

# An Unambiguous Clustering and Ranking of Software Cost Estimation Models by Modified Scott-Knott Approach

N.Padma Priya, PG Scholar, Dept of Computer Science and Engg , Sona College of Technology, Salem, India **D.Vidyabharathi,** Assistant Professor (SG), Dept of Computer Science and Engg, Sona College of Technology, Salem, India

Abstract—Software Cost Estimation can be described as the process of predicting the most realistic effort required to complete a software project. Due to the strong relationship of accurate effort estimations with many crucial project management activities, the research community has been focused on the development and application of a vast variety of methods and models trying to improve the estimation procedure. From the diversity of methods emerged the need for comparisons to determine the best model. In the existing work Scott-Knott test was used to rank and cluster the software estimation models. The test proposed by Scott Knott, a procedure of means grouping, is an effective alternative to perform procedures of multiple comparisons without ambiguity. This study aimed to propose a modification related to the partitioning and means grouping in the said procedure, to obtain results without ambiguity among treatments, organized in more homogeneous groups. In the proposed methodology, treatments that did not participate in the initial group are joined for a new analysis, which allows for a better group distribution. The proposed methodology is considered effective, aiming at the identification of elite cultivar groups for recommendation.

Index Terms— software cost estimation; software metrics; software effort estimation; management; statistical methods.

## I. INTRODUCTION

Prediction of the effort is used to complete the software project by comparing the prediction models over past historical data set. This framework is based on a multiple comparisons algorithm, to rank several cost estimation models.

Software Engineering cost model and estimation techniques are used for budgeting, trade-off, risk analysis, and project planning with control to provide software improvement investment analysis.

The estimation increases the breadth of the search for relevant studies which conduct more studies on estimation methods commonly used by the software industry and also increases the awareness of how properties of the data set impacts the results when evaluating the estimation methods.

Accuracy is measured by the Magnitude of Relative Error (MRE) and MRE to the Estimate (MER). This can be achieved by accurate cost estimation. This needs the knowledge of size specifications, source code, manuals and the rate at which the requirements are likely to change during development and also the probable number of bugs that are likely to be encounter. The capability of development team and the salary over head incase if team increases along with the tools are necessary for estimation. Best of our knowledge, the problem of simultaneous comparisons among multiple prediction models has not been studied yet in the sense that there is no statistical procedure which can identify the significant differences between a number of cost estimation methods and at the same time be able to rank and cluster them, designating the best ones. All of the issues discussed above lead us to conclude that there is an imperative need to investigate what the state of the art in statistics is before trying to derive conclusions and unstable results concerning the superiority of a prediction

Model over others for a particular data set. The answer to this problem cannot constitute a unique solution since the notion of "best" is quite subjective. In fact, a practitioner can always rank the prediction models according to a predefined accuracy measure, but the critical issue is to identify how many of them are evidently the best, in the sense that their difference from all the others is statistically significant. Hence, the research question of finding the "best" prediction technique can be restated as a problem of identifying a subset or a group of best techniques.

The aim of the paper is therefore to propose a statistical framework for comparative SCE experiments concerning multiple prediction models. It is worth mentioning that the setup of the current study was also inspired by an analogous attempt dealing with the problem of comparing classification models in Software Defect Prediction, a research area that is also closely related to the improvement of software quality.

The proposed methodology is based on the analysis of a Design of Experiment (DOE) or Experimental Design, a basic statistical tool in many applied research areas such as engineering, financial, and medical sciences. In the field of SCE it has not yet been used in a systematic manner. Generally, DOE refers to the process of planning, designing, and analyzing an experiment in order to derive valid and objective conclusions effectively and efficiently by taking into account, in a balanced and systematic manner, the ources of variation.

## Priya et al., International Journal of Advanced Research in Computer Science and Software Engineering 4(2), February - 2014, pp. 785-790

In the present study, DOE analysis is used to compare different cost prediction models by taking into account the blocking effect, i.e., the fact that they are applied repeatedly on the same training-test datasets. The proposed statistical methodology is also based on an algorithmic procedure which is able to produce nonoverlapping clusters of prediction models, homogeneous with respect to their predictive performance. For this purpose, we utilize a specific test from the generic class of multiple comparisons procedures, namely, the Scott-Knott test which ranks the models and partitions them into clusters.

