# LOCALIZATION OF SOFTWARE DEVELOPMENT EFFORT ESTIMATION BASED ON CLASSIFICATION OF PROJECTS

VAHID KHATIBI BARDSIRI

A thesis submitted in fulfilment of the
requirements for the award of the degree of
Doctor of Philosophy (Computer Science)

Faculty of Computing
Universiti Teknologi Malaysia

JUNE 2013

I declare that this thesis entitled: *"Localization of Software Development Effort Estimation based on Classification of Projects"* is the result of my own research except as cited in references. The thesis has not been accepted for any degree and is not concurrently submitted in candidature of any other degree.

Signature        :

Name            : VAHID KHATIBI BARDSIRI

Date            :

"To my beloved wife and daughter"

# ACKNOWLEDGEMENT

I would like to express my special appreciation and thanks to my supervisor Assoc.Prof.Dr. Dayang Norhayati Abang Jawawi, you have been a tremendous mentor for me. I would like to thank you for encouraging my research and for allowing me to grow as a research scientist. Your advice on both research as well as on my career have been invaluable. My sincere appreciation also extends to Assoc.Prof.Dr. Siti Zaiton mohd Hashim who has provided assistance at various occasions. Her views and tips have been really useful throughout the current research. At the same time, I also appreciate Universiti Teknologi Malaysia for International Doctorate Fellowship (IDF) award.

A special thanks to my family. Words cannot express how grateful I am to my family. Your prayer for me was what sustained me thus far.

I would also like to thank my beloved wife, Elham. Thank you for supporting me, and especially I can't thank you enough for encouraging me throughout this experience. To my beloved daughter, Bita, I would like to express my thanks for being such a good girl always cheering me up. I will keep on trusting You for my future.

# ABSTRACT

Accurate estimation of software development effort is strongly associated with the success or failure of software projects. Since the existing estimation models are confronted by special characteristics of software projects such as non-normality of project attributes, intangible product, and lack of well-defined measurement techniques, the accuracy of estimates is not convincing and the estimation models are not sufficiently adaptable to support a wide range of software projects. The main issue is that prior studies have considered software projects to be similar to other types of projects regardless of the fact that software projects are entirely different. Furthermore, an accurate estimation of development effort in software projects is unreachable in global space, thus this research proposed that the estimation process be localized based on the characteristics of the project being estimated. The aim of the study is to alleviate the inconsistency and heterogeneous nature of software projects through the ideas of localization and classification. This research focused on a widely accepted estimation model called analogy based estimation (ABE). Four different estimation models were proposed to overcome the problems of: low quality in construction of machine-learning based estimation models, difficulty of attribute weighting, limited support of software projects, and biased comparisons. Each proposed model was a hybrid model constructed using a classification method, ABE, and a complimentary method (soft computing or statistical techniques). Performance evaluations of the models were carried out using five real software project datasets and the estimates were compared to those achieved by other well-known estimation models. A comparative study was also conducted for further verification of the models. The evaluation results confirmed the ability of the proposed models to improve the accuracy of estimates, quality of training, efficiency of attribute weighting, adaptability of estimation process, correctness of comparisons, and support domain of software projects.

# ABSTRAK

Ketepatan anggaran untuk usaha pembangunan perisian adalah sangat berkait rapat dengan kejayaan atau kegagalan projek perisian. Oleh kerana model anggaran usaha yang sedia ada berhadapan dengan ciri-ciri khas projek perisian seperti ketidak-normalan atribut projek, ketidak-ketaraan produk, dan kekurangan teknik pengukuran yang jelas, ketepatan anggaran tidak begitu menyakinkan dan model anggaran tersebut tidak mudah disesuaikan untuk menyokong pelbagai projek perisian. Isu utama ialah kajian sebelum ini menganggap projek perisian adalah sama seperti jenis projek lain tanpa mengambilkira kenyataan yang projek perisian berbeza secara keseluruhannya. Tambahan pula, anggaran yang tepat pada projek perisian tidak tercapai pada ruang global, oleh itu penyelidikan ini mencadangkan proses anggaran yang disetempatkan berdasarkan ciri-ciri projek yang dianggarkan. Tujuan kajian ini adalah untuk mengurangkan ketidak-konsistenan dan tabii berheterogen projek perisian menerusi idea penyetempatan dan pengelasan. Penyelidikan ini memberi fokus kepada model anggaran penganggaran usaha yang telah diterima pakai secara meluas yang dinamakan model berasaskan analogi (ABE). Empat model anggaran yang berbeza telah dicadangkan untuk mengatasi: masalah kualiti rendah pada pembinaan model anggaran berasaskan pembelajaran-mesin, kesukaran penentuan pemberat atribut, sokongan terhad kepada projek perisian dan perbandingan berat sebelah. Setiap model yang telah dicadangkan adalah model gabungan yang dibangunkan dengan menggunakan kaedah pengelasan, ABE, dan kaedah pelengkap (pengkomputeran lembut atau teknik statistik). Penilaian-penilaian prestasi model-model telah dilakukan terhadap lima set data sebenar projek perisian dan anggaran dibandingkan dengan pencapaian model-model anggaran yang sedia ada. Kajian perbandingan telah dilaksanakan untuk pengesahan lanjut model. Keputusan penilaian mengesahkan keupayaan model yang dicadangkan memperbaiki ketepatan anggaran, kualiti latihan pembelajaran, kecekapan pemberat atribut, kebolehsesuaian proses anggaran, ketepatan perbandingan dan sokongan bidang projek perisian.

# TABLE OF CONTENTS