# MISSING DATA IMPUTATION, OPTIMIZATION AND LOCALIZED SOFTWARE DEVELOPMENT EFFORT ESTIMATION

MUHAMMAD ARIF SHAH

UNIVERSITI TEKNOLOGI MALAYSIA

## UNIVERSITI TEKNOLOGI MALAYSIA

# DECLARATION OF THESIS / UNDERGRADUATE PROJECT REPORT AND COPYRIGHT

Author's full name : MUHAMMAD ARIF SHAH

Date of Birth : 05 October 1986

Title : Missing Data Imputation, Optimization and Localized Software Development Effort Estimation

Academic Session : 2018-2019/2

I declare that this thesis is classified as:

☐ **CONFIDENTIAL** (Contains confidential information under the Official Secret Act 1972)*

☑ **RESTRICTED** (Contains restricted information as specified by the organization where research was done)*

☐ **OPEN ACCESS** I agree that my thesis to be published as online open access (full text)

1. I acknowledged that Universiti Teknologi Malaysia reserves the right as follows:

2. The thesis is the property of Universiti Teknologi Malaysia

3. The Library of Universiti Teknologi Malaysia has the right to make copies for the purpose of research only.

4. The Library has the right to make copies of the thesis for academic exchange.

Certified by:

| **SIGNATURE OF STUDENT** | **SIGNATURE OF SUPERVISOR** |
|---|---|
| MH5145682 | ASSOC. PROF. DR. DAYANG NORHAYATI ABANG JAWAWI |
| **Passport number** | **NAME OF SUPERVISOR** |
| Date: 20 JUNE 2019 | Date: 20 JUNE 2019 |

NOTES : If the thesis is CONFIDENTIAL or RESTRICTED, please attach with the letter from the organization with period and reasons for confidentiality or restriction

**Status Declaration Letter**

Librarian                                                                                       20 JUNE 2019

Perpustakaan Sultanah Zanariah

UTM, Skudai

Johor

Respected Sir/Madam,

CLASSIFICATION OF THESIS AS RESTRICTED


MISSING DATA IMPUTATION, OPTIMIZATION AND LOCALIZED
SOFTWARE DEVELOPMENT EFFORT ESTIMATION


Please be informed that the above-mentioned thesis entitled — **Missing Data Imputation, Optimization And Localized Software Development Effort Estimation** be classified as RESTRICTED for a period of three (3) years from the date of this letter. The reasons for this classification are:


(i) COMMERCIALIZATION OF RESEARCH RESULTS.
(ii) COPYRIGHT RULES RELATED TO WOS JOURNALS.


Thank you.

Sincerely yours,

ASSOC.PROF.DR. DAYANG NORHAYATI ABANG JAWAWI

N28-305-11

07- 5538870

"We hereby declare that we have read this thesis and in our opinion this thesis is sufficient in term of scope and quality for the award of the degree of Doctor of Philosophy in (Computer Science)"

| | | |
|---|---|---|
| Signature | : | _____ |
| Name of Supervisor I | : | ASSOC. PROF. DR. DAYANG NORHAYATI ABANG JAWAWI |
| Date | : | 20 JUNE 2019 |

| | | |
|---|---|---|
| Signature | : | _____ |
| Name of Supervisor II | : | DR. MOHD ADHAM BIN ISA |
| Date | : | 20 JUNE 2019 |

**BAHAGIAN A - Pengesahan Kerjasama***

Adalah disahkan bahawa projek penyelidikan tesis ini telah dilaksanakan melalui kerjasama antara _____dengan _____

Disahkan oleh:

Tandatangan :                            Tarikh :

Nama :

Jawatan :

(Cop rasmi)

*Jika penyediaan tesis atau projek melibatkan kerjasama.*

═══════════════════════════════════════════════

**BAHAGIAN B - Untuk Kegunaan Pejabat Sekolah Pengajian Siswazah**

Tesis ini telah diperiksa dan diakui oleh:

Nama dan Alamat Pcmeriksa Luar     **:**




Nama dan Alamat Pcmeriksa Dalam   **:**




Nama Penyelia Lain (jika ada)        **:**




Disahkan oleh Timbalan Pendaftar di SPS:

