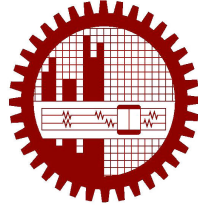


Knowledge Discovery From Academic Data Using Data Mining Technique



By

Md. Mostafizur Rahman

Student Id: 1005017

Sabid Bin Habib

Student Id: 1005102

A thesis submitted to the Department of Computer Science And Engineering
in partial fulfillment of the requirements for the degree of Bachelor of Science
in Computer Science and Engineering

Supervised by

Dr. Abu Sayed Md. Latiful Hoque

Professor

Department of Computer Science And Engineering
Bangladesh University Of Engineering And Technology

February 29, 2016

Declaration of Authorship

We, Md.Mostafizur Rahman and Sabid Bin Habib, declare that this thesis titled, "Knowledge Discovery From Academic Data Using Data Mining Technique" and the work presented in it are our own. We confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where We have consulted the published work of others, this is always clearly attributed.
- Where we have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely our own work.
- We have acknowledged all main sources of help.
- Where the thesis is based on work done by us, We have made clear exactly what was done by others and what we have contributed.

Name:

Signature and Date:

Name:

Signature and Date:

Abstract

Educational Data mining is an emerging research area in Data Mining. This area associates to discovering relevant and useful knowledge from academic records using various data mining techniques. Every educational institute wants to provide maximum facility to the students and faculties and expects maximum output from the students and faculties. It is very important to regularly monitor the students' performances for better accountability. More importantly, the format of the data records can be modified or changed to apply knowledge discovery techniques and gather some statistical output to analyze the performance of the students and faculties as well as the departments.

In this thesis work, we present a knowledge discovery process in BUET Institutional Information System student data using classification rules. We derived a technique to find the overall performance of a student and a department. Before applying classification, we preprocessed the given data set to remove unnecessary and irrelevant data and columns. After combining relevant columns and removing impurities from data, we designed our training data with our five different class labels. We apply the ID3 classification algorithm on our designed test data according to the decision tree made from our training data. We analyze the final output to gather some statistical results for better visualization of the findings. These results hope to be helpful to analyze the performance of individual students and departments and decision making for top management.

Acknowledgements

Special thanks to our very helpful supervisor ASM Latiful Hoque for his outstanding help throughout the full time. We express out gratitude to the authority of BIIS for providing the Sample data for our thesis. We also express our gratitude to our beloved friends and family members who helped us completing this thesis.

Contents

Declaration of Authorship	i
Abstract	ii
Acknowledgements	iii
1 Introduction	1
1.1 Problem Definition	1
1.2 Scope of the Work	2
1.3 Objectives	3
1.4 Thesis organization	3
2 Literature Study	5
2.1 Knowledge Discovery Steps	5
2.1.1 Understanding the Application Domain	6
2.1.2 Selection of Dataset	6
2.1.3 Data Cleaning	6
2.1.4 Data transformation	7
2.1.5 Finding interesting features in the database	7
2.1.6 Selection of data mining task	7
2.1.7 Selection of Data mining method	7
2.1.8 Data mining	8
2.1.9 Pattern evaluation	8
2.1.10 Knowledge consolidation	8
2.2 Data Mining Concepts	9

2.2.1	Database Data	9
2.2.2	Data Warehouses	10
2.2.3	Transactional Data	10
2.2.4	Other Kinds of Data	11
2.3	Preprocessing	13
2.3.1	Data Cleaning	13
2.3.2	Data Integration	15
2.3.3	Data Transformation	15
2.3.4	Data Transformation by Normalization	16
2.3.5	Data Discretization	16
	Discretization by Binning	17
	Discretization by Histogram Analysis	17
	Discretization by Cluster, Decision Tree, and Correlation Analyses	18
2.4	Classification	19
2.4.1	Basic Concept	19
2.4.2	General Approach to Classification	19
2.4.3	Decision Tree Induction	21
2.4.4	Decision Tree Induction Algorithm	21
2.4.5	Attributes Selection Measures	22
	Information Gain	24
	Gain Ratio	25
	Gini Index	26
	Other Attribute Selection Measures	27
3	Analysis of BIIS Data	28
3.1	Scope	28
3.2	Database Structure	29
3.3	Problems in Existing Structure	30

4	Preprocessing	33
4.1	Technique And Design	33
4.1.1	Data Cleaning	33
4.1.2	Data Normalization	35
4.1.3	Data Transformation	36
4.2	Results	36
5	Classification	39
5.1	Decision Tree	39
5.2	Algorithm	40
5.3	Result of Classification	41
5.4	Statistical Analysis	41
5.5	Result of Statistical Analysis	44
5.5.1	Department wise Performance	44
	EEE Department	44
	CSE Department	44
	IPE Department	46
	ME Department	46
	CE Department	48
	MME Department	48
	CHE Department	50
	NAME Department	50
	URP Department	52
	ARCH Department	52
	WRE Department	54
	Overall Department wise Performance	54
5.5.2	Impact of Gender on Performance	55
5.5.3	Impact of Hall Status on Performance	55
5.5.4	Impact of Class Attendance Marks	56

5.5.5	Impact of Class Test Marks on Performance	58
5.5.6	Impact of CGPA on Overall Performance	58
5.5.7	Impact of Credit Completion	59
6	Conclusion	60
6.1	Summary of Thesis	60
6.2	General Findings	60
6.3	Future Works	61
A	Data Samples	62
B	More Statistical Analysis Results	64
B.1	Department wise Different Grade Comparison	64
C	Details of Department Name	66

List of Figures

2.1	Knowledge discovery steps	6
2.2	Data Mining	9
2.3	Classification	20
2.4	A Decision Tree	22
5.1	A Decision Tree From Training Data	41
5.2	Performance of EEE Department	45
5.3	Performance of CSE Department	45
5.4	Performance of IPE Department	46
5.5	Performance of ME Department	47
5.6	Performance of CE Department	48
5.7	Performance of MME Department	49
5.8	Performance of CHE Department	50
5.9	Performance of NAME Department	51
5.10	Performance of URP Department	52
5.11	Performance of ARCH Department	53
5.12	Performance of WRE Department	54
5.13	Overall Department wise Performance	55
5.14	Impact of Gender on Performance	56
5.15	Impact of Hall Status on Performance	57
5.16	Impact of Attendance on Performance	57
5.17	Impact of Class test marks on Performance	58
5.18	Impact of Credit completion on Performance	59

A.1	BIIS Data	63
A.2	Processed BIIS Data	63
B.1	Different Grades for Dept	65
B.2	Different Grades for Subjects	65

List of Tables

3.1	Multiple entries for single student in the Universal Database . .	32
3.2	Blank attributes in Database	32
3.3	Rounded up value of CGPA and GPA	32
3.4	Different range of values	32
4.1	Blank attributes in the Universal Database	37
4.2	Universal database after filling up blank attributes	37
4.3	Redundant attributes in Database	37
4.4	Database after cleaning of redundant attributes	37
4.5	Database with binary values	37
4.6	Database after conversion of binary attributes to numeric value	38
4.7	Database before Data cleaning	38
4.8	Letter Grade and it's respective percentage of number	38
4.9	Database after addition of CourseCreditHour attribute	38
5.1	Training Data Sample	40
5.2	Test Data Sample	43
5.3	Final Output Sample	43
5.4	Performance of EEE department	44
5.5	Performance of CSE Department	44
5.6	Performance of IPE Department	46
5.7	Performance of ME Department	47
5.8	Performance of CE Department	48
5.9	Performance of MME Department	49

5.10 Performance of CHE Department	50
5.11 Performance of NAME Department	51
5.12 Performance of URP Department	52
5.13 Performance of ARCH Department	53
5.14 Performance of WRE Department	54
5.15 Impact of Gender on Performance	55
5.16 Impact of Hall Status on Performance	56
5.17 Impact of Hall Status on Performance	57
5.18 Impact of Class Test Marks on Performance	58

Chapter 1

Introduction

Bangladesh University of Engineering and Technology(*BUET*) is a renowned Technological University of *Bangladesh*. Brilliant students around the country get admitted here for quality education. The huge amount of BIIS data represent the states of all the students of BUET in academic education. For a large educational institute like public university which generates large volumes of data, it requires an efficient way to apply data mining techniques for obtaining knowledge on the development and performance improvement of academic activities. The knowledge acquired from BIIS can be sufficient to look for answers to such questions as: which factors determine better or worse academic performance of students? Which factors in a Department can be crucial for quality education? Concepts and techniques of data mining are essential to discover the hidden knowledge from large datasets [1] .

1.1 Problem Definition

Every year BUET enrolls almost 1000 brilliant Engineering minded student into 11 different departments from all over the country and simultaneously almost 1000 students graduate from BUET. The BIIS data holds all the academic information of individual student in detail. Statistics from BIIS data shows that academic performance of a student varies from department to department. After close observation, it has been revealed that department might not be the only

factor behind the performance of a student. There might be some other reasons such as CGPA, Hall Status, Gender, Class Test Marks, Attendance Status which affect a student's performance. Only statistical analysis is not sufficient for finding out the knowledge from the BIIS data. The hidden knowledge inside the institutional academic and personal data of students is necessary to find out the possible effects on any student's academic performance. That is why knowledge discovery and data mining from academic data is essential for educational institution like BUET to improve academic performance of students as well as reshaping the decision makings for the betterment of the institution. Data mining techniques are very effective for discovering the hidden knowledge from educational data and applying it properly for the decision making. But all the data mining techniques can not be applied directly on academic data because of complex structure. This requires rigorous preprocessing. Bringing all the relevant data in a useful scope and applying classification methods on them are other problems of this research.

1.2 Scope of the Work

The BIIS data represents the detailed academic performance of an individual student throughout the undergrad period. These students are under 5 different faculties. These five faculties are Faculty of Architecture & Planning, Faculty of Civil Engineering, Faculty of Mechanical Engineering, Faculty of Engineering and Faculty of Electrical & Electronic Engineering. There are 11 different departments under these faculties which are Department of Electrical & Electronic Engineering, Department of Computer Science & Engineering, Department of Civil Engineering, Department of Mechanical Engineering, Department of Architecture, Department of Water Resource Engineering, Department of Naval Architecture & Marine Engineering, Department of Chemical Engineering, Department of Materials & Metallurgical Engineering, Department of Urban & Regional

Planning and Department of Industrial & Production Engineering. [2] The scope of knowledge discovery from BIIS data is immense in context of undergraduate students. In this research BIIS data of 10 already graduated batches have been used which represents almost 10 thousand students maintaining the privacy of the data. No current student has been brought under this research.

1.3 Objectives

Every year, a number of brilliant minded students are admitted into BUET, but at the end of their academic period, all of them can not utilize the full potential, it affects their academic performance. The main objectives of this research study are

- To find out the effect of different factors to the performance of a student.
- To discover knowledge of students' academic performance and personal statistics through the impact of different assessment and factors e.g. Class Test, Attendance, CGPA, Credit Completion etc.
- Compare the impact of different factors e.g. Hall Status, Gender, Class Test, Attendance, Departments on the academic performances of a student

1.4 Thesis organization

We have developed a technique to discover knowledge using ID3 classification algorithm from institutional data of students who have completed their undergraduate in the department of CSE, BUET.

All the literature studies e.g., preprocessing of raw data, preliminaries of Knowledge Discovery and Data Mining, Building decision tree, ID3 algorithm and related works have been elaborated in chapter 2.

Description of the raw academic and analysis have been described in chapter 3. Basically, the raw BIIS data is not suitable for implementation of necessary procedures, the status of the data has been elaborated in chapter 3.

The existing raw data needs to be rigorously preprocessed before the implementation of classification. The step by step methodologies and the result of this preprocess has been described in chapter 4.

On chapter 5, the implementation of ID3 classification algorithm on the test data using training set has been elaborated. The result of the procedures taken has also been depicted in this section in detail.

Finally, we have illustrated the summary of the findings along with quantitative analysis. We have also encompassed the scope of the extension of this research work by illustrating some significant future works in chapter 6.

Chapter 2

Literature Study

2.1 Knowledge Discovery Steps

Knowledge Discovery is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data

[4]. With the emphasis on collecting data increasing around the world, there is an urgent need for a new generation of different techniques, methods and algorithms to assist researchers, analysts, decision makers and managers in extracting useful patterns from the rapidly growing volumes of data. These techniques and tools are the subject of the emerging field of knowledge discovery in databases. Knowledge Discovery and Data Mining is an interdisciplinary area focusing upon methodologies for extracting useful knowledge from data. Knowledge Discovery and Data Mining is an interdisciplinary area focusing upon methodologies for extracting useful knowledge from data. According to [5], Though Knowledge Discovery is used synonymously to represent data mining, both these are actually different. Some preprocessing steps before data mining and post processing steps after data mining are to be completed to transform the raw data as useful knowledge.

Knowledge Discovery is an iterative process that transforms raw data into useful information. Different steps of Knowledge Discovery in Databases are discussed in [6].

Figure 2.1 shows the general process of knowledge discovery.

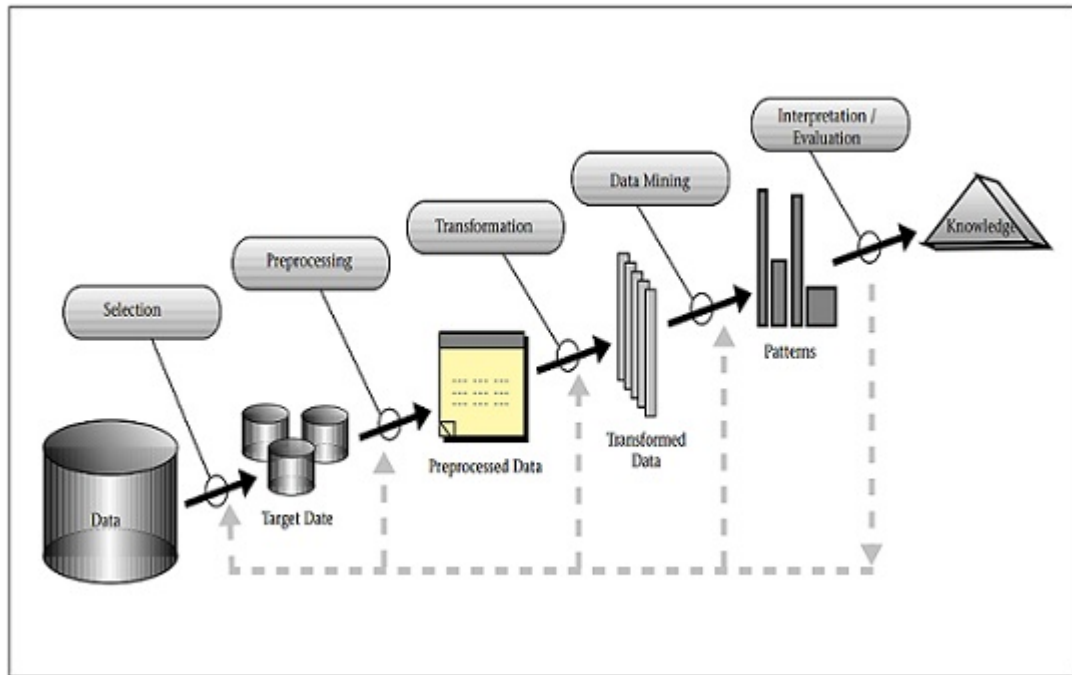


FIGURE 2.1: Knowledge discovery steps

2.1.1 Understanding the Application Domain

The first step is understanding requirements. It is needed to have a clear understanding about the application domain and your objectives. It should be also known whether the data is going to be described or information is predicted.

