# ACKNOWLEDGEMENTS

## ABSTRACT

Recent studies have shown that the significance of data-driven decision making and the notion of "Big-Data" are recognized broadly. Advances in data collection and storage technologies have made large organizations to accumulate vast volumes of data with ease in digital format. But extracting useful information from such massive volumes of data using traditional data analysis techniques is quite challenging. The various factors that impede the progress of "Big-Data" analysis in all phases during its analysis are heterogeneity, scale, complexity, timeliness and privacy problems. Also, it is important to decide which data to retain and which to discard during the data acquisition phase itself. Data Mining is a solution to overcome much of the challenges faced by the traditional data analysis techniques in such situations. Data Mining is a technology that combines traditional data analysis methods with sophisticated algorithms for processing large volumes of data. It is a process of automatically discovering novel and useful patterns that are hidden in large data repositories.

Traditional data mining techniques have been developed mainly for structured data types. But much of the data available today is not in the native structured format; for example, tweets and blogs have large pieces of unstructured text data; Images and videos are structured for storage and display, but not for semantic content and search; Web pages have data of different formats embedded in them and hence are semi-structured. Transforming such contents into a structured format for later analysis is itself a major challenge.

The objective of this thesis is to design and implement algorithms for improving subject based classification of web page and medical image data sets using data mining techniques. Classifying

a web page/a medical image into one of a pre-defined category is known as web page/medical image classification. With millions of new web pages being added each day the volume of WWW will defeat any conceivable team of human classifiers. Hence automatic classification of web pages is required. Such automated tools also help the search engines to construct web directories and hence make a relevant and quick retrieval of information for the user query. The ever increasing amounts of patient data in the form of medical images, imposes new challenges to clinical routine, such as diagnosis, treatment and monitoring. Medical data mining refers to the process of transforming raw imaging data using knowledge-based data mining algorithms into clinically relevant information. The target mining model enables a physician to spend less time in spending on the image volumes to extract the clinical information in it, while improving the diagnostic accuracy.

The data sets are classified using the content embedded in them and hence the present classification framework is subject-based. Algorithms for the various pre-processing steps namely feature extraction, feature selection and discretization are designed and implemented to improve the predictive accuracy of the classification model. Two new feature selection methods are designed and implemented in this thesis. First is a hybrid model using CFS (Correlation Based Feature Selection) and Decision Tree. Second is a novel method using Ward's minimum variance measure to identify clusters of redundant features. This research also presents two new classification models for web page and image data sets. The first model is a probabilistic method (Probabilistic Web Page Classifier and Probabilistic Medical Image Classifier). Experimental analysis is done for both binary-class and multi-class classification.

The results of the research identify that the present feature selection method helps in improving the predictive accuracy and modeling the classifier with less number of features. The performance of many of the binary-class classifiers is better in the discrete domain than in continuous domain. The predictive accuracy of the present MKNN classification framework with both web page and image data sets is better than the traditional KNN for both binary-class and multi-class classification. The performance of the present PWPC and PMIC is also better than many of the existing classifiers for both binary-class and multi-class classification. Its performance depends on the discretization algorithm that is run prior to classification.

# TABLE OF CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

WWW   –   World Wide Web

WPC   –   Web Page Classification

HTML   –   Hyper Text Markup Language

TF   –   Term Frequency

DF   –   Document Frequency

IDF   –   Inverse Document Frequency

MIC   –   Medical Image Classification

KDD   –   Knowledge Discovery in Databases

MDL   –   Minimum Description Length

ARFF   –   Attribute Relation File Format

NB   –   Naïve Bayes

ML   –   Machine Learning

KNN   –   K Nearest Neighbor

DT   –   Decision Tree

CFS   –   Correlation Based Feature Selection

URL   –   Uniform Resource Locator

PCA   –   Principal Component Analysis

UCI   –   University of Irvine

SVM   –   Support Vector Machine

EM   –   Expectation Maximization

CART   –   Classification and Regression Trees

GA   –   Genetic Algorithms

BoW     -     Bag of Words

AI     –     Artificial Intelligence

PSO     –     Particle Swarm Optimization

FP-tree     –     Frequent Pattern tree

MKNN     –     Modified K Nearest Neighbor

ROC     –     Receiver Operating Characteristics Curve

AUC     –     Area under the Curve

TP     –     True Positive

TPR     –     True Positive Rate

FP     -     False Positive

FPR     –     False Positive Rate

MLP     –     Multilayer Perceptron

RBF     –     Radial Basis Function