



**Ain Shams University**  
**Faculty of Engineering**  
Computer & Systems Engineering Department

## **Sizing and Estimation of Maintenance Software Projects**

### **Thesis**

Submitted in partial fulfillment of the Requirements for the M.Sc.  
Degree in Electrical Engineering  
(Computer and Systems Engineering)

Submitted by

**Nermeen Fawzy Ahmed Mohamed**  
B.Sc of Electrical Engineering  
Computer and Systems Engineering Department  
Faculty of Engineering, Ain Shams University, 2002

Supervised by

**Prof. Dr. Abdel-Monein Abdel-Zaher Wahdan**

**Prof. Dr. Ashraf Mohamed Salem**

Cairo 2009

# **APPROVAL SHEET**

## **Acknowledgements**

All praise be to ALLAH the High and All Mighty for everything I realize and what I don't realize. I would like to express my deepest appreciation to my thesis supervisors, Professor Abdel-Moneim Abdel-Zaher Wahdan, and Professor Ashraf Mohamed Salem. I greatly valued their support, care and guidance for me through all the work phases. Thanks also for my parents for supporting me all the way till I have become engineer, and for their help and continuous motivation. Thanks to my Kid, Omar, for his pure inspiring smile. Finally thanks to my husband for providing me with unconditional support and care. I'll never be thankful enough for him.

## ABSTRACT

The objective of this thesis is to study of the existing software cost estimation techniques, and to reveal some of the drawbacks of these techniques, and to propose a new model to get rid of these drawbacks. The proposed estimation model is a modified version of the famous COCOMO model, in which we adopted the use of fuzzy logic to override some of its deficiencies. The introduction of fuzzy logic tolerates the existence of imprecision and uncertainty in measurements throughout the whole estimation process.

As the proposed cost estimation model tolerates the imprecision in the value of the cost factors, then the model is less sensitive to slight changes between projects, and makes it give better estimation results. On the other hand the original version of the intermediate COCOMO'81 does not tolerate the imprecision of the values of the cost factors.

We validated our proposed model on the COCOMO benchmark which is composed of 4 databases of software projects. Each database contains 63 different projects. The results obtained from this validation show the advantage of our proposed model over the original COCOMO model.

### **Keywords:**

Software cost estimation, Software engineering, Fuzzy Logic.

# Contents

<b>Acknowledgements .....</b>	<b>3</b>
<b>ABSTRACT .....</b>	<b>4</b>
<b>Abbreviations .....</b>	<b>7</b>
<b>Chapter 1 .....</b>	<b>9</b>
<b>Introduction .....</b>	<b>9</b>
<b>Chapter 2 .....</b>	<b>13</b>
<b>Software Cost Estimation Techniques .....</b>	<b>13</b>
<b>2.1. Man-Month measure .....</b>	<b>13</b>
<b>2.2. Cost estimation techniques .....</b>	<b>15</b>
<b>2.2.1. Estimation based on expertise: .....</b>	<b>16</b>
<b>2.2.2. Estimation by analogy: .....</b>	<b>16</b>
<b>2.2.3. Descending and ascending estimations: .....</b>	<b>17</b>
<b>2.3. Cost estimation models .....</b>	<b>18</b>
<b>2.4. Parametric models for cost estimation .....</b>	<b>21</b>
<b>2.4.1. Linear models .....</b>	<b>21</b>
<b>2.4.2. Analytical models .....</b>	<b>27</b>
<b>2.4.3. Limitations of the parametric models .....</b>	<b>31</b>
<b>2.5. Non-Parametric models for cost estimation .....</b>	<b>34</b>
<b>2.5.1. Estimation Models based on Neural Networks: .....</b>	<b>34</b>
<b>2.5.2. Estimation Models based on Regression Trees: .....</b>	<b>37</b>
<b>2.5.3. Estimation Models using the genetic algorithms: .....</b>	<b>41</b>
<b>2.6. Validation of Cost Estimation Models .....</b>	<b>47</b>
<b>Chapter 3 .....</b>	<b>53</b>
<b>Fuzzy Logic .....</b>	<b>53</b>
<b>3.1 Introduction .....</b>	<b>54</b>
<b>3.2 Fuzzy sets .....</b>	<b>55</b>
<b>3.2.1. Definitions and terminology .....</b>	<b>55</b>
<b>3.2.2. Operations on fuzzy sets .....</b>	<b>57</b>
<b>3.2.3. Properties of fuzzy sets .....</b>	<b>59</b>
<b>3.3 Linguistic variables and numerical values .....</b>	<b>62</b>
<b>3.4 Fuzzy logic, probability theory, and measurement theory .....</b>	<b>63</b>
<b>3.5 Fuzzy reasoning .....</b>	<b>64</b>
<b>3.5.1. Fuzzy propositions .....</b>	<b>67</b>
<b>3.5.2. Fuzzy rules .....</b>	<b>68</b>

