IN THE NAME OF GOD

# TEXT CLASSIFICATION BASED ON PROBABILISTIC
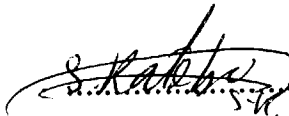
# LEARNING MODELS

BY·
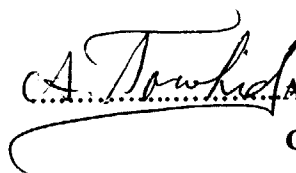**HEYDAR DAVOODI**

**THESIS**

SUBMITTED TO THE SCHOOL OF GRADUATE STUDIES IN
PARTIAL FULFILMENT OF THE REQUIREMENTS FOR
THE DEGREE OF MASTER OF SCIENCE (M.Sc)

IN
COMPUTER ENGINEERING-
ARTIFICIAL INTELLIGENCE AND ROBOTICS
SHIRAZ UNIVERSITY
SHIRAZ, IRAN

EVALUATED AND APPROVED BY THE THESIS COMMITTEE AS: **EXCELLENT**

....S. D. KATEBI, Ph.D., PROF. OF COMPUTER
ENGINEERING (CHAIRMAN)

....M. ZOLGHADRI JAHROMI, Ph.D., ASSISTANT PROF.
OF COMPUTER ENGINEERING (CHAIRMAN)

A. TOWHIDI, Ph.D., ASSISTANT PROF. OF
COMPUTER ENGINEERING

**April 2001**

٢٨٢٧٧

# ACKNOWLEDGEMENT

۳۸۲ ۱۷

*Dedicated to*

*my mother*

*for the first lullaby*

# ABSRACT

## Text Classification Based on Probabilistic
## Learning Models

By

**Heydar Davoodi**

As the amount of information available to human increases, the role of automatic information organization becomes more important. The area of classifying the text documents into several groups has been the subject of intensive researches during the last decade.

Classification is an important method for data analysis and there are several techniques in artificial intelligence and pattern recognition which have been proposed for this task, but these methods can not be applied efficiently in text classification task because in this problem we will encounter a high number of features. This thesis aims to investigate the concepts that must be considered in text classification task, like: feature extraction document representation, feature subset selection and proper machine learning algorithms for this task. In feature subset selection we have tried to show the relation between probabilistic classification and criterion which is used in our work. Also, we

IV

completely considered a probabilistic framework for text classification. In this framework the probability of belonging one document to all classes will be estimated. In this method we use Bayesian networks as an effective and efficient way of saving a joint probability distribution of variables. Then we have proposed architecture for building a knowledge model using Bayesian Networks. By measuring the performance and comparing the results with one classic algorithm, which commonly used in Information Retrieval, we have shown that these proposed approaches are efficient. We used Reuters data for learning and test and experiments have been carried out in 5 categories of these data. Finally we have shown that inference is also efficient in the proposed models.

# TABLE OF CONTENTS

**Content**                                             **Page**

**Content**                                                                **Page**

**CHAPTER 3:   MACHINE   LEARNING   FOR   TEXT**

**Content**                                                **Page**

**Content**                                                      **Page**

# LIST OF TABLES

XII

# LIST OF FIGURES