population consists of a number of probable decision trees.

8.3 Advantages

Such a hybrid learning system will identify significantly better feature subsets than those produced by existing methods for two reasons.

1. The power of Genetic algorithms is being exploiting to efficiently explore the non-linear interactions of a given set of features.
2. By using ID3 in the evaluation loop, an efficient mechanism for directly measuring classification accuracy is present

Chapter 9

Our State-of-the-art Object-Oriented System

Object-oriented programming (OOP) is a programming paradigm that represents concepts as "objects" that have data fields (*attributes that describe the object*) and associated procedures known as methods. Objects, which are instances of classes, are used to interact with one another to design applications and computer programs.

Had it not been for the presence of the OOP paradigm, our efforts in this project would have gone in vain, and we do not use that term lightly. Code management was a whole lot easier when compared to our past experience with procedural programming. In this chapter, we describe the different packages that were created by us to efficiently manage our code. Know first that our application consists of four diverse packages, namely:

1. **org.ck.sample**
2. **org.ck.dt**
3. **org.ck.ga**
4. **org.ck.gui**

We describe each package in detail, in the following sections. We follow a bottom-up methodology for explaining the layout of the classes.

9.1 org.ck.sample

This package allows us to efficiently manage and encapsulate the details of the data samples, provided by users, which are required for analysis. It is this class that allows our application to accept generic data sets. It consists of five classes:

1. **DataHolder** - This class keeps track of names of files that contain training and testing samples; lists of features; their corresponding classification values; and Probability values. It provides this information, when required, to the front-end or back-end of our application. To make a long story short, this class acts like a middleman between the back-end and front-end of our application.
2. **Feature** - This class stores a mapping between a feature name and a feature value.
3. **Sample** - This class stores all the features and the corresponding classification value for one training/test sample only.
4. **SampleCollection** - A SampleCollection, as the name suggests, is a collection of samples. In essence, this class reads the sample data from a file (using information provided by a DataHolder object) and initializes all the necessary data structures to store the data values for Classification analysis.
5. **SampleSplitter** - This class contains methods that operate on a given SampleCollection, in order to split it into two new SampleCollections, based on a given Feature. It also calculates the information gain (*as described in Chapter 3*) of the given split operation.

9.2 org.ck.dt

This package allows us to efficiently manage and encapsulate methods for all stages of Decision Tree Learning.

1. **DecisionTreeNode** - A DecisionTreeNode is a structure which may have either:

- (a) a classification value, if it happens to be a leaf node
- (b) a list of one or more children `DecisionTreeNode`s, if it happens to be a decision node, i.e., an internal node.

Know that the types of these nodes is defined by the decision tree that is constructed, and a node has at most one parent.

2. **DecisionTreeConstructor** - This class takes a `SampleCollection` (*containing training samples*) as input, builds a decision tree (*as described in Chapter 3*) and stores the root `DecisionTreeNode` of the decision tree. In essence, a `DecisionTreeConstructor` consists of a number of `DecisionTreeNode`s.
3. **DecisionTreeClassifier** - This class keeps track of the measurements of the `DecisionTree` constructed by a `DecisionTreeConstructor` object. It keeps track of the training as well as the test `SampleCollection`, and runs each sample through the `Decision Tree` that was constructed, to find out its classification accuracy, which it stores and retrieves, when required.
4. **Discretizer** - This class provides implementations of algorithms used for discretization. As mentioned earlier, decision trees work with discretized values, and if continuous-valued features are present, they have to be discretized. The `Discretizer` class contains two algorithms for discretization:
 - (a) A naive discretizer that discretizes data based on the median, with those values below the median being set to 0 and those values above the median being set to 1.
 - (b) An Equal-Binning Discretizer that discretizes the values of certain feature of a collection of samples, by putting each value into particular bins. After discretization, the values can be any integer between 0 and `binSize` (inclusive).

9.3 org.ck.ga

This package takes care of all operations of the Genetic algorithm that was mentioned earlier.