2.1.2 Selection of Dataset

Data mining is done on current or past records. Thus, a data set or subset of data should be selected, in other words data samples, on which you need to perform data analysis and get useful knowledge. There should be enough quantity of data to perform data mining.

2.1.3 Data Cleaning

Data cleaning is the step where noise and irrelevant data are removed from the large data set. This is a very important preprocessing step because the outcome would be dependent on the quality of selected data. As part of data

cleaning, duplicate records might have to be removed, logically correct values for missing records might have to be entered, unnecessary data fields might have to be removed, data format standardized, update data in a timely manner and so on.

2.1.4 Data transformation

With the help of dimensionality reduction or transformation methods, the number of effective variables is reduced and only useful features are selected to depict data more efficiently based on the goal of the task. In short, data is transformed into appropriate form making it ready for data mining step.

2.1.5 Finding interesting features in the database

This step is extremely important in the field of International Studies. Researchers and practitioners with different backgrounds and different languages may work on a given database, getting different results. Each group may consider different attributes in doing so.

2.1.6 Selection of data mining task

Based on the objective of data mining, appropriate task is selected. Some common data mining tasks are classification, clustering, association rule discovery, sequential pattern discovery, regression and deviation detection. You can choose any of these tasks based on whether you need to predict information or describe information.

2.1.7 Selection of Data mining method

Appropriate method(s) is to be selected for looking for patterns from the data. You need to decide the model and parameters that might be appropriate for

the method. Some popular data mining methods are decision trees and rules, relational learning models, example based methods etc.

2.1.8 Data mining

Data mining is the actual search for patterns from the data available using the selected data mining method.

2.1.9 Pattern evaluation

This is a post processing step in KDD which interprets mined patterns and relationships. If the pattern evaluated is not useful, then the process might again start from any of the previous steps, thus making KDD an iterative process.

2.1.10 Knowledge consolidation

This is the final step in Knowledge Discovery. The knowledge discovered is consolidated and represented to the user in a simple and easy to understand format. Mostly, visualization techniques are being used to make users understand and interpret information.

Though these are the main steps in any Knowledge Discovery process, some of the steps could be done combined during the actual process. For example, considering the convenience, data selection and data transformation can be combined together. Even after presenting knowledge to the user, new data can be added to the data set or mining can be further refined or a different data mining method can be chosen to get more accurate results. Thus, Knowledge Discovery is completely an iterative process.

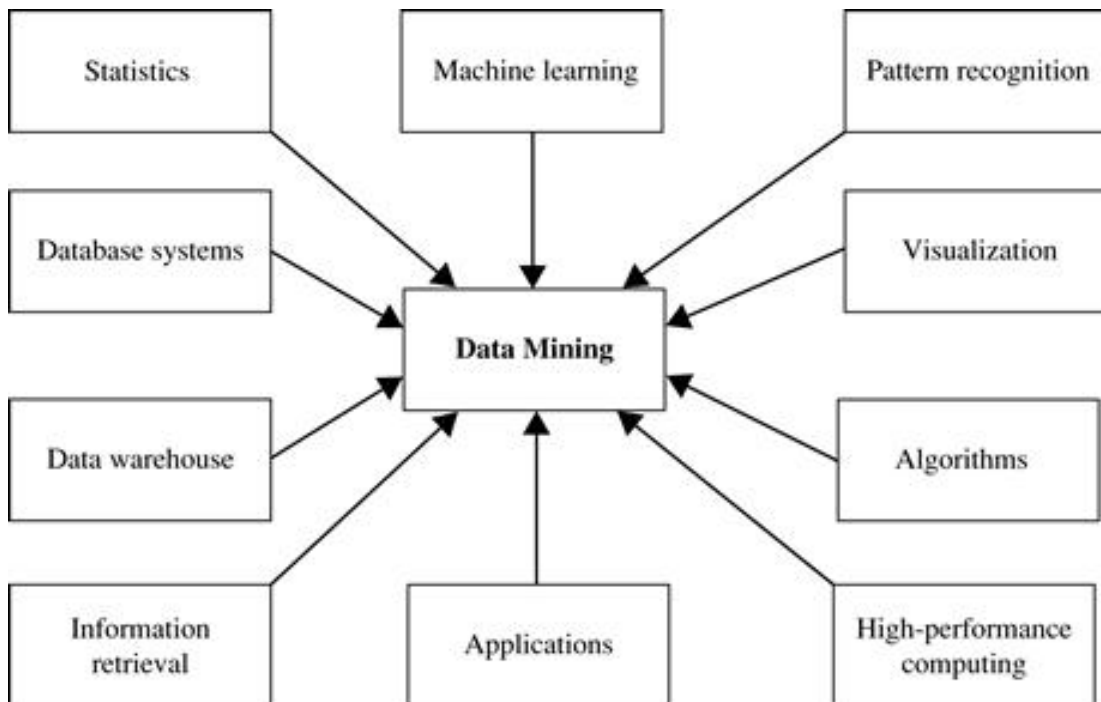


FIGURE 2.2: Data Mining Concept

2.2 Data Mining Concepts

Data mining is the process of discovering interesting patterns and knowledge from large amounts of data. The data sources can include databases, data warehouses, the Web, other information repositories, or data that are streamed into the system dynamically.

Figure 2.2 shows the general parts of data mining. As a general technology, data mining can be applied to any kind of data as long as the data are meaningful for a target application. The most basic forms of data for mining applications are database data, data warehouse data, and transactional data.

2.2.1 Database Data

A database system, also called a database management system (DBMS), consists of a collection of interrelated data, known as a database, and a set of software programs to manage and access the data. The software programs provide mechanisms for defining database structures and data storage; for specifying

and managing concurrent, shared, or distributed data access; and for ensuring consistency and security of the information stored despite system crashes or attempts at unauthorized access.

A relational database is a collection of tables, each of which is assigned a unique name. Each table consists of a set of attributes (columns or fields) and usually stores a large set of tuples (records or rows). Each tuple in a relational table represents an object identified by a unique key and described by a set of attribute values. A semantic data model, such as an entity-relationship (ER) data model, is often constructed for relational databases. An ER data model represents the database as a set of entities and their relationships.

2.2.2 Data Warehouses

A data warehouse is usually modeled by a multidimensional data structure, called a data cube, in which each dimension corresponds to an attribute or a set of attributes in the schema, and each cell stores the value of some aggregate measure such as count or sum. A data cube provides a multidimensional view of data and allows the precomputation and fast access of summarized data.

2.2.3 Transactional Data

In general, each record in a transactional database captures a transaction, such as a customer's purchase, a flight booking, or a user's clicks on a web page. A transaction typically includes a unique transaction identity number (trans_ID) and a list of the items making up the transaction, such as the items purchased in the transaction. A transactional database may have additional tables, which contain other information related to the transactions, such as item description, information about the salesperson or the branch, and so on.

2.2.4 Other Kinds of Data

Besides relational database data, data warehouse data, and transaction data, there are many other kinds of data that have versatile forms and structures and rather different semantic meanings. Such kinds of data can be seen in many applications: time-related or sequence data (e.g., historical records, stock exchange data, and timeseries and biological sequence data), data streams (e.g., video surveillance and sensor data, which are continuously transmitted), spatial data (e.g., maps), engineering design data (e.g., the design of buildings, system components, or integrated circuits), hypertext and multimedia data (including text, image, video, and audio data), graph and networked data (e.g., social and information networks), and the Web (a huge, widely distributed information repository made available by the Internet). These applications bring about new challenges, like how to handle data carrying special structures (e.g., sequences, trees, graphs, and networks) and specific semantics (such as ordering, image, audio and video contents, and connectivity), and how to mine patterns that carry rich structures and semantics.

Various kinds of knowledge can be mined from these kinds of data. Here, we list just a few. Regarding temporal data, for instance, we can mine banking data for changing trends, which may aid in the scheduling of bank tellers according to the volume of customer traffic. Stock exchange data can be mined to uncover trends that could help you plan investment strategies.

We could mine computer network data streams to detect intrusions based on the anomaly of message flows, which may be discovered by clustering, dynamic construction of stream models or by comparing the current frequent patterns with those at a previous time. With spatial data, we may look for patterns that describe changes in metropolitan poverty rates based on city distances from major highways. The relationships among a set of spatial objects can be examined in order to discover which subsets of objects are spatially

autocorrelated or associated. By mining text data, such as literature on data mining from the past ten years, we can identify the evolution of hot topics in the field. By mining user comments on products (which are often submitted as short text messages), we can assess customer sentiments and understand how well a product is embraced by a market. From multimedia data, we can mine images to identify objects and classify them by assigning semantic labels or tags. By mining video data of a hockey game, we can detect video sequences corresponding to goals. Web mining can help us learn about the distribution of information on the WWW in general, characterize and classify web pages, and uncover web dynamics and the association and other relationships among different web pages, users, communities, and web-based activities.

It is important to keep in mind that, in many applications, multiple types of data are present. For example, in web mining, there often exist text data and multimedia data (e.g., pictures and videos) on web pages, graph data like web graphs, and map data on some web sites. In bioinformatics, genomic sequences, biological networks, and 3-D spatial structures of genomes may co-exist for certain biological objects. Mining multiple data sources of complex data often leads to fruitful findings due to the mutual enhancement and consolidation of such multiple sources. On the other hand, it is also challenging because of the difficulties in data cleaning and data integration, as well as the complex interactions among the multiple sources of such data.

While such data require sophisticated facilities for efficient storage, retrieval, and updating, they also provide fertile ground and raise challenging research and implementation issues for data mining. Data mining on such data is an advanced topic. The methods involved are extensions of the basic techniques presented in this book.

2.3 Preprocessing

Data preprocessing describes any type of processing performed on raw data to prepare it for another processing procedure. Commonly used as a preliminary data mining practice, data preprocessing transforms the data into a format that will be more easily and effectively processed for the purpose of the user.

Raw data is highly susceptible to noise, missing values, and inconsistency. The quality of data affects the data mining results. In order to help improve the quality of the data and, consequently, of the mining results raw data is pre-processed so as to improve the efficiency and ease of the mining process. Data preprocessing is one of the most critical steps in a data mining process which deals with the preparation and transformation of the initial dataset. Data preprocessing methods are divided into following categories

2.3.1 Data Cleaning

Data that is to be analyzed by data mining techniques can be incomplete (lacking attribute values or certain attribute of interest, or containing only aggregate data), noisy (containing errors, or outlier values which deviate from the expected), and inconsistent (e.g., containing discrepancies in the department codes used to categorize items). [7] Incomplete, noisy and inconsistent data are common-place properties of large, real-world databases and data warehouses.

Therefore, a useful preprocessing step is to run data through some data cleaning routines. *Missing Values*: If there are tuples that have no recorded value for several attributes, then missing value can be filled in for attributes by the methods below

1. Ignore tuple if the class label is missing
2. Fill in missing values manually
3. Use global constant to fill in missing values

4. Use the most probable value to fill in the missing value

Inconsistent Data: There may be inconsistencies in the data recorded for some transactions. Some data inconsistencies may be corrected manually using external references. For example, errors made at data entry may be corrected by performing a paper trace. This may be coupled with routine designed to help correct the inconsistent use of codes. *Conversion: Nominal to Numeric* Sometimes it is more efficient for some programs to calculate numeric values than nominal ones. There are different strategies for conversions

- Binary to Numeric: E.g. Gender=M,F. Convert field with 0,1 values

Gender=M -> Gender=0

Gender=F -> Gender=1

- Ordered to Numeric: Ordered attributes (e.g. Grade) can be converted to numbers preserving natural order to allow comparisons,
e.g. A -> 3.75; B -> 3

2.3.2 Data Integration

Data analysis task can involve data integration, which involves combining data residing in different sources and providing users with a unified view of these data.[8] In case there are tuples which represent a single instance, those tuples can be combined using various methods. In this case there can be problems like attributes not matching or attributes missing. Using various methods, these problems can be overcome and data integration is implemented.

2.3.3 Data Transformation

In data transformation, the data are transformed or consolidated into forms appropriate for mining. Strategies for data transformation include the following

- **Smoothing**, which works to remove noise from the data. Techniques include binning, regression, and clustering.
- **Attribute construction** (or feature construction), where new attributes are constructed and added from the given set of attributes to help the mining process.
- **Aggregation**, where summary or aggregation operations are applied to the data. For example, the daily sales data may be aggregated so as to compute monthly and annual total amounts. This step is typically used in constructing a data cube for data analysis at multiple abstraction levels.
- **Normalization**, where the attribute data are scaled so as to fall within a smaller range, such as 1.0 to 1.0, or 0.0 to 1.0.
- **Discretization**, where the raw values of a numeric attribute (e.g., age) are replaced by interval labels (e.g., 0 to 10, 11 to 20, etc.) or conceptual

labels (e.g., youth, adult, senior). The labels, in turn, can be recursively organized into higher-level concepts, resulting in a concept hierarchy for the numeric attribute.

- **Concept hierarchy generation for nominal data**, where attributes such as street can be generalized to higher-level concepts, like city or country. Many hierarchies for nominal attributes are implicit within the database schema and can be automatically defined at the schema definition level.

2.3.4 Data Transformation by Normalization

The measurement unit used can affect the data analysis. For example, changing measurement units from meters to inches for height, or from kilograms to pounds for weight, may lead to very different results. In general, expressing an attribute in smaller units will lead to a larger range for that attribute, and thus tend to give such an attribute greater effect or "weight". To help avoid dependence on the choice of measurement units, the data should be normalized or standardized. This involves transforming the data to fall within a smaller or common range such as $[1, 2]$ or $[0.0, 1.0]$.

There are many methods for data normalization. Common methods are

- Minmax Normalization
- zscore Normalization
- Normalization by Decimal scaling

2.3.5 Data Discretization

Data discretization transforms numeric data by mapping values to interval or concept labels. Such methods can be used to automatically generate concept hierarchies for the data, which allows for mining at multiple levels of granularity. Discretization techniques include binning, histogram analysis, cluster

analysis, decision tree analysis, and correlation analysis. For nominal data, **concept hierarchies** may be generated based on schema definitions as well as the number of distinct values per attribute.

Discretization by Binning

Binning is a top-down splitting technique based on a specified number of bins. For example, attribute values can be discretized by applying equal-width or equal-frequency binning, and then replacing each bin value by the bin mean or median, as in smoothing by bin means or smoothing by bin medians, respectively. These techniques can be applied recursively to the resulting partitions to generate concept hierarchies.

Binning does not use class information and is therefore an unsupervised discretization technique. It is sensitive to the user-specified number of bins, as well as the presence of outliers.

Discretization by Histogram Analysis

Like binning, histogram analysis is an unsupervised discretization technique because it does not use class information. A histogram partitions the values of an attribute, A , into disjoint ranges called buckets or bins. Various partitioning rules can be used to define histograms. In an equal-width histogram, for example, the values are partitioned into equal-size partitions or ranges (e.g., earlier in Figure 3.8 for price, where each bucket has a width of \$10). With an equal-frequency histogram, the values are partitioned so that, ideally, each partition contains the same number of data tuples. The histogram analysis algorithm can be applied recursively to each partition in order to automatically generate a multilevel concept hierarchy, with the procedure terminating once a prespecified number of concept levels has been reached. A minimum interval size can also be used per level to control the recursive procedure. This specifies the minimum width of a partition, or the minimum number of values for each

partition at each level. Histograms can also be partitioned based on cluster analysis of the data distribution, as described next.