3.5.3. Fuzzy implications .....	69
3.5.4. Aggregation of fuzzy rules: .....	70
<b>Chapter 4 .....</b>	<b>72</b>
<b>A mixed Fuzzy-COCOMO Approach for Software Cost Estimation ...</b>	<b>72</b>
4.1 Introduction .....	72
4.2 Objectives of software measurement .....	73
4.3 Linguistic values in software measurement: .....	75
4.4 Managing the uncertainties of the COCOMO'81 model:.....	76
4.4.1. The COCOMO'81 model: .....	76
4.4.2. The Intermediate COCOMO'81 model: .....	78
4.4.3. The deficiency of the Intermediate COCOMO'81 model: ...	80
4.4.4. Fuzzy set representation of the linguistic values: .....	82
4.4.4.1 DATA (Size of the database) .....	83
4.4.4.2 TIME (Hardware utilization constraint) .....	83
4.4.4.3 STOR (Main storage constraint) .....	84
4.4.4.4 ACAP (Analyst Capability).....	85
4.4.4.5 PCAP (Programmer Capability).....	86
4.4.4.6 AEXP (Application Experience) .....	87
4.4.4.7 SCED (Required Development Schedule) .....	87
4.4.4.8 VEXP (Virtual machine experience).....	88
4.4.4.9 LEXP (Programming language experience).....	89
4.4.4.10 TURN (Computer Turnaround Time) .....	90
4.4.4.11 VIRT (Virtual machine volatility).....	91
4.4.5. Application of the proposed technique:.....	92
<b>Chapter 5 .....</b>	<b>98</b>
<b>Experimental Results .....</b>	<b>98</b>
5.1. Evaluation of the precision of the proposed approach.....	98
<b>Chapter 6 .....</b>	<b>104</b>
<b>Conclusions and Suggestions for Future Work .....</b>	<b>104</b>
6.1. Conclusions .....	104
6.2. Suggestions for future work .....	106
<b>References .....</b>	<b>107</b>

## Abbreviations

ACAP:	<u>A</u> nalyst <u>C</u> APability
AEXP:	<u>A</u> pplication <u>E</u> XPerience
AKDSI:	<u>A</u> ddjusted <u>K</u> ilo <u>D</u> elivered <u>S</u> ource <u>I</u> nstructions
ASE:	<u>A</u> verage <u>S</u> quared <u>E</u> rror
CASE:	<u>C</u> omputer <u>A</u> ided <u>S</u> oftware <u>E</u> ngineering
CBR:	<u>C</u> ase <u>B</u> ased <u>R</u> easoning
COCOMO:	<u>C</u> ONstructive <u>C</u> Ost <u>M</u> ODEl
COD:	<u>C</u> oefficient <u>O</u> f <u>D</u> etermination
CPLX:	Product <u>C</u> om <u>P</u> Le <u>X</u> ity
CTP:	<u>C</u> omputational <u>T</u> heory of <u>P</u> erceptions
CW:	<u>C</u> omputation with <u>W</u> ords
DSI:	<u>D</u> elivered <u>S</u> ource <u>I</u> nstructions
FP:	<u>F</u> unction <u>P</u> oint
GSC:	<u>G</u> eneral <u>S</u> ystem <u>C</u> haracteristics
GP:	<u>G</u> enetic <u>P</u> rogramming
IDS:	<u>I</u> nitial <u>D</u> ata <u>S</u> et
INQ:	<u>I</u> NQuiry
ISBSG:	<u>I</u> nternational <u>S</u> oftware <u>B</u> enchmarking <u>S</u> tandards <u>G</u> roup
KDSI:	<u>K</u> ilo <u>D</u> elivered <u>S</u> ource <u>I</u> nstructions
KLOC:	<u>K</u> ilo <u>L</u> ine <u>O</u> f <u>C</u> ode
LEXP:	Programming <u>L</u> anguage <u>E</u> XPerience
MIQ:	<u>M</u> achine <u>I</u> ntelligent <u>Q</u> uotient
MM:	<u>M</u> an- <u>M</u> onth
MMRE:	<u>M</u> ean <u>M</u> agnitude <u>R</u> elative <u>E</u> rror
MODP:	<u>M</u> ODern Programming <u>P</u> ractices
MP:	<u>M</u> ultilayer <u>P</u> erceptron

