# Universiti Teknologi MARA

# Mining the DNA Sequence for Race Classification

## Siti Nurihan Hj.Ariffin

Thesis submitted in fulfillment of the requirements for
**Bachelor of Science (Hons) Information System Engineering**
**Faculty of Information Technology And**
**Quantitative Science**

May 2008

# ACKNOWLEDGEMENT

# ABSTRACT

Bioinformatics is a new field that arises in Malaysia. It offers many benefits especially in the genes and genomic research. By using the computers technology, it helps the researchers to understand the large amount of data in the bioinformatics. By implementing data mining techniques, the researchers are able to analyze the data from different perspective and summarizing the data into useful information. This research focusing on identifying the data mining methods/ techniques for the race classification and predict the races based on their DNA sequence using the appropriate technique. The comparison between classification techniques was done in order to get the best technique for classification. By implementing the existing tool of neural network, the race classification was achieved. From this research, it helps in identifying the techniques for classification.

# TABLE OF CONTENTS

**TITLE**                                                                    **PAGE**

## CHAPTER ONE: INTRODUCTION

## CHAPTER TWO: LITERATURE REVIEW

## CHAPTER THREE: RESEARCH APPROACH AND METHODOLOGY

## CHAPTER FOUR: ANALYSIS AND FINDINGS