

Selecting a Standard Set of Attributes for Cost Estimation of Software Projects

Arvind Sharma

Assistant Professor

PG Deptt. Of CS&IT, DAV COLLEGE, AMRITSAR

Abstract— the aim of the software engineering is to enhance projects that produce the needed results within limited schedule and budget. So that, software effort estimation becomes a valuable manner since it limits the problems of overestimate and underestimate for the software.

Software cost estimation is the process of predicting the effort required to develop a software system. There are many estimation models over the last decade, and in this paper, we use six public cost estimation data sets that we obtained from promise repository. We perform regression analysis over these data sets and perform a feature selection in order to get the most effective attribute to the effort. Finally, we analyze and compare the results obtained from each data set to build a framework for the standard set of metrics that we suggest each cost estimation data set must contain.

Keywords- cost estimation, feature selection;ols regression, metrics

I. INTRODUCTION

In recent years, software has become the most expensive component of computer system projects; for that, software cost estimation becomes critical and significant to both developers and customers by predicting the effort required to improve a software system. It is essential for both developers and customers to get accurate software cost estimation by accurately estimating the new project cost. Project managers can provide the customers with an accurate deadline for their projects and debate some of the contract negotiation's issues. So, customers can expect actual development costs to be in line with estimated cost. As well as, these estimations can be used by developers to generate reports and proposals in order to determine what resources are needed to commit to the project. These resources will be used as a result of the prioritizing development projects with respect to an overall business plan. Accurate software cost estimation makes project

management easier to be managed and controlled as resources are better matched to the real needs [9].

Many researchers in software engineering field have in depth studied how to predict the software project cost which is important for the project managers and software development organizations. Cost estimation, or what is called "effort prediction" is the process of estimating the cost of the software system development. This estimation can generally be estimated through three methods: experts' judgment, algorithmic model and analogy- based method.

Several techniques have been proposed in the past decades in order to make an accurate cost estimation for the projects and then avoid the overruns in the budget and increase the organization efficiency as a result of improving software investment analysis [1]. A key factor in selecting cost estimation model is the accuracy of its metrics since these models depend on their metrics which act as an input to the model. Metric can be defined as a "quantitative measure of the degree to which a system, component, or process possess a given attribute in order to produce a reliable assessment of these attributes in the real problem"[1].

There are too many data sets in the software cost estimation area with different attributes. That makes it difficult and annoying to a person who wants to analyze and correlate the data set in order to use the best and the strongest one in his analysis, and make cost estimation for their projects. So that, there is a necessary need for a valuable standard model. This suggestion can be constructed by obtaining the similarities among those available datasets. This can be applied by first understanding the different attributes for each dataset and finding the effective set of metrics which directly affects effort estimation using regression. After that, we make a comparison among different datasets in order to build a framework which consists of a common set of metrics which serves the user with best prediction, and improves the accuracy estimate of effort required to build a software system.

II. BACKGROUND

A- Cost estimation

Projects managers need to determine the cost of their ordered projects in order to manage their budgets and determine the deadline of the projects as well as improving the overall quality of any new projects. So, project managers need an effective model which helps them in accurate cost estimation according to their current project requirements.

Heemstra [6] discussed many questions related to software cost estimation; those questions are emphases on the causes for the excesses of the budgets and ruled periods. He explained the preconditions for the estimation and the way for this process. The benefit the software project management can obtain from the used models was explained

by viewing the strengths and weaknesses of cost estimation models.

B- Cost estimation methods

Cost estimation methods can be grouped under three categories: expert judgment, analogy- based estimation [10] and algorithmic estimation [9]. In this paper, we are interested in Algorithmic estimation, where the algorithmic model estimates the software cost through some formulas, depending mainly on the size of the project which is measured in terms of Function Point, Object Point and Line of Code (LOC).

In addition to the size of the project, there are several variables participated in the algorithmic model function such as:

$$\text{Effort} = \text{function}(\text{var}_1, \text{var}_2, \text{var}_3 \dots \text{var}_n)$$

Where effort is a cost estimation measure that is usually measured by (person-month), function refers to the function form, and $(\text{var}_1, \text{var}_2, \text{var}_3 \dots \text{var}_n)$ refers to the cost factors.

In the COCOMO 2 model, Boehm proposed a set of cost factors which are classified into four groups as:-

- Product factor which refers to the software product features like product complexity, database size used...etc.
- Computer factor which refers to the computer characteristics.
- Personal factor which refers to the development staff capability.
- Project factor which and the project environment refers to the project work process.

Through the algorithmic model, some mathematical formula forms can be used such as [9]:

* Linear models.

* Multiplicative models: They refer to the form:

$$(\text{Effort} = a_0 \prod a_i^{\text{var}_i}) \quad (1)$$

Where var_i refers to the project factor, $a_0, a_1 \dots$ are the coefficients extracted from experimental calculations.

* Power function model: They refer to the form $(\text{effort} = a * s^b)$. (2)

Where s is the software size usually measured by lines of code (LOC), a and b are the coefficients extracted from experimental calculations of the data set.