MRE:	<u>M</u> agnitude <u>R</u> elative <u>E</u> rror
MSE:	<u>M</u> ean <u>S</u> quare <u>E</u> rror
MU:	<u>M</u> emory <u>U</u> tilization
PCAP:	<u>P</u> rogrammer <u>C</u> APability
RELY:	Required Software <u>R</u> ELiability
SCED:	Required Development <u>S</u> ChEDule
SDMRE:	<u>S</u> tandards <u>D</u> eviation of <u>M</u> agnitude <u>R</u> elative <u>E</u> rror
STOR:	Main <u>S</u> TORage Constraint
TDS:	<u>T</u> erminal <u>D</u> ata <u>S</u> et
TE:	<u>T</u> otal <u>E</u> ffort of all the phases of a project
TURN:	Computer <u>T</u> URNaround Time
VEXP:	<u>V</u> irtual Machine <u>E</u> XPerience
VIRT:	<u>V</u> IRTual Machine Instability



# **Chapter 1**

## **Introduction**

The good management of the software development process is based on a good estimation of its cost. Without a good estimate, managers and other persons responsible for a software project will not be able to assign the budgets and the competences necessary for the project development. Thus, managers depending on the results of estimation could make serious errors by assigning inadequate time plans or insufficient budgets. These errors could be also reflected on the users. Extending the development time of a project would involve a delay in the delivery and installation of the software. The users would thus fall into economical and organizational problems, especially if they have already made up their planning and set up new structures by hiring new employees and arranging buildings for example.

Although precise estimates are difficult to establish, it is well known that software development costs represent a significant part of the project. In many software projects the development cost exceeds 80% of the total cost of the project [BOE81]. The growth of software costs as well as the difficulties that we meet to estimate them and control them still constitutes a major concern for project managers and researchers in this domain. Many problems are faced such as the increasing demands of users with their wide diversity, the continuously increasing costs, the effect of deadlines extensions, ... etc.

Due to the importance of finding a reliable estimate of the software development cost, this subject became the target of many research works. A close cooperation between the research and the industry was important in order

to develop and to validate estimation approaches. Cost estimation comprises the study of the three following attributes: the *effort*, the *time* and the *cost*.

Several techniques were proposed in order to formalize the cost estimation process [BOE81]. The techniques that are based on models are very well documented in the literature. These techniques can be classified into two categories [SHE97]:

- Parametric models
- Non parametric models

The parametric models are essentially based on a central equation expressing the effort as a function of certain attributes that affect the cost [HAL77] [WAL77] [PUT78] [BOE81] [BOE00]. These attributes are called *cost drivers*. Other parametric models are based on the calculation of function points [ALB83] [ABR96] [MAT94].

Non parametric models make use of artificial intelligence techniques such as analogy reasoning, genetic algorithms, rule based systems, and decision trees. These models were utilized to get rid of some of the inconveniences of parametric models. Non parametric models don't necessitate any specific form to express the cost in terms of cost factors. Examples of non parametric models are found in [HAM00], [POR90a], [POR90b], [MEN02], [SRI95], [SER95] [HUG96], [SHE97], [WIT97], [KAD00], and [BUR01].

Each type of the estimation models has its advantages and its disadvantages, and no one could prove that one type performs better than the other in all the situations. Actually the performance of a model depends on several characteristics of a given database of historical projects, such as their size, the number of attributes describing the projects, and the presence of extreme

points [SHE01]. The estimation using analogy and the Case-Based Reasoning (CBR) are promising techniques that are well adapted to the problem of cost estimation. These methods can be used when we have little knowledge about the problem to be solved, and for which the optimal solution is not known a priori. The estimation by analogy consists of using the of historical data to provide a solution of the studied problem [KAD00], [HAM00].

The thesis is composed of six chapters. In chapter 2 we present a survey on the different techniques used for software cost estimation. We emphasize more on the estimation techniques which are based on estimation models, being parametric or non-parametric. For each model we present the basic concepts, the advantages, the disadvantages, and some examples.