1. **Genome** - This class takes as input, a `SampleCollection`, and initializes a chromosome with random values for the presence/absence of features, as defined in Chapter 5. It keeps track of this chromosome, and provides methods to manipulate this chromosome; to calculate the fitness score of this chromosome; and to throw an exception when the fitness value threshold has been crossed or when the best solution has been discovered. It also provides facilities to switch between a chromosome and the corresponding optimal decision tree to which it is bound.
2. **Population** - As defined earlier, a population is a collection of genomes. And this is exactly what this class is. Initially, the `Population` class randomly initializes a large number of genomes, of which it keeps track. It provides methods such as roulette selection, reproduction, crossover, and mutation to operate on the population and discover the best genome, and hence, the best decision tree with the appropriate feature subset.
3. **OptimalScoreException** - This class is responsible for catching the best genome as soon as it is discovered, since the best genome should never be allowed to escape. It should be caught and nurtured for future use.

9.4 org.ck.gui

As you've probably guessed by now, this package handles the Graphical User Interface of our application, with all its bells and whistles. We made use of the Standard Widget Toolkit for the GUI of our application. This package consists of the following classes:

1. **WelcomeWindow** - This class takes care of drawing the window that appears when our application is first switched on, and obviously, its name should be Wel-

comeWindow, nothing more, nothing less. It displays a list of clickable options, namely

- (a) Train Decision Tree
- (b) Classify Data Sets
- (c) View on Github
- (d) Exit Application

We have organized the code in this package in such a way that all the options (*except for the last one - "Exit Application"*), correspond to a different class which handles the creation of the corresponding window.

2. **MainWindow** - This class manages the window that is opened when a user clicks the *Train Decision Tree* option in the Welcome Window. In this window, the user can select the appropriate options required to construct a decision tree using our Hybrid DT/GA Classifier. By the way, did we mention that a constructed decision tree can be saved for later usage?
3. **ClassifyWindow** - This class manages the window that is opened when a user clicks the *Classify Data Sets* option in the Welcome Window. A user is provided with an interface to select a saved decision tree, and classify new samples based on it. It really saves a lot of time in this fast-paced world of ours.
4. **BrowserWindow** - This class manages the window that is opened when a user clicks the *View on Github* option in the Welcome Window. In order to see and verify whether our code is original or not, users can see the online repository of our code (including its version history) on Github, in this window. Verification couldn't have been more easier.
5. **Constants** - This interface (mark my words, this is not a class) contains a list of constants used by all the classes in all the packages. This interface really makes

updating our software and meddling with various values much easier, like never before.

Chapter 10

CONCLUSION AND FUTURE WORK

10.1 Summary

In this mini project, we were able to successfully implement and test the performance of Decision Tree-based classifiers. The Decision Tree classifier was optimized using a Genetic Algorithm to select a subset of the features that were to be used in constructing an optimal decision tree.

Although our program works with generic data samples, it must be noted that when we started this project, our main intention was to classify ground water samples into two classes, namely Potable and Non-Potable Water. However, thanks to the miracle of Object-Oriented Programming Concepts, we were able to extend our application, which was developed in Java and Java SWT. We were able to extend this application to work with any generic samples. Two other samples/ Classification problems were addressed:

- Diagnosing whether a Horse has colic or is healthy, based on its Blood Sample Data.
- Classifying/Determining the quality of a wine based on Data Samples containing its quality parameters

The hybrid GA /decision tree algorithm needs to be tested further to realize its true

potential. Clearly more work needs to be done. The test results show that the Decision Trees constructed using the Genetic algorithm-based feature selector, were more efficient and accurate in classifying the data than the Decision Trees constructed by selecting features manually.

10.2 Limitations

These are a few limitations of this application.

1. The application uses about 800MB-1GB of RAM.
2. The GA feature selector takes a considerable amount time to optimize the Decision Tree depending on the size of the training and testing samples.
3. The GA optimizer may take a long time to converge or may not converge at all, which usually results in the application crashing, due to high memory usage.
4. The Decision Tree classifier with GA-based feature selection requires the use of accurate as well as a large number of training and testing samples. The efficiency of the Decision Tree constructed is solely based on the input training and testing samples.