COCOMO (Constructive Cost Model), proposed by Boehm, is an example of the algorithmic model, which is a software package that helps assisting projects managers in planning a cost estimation of a software development project through an interactive interface. Furthermore, COCOMO is a highly subjective to the users input variables by noting the equations coefficients.

Korte and Port [8] implemented a standard and easily-analyzed statistical methods –standard error and bootstrapping –to several COCOMO 1 model research findings. The primary focus is on the confidence obtained from the findings that are experimentally based on estimators for error distribution parameters. As a result, this will reduce the contradictory and the lack of confidence in several published cost estimation research results based on Precisely MMRE and PRED comparisons such as model selection.

C - Feature subset selection

Feature selection which is referred to as subset selection is a pre- processing procedure used in machine learning where a subset of the features obtained from the data is chosen for implementation of a learning algorithm. The best subset includes the minimum number of diminutions which increase the accuracy by eliminating the inefficient dimensions [11].

Das and Kempe [5] demonstrated the problem of selecting a subset of k random variables to perceive that will produce the best linear prediction of different important variable, taking into account the pair wise correlations between the other variables and the predictor variable.

Azzeh, Neagu and Cowling [2] checked the influence of using feature subset selection algorithms in enhancing the accuracy of analogy software effort estimation models. They validated their works using two established data sets (ISBSG and Desharnais) use MMRE as evaluation criteria for all feature subset selection algorithms. They concluded that the employment of a fuzzy feature subset selection algorithm in analogy software effort estimation can give a valuable result. Menzies, Port and Boehm [4] discovered that COCOMO's estimates can be enhanced by using WRAPPER which is a feature subset selection method improved by the data mining industry. The results showed that the features subset selection always enhance the PRED(30) values without rising variance when applying on different data sets and as a result this will improve the strength of the COCOMO's prediction.

Kirsopp, Shepperd and Hart [7] explained the employment of using search techniques to aid the enhancement of case – based reasoning (CBR) system used in the software effort prediction. They checked the use of random searching, hill climbing and forward sequential selection (FSS) in order to get the optimal feature subsets that affect the effort prediction. They concluded that using a random search is better than using all the features. However using hill climbing and FSS can give better results than random search. They suggested using some form of heuristic –based initialization that can help in gaining improved results.

III. The Methodology

Software cost estimation can be considered as an empirical process that can be used to calculate the effort and the

(adj) – Maxwell

For cocomonasa data set, the model is displayed in Figure 6 that shows the selected attributes with R square (adj) value

Vars	R-Sq	R-Sq (adj)	Mallows	C-p	S	R	D	C	T	S	V	T	A	A	P	F	V	L	M	T	S	L	
						Y	A	X	E	R	T	N	P	P	P	P	P	P	P	P	P	P	P
1	84.2	84.0	36.5	262.97																			
1	85.6	85.6	424.7	641.76																			
2	86.1	85.6	27.9	249.64																			
2	85.8	85.3	29.2	251.66																			
3	88.6	88.0	14.8	227.46																			
3	87.6	86.9	20.9	237.62																			
4	89.7	89.0	10.3	218.03																			
4	89.5	88.7	11.7	220.59																			
5	90.5	89.6	7.6	211.37																			
5	90.5	89.6	8.1	212.19																			
6	91.1	90.1	6.3	206.95																			
6	91.0	90.0	6.9	208.00																			
7	91.5	90.3	6.1	204.41																			
7	91.4	90.3	6.3	204.87																			
8	91.8	90.5	6.2	202.51																			
8	91.7	90.4	6.5	203.03																			
9	92.3	90.9	5.4	198.42																			
9	92.0	90.5	6.7	201.28																			
10	92.5	90.9	6.2	197.77																			
10	92.3	90.8	6.9	199.41																			
11	92.6	90.9	7.5	198.36																			
11	92.5	90.8	8.0	199.51																			
12	92.6	90.7	9.4	200.14																			
12	92.6	90.7	9.5	200.34																			
13	92.6	90.6	11.2	201.79																			
13	92.6	90.5	11.4	202.23																			
14	92.7	90.4	13.0	203.67																			
14	92.6	90.4	13.2	204.01																			
15	92.7	90.2	15.0	205.89																			
15	92.7	90.2	15.0	205.97																			
16	92.7	90.0	17.0	208.27																			

is 90.9% and R square value equal to 92.3%.