The third chapter presents the role of the fuzzy logic in the development of intelligent systems. This chapter discusses the advantages of using fuzzy logic for the presentation and the treatment of imprecision and the uncertainty in intelligent systems. It then presents the basic concepts of fuzzy logic such as the fuzzy sets, fuzzy propositions, fuzzy rules, fuzzy implication and aggregation of fuzzy rules.

In the fourth chapter, we discuss the massive use of linguistic variables by most of the cost estimation models. These models adopt a classic representation of these linguistic values. We illustrate the inconveniences of such representation, and we suggest the use of fuzzy representation. In this chapter we apply the fuzzy logic for the representation, and for the treatment of the six linguistic values used by the COCOMO model. We introduce new cost factors  $Fuzzy\_C_{ij}$ , which are calculated from the classical factors  $C_{ij}$  after modifying them using the membership functions associated with the different fuzzy sets.

In chapter 5 we evaluate the precision of the propose algorithm. The evaluation is made by comparing the estimations performed by the new approach with the results obtained from the original version of the intermediate COCOMO'81 model.

Finally in chapter 6, we present our conclusions and our suggestion for further studies.

## Chapter 2

### Software Cost Estimation Techniques

In this chapter, we give a literature review of the various cost estimation techniques in software development. We briefly present those techniques based on the *expertise*, the *analogy*, and the *descending* and *ascending* strategies. We then present in full details the technique based on modeling. This technique expresses the development effort as a function of the factors affecting the cost. For this technique there are two modeling approaches: the *parametric* approach, that uses statistical methods, and the *non-parametric* approach, that uses the artificial intelligence methods, such as the decision trees and the genetic algorithms. We terminate this chapter by summarizing the results obtained from some historic projects.

#### 2.1. Man-Month measure

The development effort of software represents the number of persons as well as the number of months necessary to complete the project. For example, if 10 persons worked for 5 months in a project, the development effort of this project is then 50 Man-Months. Consequently, if we know the monthly salary of the persons working in the project, then the project cost is simply the sum of the salaries multiplied by the number of months of work. This implies that the three attributes: *effort*, *time* and *cost*, are closely dependent.

Actually the major part of the development costs corresponds to the salaries. Thus, the cost of a software project is directly proportional to the number of Man-Months necessary to finish of the project. There are obviously, other costs, for example the hardware, the transportation, and staff training, but those

are less difficult to obtain. These other costs are not based on the productivity of the personnel, or the software size to be developed.

Due to the fact that the software development cost is determined by the number of persons and months of work, i.e. the development effort, the two terminologies of *cost estimation* and *effort estimation* are considered equivalent in software engineering. In general, most of the approaches of software cost estimation determine, in a first phase, the development effort, and then, they evaluate the cost by using the salaries of the involved personnel. The Man-Month unit, often used to measure the effort of software development, was the subject of a relevant criticism in [BRO95]. He showed that assigning more programmers to a project running behind schedule will make it even later, due to the time required for the new programmers to learn about the project, and the increased communications overhead. When  $N$  people have to communicate among themselves (without a hierarchy), as  $N$  increases, their output decreases and can even become negative (i.e. the total work remaining at the end of a day is greater than the total work that had been remaining at the beginning of that day, such as when many bugs are created).

One of the results shown in [BRO95] was that persons and months are not interchangeable. It is not always true that we can reduce the development time by increasing the number of persons working on the project. This assumption remains valid in certain traditional industries where a task can be easily divided and implemented by a group of persons without intercommunication between them. Thus, in such cases, the addition of the personnel in the project makes it possible to reduce its development time. However, in software this situation is rarely encountered; the activities of software development are often

not divisible. Sometimes tasks can be partitioned but with very high rate of communication between the involved persons.

## 2.2. Cost estimation techniques

In the literature there are a plenty of cost estimation techniques. [BOE81] presents seven of these techniques that are widely used nowadays. These seven techniques are:

- Estimation based on expertise
- Estimation by analogy
- Estimation based on Parkinson law (a project will cost what you decided to spend)
- Estimation based on the customer price (a project will cost what the customer is willing to spend)
- The descending estimate (Top-down)
- The ascending estimate (Bottom-up)
- Estimation based on modeling

However, according to [FEN97] the techniques based on *Parkinson law* and the *customer price* cannot be considered as techniques for cost estimation. The method based on Parkinson law expresses the goal not an estimate. The technique based on Parkinson law is a constraint with the project which can determine its feasibility.

Each one of the above techniques has its advantages and disadvantages, but according to [BOE81] none of these techniques can be used exclusively. In the following we briefly present these techniques.