10.3 Future enhancements

Some of the future enhancements are :

1. The application could be made more responsive by using Threads and Parallel/- Cloud Computing
2. The Decision Tree Classifier of this application could be optimized using Neural Networks which are more efficient than Genetic Algorithms.
3. An interesting extension to be explored is the possibility of additional feedback from ID3 concerning the evaluation of a feature set.

Bibliography

- [1] Genetic Algorithms -http://en.wikipedia.org/wiki/Genetic_algorithm
- [2] Genetic Algorithm for constructing DT - <http://www.jprr.org/index.php/jprr/article/viewFile/44/25>
- [3] Decision Trees - <http://web.cecs.pdx.edu/~mm/MachineLearningWinter2010/pdfslides/DecisionTrees.pdf>
- [4] Project brief for the DT using Horse data sets - <https://cs.uwaterloo.ca/~ppoupart/teaching/cs486-spring06/assignments/asst4/asst4.pdf>
- [5] Supervised and Unsupervised Discretization of Continuous Features - <http://robotics.stanford.edu/users/sahami/papers-dir/disc.pdf>
- [6] Hybrid learning using Genetic Algorithms and Decision Trees for pattern classification - <http://cs.gmu.edu/~eclab/papers/ijcai95.pdf>
- [7] Kardi Tutorials on Decision Trees- <http://people.revoledu.com/kardi/tutorial/DecisionTree/index.html>