The proposed statistical framework is applied on a relatively large-scale set of 11 methods over six public domain datasets from the PROMISE repository and the International Software Benchmarking Standards Group (ISBSG) . Finally, in order to address the disagreement on the performance measures, we apply the whole analysis on three functions of error that measure different important aspects of prediction techniques: accuracy, bias, and spread of estimates.

## **II.** LITERATURE SURVEY

The software cost estimation researchers [2] supports with a list of journal papers with relevant historical papers. The First task is to estimate the purpose to include the particular journal paper for estimating the effort and cost. The second task is to identify the relevant papers. The third task is to classify the papers on their properties with respect to their estimation topics, estimation approach, research approach and result analysis to valid the threads. This review increases the breath of search. Complete search gives study on estimation methods by software industry and predicts the awareness about result impacting. This directs and supports the future estimation research. The other review principally aims to introduce software estimation researchers to the variety of formal estimation models.

The software cost estimation [3] includes the following approaches such as model based, expertise based, learning oriented, dynamic based, regression based and composite cocomo. These approaches capture the knowledge and experience by domain of interest. It estimates the effort hours, staff size and deployment, portfolio impact, risk, maintenance, schedule and hardware resource requirement.

The resembling techniques provide best accuracy measure without its superiority. Simulation tools [4] include regression and estimation by analogy. The accuracy measures

(I)Magnitude of Relative Error (MRE)  $y_A = actual \quad y_E = estimated$ MRE= $\frac{|y_A - y_E|}{y_A}$ 

(ii)Magnitude of Relative Error to Estimate (MER).  $MER = \frac{|y_A - y_E|}{y_E}$ 

Known data set validates the result by simulation method, traditional parametric and non parametric procedures. The datasets are tested by parametric and non parametric paired sample tests, bootstrap confidence intervals, and permutation tests and finish test.

Mean Magnitude of Relative Error (MMRE) accurately select the best model [5] and predict the effort from size in 3 ways. The ways are

OL: Ordinary Least squares on the raw data,

MR: Median Regression technique on raw data,

LNOLS: Ordinary Least squares Regression.

The problem of relying on within company [6] involves time to accumulate data, technologies will get change and the data will be collected in consistent manner. This review gives the complete analysis which project used to construct each model: accuracy measured, cross validation methods, fully defined methods, good comparison method. Mean Magnitude of Relative Error [7] predicts the software performance

 $MMRE = \frac{|y - y^{\wedge}|}{y}$ y=actual, y^ =prediction

It selects the model that is closest to the true model most of the time. MMRE is preferred which can be applied with ease to compare a linear regression model with a non linear arbitrary function estimator.

## III. PROPOSED SYSTEM

I. User Interface: A graphical user interface has been developed which provide user with some predefined options as well as some options are provided where user can input in plain English. Predefined options are provided in cases where a numeric value is needed, otherwise natural language has been used for both questions as well as answers the next question displays on the basis of previous response from the user. Thus an intelligent interaction occurs between user and computer.

ii. Natural Language Processor: NLP has been used to translate user response and query to specific rules and vice versa. It simply acts as an interface between User Interface and Inference Engine.



iii. Inference Engine: The basic objective of Inference Engine is to access knowledge Base on the basis of input parameters, supplied by the user. The developed Inference Engine is level 2-Type engine which not only provides basic reasoning but explanation facility has also been added that reproduces the logic to reach its conclusion. In order to reach a conclusion and offer an expert advice to the user, reasoning of the engine has been further strengthened by adding a database of static information. This database contains static information needed for calculation like effort adjustment factors in COCOMO.

iv. Knowledge Base: As the objective of the system is to effort estimation for different types of software development including variation of technology used as well as methodology followed, therefore four sets of rules have been incorporated in the knowledge base to support software effort determination for:

I. Line of Code base software, II. Component base software

a. Define the Public-domain data set

Public-domain data set with different characteristics are used in order to address the inherent problem of prediction systems, which means their high dependency on the types of data. Alternative error functions measuring different important aspects of error are studied. The repositories contain data from a wide range of projects are in the public domain.

b. Candidate prediction methods

The candidate methods can be grouped into three main categories that are [1] regression-based models, analogybased models, and machine learning methods. All these models are well-established methods, they are applied in SCE.