Figure 6 Models building with different values of R square (adj) – cocomonasa

After applying the regression method, we got the following equation:

$$\text{Effort} = -2254 - 724 \text{ rely} + 910 \text{ cplx} + 1261 \text{ time} - 543 \text{ stor} - 702 \text{ virt} + 1991 \text{ acap} - 1135 \text{ aexp} + 950 \text{ pcap} + 7.05 \text{ loc.} \quad (9)$$

As to nasa93, we have the following equation for regression. Effort = 4771-1174 pr1- 1217 pr2-1615 pr3 – 2219 pr4 -802 pr5 -2368 pr6 -1712 pr7 -625 C1+ 1545 C2– 404 C3+ 108 C4-419 C5+ 159 C6+ 1210 C7– 307 C8+539 C9– 801 C10- 438 C11+ 474 C12– 72 C13-363 center + 1242 time +6.24 equivphyskloc – 2164 virt -2030 modp. (10)

It is noteworthy to mention that we perform some operations on the variables that are text in order to make the regression available; here we have the project name and cat2 that are indicator variables. The following table shows the indication of the used attributes.

Table 2 Indicator variables

Attribute	category
Pr1	De
Pr2	erb
Pr3	gal
Pr4	hst
Pr5	slp
Pr6	spl
Pr7	x
C1	Application groubd
C2	avionics
C3	avionicsmonitoring
C4	batchdataprocessing
C5	communications
C6	datacapture
C7	launchprocessing
C8	missionplanning
C9	monitor_control
C10	operatingsystem
C11	realdataprocessing
C12	science
C13	simulation

V. STANDARD SET OF ATTRIBUTES

Loc appears in different data sets with different names such as loc, ksloc and equivphyskloc. We have some attributes with different names such as T13 staff application knowledge is the same with aexp, and T15 (staff team skills) is the same to team.

Table 3 Resulted attributes from regression analysis

Data set	Attributes
Cocomo81	Rely, data, acap, aexp, modp, loc
Cocomo_sdr	Pmat, flex, time, pvol, prec, resl, loc
cocomonasa	Rely, cplx, time, stor, virt, acap, aexp, pcap, loc
kemerer	Hardware, duration, adjfp
Maxwell	Source, telonuse, T02, T13, T14, T15, duration, size
Nasa93	Program name, cat2, time, virt, modp, equivphyskloc

Table 4 Framework for the standard set of metrics used for effort estimation

Category	Standard attributes
Product Factors	Rely data cplx duration pmat flex prec
Platform Factors	Hardware pvol Time stor virt Source T02 (development environment adequacy)
Personnel Factors	acap aexp pcap modp team T14 (Staff tool skills)
Project Factors	Loc adjfp Size project name Cat2 telonuse resl

VI. CONCLUSION

Cost estimation plays a significant role in the software development since it attempts to avoid software overestimate or underestimate for the allocated budget. Various cost estimation models exist, but the effectiveness of each depends basically on its constituent attributes. In this paper, we make a feature subset selection over public selected data sets: cocomo81, cocomo_sdr, cocomonasa, Kemmerer, Maxwell and nasa93. We use regression analysis to obtain the best model built using minimum number of attributes. The results from the selected attributes are compared in order to build a framework of the standard set of metrics for effort estimation. This framework consists of attributes that may be classified into four categories: product factors, platform factors, personnel factors and project factors.

REFERENCES

[1] Ayyildiz, M., Kalipsiz, O., & Yavuz, S. (2006). A Metric-Set and Model Suggestion for Better Software Project Cost Estimation. World academy of science, Engineering and Technology, 23,167-172.
 [2] Azzeh, M. Neagu, D. & Cowling. (2008). improving analogy software effort estimation using a fuzzy feature subset selection algorithm, promise 08 proceeding of the 4th international workshop on predictor models in software engineering, 71-78.
 [3] Briand, L., Emam, K., Surmann, D., Wiecek, I., & Maxwell, K. (1999). an assessment and comparison of common software cost estimation modeling techniques, ICSE '99 proceedings of the 21st international conference on

software engineering,313-322.

- [4] Chen, Z.Menzies, T.Port, D. &Boehm, B. (2005).feature subset selection can improve software cost estimation accuracy. PROMISE '05 proceedings of the 2005 workshop on predictor models in software engineering, 1-6.
- [5] Das, A.Kempe, D. (2008).Algorithms for subset selection in linear regression. STOC 08 proceedings of the 40th annual ACM symposium on theory of computing, 1-10.
- [6] Heemstra, F.J. (1992). Software cost estimation. Information and software technology, volume 34, issue 10, 627-639.
- [7] Kirsopp, C.Shepperd, M.Hart, J. (2002).search heuristic, case –based reasoning and software project effort prediction.GECCO 02 proceedings of the genetic and evolutionary computation conference, 1-8.
- [8] Korte, M. & Port, D. (2008) confidence in software cost estimation results based on -MRE and PRED. PROMISE'08 proceedings of the 4th international workshop on predictor models in software engineering, 63-70.
- [9] Leung, H., & Fan, Z. (2002). Software Cost Estimation. Software engineering and knowledge engineering, 1-14.
- [10] Liyan-Fu. (2009).Improvement and Implementation of Analogy Based Method for Software Project Cost Estimation (doctoral thesis).Retrieved from <http://scholarbank.nus.edu.sg/bitstream/handle/10635/17724/LIYF.pdf>.
- [11] Sewell, M. (2007). Feature selection. Available on <http://machine-learning.Martinsewell.com>.