Appendices

Appendix A : Source Code

Appendix B : Screen Shots

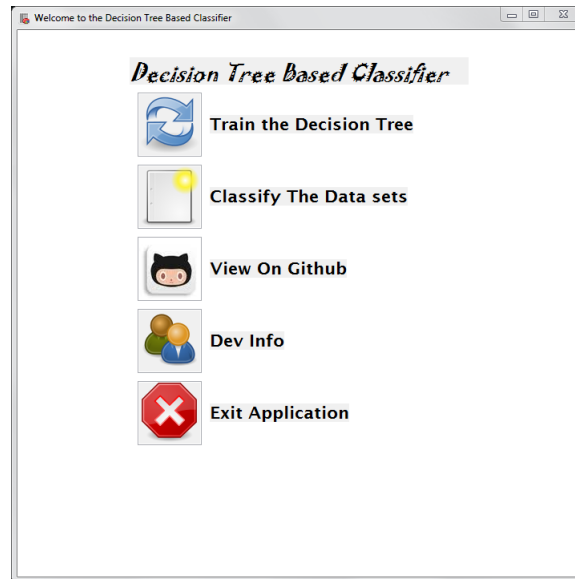


Figure 10.1: Application Window - Welcome Screen

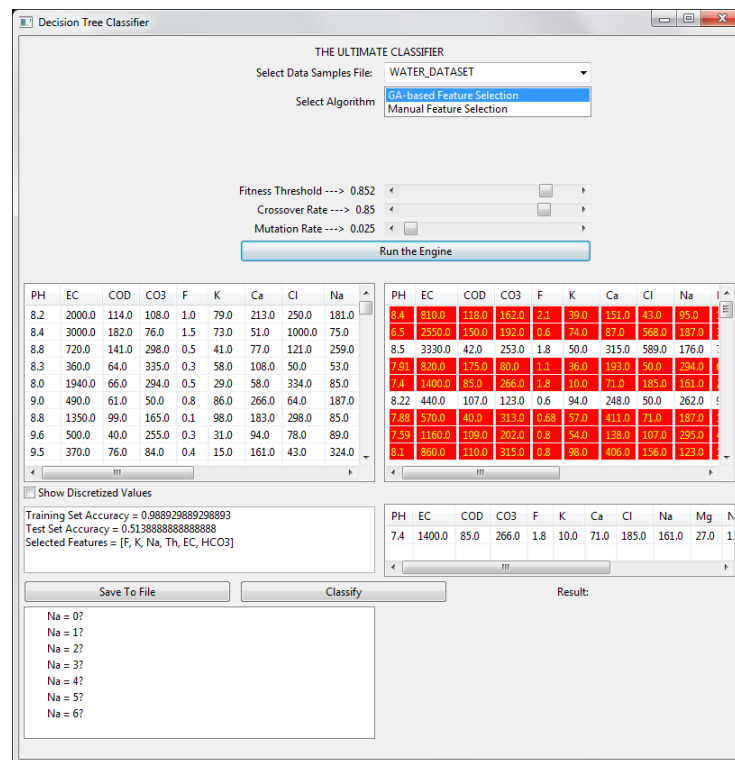


Figure 10.2: Decision Tree Constructor Window.

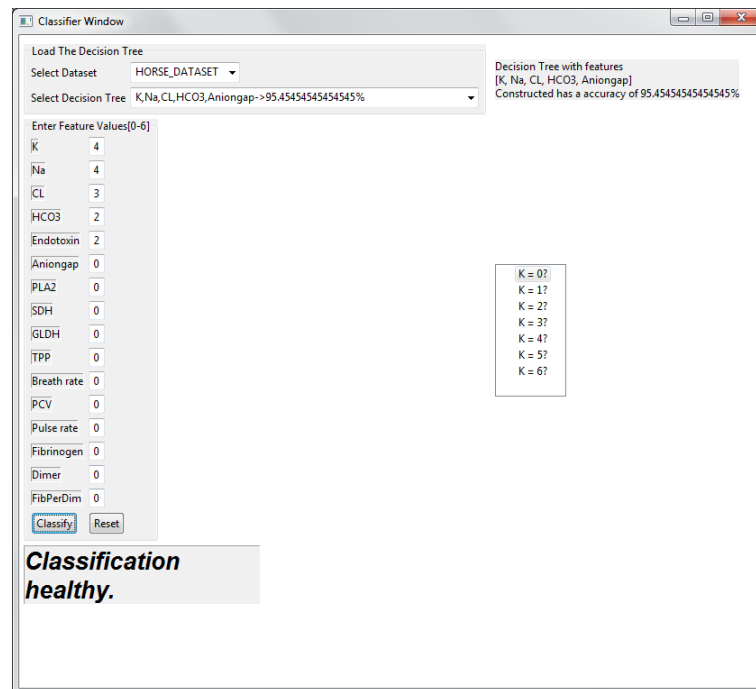


Figure 10.3: Decision Tree Classifier Window.

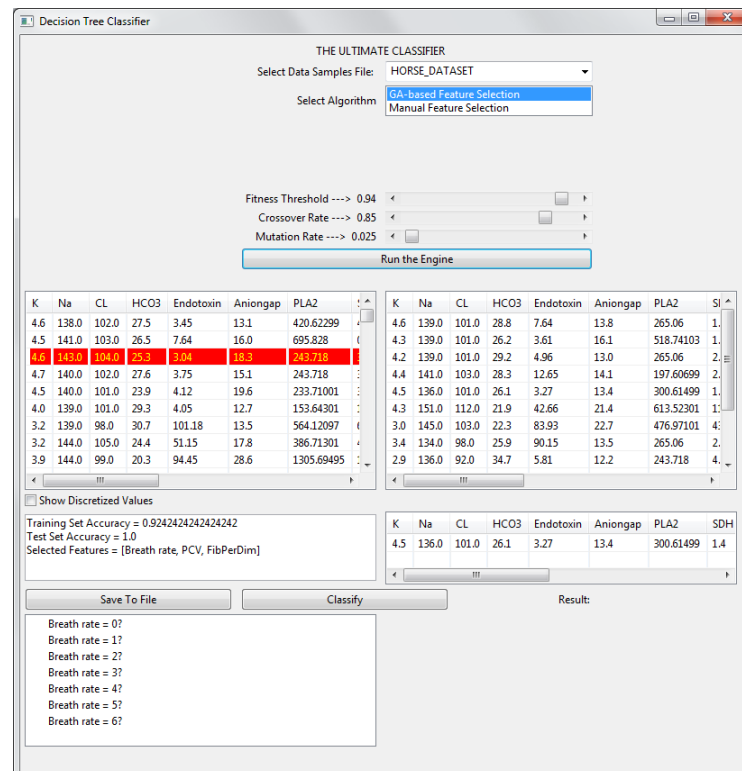


Figure 10.4: Decision Tree Construction with GA based Feature Selector.