Discretization by Cluster, Decision Tree, and Correlation Analyses

Clustering, decision tree analysis, and correlation analysis can be used for data discretization. We briefly study each of these approaches.

Cluster analysis is a popular data discretization method. A clustering algorithm can be applied to discretize a numeric attribute, A , by partitioning the values of A into clusters or groups. A into consideration, as well as the closeness of data points, and therefore is able to produce high-quality discretization results.

Clustering can be used to generate a concept hierarchy for A by following either a top-down splitting strategy or a bottom-up merging strategy, where each cluster forms a node of the concept hierarchy. In the former, each initial cluster or partition may be further decomposed into several subclusters, forming a lower level of the hierarchy. In the latter, clusters are formed by repeatedly grouping neighboring clusters in order to form higherlevel concepts.

Measures of correlation can be used for discretization. ChiMerge is a χ^2 -based discretization method. The discretization methods that we have studied up to this point have all employed a top-down, splitting strategy. This contrasts with ChiMerge, which employs a bottom-up approach by finding the best neighboring intervals and then merging them to form larger intervals, recursively. As with decision tree analysis, ChiMerge is supervised in that it uses class information. The basic notion is that for accurate discretization, the relative class frequencies should be fairly consistent within an interval. Therefore, if two adjacent intervals have a very similar distribution of classes, then the intervals can be merged. Otherwise, they should remain separate.

2.4 Classification

2.4.1 Basic Concept

Classification is a form of data analysis that extracts models describing important data classes. Such models, called classifiers, predict categorical (discrete, unordered) class labels. For example, we can build a classification model to categorize bank loan applications as either safe or risky. Such analysis can help provide us with a better understanding of the data at large. Many classification methods have been proposed by researchers in machine learning, pattern recognition, and statistics. Most algorithms are memory resident, typically assuming a small data size. Recent data mining research has built on such work, developing scalable classification and prediction techniques capable of handling large amounts of disk-resident data. Classification has numerous applications, including fraud detection, target marketing, performance prediction, manufacturing, and medical diagnosis.

The data analysis task is classification, where a model or classifier is constructed to predict class (categorical) labels. This model is a predictor. Regression analysis is a statistical methodology that is most often used for numeric prediction; hence the two terms tend to be used synonymously, although other methods for numeric prediction exist. Classification and numeric prediction are the two major types of prediction problems.

2.4.2 General Approach to Classification

Data classification is a two-step process, consisting of a learning step (where a classification model is constructed) and a classification step (where the model is used to predict class labels for given data).

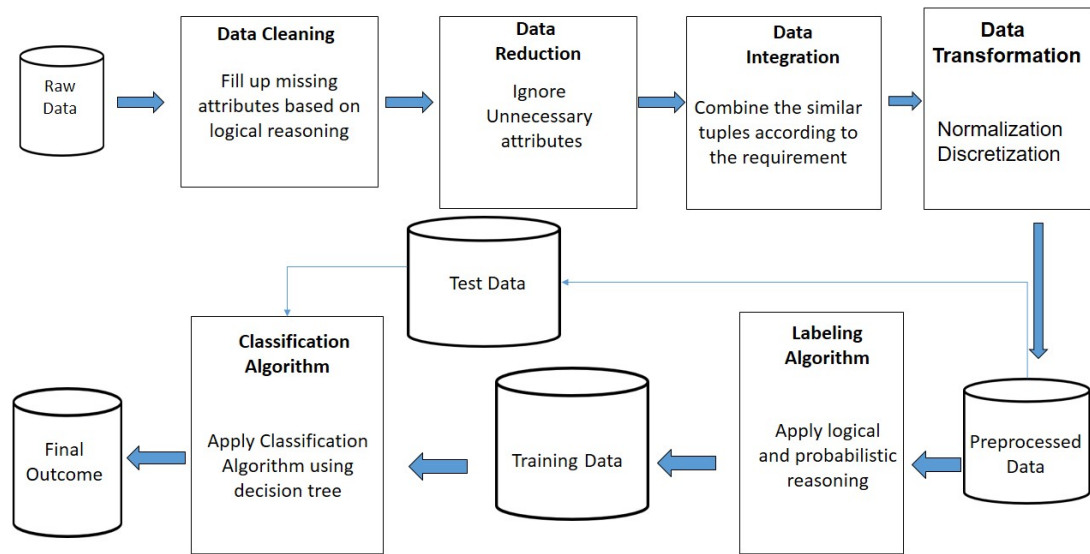


FIGURE 2.3: General Process of Classification

In the first step, a classifier is built describing a predetermined set of data classes or concepts. This is the learning step (or training phase), where a classification algorithm builds the classifier by analyzing or "learning from" a training set made up of database tuples and their associated class labels. A tuple, X , is represented by an n -dimensional attribute vector, $X = (x_1, x_2, \dots, x_n)$ depicting n measurements made on the tuple from n database attributes, respectively, A_1, A_2, \dots, A_n .¹ Each tuple, X , is assumed to belong to a predefined class as determined by another database attribute called the class label attribute. The class label attribute is discrete-valued and unordered. It is categorical (or nominal) in that each value serves as a category or class. The individual tuples making up the training set are referred to as training tuples and are randomly sampled from the database under analysis. In the context of classification, data tuples can be referred to as samples, examples, instances, data points, or objects.²

¹Each attribute represents a "feature" of X . Hence, the pattern recognition literature uses the term feature vector rather than attribute vector.

²In the machine learning literature, training tuples are commonly referred to as training samples. Throughout this text, we prefer to use the term tuples instead of samples.

Figure 2.3 shows the general procedure of classification.

2.4.3 Decision Tree Induction

Decision tree induction is the learning of decision trees from class-labeled training tuples. A decision tree is a flowchart-like tree structure, where each internal node (nonleaf node) denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node (or terminal node) holds a class label. The topmost node in a tree is the root node. The decision tree in Figure 2.4 is a tree for the concept `buy_computer` that indicates whether a customer at a company is likely to buy a computer or not. Each internal node represents a test on an attribute. Each leaf node represents a class. The benefits of having a decision tree are

- It does not require any domain knowledge.
- is easy to comprehend.
- The learning and classification steps of a decision tree are simple and fast.

2.4.4 Decision Tree Induction Algorithm

A machine researcher named J. Ross Quinlan in 1980 developed a decision tree algorithm known as ID3 (Iterative Dichotomiser). Later, he presented C4.5, which was the successor of ID3. ID3 and C4.5 adopt a greedy approach. In this algorithm, there is no backtracking; the trees are constructed in a top-down recursive divide-and-conquer manner.

Algorithm 1 is the general approach to build a decision tree as described in [3].

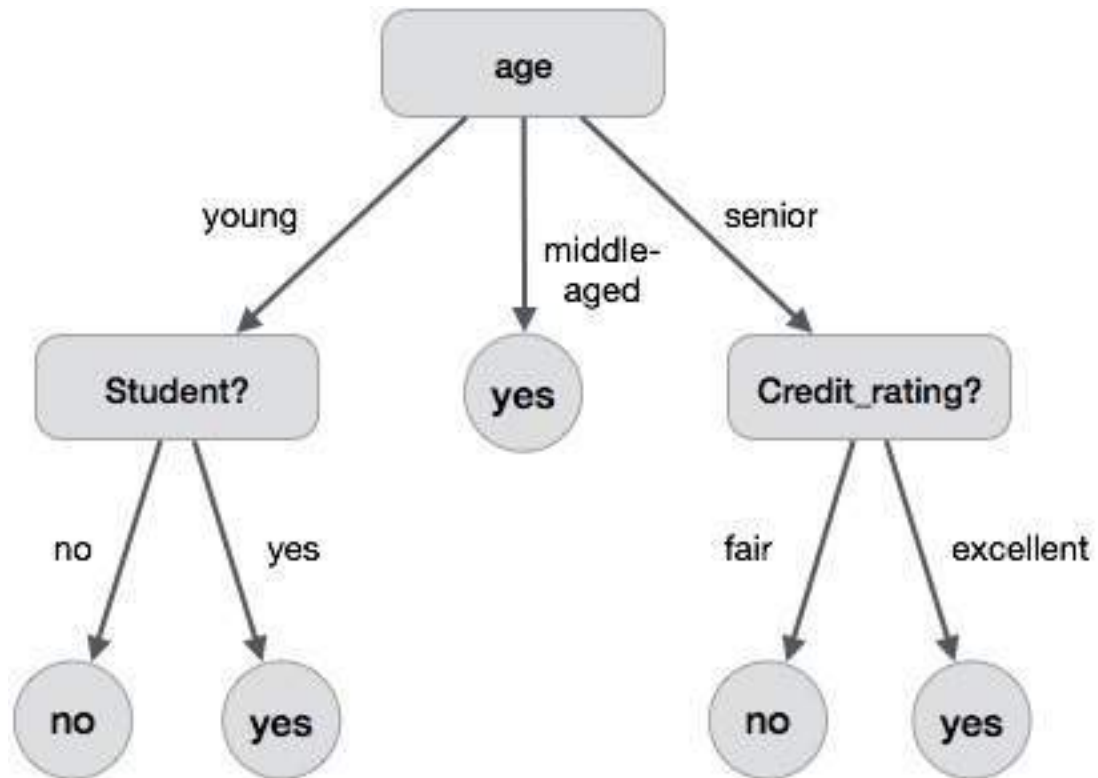


FIGURE 2.4: A Decision Tree

2.4.5 Attributes Selection Measures

An attribute selection measure is a heuristic for selecting the splitting criterion that "best" separates a given data partition, D , of class-labeled training tuples into individual classes. If we were to split D into smaller partitions according to the outcomes of the splitting criterion, ideally each partition would be pure (i.e., all the tuples that fall into a given partition would belong to the same class). Conceptually, the "best" splitting criterion is the one that most closely results in such a scenario. Attribute selection measures are also known as splitting rules because they determine how the tuples at a given node are to be split.

Input :

- Data partition, D , which is a set of training tuples and their associated class labels;
- *attribute_list*, the set of candidate attributes;
- *Attribute_selection_method*.

Output: A decision tree.

```

1: procedure
2:   create a node  $N$ ;
   if tuples in  $D$  are all of the same class,  $C$ , then
     return  $N$  as leaf node labeled with the class  $C$  ;
   end
   if attribute_list is empty then
     return  $N$  as leaf node labeled with the majority class in  $D$  ;
   end
3:   apply Attribute_selection_method( $D$ , attribute_list) to find best
   splitting criterion ;
4:   label node  $N$  with splitting_criterion ;
   if splitting_criterion is discrete valued and multiway splits allowed then
     attribute_list  $\leftarrow$  attribute_list - splitting_attribute;
   end
   for each outcome  $j$  of splitting_criterion do
5:     let  $D_j$  be the set of data tuples in  $D$  satisfying outcome  $j$  ;
     if  $D_j$  is empty then
       attach a leaf labeled with the majority class in  $D$  to node  $N$  ;
     end
     else
       attach a leaf labeled with the majority class in  $D$  to node  $N$  ;
     end
   end
6:   return  $N$  ;
7: end procedure

```

Algorithm 1: Generate_decision_tree

The attribute selection measure provides a ranking for each attribute describing the given training tuples. The attribute having the best score for the measure³ is chosen as the splitting attribute for the given tuples. If the splitting attribute is continuous-valued or if we are restricted to binary trees, then, respectively, either a split point or a splitting subset must also be determined

³Depending on the measure, either the highest or lowest score is chosen as the best (i.e., some measures strive to maximize while others strive to minimize).

as part of the splitting criterion. The tree node created for partition D is labeled with the splitting criterion, branches are grown for each outcome of the criterion, and the tuples are partitioned accordingly. Three popular attribute selection measures are

- Information gain
- Gain ration
- Gini index

The notation used herein is as follows. Let D , the data partition, be a training set of class labeled tuples. Suppose the class label attribute has m distinct values defining m distinct classes, C_i (for $i = 1, 2, \dots$). Let be the set of tuples of class C_i in D . Let $|D|$ and $|C_{i,D}|$ denote the number of tuples in D and $C_{i,D}$ respectively.

Information Gain

ID3 uses information gain as its attribute selection measure. This measure is based on pioneering work by Claude Shannon on information theory, which studied the value or "information content" of messages. Let node N represent or hold the tuples of partition D . The attribute with the highest information gain is chosen as the splitting attribute for node N . This attribute minimizes the information needed to classify the tuples in the resulting partitions and reflects the least randomness or "impurity" in these partitions. Such an approach minimizes the expected number of tests needed to classify a given tuple and guarantees that a simple (but not necessarily the simplest) tree is found.

The expected information needed to classify a tuple in D is given by

$$Info(D) = - \sum (p_i * \log(p_i))$$

where p_i is the nonzero probability that an arbitrary tuple in D belongs to class C_i and is estimated by $|C_{i,D}|/|D|$. A log function to the base 2 is used, because the information is encoded in bits. $Info(D)$ is just the average amount of information needed to identify the class label of a tuple in D . Note that, at this point, the information we have is based solely on the proportions of tuples of each class. $Info(D)$ is also known as the entropy of D . How much more information would we still need (after the partitioning) to arrive at an exact classification? This amount is measured by

$$Info_A(D) = \sum \frac{|D_j|}{|D|} * Info(D_j)$$

Information gain is defined as the difference between the original information requirement (i.e., based on just the proportion of classes) and the new requirement (i.e., obtained after partitioning on A). That is,

$$Gain(A) = Info(D) - Info_A(D)$$

Gain Ratio

The information gain measure is biased toward tests with many outcomes. That is, it prefers to select attributes having a large number of values. For example, consider an attribute that acts as a unique identifier such as *product_ID*. A split on *product_ID* would result in a large number of partitions (as many as there are values), each one containing just one tuple. Because each partition is pure, the information required to classify data set D based on this partitioning would be $Info_{product_ID}(D) = 0$. Therefore, the information gained by partitioning on this attribute is maximal. Clearly, such a partitioning is useless for classification.

C4.5, a successor of ID3, uses an extension to information gain known as gain ratio, which attempts to overcome this bias. It applies a kind of normalization to information gain using a "split information" value defined analogously with $Info(D)$ as

$$SplitInfo_A(D) = - \sum \left(\frac{|D_j|}{|D|} * \log_2 \left(\frac{|D_j|}{|D|} \right) \right)$$

This value represents the potential information generated by splitting the training data set, D , into v partitions, corresponding to the v outcomes of a test on attribute A . Note that, for each outcome, it considers the number of tuples having that outcome with respect to the total number of tuples in D . It differs from information gain, which measures the information with respect to classification that is acquired based on the same partitioning. The gain ratio is defined as

$$GainRatio = \frac{Gain(A)}{SplitInfo_A(D)}$$

Gini Index

The Gini index is used in CART. Using the notation previously described, the Gini index measures the impurity of D , a data partition or set of training tuples, as

$$Gini(D) = 1 - \sum (p_i^2)$$

where p_i is the probability that a tuple in D belongs to class C_i and is estimated by $|C_{i,D}|/|D|$. The sum is computed over m classes.

When considering a binary split, we compute a weighted sum of the impurity of each resulting partition. For example, if a binary split on A partitions D into D_1 and D_2 , the Gini index of D given that partitioning is

$$Gini_A(D) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2).$$

The reduction in impurity that would be incurred by a binary split on a discrete- or continuous-valued attribute A is

$$\delta Gini(A) = Gini(D) - Gini_A(D).$$

Other Attribute Selection Measures

Many other attribute selection measures have been proposed. **CHAID**, a decision tree algorithm that is popular in marketing, uses an attribute selection measure that is based on the statistical χ^2 test for independence. Other measures include **C-SEP** (which performs better than information gain and the Gini index in certain cases) and G-statistic (an information theoretic measure that is a close approximation to χ^2 distribution).

