Computer Science Department
Faculty of Computer & Information Sciences
Ain Shams University

# Computational Intelligence Techniques for Big Data Analytics

A Thesis Submitted to Computer Science Department,
Faculty of Computer & Information Sciences
Ain Shams University, Cairo, Egypt

In Partial Fulfilment of the Requirements for
the Degree of Doctor of Philosophy in Computer Science

## By

Mahmoud Ibrahim Elbattah
Doctoral Researcher
Ain Shams University

## Under Supervision of

Prof. Dr. Abdel-Badeeh Mohamed Salem
Computer Science Department
Faculty of Computer & Information Sciences
Ain Shams University

Prof. Dr. Mohamed Roushdy
Computer Science Department
Faculty of Computer & Information Sciences
Ain Shams University

Prof. Dr. Mostafa Aref
Computer Science Department
Faculty of Computer & Information Sciences
Ain Shams University

April 2017

# Declaration

This is to certify that this work has not been accepted in substance for any academic degree and is not being concurrently submitted in candidature for any other degrees.

Any other portions of this thesis for which the author is indebted to other sources are mentioned and explicit references are given.

Scholar Name: Mahmoud Elbattah

Signature:

*Elbattah*

# Acknowledgements

This work would not have been possible without all who have helped me along the way. First, the author would like to dedicate this work to his parents for their unbreakable faith and moral support over long years. Their constant encouragement helped the author to continue to be persistent and diligent in his work.

Second, the author would like to express his sincere gratitude to the thesis advisors Prof. M.Roushdy, Prof. M.Aref, and Prof A.Salem for their guidance, advice, criticisms, encouragements and insight throughout the research development. In particular, Prof. A.Salem showed immense understanding and patience for the research project during difficult times.

Furthermore, the author would like to thank Prof. Dima Shepelyansky (Université Paul Sabatier Toulouse, France) who organised the Summer School for Advanced Sciences of Luchon in July 2014. The author believes that the knowledge gained from that summer school provided useful insights to address the research problem from a different angle. In particular, the lectures and tutorials provided by Prof. Andreas Kaltenbrunner (Technology Centre of Catalonia, Spain) on network analysis provided a solid foundation to develop an important part of the research project.

Next, the author wishes to thank the providers of the MOOCs (massive open online course) such as Coursera.org and edX.org. The author cannot be more fortunate to have gained a learning experience from such a supporting and knowledgeable community. Specifically, the author is grateful to Prof. Andrew Ng for his comprehensive course on Machine Learning provided on Coursera.org. Special thanks also to Pro.f Yaser S. Abu-Mostafa (California Institute of Technology) who provided the wonderful course "Learning from Data".

Equally important, the author would like to give special thanks to Dr. Steve Elston (Quantia Analytics) who delivered the "DAT203x Data Science and Machine Learning Essentials" module on edX. The author can confidently confirm that Dr Steve provided him with valuable technical knowledge on handling Big Data processing tasks using the Azure Machine Learning Studio.

# Abstract

The world of Big Data continues to expand, and forge the landscape of decision making and analytics. Datasets are rapidly growing in size and complexity, and there is a pressing need to develop solutions to harness this deluge of data for producing useful insights. This study addresses the tasks in relation to storing, querying, analysing and visualising Big Data from a graph-based perspective. Through the study, datasets extracted from the immense knowledgebase of Freebase are utilised. Initially, a web-based tool for data visualisation was developed, named as FreebaseViz, for visually exploring the schema of Freebase data. The visualisation design is built upon node-link network layouts, which can facilitate exploring connectivity, visual search and analysis, and visualising patterns underlying the schema graph. FreebaseViz is claimed to enable users to interact with the schema visualisations to filter and drill into lower levels of detail, and highlight subsets of the schema graph. In addition, a graph database-oriented approach is embraced in a further bid to boost the visualisation query-ability using graph-based query operations.

Subsequently, the study conducted a graph-driven methodology for the analysis and visualisation of Freebase complex schema. Specifically, our methodology utilised Freebase schema objects in order to construct a directed weighted graph. The schema graph is employed to perform a modularity-based analysis in order to detect communities underlying Freebase data. In light of that, the detected communities were effectively used for the purpose of revealing unobserved or implicit domain relationships.

In terms of storing and querying large-scale datasets, a graph database-oriented approach is proposed, which considered Freebase data as a large graph. The proposed approach endeavoured to address the limitations encountered within traditional relational models. Furthermore, scalability and query efficiency of the approach are verified based on empirical experiments using a subset of Freebase data that comprised a large-scale graph consisting of more than 500K nodes, and 2M edges

Furthermore, the study addresses the problem of entity clustering within large-scale knowledge graphs with application to the knowledgebase of Freebase. Particularly, the clustering task is approached form a mere graph-driven perspective. Entities were aimed to be clustered based on structural similarity within a knowledge graph. In this manner, entities were clustered in an unsupervised fashion by matching their link-based structure rather than relational attributes.

Eventually, the study aimed to develop an approach for estimating the consistency of knowledgebase triples. The proposed approach was based on utilising machine learning in order to learn the graph-based patterns of the triples. Specifically, the study investigated the feasibility of training a model to learn triples patterns in terms of subject-predicate-object. The validity of the method was experimented using a relatively large-scale subset of Freebase data. The dataset incorporated about 10M triples, which contained 6M true patterns and 4M false patterns randomly generated. The study availed of the cloud platform of Microsoft Azure in order to conduct the large-scale machine learning experiments efficiently. On top of the Azure platform, an Apache Spark cluster was deployed to realise a distributed computing environment. The classifier model evidently demonstrated a relatively high accuracy. Broadly, the study endeavoured to present and emphasise the appropriateness of graph-based methods for dealing with Big Data scenarios in terms of storage, querying, visualisation, and predictive analytics.

# Table of Contents

# Table of Contents (cont'd)

# Table of Contents (cont'd)

# List of Figures

# List of Figures (cont'd)

# List of Tables