An alternative prediction technique was also based on the conclusions of a systematic review on SCE studies. Jorgensen and Sheppard [4] pointed out that the regression-based models dominate since half of all studies deal with the problem of fitting, improvement of a regression model. Furthermore, the researchers' interest for the analogy-based techniques [9] is steadily increased during the end of the decade. At last, the distribution of estimation methods also reveals that the proportion of machine learning techniques (Classification and Regression Trees and Neural Networks) presents an increasing trend.

It is obvious that the prediction techniques used in our experimentation is to tuning of certain parameters in order to build meaningful correct models. Consider example, the ratio-scaled variables of regression-based models[1] are checked in order to investigate whether the normality assumption is satisfied, and also the nominal and ordinal variables are replaced with new dummy variables and then a stepwise procedure is adopted to extract the most significant independent variables. In analogy based methods [8], the dissimilarity measure taking into account various types of variables, for the selection of the best number of the "neighbor" projects is determined through the leave-one-out cross-validation

## Priya et al., International Journal of Advanced Research in Computer Science and Software Engineering 4(2), February - 2014, pp. 785-790

procedure. Regarding neural network models specifies the number of nodes for hidden layers .In R Miner, the NN hyper parameter H [2] is optimized using a grid search with a backward selection algorithm, to avoid over fitting, there is an internal k-fold process is used. Thus the best parameter is selected with; the model is retrained with all training data.

CART model is concerned; utilize the Recursive Partitioning algorithm [1] as implemented in S-PLUS in which the model is fitted using binary recursive partitioning whereby the data are successively split along coordinate axes of the predictor variables so that the split which maximally distinguishes the response variable at any node in the left and the right branches which are to be selected. This splitting continues until nodes are pure or data are too sparse, to the recommendations of S-PLUS manual .Finally, for the case of the Naive Bayes classifier methodology [8] computes the conditional a-posterior probabilities of the dependent variable given the independent predictors using the Bayes rule,

#### c. Method comparison results

#### K-fold cross-validation with Design of Experiment

DOE [1] constitutes an entire branch area in statistics involving fundamental concepts that have to be specified and controlled in advance. The basic element of a DOE [1] is the experimental unit, which is the "object" on which the researcher wishes to measure a response variable. The purpose is to study the effect of one or more factors (categorical variables) on the response variable. The different categories of a factor are known as levels or treatments [1]. In the experimental setup [6], the predictive performance of each competitive model is evaluated through a k-fold cross-validation approach in which the original data set is randomly partitioned into k sub samples of equal size. During a repeated procedure, each one of the sub samples is considered as the validation sample (test set) and the remaining k -1 sub samples as the training sets used for fitting the models.

#### Repeated Measures Design similarly to the Randomized Complete Block Design (RCBD)

The RCBD [1] incorporates an additional factor takes into account the grouping of similar experimental units. The incorporation of this extra factor is considered advantageous in order to identify true differences between treatments or, equivalently [1], the true treatment effect. Indeed, when different treatments are applied to similar (or the same) experimental units which form, in any sense, a block, there is a source of variation between blocks which cannot be explained by the difference between treatments [1]. This source of variation is represented by the block factor that is considered in the analysis. In our context, the splitting of data [4] into different training-test pair's represents the blocking factor, i.e., each block is a specific pair of training-test subsets, where all models are applied and validated.

#### d. Principles of cluster analysis

Scott-Knott procedures are also presented in a graphical manner for two cases. The diagram [1] plots comparative models (x-axis) against the transformed mean errors (y-axis), whereby all methods are sorted according to their ranks. The vertical dashed lines indicate which models give statistically different results and thus are clustered into homogeneous group. The Scott- Knott algorithm resulted in four homogeneous groups of models with similar performances. Each small vertical solid line represents the prediction performance of the competitive models depicts the mean value of the transformed error function [1]. It is clearly inferred from the results of Scott-Knott tests, where the analogy-based techniques are clustered together in the same group of methods for all experiments [1]. Statistical methodology is also based on an algorithmic procedure which is able to produce non-overlapping clusters of prediction models and homogeneous with their predictive performance. It utilizes a specific test of Scott-Knott test which ranks the models and partitions them into clusters. The clustering refers to the treatments being compared and not to the individual cases, while the criterion for clustering together treatments is the statistical significance of differences between their mean values.