Attribute selection measures based on the Minimum Description Length (MDL) principle have the least bias toward multivalued attributes. MDL-based measures use encoding techniques to define the "best" decision tree as the one that requires the fewest number of bits to both (1) encode the tree and (2) encode the exceptions to the tree (i.e., cases that are not correctly classified by the tree). Its main idea is that the simplest of solutions is preferred.

Other attribute selection measures consider multivariate splits (i.e., where the partitioning of tuples is based on a combination of attributes, rather than on a single attribute). The CART system, for example, can find multivariate splits based on a linear combination of attributes. Multivariate splits are a form of attribute (or feature) construction, where new attributes are created based on the existing ones.[3]

Chapter 3

Analysis of BIIS Data

3.1 Scope

BUET offers educational facilities for more than 5000 students in an academic session. These students are under departments. Number of students in a respective department is different according to the department facility. Department of Electrical & Electronic Engineering, Department of Civil Engineering, Department of Mechanical Engineering can afford 195 students per batch. Department of Computer Science & Engineering has 120 students. Department of Architecture afford 55 students, Department of Water Resource Engineering, Department of Materials & Metallurgical Engineering, Department of Urban & Regional Planning and Department of Industrial & Production Engineering hold 30 students each. Department of Naval Architecture & Marine Engineering and Department of Chemical Engineering has 60 students each [2]. These numbers have been changed throughout year by year but it represent the period on which this research has been more focused. In this research BIIS data of 10 already graduated batches have been used which represents almost 10 thousand students. Academic performance of no current student has been brought under this research.

3.2 Database Structure

There are 8 parts of BIIS data and each data sheet has 65,000 tuples. Each tuple represents the performance of a student in a particular course. The attributes of the tuple are described below

- **Serial** : It is the serial number of a particular student. It is not actually the student id, but it's a unique constraint and primary key of the database table.
- **Department** : Respective department of the student.
- **HallStatus** : Indicates whether the student is hall resident or attached.
- **District** : Home district of the student.
- **Thana** : Thana of the student.
- **Gender** : Gender of the student.
- **PlaceOfBirth** : Birth place of the student.
- **StartingOfAcademicYear** : Indicates on which academic year the student has started undergrad.
- **AdmissionDate** : Admission day at BUET.
- **CourseName** : The particular course taken by student.
- **LetterGrade** : Lettergrade of the respective course.
- **ClassAttendanceMark** : It represents the attendance mark of the student in respective course. The attendance mark is allocated 10% of total course mark. E.g., For 3 credit courses, it's marked out of 30, for 4 credit course it's marked out of 40.

- **ClassTestMark** : Class test mark is allocated 20% of the total number of a course. E.g., For 3 credit course, highest ct mark is 60, for 4 credit course it's 80.
- **PartAMark** : For theory courses the term final examination holds 70% weight of total course. This huge exam is divided into 2 parts. Generally internal examiner handles Part A. E.g., For 3 credit course, part A holds 105 marks, for 4 credit courses it's 140.
- **PartBMark** : Same as part A, generally external examiner handles part B.
- **TotalNumber** : Total mark obtained by a student in a particular course..
- **LevelName** : The level respective course has been taken.
- **TermName** : The term respective course has been taken.
- **GPA** : GPA of student in particular term.
- **CGPA** : CGPA upto the respective term.
- **TermCount** : Term count of the student.
- **CreditHourEarned** : Credit hour earned by the student in the respective term.
- **TotalCreditHourCompleted** : Total credit hour earned by the student upto respective term.

3.3 Problems in Existing Structure

The universal database holds raw academic data of the students. There are plenty of disturbance in this database. They are described below :

- The first and most important problem is, all the eight data sheets don't have same attributes among them. For example, Sheet 1 doesn't have a

CourseName column whereas Sheet 2 has one. There are many more incidents where attributes don't match.

- One student enrolls into a number of courses in academic career which can vary from 65 to 80. So, there are multiple tuple for one student which is not suitable for classification. Determining status from multiple tuples is very difficult. Table 3.1 is an example of this.
- The database is in a format of full outer join. So, there are plenty of blank attributes. For example, in sessional courses, there is no Part A or Part B. So, *PartAMark* or *PartBMark* remains void for sessional courses. Besides, there are also some special course for which maximum attributes remain blank.
- There are some attributes, for which there is no value was really inserted ever. In some datasheets *PlaceOfBirth*, *District*, *Thana* are blank.
- In some data sheets, *CGPA* & *GPA* have rounded value. But, for accurate results, it is very important to get the accurate *CGPA* & *GPA* upto two floating points at least.
- There is no *CourseCredit* attribute in the database. But for proper reasoning, *CourseCredit* is very important.
- Some attributes like *TotalCreditHourCompleted*, *AttendanceMarks*, *ClassTestMark* don't have the same range of value for all the tuples. So, any reasoning from these values is very difficult.
- There are many confusing attributes like *TermCount*, *AdmissionDate* which don't actually make sense.
- Reasoning for Theory and Sessional courses are entirely different, but this database don't indicate whether a course is Theory course or not.

TABLE 3.1: Multiple entries for single student in the Universal Database

Serial	Department	...	CourseName	LetterGrade	...
2	CE	...	PHY143	A-	...
2	CE	...	MATH141	B+	...
...

TABLE 3.2: Blank attributes in Database

Serial	Department	...	CourseName	PartAMark	PartBMark	...
65	CSE	...	CSE100	-	-	...
65	CSE	...	MATH141	93	112	...
...

TABLE 3.3: Rounded up value of CGPA and GPA

Serial	Department	...	GPA	CGPA	...
4478	NAME	...	4	3	...
4478	NAME	...	3	3	...
...

TABLE 3.4: Different range of values

Serial	Department	...	Attendance Mark	ClassTest Mark	Total CreditHour Completed	...
65	CSE	...	27	44	160	...
120	ARCH	...	40	71	190	...
...

Chapter 4

Preprocessing

4.1 Technique And Design

The universal database contains raw data of academic performance of the students. There are plenty of problems in this Database which has been discussed in *Chapter 3*. A rigorous preprocessing technique is highly needed for turn the database into a suitable one for the classification implementation. The classification algorithm demands a clean and quality dataset without which it can't proceed efficiently. In this procedure, there were several steps taken for the data preprocessing.

4.1.1 Data Cleaning

One of the main problems with the universal database was, there are several attributes where no value was not inserted at all. In result, there are plenty of blank in the data sheet. As the first step of cleaning, these missing values were filled up using intuition and global constant. For example, data sheet 1 had some blank attributes like the Table 4.1. The *PartAMark* & *PartBMark* fields are blank because the tuples are for sessional courses and only theory courses have Part A and Part B factor.

In this case, these two fields have been filled up with 0. After the process, table looked like Table 4.2.

There are several attributes in the dataset which are not really needed for the classification implementation. As we are going to implement the excellency of a student in academic career, attributes like *District*, *Thana*, *PlaceOf-Birth*, *AdmissionDate*, *StartingDate* are not going to affect the classification algorithm anyway. So, all the database has been got rid of all the redundant attributes. Before this procedure, the database was like Table 4.3. After data cleaning process, the redundant attributes have been taken care of. Table 4.4 is the outcome.

The universal database is filled with plenty of inconsistent data. These data types are not compatible with the procedure of ID3 classification. Some steps were taken to change these inconsistent data type to a consistent one.

The *Gender* and *HallStatus* fields are filled with binary value. *Gender* is represented as *Male* or *Female* and *HallStatus* is represented with *Resident* or *Attached*. The values are converted to numeric from binary so that the data set becomes compatible with the ID3 procedure.

A crucial fault in the universal database is, it doesn't have any *CourseCreditHour* attribute and some data sheets don't have the correct value of *CGPA* as well. The *CGPA* field is rounded in some data sheets. Which does not represent the accurate value. But for proper calculation, accurate *CGPA* and *GPA* is must. In this case, with the help of other attributes of the database, *CourseCreditHour*, *GPA* and *CGPA* is calculated.

At first, *Numberpercentage* is calculated from *LetterGrade* attribute. Then, the *CourseCreditHour* is calculated from the following equation.

$$CourseCreditHour = \frac{TotalNumber}{NumberPercentage} \quad (4.1)$$

After the modification in *CourseCreditHour*, the database looks like Table 4.9.

In the ID3 implementation, the final *CGPA* of student is used only. It is calculated using the *CourseCreditHour* and *GPA* attribute of all courses the student

is enrolled to.

$$CGPA = \frac{\sum (CourseCreditHour * CourseGPA)}{\sum CourseCreditHour} \quad (4.2)$$

In this phase, one student has only one tuple in the database.

4.1.2 Data Normalization

We used normalization by decimal scaling normalizes by moving the decimal point of values of attribute A. The number of decimal points moved depends on the maximum absolute value of A. A value, v_i , of A is normalized to v'_i by computing

$$v'_i = \frac{v_i}{10^j}$$

where j is the smallest integer such that $\max(|v'_i|) < 1$.

For example, we normalized class attendance mark by calculating marks for each credit. Maximum 10 marks are given for a theory course. So, we normalized the data by calculating attendance for each credit than dividing it by the total theory credit completed by the student.

4.1.3 Data Transformation

We used numeric data format for our algorithm. So we converted string values to a corresponding integer or numeric values. For example: We gave each department unique numeric id instead of using their name code. 'EEE' is replaced by 1, 'CSE' is replaced by 2 so on. Similarly for 'Hall Status' we replaced 'Resident' with 1 and 'Attached' with 0. Similarly for 'Gender' female is indicated by 0 and male is indicated by 1.

4.2 Results

After implementing some crucial data preprocessing procedures, the modified database is at last suitable for the next step, implementation of ID3 algorithm using decision tree. The final test data has only the relevant attributes in compatible format. The attributes are :

- Serial
- Department
- Hall Status
- Gender
- ClassTestMark
- AttendanceMark
- CGPA
- TotalCourseCompleted

TABLE 4.1: Blank attributes in the Universal Database

Serial	Department	...	CourseName	PartAMark	PartBMark	...
65	CSE	...	CSE100	-	-	...
65	CSE	...	CSE210	-	-	...
...

TABLE 4.2: Universal database after filling up blank attributes

Serial	Department	...	CourseName	PartAMark	PartBMark	...
65	CSE	...	CSE100	0	0	...
65	CSE	...	CSE210	0	0	...
...

TABLE 4.3: Redundant attributes in Database

Serial	Department	District	Thana	CourseName	LetterGrade	...
4478	NAME	Dhaka	Dhamrai	MATH141	A	...
14251	MME	Satkhira	koloroa	PHY143	B+	...
...

TABLE 4.4: Database after cleaning of redundant attributes

Serial	Department	CourseName	LetterGrade	...
4478	NAME	MATH141	A	...
14251	MME	PHY143	B+	...
...

TABLE 4.5: Database with binary values

Serial	Department	Gender	HallStatus	...
4478	NAME	Male	Attached	...
14251	MME	Male	Resident	...
...

TABLE 4.6: Database after conversion of binary attributes to numeric value

Serial	Department	Gender	HallStatus	...
4478	NAME	0	1	...
14251	MME	0	0	...
...

TABLE 4.7: Database before Data cleaning

Serial	Department	CourseName	GPA	CGPA	...
4478	NAME	MATH141	3	3	...
14251	MME	PHY143	3	3	...
...

TABLE 4.8: Letter Grade and it's respective percentage of number

LetterGrade	NumberPercentage
A+	80%-100%
A	75%-79%
A-	70%-74%
B+	65%-69%
B	60%-64%
B-	55%-59%
C+	50%-54%
C	45%-49%
D	40%-44%

TABLE 4.9: Database after addition of CourseCreditHour attribute

Serial	Department	CourseName	GPA	CourseCreditHour	...
4478	NAME	MATH141	3.5	4	...
14251	MME	PHY143	3.25	3	...
...

Chapter 5

Classification

There are five class labels in our classification model.They are

- Excellent
- Good
- Moderate
- Poor
- Very Poor

Applying ID3 classification algorithm a decision tree was created using training data and this knowledge in decision tree was used to find the class label of test data set.

5.1 Decision Tree

The training data was prepared after data preprocessing as described in Chapter 4.The final attributes in the training data are

- Student Id
- Department
- Hall Status

TABLE 5.1: Training Data Sample

SID	Dept	Hall	Gender	Attendance	CT	Cgpa	Credit	Status
4650	2	0	1	1	0.83	3.72	1	excellent
4755	9	1	1	0.82	0.71	3.39	1	poor
4769	1	1	1	0.97	0.83	3.76	0.98	moderate
4975	3	1	0	0.99	0.76	3.50	1	good
5064	10	0	1	0.33	0.33	2.61	0.75	very poor

- Gender
- Attendance marks
- Class test marks
- Earned CGPA
- Completed Credit
- Final status according to our reasoning as Status

A sample of the training data set looks like Table 5.1.

A hypothetical decision tree can be derived from the training data as depicted in Figure 5.1.

5.2 Algorithm

The algorithm used is ID3 Algorithm. So, Information gain as described in Section 2.4.5 is used as splitting criterion.

We used the structure of the algorithm described in Algorithm 1 in our implementation. We used Java programming language for implementation.

The pseudocode for the final implementation of the algorithm is shown in Algorithm 2.

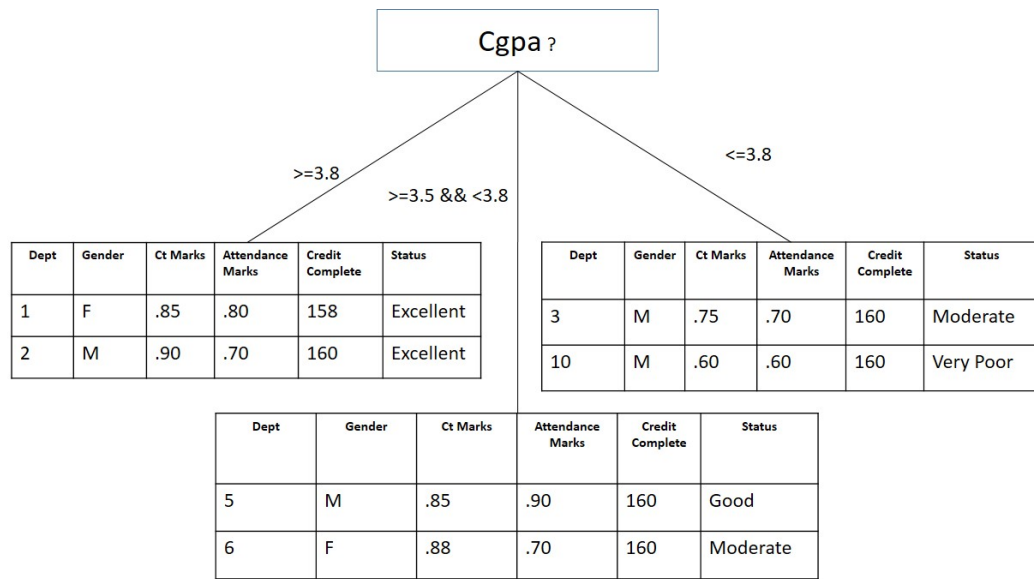


FIGURE 5.1: A Decision Tree From Training Data

5.3 Result of Classification

Test data set before applying algorithm looks like Table 5.2.