Tandatangan    :                            Tarikh : 15JULAI 2018

Nama           :

MISSING DATA IMPUTATION, OPTIMIZATION AND LOCALIZED
SOFTWARE DEVELOPMENT EFFORT ESTIMATION

MUHAMMAD ARIF SHAH

A thesis submitted in fulfilment of the
requirements for the award of the degree of
Doctor of Philosophy

School of Computing
Faculty of Engineering
Universiti Teknologi Malaysia

JUNE 2019

# DECLARATION

I declare that this thesis entitled *"Missing Data Imputation, Optimization and Localized Software Development Effort Estimation"* is the result of my own research except as cited in the references. The thesis has not been accepted for any degree and is not concurrently submitted in the candidature of any other degree.

Signature       :  ....................................................

Name           :  Muhammad Arif Shah

Date             :  20 JUNE 2019

# DEDICATION

This thesis is dedicated to my father, who taught me that the best kind of knowledge to have is that which is learned for its own sake. It is also dedicated to my mother, who taught me that even the largest task can be accomplished if it is done one step at a time.

The sacrifices of my wife (Shahida Arif) and Daughter (Eshal Arif) can never be disregarded, therefore I dedicate my work to them too.

# ACKNOWLEDGMENT

In preparing this thesis, I was in contact with many people, researchers, academicians, and practitioners. They have contributed towards my understanding and thoughts. In particular, I wish to express my sincere appreciation to my main thesis supervisor, **Assoc. Prof. Dr. Dayang Norhayati Abang Jawawi,** for encouragement, guidance, criticism, and friendship. I am also very thankful to my co-supervisor **Dr. Mohd Adham** Bin Isa for his guidance, advices, and motivation. Without their continued support and interest, this thesis would not have been the same as presented here.

My fellow postgraduate student should also be recognized for their support. My sincere appreciation also extends to all my colleagues and others who have provided assistance at various occasions. Their views and tips have been useful indeed. Unfortunately, it is not possible to list all of them in this limited space. I am also grateful to all my family members.

# ABSTRACT

The accuracy of software development effort estimation is one of the vital factors that leads to successful or failed projects. Most of the current estimation models are not adaptable enough according to the nature of software projects due to their special characteristics, such as intangibility and non-normality of project attributes. Localization-based estimation models solve the issues, but are unable to compare software projects, based on Language Type for accurate effort estimation. Secondly, these models use Analogy-Based Estimation (ABE), which completely depends upon past projects and any missing values that may cause unrealistic estimation results. This study extended the domain of localization to estimate the development effort according to the nature of software and introduced accurate missing data imputation techniques to prevent losing the most similar project. This study focused on ABE, which is a widely accepted non-algorithmic model incorporated in localized estimation. Five estimation models such as Localized Analogy Based Estimation (LABE), artificial Bee colony guided Analogy Based Estimation (BABE), Localized BABE (LBABE), Imputation and Optimization based Effort Estimation (ImOEE), Localized Imputation and Optimization based Effort Estimation (LImOEE) and three missing data imputation techniques such as Median Imputation of the Nearest Neighbours (MINN), Localized Imputation Technique (LIT), and Identical Project based Imputation (IPI) were developed in this study for accurate and unbiased development effort estimation. The techniques accurately filled the missing information in the past projects and the models dealt with the attribute weight optimization, project attribute selection in local space and model comparison. Imputation techniques consist of distance calculation and impute value calibrations with localization. The models are commonly composed by introducing the missing data imputation and soft computing techniques, and ABE. In this research, the models were evaluated using six real datasets. The results were compared with prominent estimation models. A comparative study of the developed models was performed to further validate the accuracy of the results. LImOEE model outperformed the other developed models. LImOEE showed 51%, 25%, 11%, and 31% improvements on Mean Magnitude of Relative Error (MMRE), Percentage of Prediction (PRED), Standard Accuracy (SA) and Effect Size ($\Delta$) respectively for the BABE model. For, LBABE, it showed 37% improvement on MMRE, 12% improvement on PRED (0.25), 5% improvement on SA and 18% improvement on $\Delta$. For the ImOEE model, it showed 26%, 8 %, 9% and 28% improvements on MMRE, PRED, SA and $\Delta$ respectively. The results revealed that, accurately imputing the missing data for ABE, optimizing the attribute weights, and extending the scope of localization to the important attributes have significantly improved the accuracy of software development effort estimation.

