# Cairo University

## Faculty of Graduate Studies for Statistical Research
## Department of Information Systems and Technology

# Agent based system for distributed data mining

**A thesis submitted in partial fulfillment of the requirements**
**For The M.SC Degree in Information System and Technology**

**Submitted by**
**Doaa Mostafa Hussein Ahmed**

**Supervised by**

### Prof. Dr. Hesham Ahmed Hefny
Vice-Dean for Graduate Studies
Faculty of Graduate Studies for Statistical Research
Cairo University, Egypt

### Dr. Ammar Mohammed Ammar
Assistant professor at computer Science Department
Faculty of Graduate Studies for Statistical Research
Cairo University, Egypt

### Dr. Maryam Mohey Eldin Hazman
Senior Researcher
Agriculture Research Center, Egypt

# 2021

# ABSTRACT

Data mining technology has appeared to discover patterns and trends from large quantities of data. Distributed Data Mining (DDM) is emerged from the need of mining over decentralized data sources. When using a batch approach for distributed data is complex and expensive. It requires techniques to improve performance and reduce complexity in a proper way. Multi-Agent System (MAS) is one of such techniques which handle the complexity and the distribution of data by efficient way.

This research advances the understanding of a multi-agent approach to data mining of large datasets. An agent mining architecture called ADDM (Agent based Distributed Mining) is developed for the purpose of building accurate and transparent cluster and improving the efficiency of mining a large dataset by using MAS. In the proposed architecture, several agents are distributed over local servers. The novelty of our work is to select the best clustering result on each local server side and use MAS system to manage resources to impose the processing power of distributed infrastructures and reduce time and cost.

The ADDM approach provides innovative distributed data mining model, with great research and commercial prospect for distributing mining across multiple agents and different data sources. This thesis work to value greatly the idea that combining data mining and multi-agent approaches in large scale data mining applications.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ALGORITHMS