After applying algorithm as described in Algorithm 2 the results as shown in Table 5.3 are found.

5.4 Statistical Analysis

We analyzed the final results and found some relevant statistical results. The categories are

- Department wise performance
- Impact of gender on performance
- Impact of hall status on performance
- Impact of classtest marks on performance
- Impact of attendance marks on performance

Input :

- Data partition, D , which is a set of training tuples and their associated class labels ;
- $attribute_list(SID, Dept, Hall, Attendance, CT, Cgpa, CrComplete)$, the set of candidate attributes;
- $Attribute_selection_method$: Information gain with majority voting.

Output: A decision tree.

```

1: procedure
2:   create a node  $N$ ;
      if tuples in  $D$  are all of the same class,  $C$ , then
        return  $N$  as leaf node labeled with the class  $C$  ;
      end
      if  $attribute\_list$  is empty then
        return  $N$  as leaf node labeled with the majority class in  $D$  ;
      end
3:   apply  $Attribute\_selection\_method(D, attribute\_list)$  to find best
      splitting criterion ;
4:   label node  $N$  with  $splitting\_criterion$  ;
      if  $splitting\_criterion$  is  $Dept$  or  $Hall\_Status$  or  $Gender$  then
5:     split according to the discrete values of the attribute;
6:    $attribute\_list \leftarrow attribute\_list - splitting\_attribute$ ;
      end
      else
         $\triangleright$   $splitting\_criterion$  is  $Cgpa$  or  $Attendance$  or  $ClassTest$  or
        CreditCompleted.
7:
8:     split 3 ways according to the values of the attribute for example
        for  $Cgpa$  divide at 3.8 and 3.5 into 3 parts;
9:    $attribute\_list \leftarrow attribute\_list - splitting\_attribute$ ;
      end
      for each outcome  $j$  of  $splitting\_criterion$  do
10:    let  $D_j$  be the set of data tuples in  $D$  satisfying outcome  $j$  ;
      if  $D_j$  is empty then
        attach a leaf labeled with the majority class in  $D$  to node  $N$  ;
      end
      else
        attach a leaf labeled with the majority class in  $D$  to node  $N$  ;
      end
      end
11:   return  $N$  ;
12: end procedure

```

Algorithm 2: Generate_decision_tree_For_BIIS_Data

TABLE 5.2: Test Data Sample

SID	Dept	Hall	Gender	Attendance	CT	Cgpa	Credit
4487	9	1	1	0.98	0.77	3.64	1
4488	9	0	0	0.72	0.68	3.18	1
4489	9	0	0	0.93	0.69	3.49	1
4490	9	1	1	0.61	0.59	3.02	0.96
4492	11	1	1	0.99	0.72	3.06	0.99
4493	11	1	1	0.71	0.68	3.07	0.91
4488	9	0	0	0.72	0.68	3.18	1
4489	9	0	0	0.93	0.69	3.49	1
4490	9	1	1	0.61	0.59	3.02	0.96
4492	11	1	1	0.99	0.72	3.06	0.99
4493	11	1	1	0.71	0.68	3.07	0.91
4494	11	1	0	0.89	0.78	3.25	0.99
4496	11	0	1	0.98	0.73	3.22	0.99

TABLE 5.3: Final Output Sample

SID	Dept	Hall	Gender	Attendance	CT	Cgpa	Credit	Status
4487	9	1	1	0.98	0.77	3.64	1	good
4488	9	0	0	0.72	0.68	3.18	1	moderate
4489	9	0	0	0.93	0.69	3.49	1	good
4490	9	1	1	0.61	0.59	3.02	0.96	poor
4492	11	1	1	0.99	0.72	3.06	0.99	moderate
4493	11	1	1	0.71	0.68	3.07	0.91	poor
4488	9	0	0	0.72	0.68	3.18	1	moderate
4489	9	0	0	0.93	0.69	3.49	1	good
4490	9	1	1	0.61	0.59	3.02	0.96	very poor
4492	11	1	1	0.99	0.72	3.06	0.99	poor
4493	11	1	1	0.71	0.68	3.07	0.91	very poor
4494	11	1	0	0.89	0.78	3.25	0.99	moderate
4496	11	0	1	0.98	0.73	3.22	0.99	moderate

TABLE 5.4: Performance of EEE department

Class Label	Percent
Excellent	18%
Good	42%
Moderate	32%
Poor	5%
Very Poor	3%

TABLE 5.5: Performance of CSE Department

Class Label	Percent
Excellent	21%
Good	32%
Moderate	12%
Poor	30%
Very Poor	5%

- Impact of cgpa on performance
- Impact of credit completion on performance

The details of the findings are discussed in Section 5.5

5.5 Result of Statistical Analysis

5.5.1 Department wise Performance

EEE Department

The overall performance of EEE department is shown in Figure 5.2. According to our classifier the percentage of each class label of EEE department is shown in Table 5.4

CSE Department

The overall performance of CSE department is shown in Figure 5.3. According to our classifier the percentage of each class label of CSE department is shown in Table 5.5

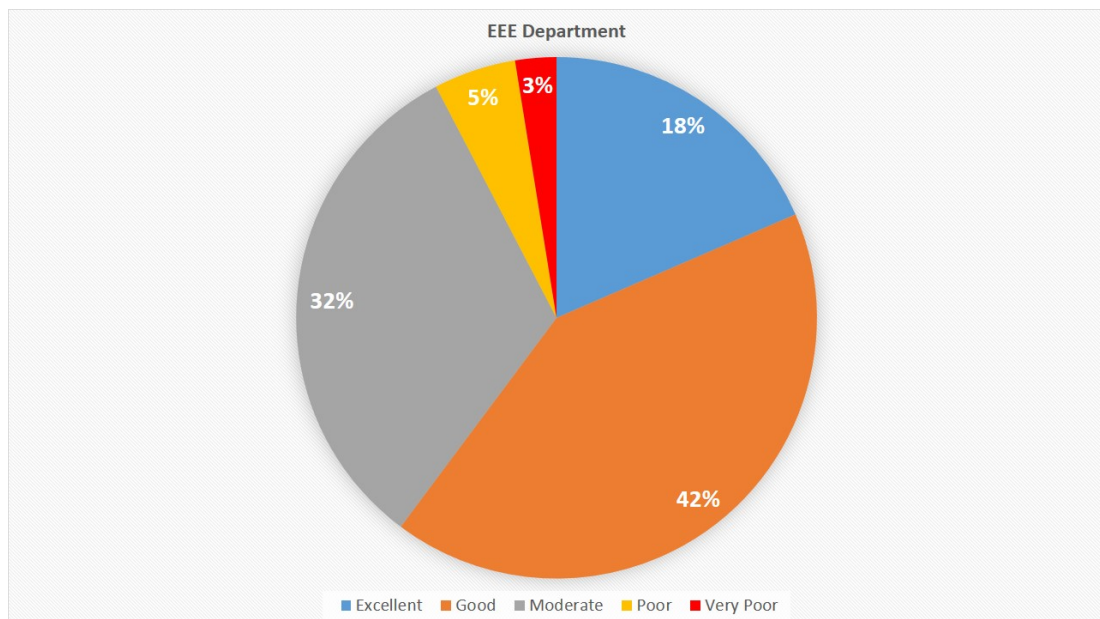


FIGURE 5.2: Performance of EEE Department

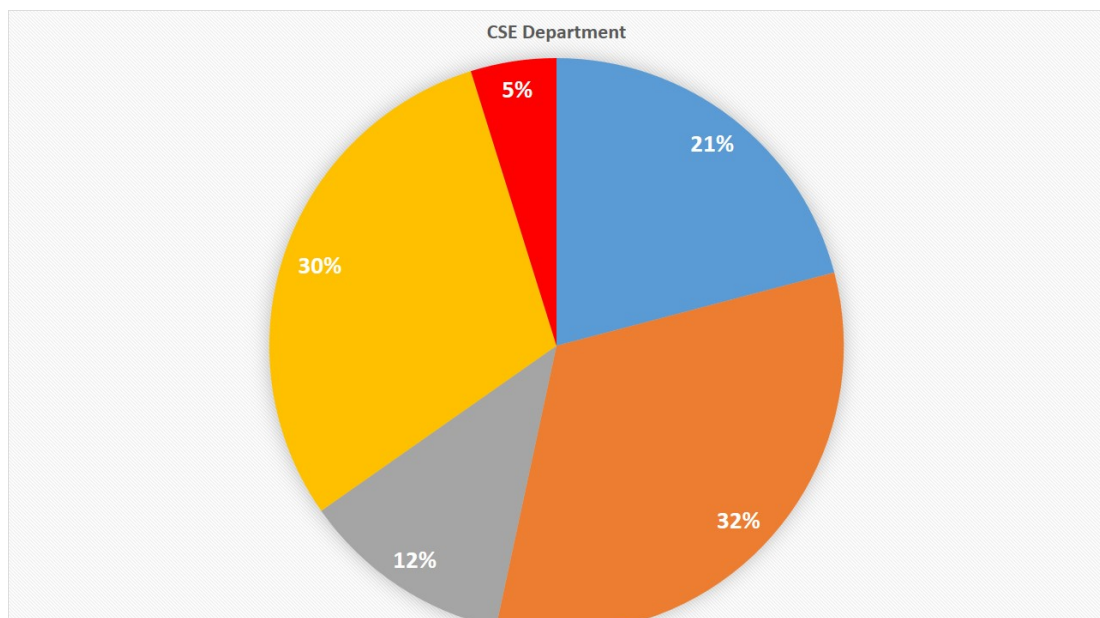


FIGURE 5.3: Performance of CSE Department

TABLE 5.6: Performance of IPE Department

Class Label	Percent
Excellent	13%
Good	33%
Moderate	17%
Poor	23%
Very Poor	14%

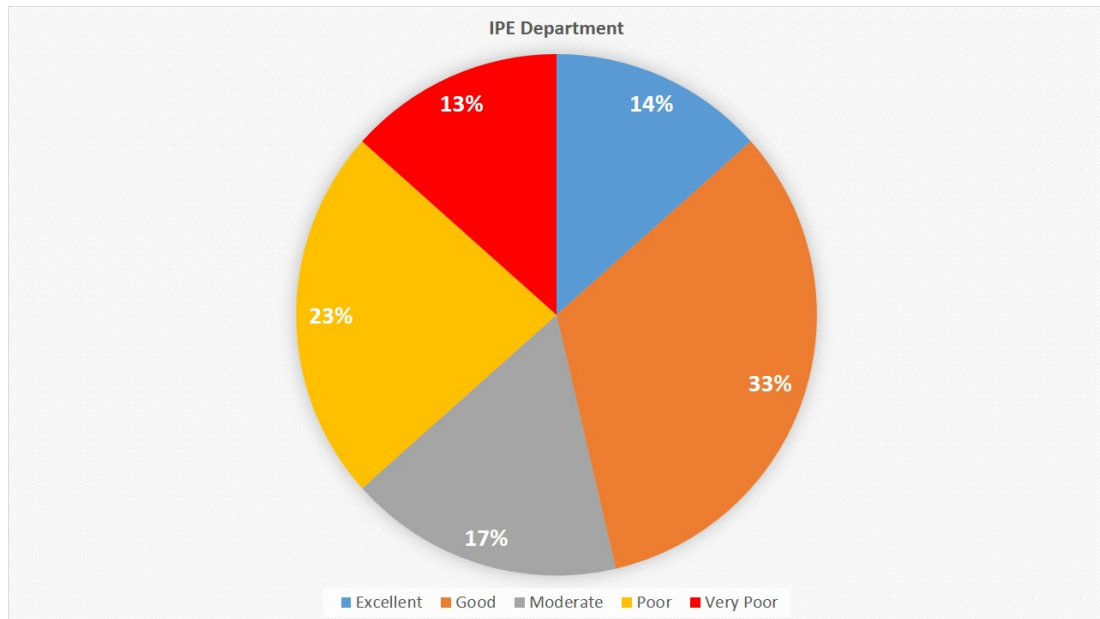


FIGURE 5.4: Performance of IPE Department

IPE Department

The overall performance of IPE department is shown in Figure 5.4. According to our classifier the percentage of each class label of IPE department is shown in Table 5.6

ME Department

The overall performance of ME department is shown in Figure 5.5. According to our classifier the percentage of each class label of ME department is shown in Table 5.7

TABLE 5.7: Performance of ME Department

Class Label	Percent
Excellent	5%
Good	42%
Moderate	37%
Poor	6%
Very Poor	10%

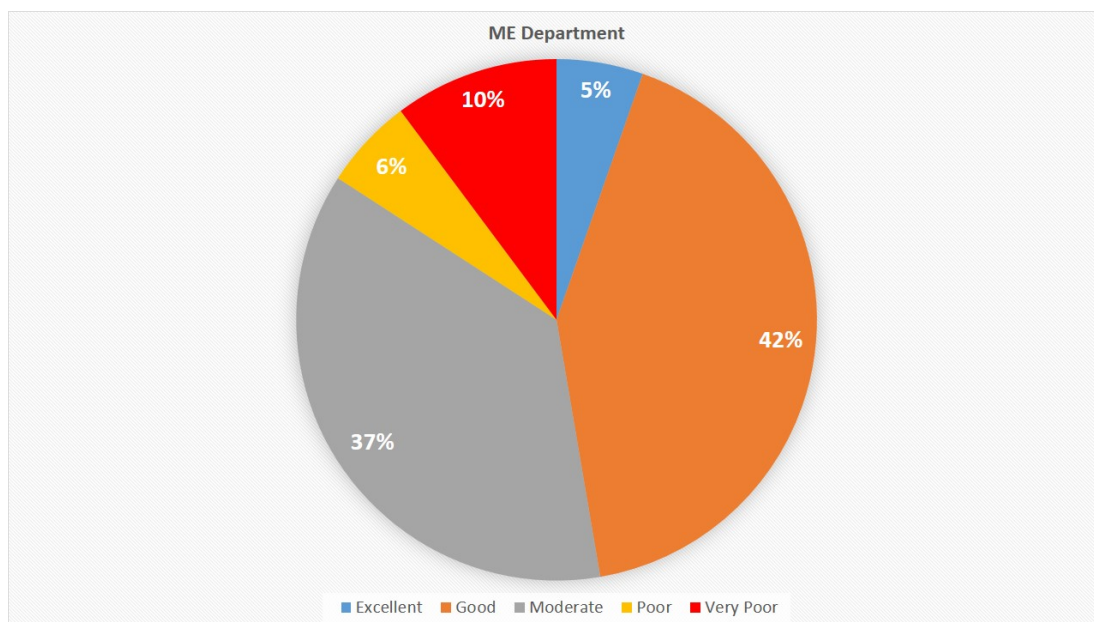


FIGURE 5.5: Performance of ME Department

TABLE 5.8: Performance of CE Department

Class Label	Percent
Excellent	5%
Good	27%
Moderate	47%
Poor	15%
Very Poor	6%

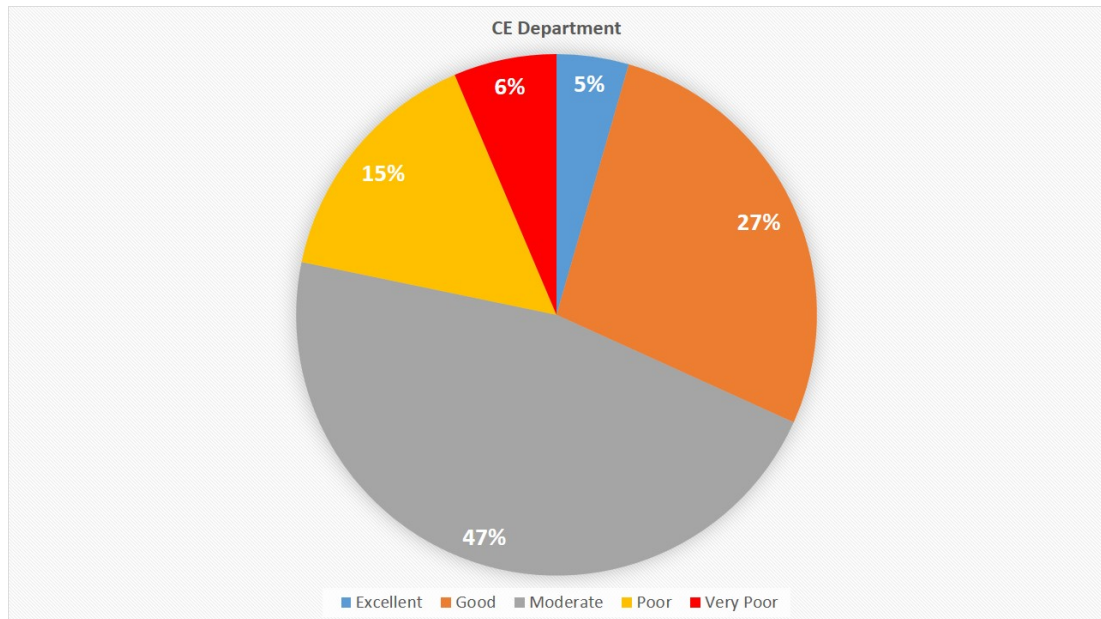


FIGURE 5.6: Performance of CE Department

CE Department

The overall performance of CE department is shown in Figure 5.6. According to our classifier the percentage of each class label of CE department is shown in Table 5.8