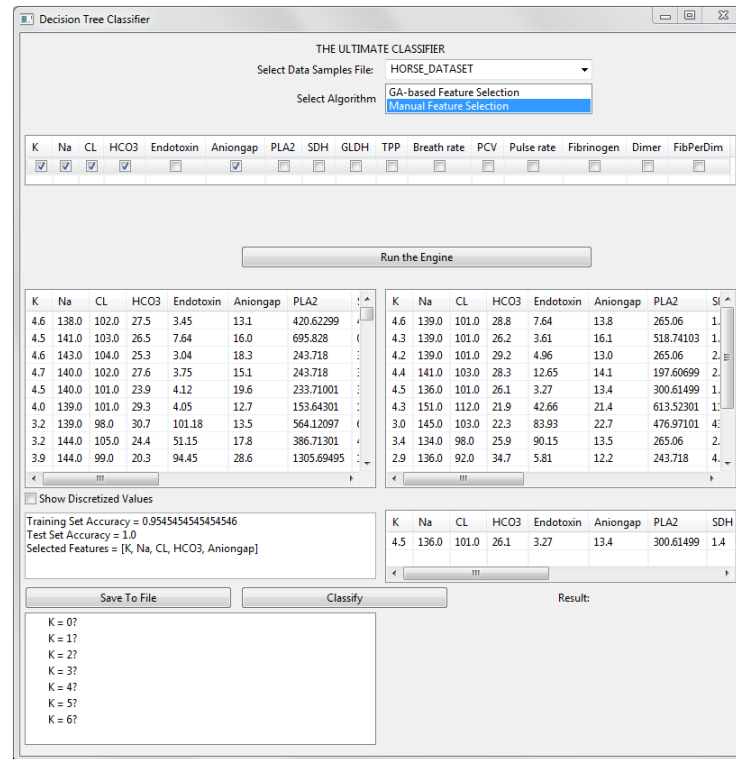


Figure 10.5: Decision Tree Construction with manual feature selection.

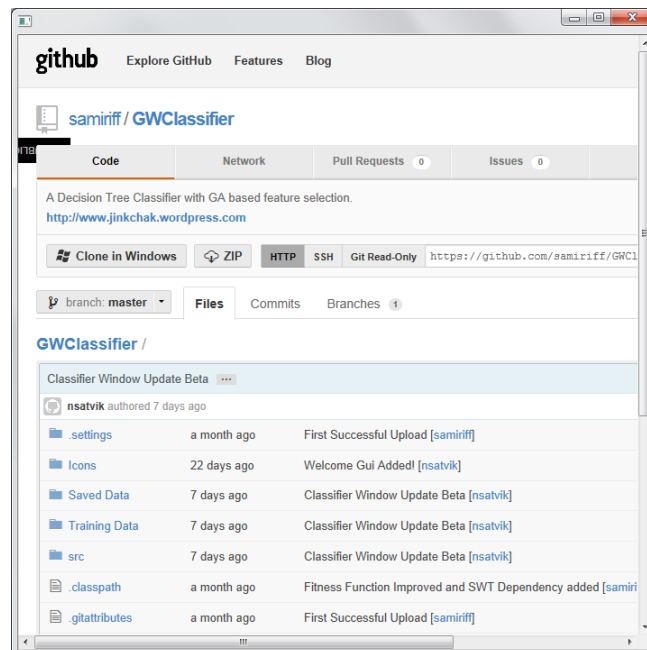


Figure 10.6: The project source code on github public repository.