# 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. I would like to take this opportunity to thank them all.

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

I would like to thank **Dr.G.Shobha**, Head of Department, Computer Science & Engineering, R.V.C.E, Bengaluru, for her valuable suggestions and expert advice.

I deeply express my sincere gratitude to my guide Mrs. **Shanta Rangaswamy,** Asst. Prof**,** Department of CSE, R.V.C.E, Bengaluru, for her able guidance, regular source of encouragement and assistance throughout this project

I thank my Parents, and all the Faculty members of Department of Computer Science & Engineering for their constant support and encouragement.

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

# Abstract

Innovation in the field of green technology by analyzing and developing efficient algorithms to reduce CPU utilization is the need of the hour. Managing the vast amount of online information is an important step towards this need. In the field of machine learning, binary classification algorithms are one such technique to process data. This data can have multiple sources and representations which influence the choice of classification technique. Data from web pages is unstructured in nature, and therefore requires many different processes and stages in classification. The aim of this project is two-fold: first, to compare important binary classification algorithms' performance in classifying samples of a large dataset from dynamic sources in the domain of online course websites, to identify valid online courses and secondly to implement the most precise algorithm in building an offline repository of all these courses for a user to access through an interactive front end.

The project involves multiple techniques such as crawling, URL filtering, metadata analysis, text extraction and processing while incorporating novel features such as social tagging from social network websites such as Delicious. A sample dataset covering all corner cases and representative examples is used for training a model classifier, which is then trained and tested using various methods like bootstrapping and cross validation. Finally the efficiency of the model is compared across five algorithms : Naive Bayes, K Nearest Neighbours, Decision Trees, Logistic Regression and Random Forest on datasets of increasing size and the most efficient one is used for classifying all the data collected from crawling. The relevant courses are put into an offline database made available to the users by an interactive and fresh user interface.

The parameter used to evaluate the algorithms was the accuracy of predictions on the test set, namely precision. It was observed that Random Forest algorithm outperforms Naive Bayes, K Nearest Neighbours, Decision Tree and Logistic Regression on this parameter. The accuracy observed upon evaluation of the test set for this classifier is 92.7%. Over 30,000 URLs were crawled and pre-processed using this classifier to create a repository of over 10,000 valid online courses.

# Table of Contents

# List of abbreviations

| | | |
|---|---|---|
| AI | : | Artificial Intelligence |
| ANFIS | : | Adaptive Neuro Fuzzy Inference System |
| FoCUS | : | Crawler Under Supervision |
| DOM | : | Document Object Module |
| NE | : | Named Entities |
| HMM | : | Hidden Markov Model |
| ME | : | Maximum Entropy |
| TBL | : | Transformation Based error-driven Learning |
| SVM | : | Support Vector Machine |
| JDBC | : | Java Database Connection |
| UDP | : | User Datagram Protocol |
| TCP | : | Transmission Control Protocol |
| OS | : | Operating System |
| HCL | : | Hardware Compatibility List |
| CPU | : | Central Processing Unit |
| MIPS | : | Million Instructions Per Second |
| GUI | : | Graphical User Interface |
| CSV | : | Common Separated Value |
| DBF | : | Database File |
| SQL | : | Structured Query Language |
| SDLC | : | Software Development Life Cycle |
| DFD | : | Data Flow Diagram |
| CLD | : | Context Level Diagram |
| API | : | Application Programming Interface |
| REST | : | Representational state transfer |

# List of Tables

# List of Figures