MME Department

The overall performance of MME department is shown in Figure 5.7. According to our classifier the percentage of each class label of MME department is shown in Table 5.9

TABLE 5.9: Performance of MME Department

Class Label	Percent
Excellent	16%
Good	37%
Moderate	22%
Poor	14%
Very Poor	11%

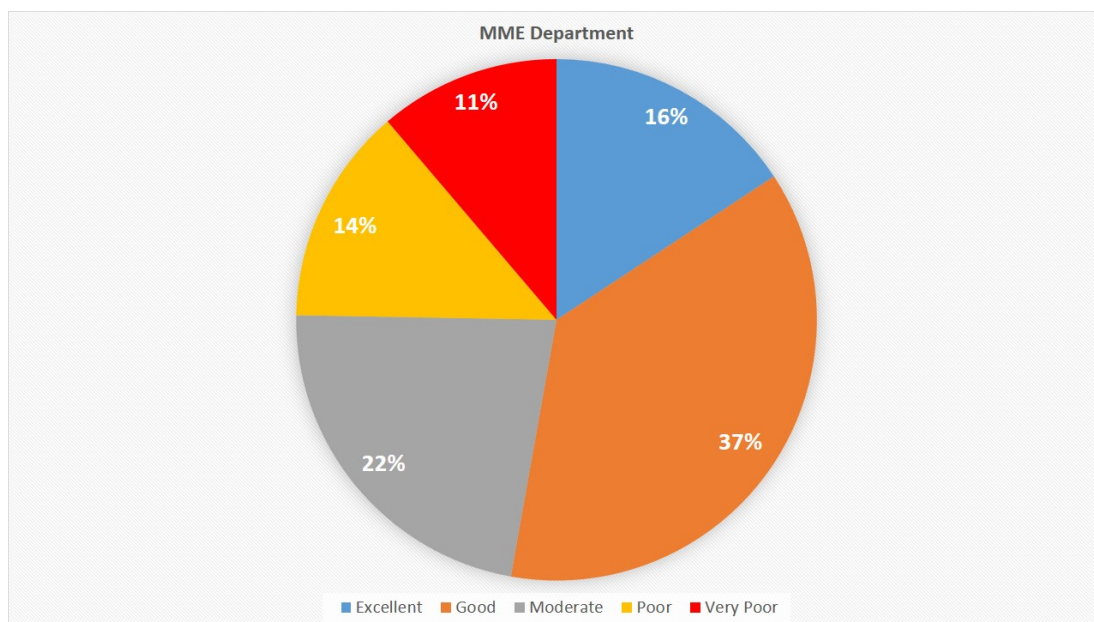


FIGURE 5.7: Performance of MME Department

TABLE 5.10: Performance of CHE Department

Class Label	Percent
Excellent	5%
Good	62%
Moderate	13%
Poor	8%
Very Poor	12%

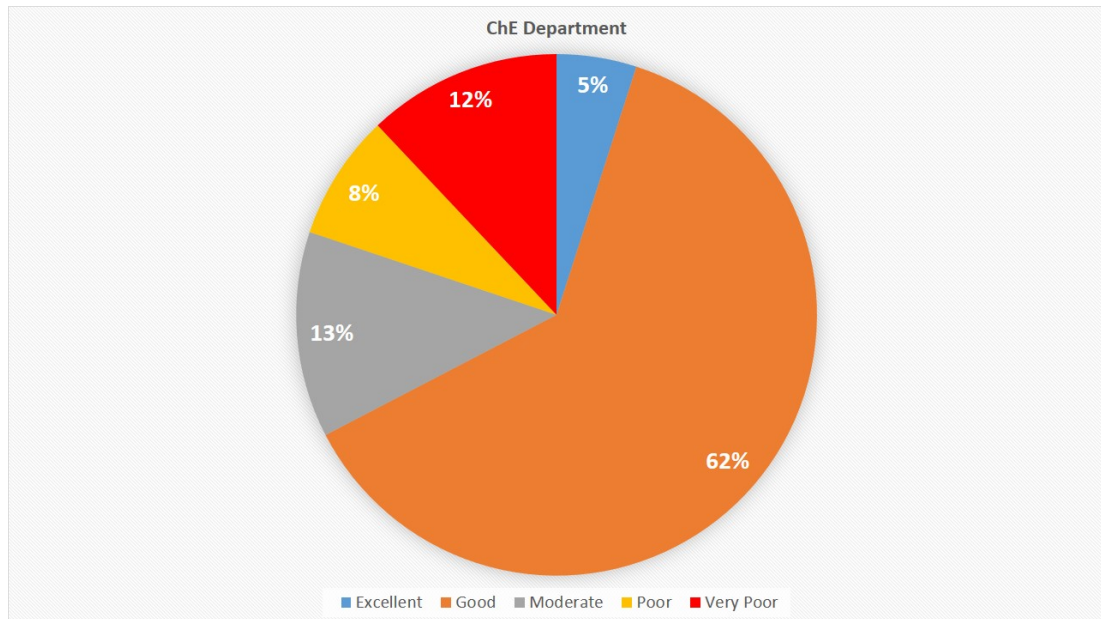


FIGURE 5.8: Performance of CHE Department

CHE Department

The overall performance of CHE department is shown in Figure 5.8. According to our classifier the percentage of each class label of CHE department is shown in Table 5.10

NAME Department

The overall performance of NAME department is shown in Figure 5.9. According to our classifier the percentage of each class label of NAME department is shown in Table 5.11

TABLE 5.11: Performance of NAME Department

Class Label	Percent
Excellent	8%
Good	48%
Moderate	27%
Poor	10%
Very Poor	7%

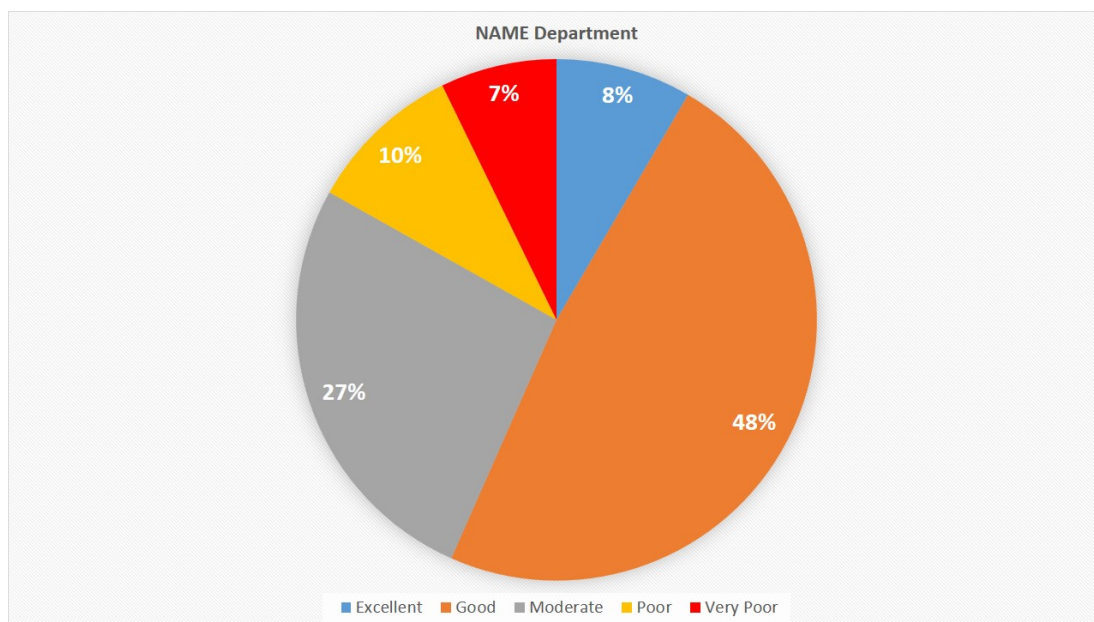


FIGURE 5.9: Performance of NAME Department

TABLE 5.12: Performance of URP Department

Class Label	Percent
Excellent	9%
Good	40%
Moderate	26%
Poor	11%
Very Poor	14%

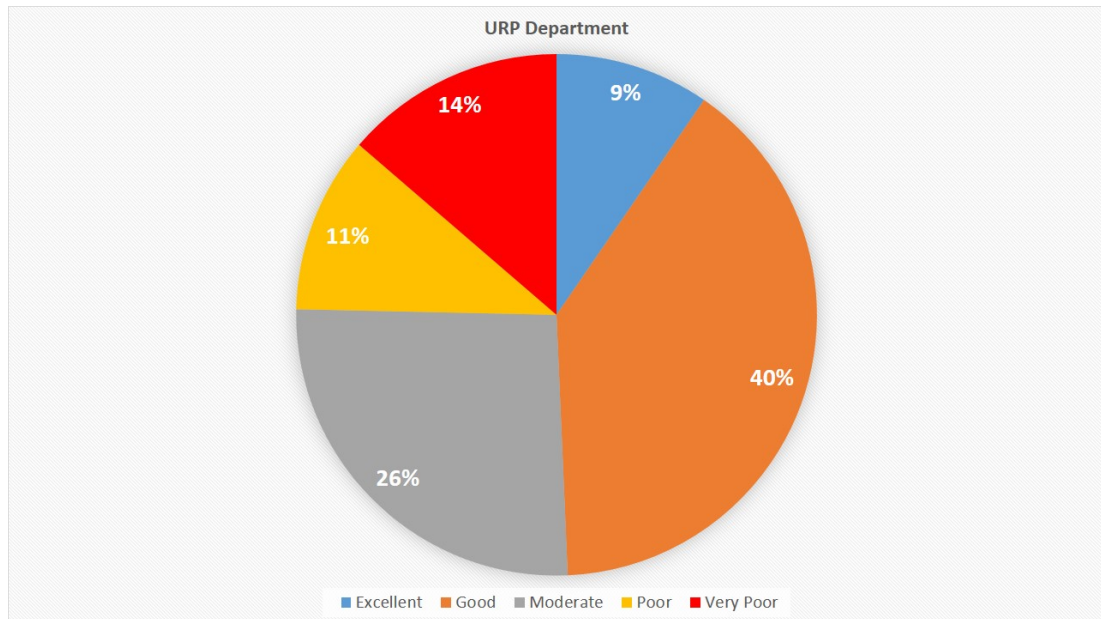


FIGURE 5.10: Performance of URP Department

URP Department

The overall performance of URP department is shown in Figure 5.10. According to our classifier the percentage of each class label of URP department is shown in Table 5.12

ARCH Department

The overall performance of ARCH department is shown in Figure 5.11. According to our classifier the percentage of each class label of ARCH department is shown in Table 5.13

TABLE 5.13: Performance of ARCH Department

Class Label	Percent
Excellent	2%
Good	8%
Moderate	32%
Poor	16%
Very Poor	39%

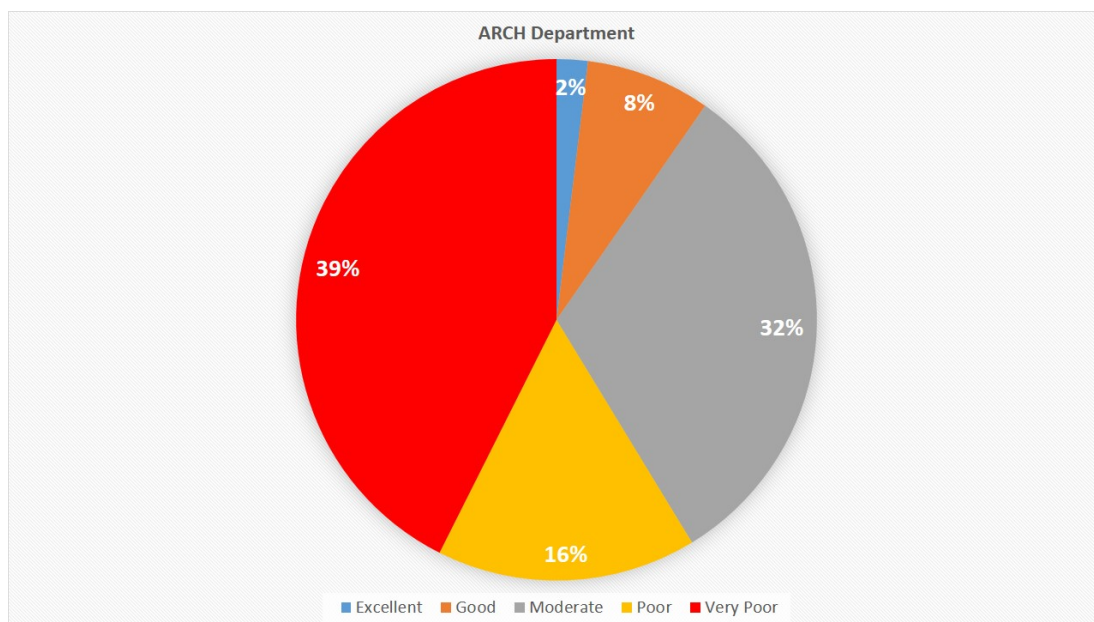


FIGURE 5.11: Performance of ARCH Department

TABLE 5.14: Performance of WRE Department

Class Label	Percent
Excellent	8%
Good	30%
Moderate	35%
Poor	12%
Very Poor	15%

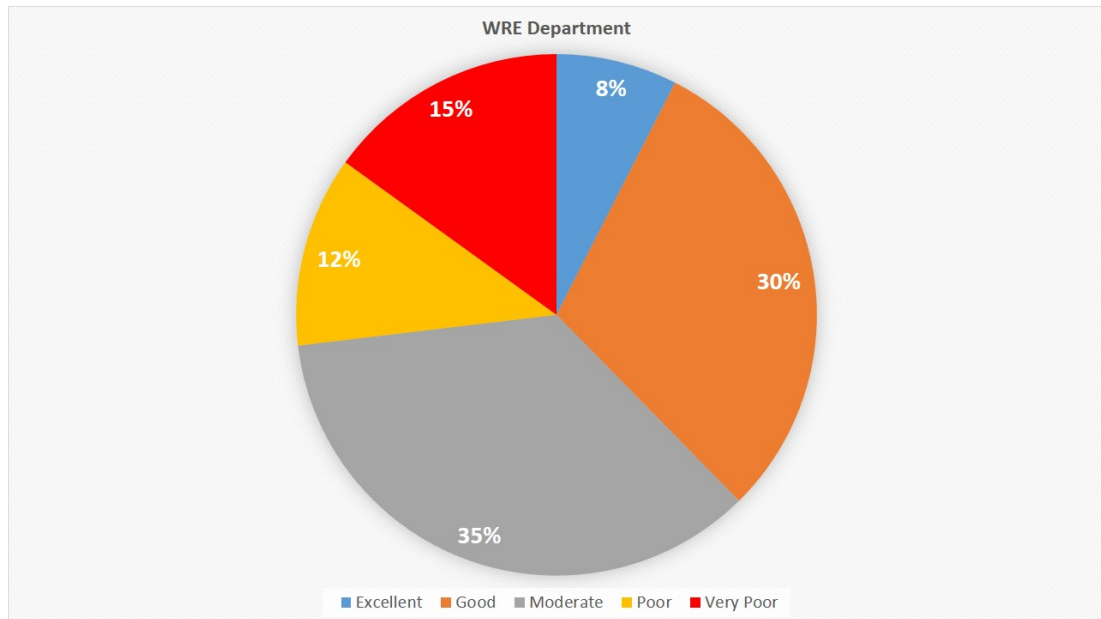


FIGURE 5.12: Performance of WRE Department

WRE Department

The overall performance of WRE department is shown in Figure 5.12. According to our classifier the percentage of each class label of WRE department is shown in Table 5.14

Overall Department wise Performance

Figure 5.13 shows the overall department wise performance according to our classifier. The amounts are calculated in percentage for better comparison.