# ABSTRAK

Ketepatan anggaran pembangunan perisian adalah salah satu faktor penting yang menyumbang kepada kejayaan dan kegagalan sesebuah projek. Kebanyakan model angaran terkini tidak sesuai berdasarkan sifat sebenar sesebuah projek berdasarkan ciri-ciri khas seperti atribut projek yang tidak ketara dan tidak normal. Model anggaran berasaskan tempatan telah menyelesaikan isu ini, tetapi tidak dapat membandingkan projek perisian, berdasar kepada jenis bahasa, untuk anggaran usaha yang tepat. Kedua, model ini menggunakan Anggaran Berasaskan-Analogi (ABE) yang keseluruhannya bergantung kepada projek lepas dan sebarang nilai yang hilang yang mungkin menyebabkan hasil anggaran yang tidak realistik. Kajian ini memperluaskan domain tempatan untuk menganggarkan usaha pembangunan berdasarkan sifat semulajadi perisian dan memperkenalkan teknik mengenal pasti taksiran data hilang untuk mengelakkan dari kehilangan kebanyakan projek yang sama. Kajian ini memberi tumpuan kepada ABE, yang merupakan model bukan beralgoritma dan berkerjasama dengan angaran tempatan. Lima model anggaran seperti *Localized Analogy Based Estimation* (LABE), *artificial Bee colony guided Analogy Based Estimation* (BABE), *Localized BABE (LBABE), Imputation and Optimization based Effort Estimation* (ImOEE), *Localized Imputation and Optimization based Effort Estimation* (LImOEE) dan tiga teknik taksiran data hilang seperti *Median Imputation of the Nearest Neighbours* (MINN), *Localized Imputation Technique* (LIT), and *Identical Project based Imputation* (IPI) dicadangkan dalam kajian ini untuk ketepatan dan anggaran perkembangan usaha yang tidak berat sebelah. Teknik yang betul memenuhi taksiran data yang hilang dalam projek lepas dan model yang mengendalikan atribut mengurangkan berat, pemilihan atribut tempatan dan perbandingan model. Teknik taksiran terdiri daripada pengiraan jarak dan menaksir nilai dengan penempatan. Model ini biasanya dibuat dengan memperkenalkan teknik menaksir data hilang dan teknik komputeran lembut, ABE. Dalam kajian ini model telah dinilai menggunakan enam data set asal. Hasil kajian telah dibandingkan dengan model anggaran yang terkenal. Kajian perbandingan untuk model yang dibangunkan telah dilakukan untuk mengesahkan ketepatan hasil kajian. Model LImOEE telah mengatasi model lain yang telah dibangunkan. LImOEE menunjukkan peningkatan sebanyak 51%, 25%, 11% dan 31% untuk *Mean Magnitude of Relative Error* (MMRE), *Percentage of Prediction* (PRED), *Standard Accuracy* (SA) dan *Effect Size* ($\Delta$) untuk model artificial Bee colony guided Analogy Based Estimation (BABE). Untuk LBABE, ia menunjukkan peningkatan 37% ke atas MMRE, peningkatan 12% ke atas PRED (0.25), peningkatan 5% ke atas SA dan peningkatan 18% ke atas $\Delta$. Untuk model ImOEE, ia menunjukkan peningkatan sebanyak 26%, 8%, 9% dan 28% ke atas MMRE, PRED, SA dan $\Delta$. Hasil kajian menunjukkan bahawa ketepatan taksiran data hilang untuk ABE, mengoptimumkan berat atribut dan memperluaskan skop tempatan kepada atribut penting telah meningkatkan ketepatan anggaran pembangunan perisian.

# TABLE OF CONTENTS

| TITLE | PAGE |
|---|---|