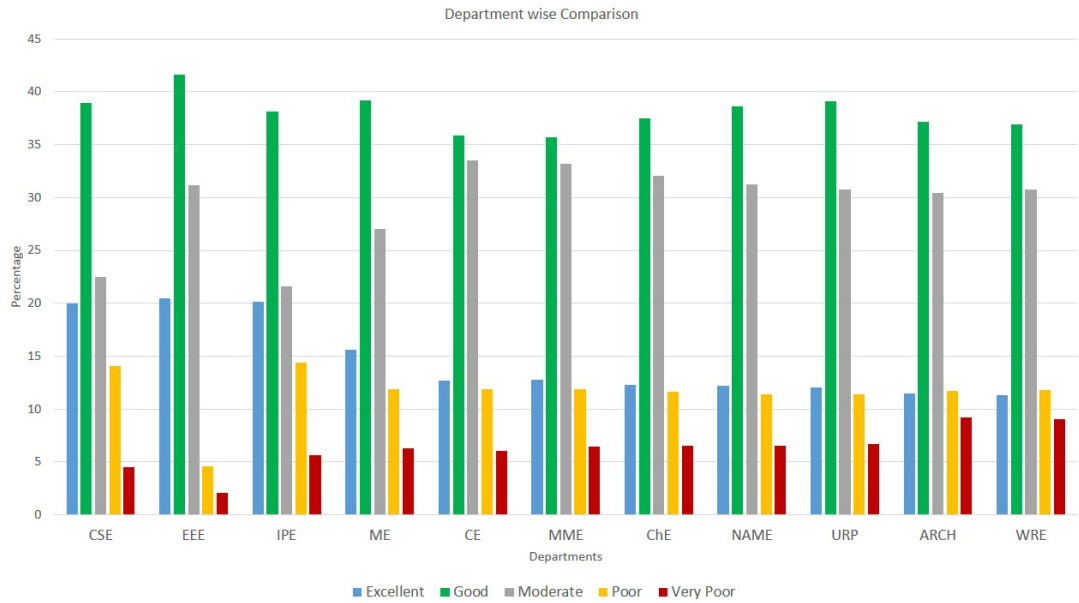


FIGURE 5.13: Overall Department wise Performance

TABLE 5.15: Impact of Gender on Performance

Class Label	Feale	Male
Excellent	14.56%	10.55%
Good	34.13%	37.59%
Moderate	29.13%	31.15%
Poor	13.26 %	11.48%
Very Poor	8.91%	9.11%

5.5.2 Impact of Gender on Performance

Male and female students have different proportion of class label according to our classifier. Female students have higher percentage of their number in top(excellent) class. In the lowest class(very poor) the percentages are almost equal. Figure 5.14 illustrate this phenomena. Table 5.15 is the numerical data table for this section.

5.5.3 Impact of Hall Status on Performance

Hall status(Resident or Attached) of a student have a profound impact on performance. Attached students have higher percentage of their number in higher

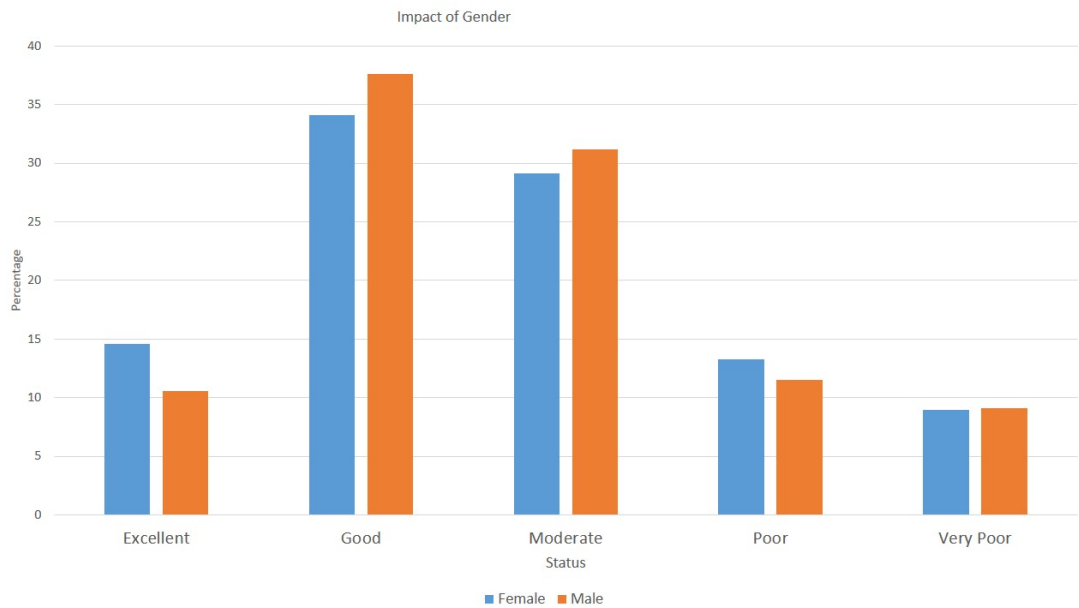


FIGURE 5.14: Impact of Gender on Performance

TABLE 5.16: Impact of Hall Status on Performance

Class Label	Resident	Attached
Excellent	7.88%	17.13%
Good	35.12%	39.97%
Moderate	35.38%	22.95%
Poor	12.65%	10.41%
Very Poor	8.94%	9.29%

classes(excellent,good). In the lower classes the percentages are almost same. Residents are dominant in 'moderate' class. Figure 5.15 illustrate this phenomena. Table 5.16 is the numerical data table for this section.

5.5.4 Impact of Class Attendance Marks

Class attendance has a profound impact on a student's performance. Students with high class test marks are tend to be in the higher label. For example students with greater than 90% attendance is dominant in top class(excellent). On the other hand student with less than 60% attendance has a very high probability of being in the 'very poor' class. Figure 5.16 describes this phenomena. Table 5.17 is the numerical data table for this section.

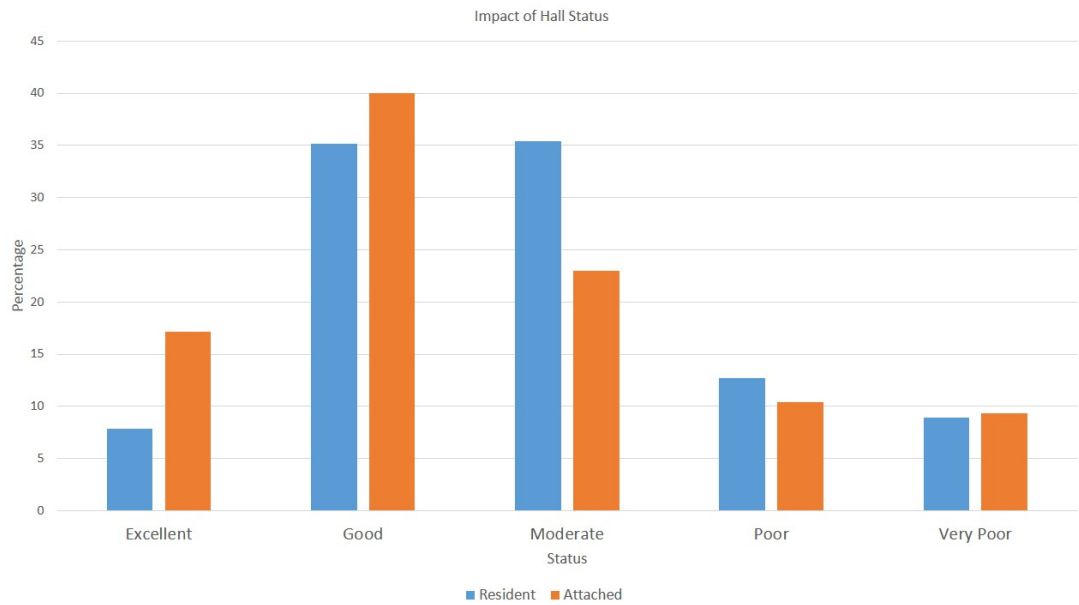


FIGURE 5.15: Impact of Hall Status on Performance

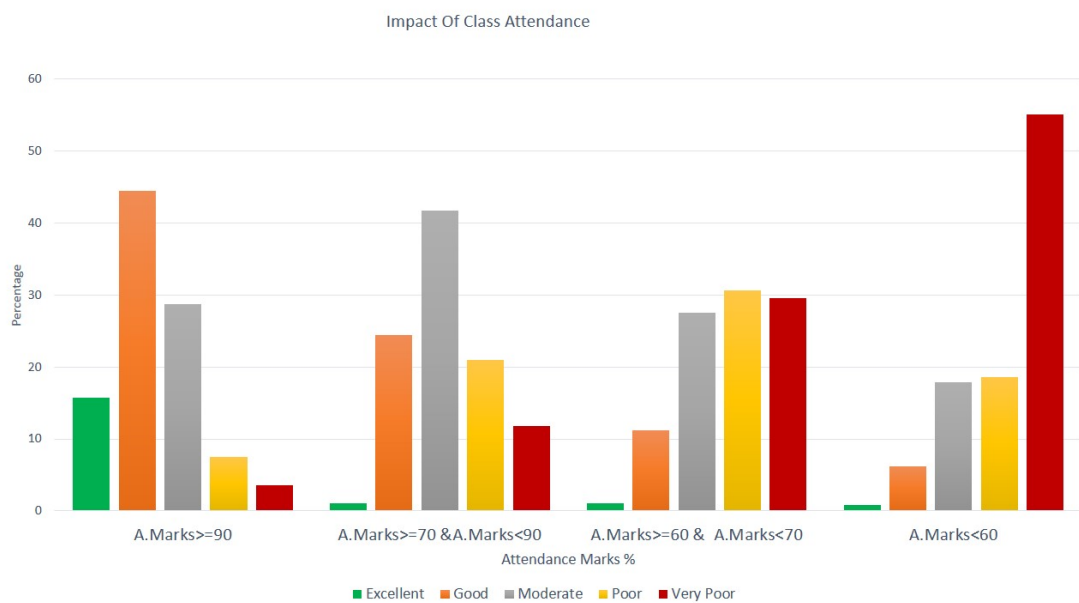


FIGURE 5.16: Impact of Attendance on Performance

TABLE 5.17: Impact of Hall Status on Performance

Class Label	>90	>70 and <90	>60 and <70	<60
Excellent	15.73%	1.01%	1.02%	0.77%
Good	44.41%	24.43%	11.22%	6.20%
Moderate	28.74%	41.75%	27.55%	17.82%
Poor	7.54%	20.97%	30.61%	18.60%
Very Poor	3.56%	11.81%	29.59%	55.03%

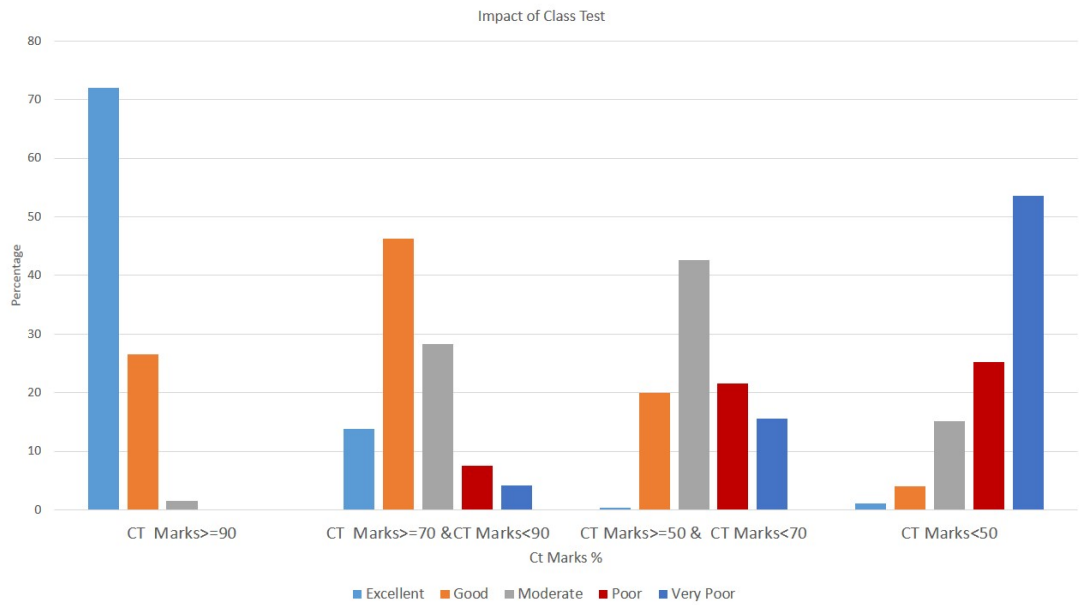


FIGURE 5.17: Impact of Class test marks on Performance

TABLE 5.18: Impact of Class Test Marks on Performance

Class Label	>90	>70 and <90	>50 and <70	<50
Excellent	73%	14.1%	1%	1%
Good	26%	47%	20%	4%
Moderate	2%	29%	43%	15%
Poor	0%	8%	22%	26%
Very Poor	0%	5%	16%	54%

5.5.5 Impact of Class Test Marks on Performance

Class test marks plays a very important role on overall performance of a student. Figure 5.17 shows the result. Table 5.18 is the data table.

5.5.6 Impact of CGPA on Overall Performance

It is quite obvious that cgpa is the most important part of a student's overall academic performance. Our classifier also agrees with this fact. About 99% of the students getting cgpa higher than 3.8 falls in the 'Excellent' class. Those with lower cgpa falls in lower classes.

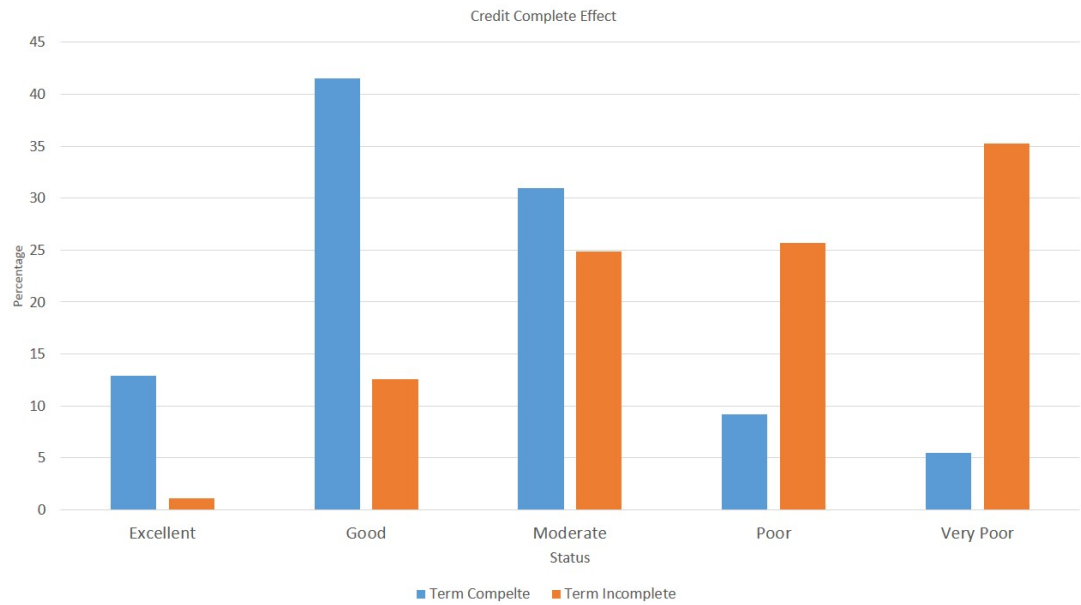


FIGURE 5.18: Impact of Credit completion on Performance

5.5.7 Impact of Credit Completion

Completing all the required credit within the regular period of time plays a vital role in overall academic performance. Figure 5.18 shows the result.

More statistical analysis results are discussed in Appendix B.

Chapter 6

Conclusion

6.1 Summary of Thesis

We used classification technique to discover knowledge from biis student data in this thesis project. All programming and necessary statistical analysis was done using java programming. The raw data set was full of impurities. So, we have to do an extensive amount of preprocessing on the data before applying it in our algorithm. Training data set was prepared very carefully and with logistic reasoning. The results of the classier was statistically analyzed and was stored in tabular form and graphical form for better visualization of the findings. The accuracy of the classifier is about 76%. The low accuracy is due to the variations and impurities of the input data.

6.2 General Findings

The Classifier has five class labels. They are excellent,good,moderate,poor,very poor. After completing statistical analysis some major findings are

- EEE and CSE departments have a higher percentage in top levels. ARCH department has a higher percentage at the lowest level.
- On average female students have better academic performance than male students.

- Attached students dominate in the higher classes over the resident students.
- Class Attendance and class attendance is very important for a student's overall performance. Students with good attendance trends to be in higher classes.
- Class test marks are also very important for a student. Students with high percentage of class test marks trend to be in the higher classes. On the other hand students with lower percentage of class test marks fall in lower classes.

6.3 Future Works

We used only classification rules to discover knowledge from data. In future we want to combine the clustering algorithm with this one to get more accurate result.

The current accuracy is low. So we want to improve the accuracy by applying accuracy improvement algorithm.

We want to use this project to make a complete system for student evaluation.

Appendix A

Data Samples

Initial dataset looked like Figure A.1. After processing and eliminating redundancy the data set looks like A.2.

Credit column is added as there was no valid credit column in the given data set. Credit was calculated from the obtained total marks and grade.

Similarly Cgpa column was also created as there was no valid cgpa column in majority of the given data sets.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W		
1	SERIAL	DEPARTM	HALLSTAT	DISTRIK	THANA	Gender	PLACE	STARTING	ADMISSIO	ADMISSIO	LETTER	CLASSATT	CLASSTES	PART	PART	TOTAL	LEVELN	TERMNAN	GPA	CGPA	TERMCOU	CRHREAR	TOTCRHR	COMPL	
2		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						10	1	1	2.2	2.56	6	5	7		
3		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						43	1	1	0	3.75	-1	0	2		
4		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0				Abs	1	1	0	2.64	11	0	30		
5		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						Abs	1	1	3.5	2.68	14	2	31		
6		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						Abs	1	2	0	2.64	12	0	30		
7		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						Abs	1	2	2	2.62	15	3	34		
8		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0	0	30	25	55	1	2	0	2.64	12	0	30		
9		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0	0	37	16	53	1	2	2	2.62	15	3	34		
10		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0	Abs	29	7	36	1	1	3.5	2.68	14	2	31		
11		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0	Abs	Abs	Abs	0	1	1	0	2.64	11	0	30		
12		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		28	9	12	8	57	1	1	2.2	2.56	6	5	7		
13		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		40	12	0	0	52	1	1	0	3.75	-1	0	2		
14		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						0	1	2	2.44	2.5	7	6	13		
15		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		0	14	Abs	Abs	14	1	2	0	2.64	12	0	30		
16		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		28	11	11	14	64	1	2	2.44	2.5	7	6	13		
17		2 CE	Resident		M		2003	1/1/2003	1/2/2003								1	2			16				
18		2 CE	Resident		M		2003	1/1/2003	1/2/2003	B-		40		67	61	62	230	1	2	2.78	2.64	8	12	25	
19		2 CE	Resident		M		2003	1/1/2003	1/2/2003	B-						82	1	2	2.78	2.64	8	12	25		
20		2 CE	Resident		M		2003	1/1/2003	1/2/2003	B-		27		28	49	68	172	1	2	2.78	2.64	8	12	25	
22		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F		10	10	18	10	48	1	1	2.65	2.64	9	5	30		
23		2 CE	Resident		M		2003	1/1/2003	1/2/2003	F						0	1	1	2.65	2.64	9	5	30		
24		2 CE	Resident		M		2003	1/1/2003	1/2/2003	B-						84	1	1	2.65	2.64	9	5	30		
25		2 CE	Resident		M		2003	1/1/2003	1/2/2003	A-						106	1	1	3.5	2.68	14	2	31		

FIGURE A.1: BIIS Data

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	SERIAL	DEPARTM	HALLSTAT	Gender	STARTING	ADMISSIO	ADMISSIO	LETTE	CLASSATT	CLASS	PARTAMA	PARTBM	TOTALN	LEVELN	TERMNAN	GPA	CGPA	TERMCOU	CRHREAR	TOTCRHR	COURSE	CREDIT
2	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	10	1	1	2.2	2.56	6	5	7	0.75	
3	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	43	1	1	0	3.75	-1	0	2	1.5	
4	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	INF	NA	Abs	1	1	0	2.64	11	0	30	0	
5	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	Abs	1	1	3.5	2.68	14	2	31	0	
6	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	Abs	1	2	0	2.64	12	0	30	0	
7	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	Abs	1	2	2	2.62	15	3	34	0	
8	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	0	30	25	55	1	2	0	2.64	12	0	30	1.5
9	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	0	37	16	53	1	2	2	2.62	15	3	34	1.5
10	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	Abs	29	7	36	1	1	3.5	2.68	14	2	31	1.5
11	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	Abs	Abs	Abs	0	1	1	0	2.64	11	0	30	0
12	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		28	9	12	8	57	1	1	2.2	2.56	6	5	7	1.5
13	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		40	12	0	0	52	1	1	0	3.75	-1	0	2	1.5
14	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	0	0	1	2	2.44	2.5	7	6	13	0
15	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	14	Abs	Abs	14	1	2	0	2.64	12	0	30	0.75
16	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		28	11	11	14	64	1	2	2.44	2.5	7	6	13	1.5
17	2	CE	Resident	M	2003	1/1/2003	1/2/2003	NR	INF	INF	NA	NA	NR	1	2	NR	NR	16	NR	NR	0	
18	2	CE	Resident	M	2003	1/1/2003	1/2/2003	B	INF	INF	NA	NA	90	1	2	2.78	2.64	8	12	25	1.5	
19	2	CE	Resident	M	2003	1/1/2003	1/2/2003	B-		40	67	61	62	230	1	2	2.78	2.64	8	12	25	0
20	2	CE	Resident	M	2003	1/1/2003	1/2/2003	B-	INF	INF	NA	NA	82	1	2	2.78	2.64	8	12	25	1.5	
21	2	CE	Resident	M	2003	1/1/2003	1/2/2003	B-		27	28	49	68	172	1	2	2.78	2.64	8	12	25	3
22	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		10	10	18	10	48	1	1	2.65	2.64	9	5	30	1.5
23	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F	INF	INF	NA	NA	0	0	1	1	2.65	2.64	9	5	30	0
24	2	CE	Resident	M	2003	1/1/2003	1/2/2003	B-	INF	INF	NA	NA	84	1	1	2.65	2.64	9	5	30	1.5	
25	2	CE	Resident	M	2003	1/1/2003	1/2/2003	A-	INF	INF	NA	NA	106	1	1	3.5	2.68	14	2	31	1.5	
26	2	CE	Resident	M	2003	1/1/2003	1/2/2003	F		0	0	Abs	Abs	0	1	1	3.5	2.68	14	2	31	0

FIGURE A.2: Processed BIIS Data

Appendix B

More Statistical Analysis Results

We also derived some statistical result from our processed data. This results is not related to our classifier. Some of the results are discussed below

B.1 Department wise Different Grade Comparison

Figure B.1 shows the frequency of different letter grades for different department. For example students of EEE department get A+ more than other grades. On the other hand students of Architecture department get most percentage of F grade among other grades.

Figure B.2 shows the frequency of grades in different departmental subjects. As we can see students get most A+ in EEE subjects and get most F grade in MATH subjects.

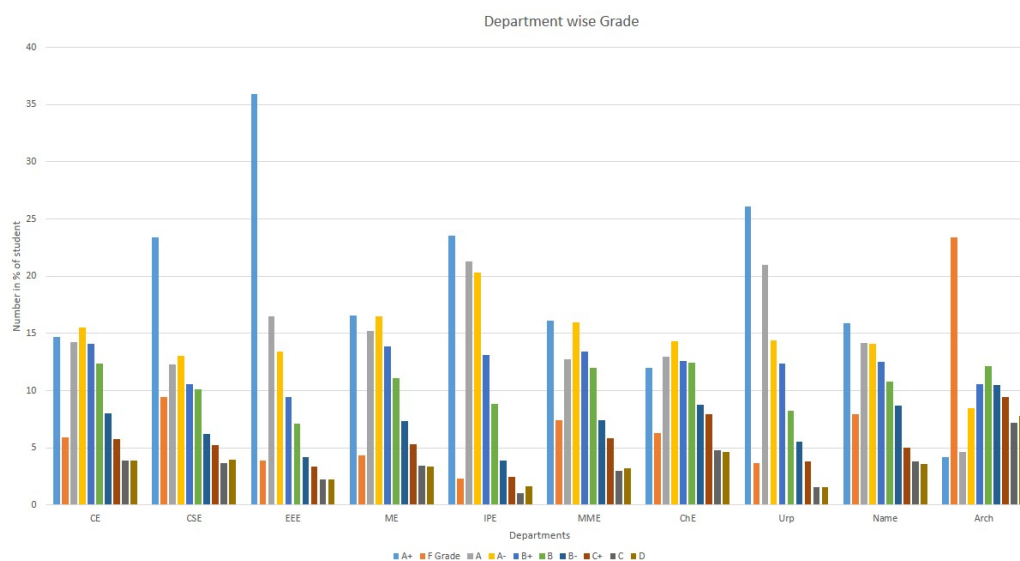


FIGURE B.1: Different Grades for Dept

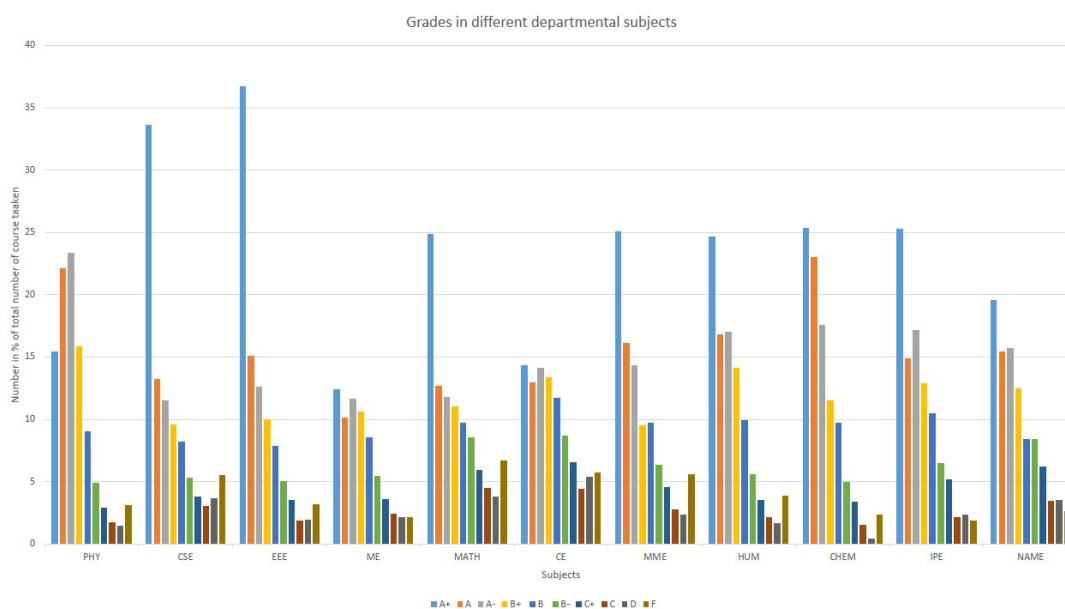


FIGURE B.2: Different Grades for Subjects

Appendix C

Details of Department Name

EEE = Electrical and Electronics Engineering .

CSE = Computer Science and Engineering.

IPE = Industrial Production Engineering.

ME = Mechanical Engineering.

NAME = Naval Architecture and Marine Engineering.

ARCH = Architecture Department.

URP = Urban and Regional Planning.

WRE = Water Research Engineering.

CE = Civil Engineering.

ChE = Chemical Engineering.

Bibliography

- [1] Han J., Kamber M. and Pei J. *Data Mining: Concepts and Techniques*. Morgan Kauffman Publishers, San Francisco, 2011.
- [2] Bangladesh University of Engineering and Technology, General Information on: <http://www.buet.ac.bd>
- [3] Data Mining Concepts and Techniques By Jiawei Han, Micheline Kamber. Jian Pei.
- [4] Osama Fayyad et al., 1996
- [5] http://researcher.watson.ibm.com/researcher/view_group.php?id=144
- [6] <http://www.usc.edu/dept/ancntr/Paris-in-LA/Analysis/discovery.html>
- [7] http://www.cs.ccsu.edu/~markov/ccsu_courses/datamining-3.html
- [8] Maurizio Lenzerini (2002). "Data Integration: A Theoretical Perspective" PODS 2002. pp. 233–246.