#### e. Performance evaluation

In order to address the disagreement on the performance measures [1], we apply the whole analysis on three functions of error that measure different important aspects of prediction techniques: accuracy, bias, and spread of estimates.

## IV. RESULTS AND DISCUSSIONS

#### SCOTT-KNOTT ALGORITHM:

The Scott-knott algorithm is utilized in cluster analysis to segregate the group of data into separate clusters. The procedure for the Scott-knott algorithm as follows:

**Step 1**: sort the means of the error measures  $e_i^-$ , j=1, d for each model in ascending order.

$$e_{(1)}^- \le e_{(2)}^- \le \dots \le e_{(d)}^-$$
 (1)

**Step 2**: For each  $e_{(j)}^-$ ,  $j = 1, \dots, d-1$ , separate the group of all ordered means E into two subgroups  $E_1 = \{e_{(1)}^-, \dots, e_{(j)}^-\}$  and  $E_2 = \{e_{(j+1)}^-, \dots, e_{(d)}^-\}$  and compute the between groups sum of squares:

$$G_{j} = k(|E_{1}|)(e_{E_{1}}^{-} - e_{E}^{-})^{2} + |E_{2}|(e_{E_{2}}^{-} - e_{E}^{-})^{2})$$
(2)

Where  $|E_1|, |E_2|$  are the cardinalities of the two subgroups and  $\overline{e_1}, \overline{e_E}, \overline{e_E_2}$  are the means of groups.

$$e_{\rm E}^{-} = \frac{1}{d} \sum_{j=1}^{d} e_{(j)}^{-}$$
(3)  
$$e_{\rm E_{1}}^{-} = \frac{1}{|{\rm E}_{1}|} \sum_{j \in {\rm E}_{1}} e_{(j)}^{-}$$
(4)

$$e_{E_2}^- = \frac{1}{|E_2|} \sum_{j \in E_2} e_{(j)}^-$$
(5)  
**Step 3:** Find the partition that maximizes the value of the sum of squares:  

$$G_i = \max\{G_i, j = 1, ..., d\}$$
(6)

Step 4: compute

$$\lambda = \frac{\pi}{2(\pi - 2)} \frac{G_j}{s^2}$$
(7)  
Distribution is computed by

$$V = \frac{k}{(\pi - 2)}$$

(8)

**Step 5:** If  $\lambda > X_v^2$  then the same test is applied to each group separately.

 $X_v^2$ 

If  $\lambda < X_v^2$  then all means belongs to the same homogeneous group.

The methodology proposes an alteration in the way of partitioning groups. The process begins with the formation of groups that maximize the sum of squares, based on the same concept as the Scott-Knott test. Two groups were formed first (one with models 9 and 10 and the second with 1, 2, 3, 4, 5, 6, 7 and 8).

Upon the formation of these groups, the second group was discarded and the possible partitions in the first group performed, resulting in two new subgroups (one with model 10 and the other with 9). This second subgroup was also discarded. Consequently, model 10 represented a group, the first formed group.

A new grouping analysis was performed with all previously discarded models (1, 2, 3, 4, 5, 6, 7, 8, and 9), which divided the models in two groups (group one 6, 7, 8 and 9 and group two 1, 2, 3, 4 and 5). Once again, the models of the second group were discarded and new possible partitions sought in the first group. No possibility of forming new subgroups was verified. Consequently, the models 6, 7, 8 and 9 represent the second group. As the procedure continues, new analyses are carried out with the previously discarded models (1, 2, 3, 4 and 5), until all models are grouped. Summing up, the new procedure consists in the removal of the models that form a new group and in the performance of new analyses with the remaining models, so that at each step a new group is formed while the number of remaining models.



More precisely, Absolute Error (AE) is used in order to evaluate the accuracy of models, whereas error ratio z has been adopted as a measure of bias accounting for underestimations (z < 1) or overestimations (z > 1) with an optimum value of 1. The most widely known MRE indicator was also used since, it provides a measure of the spread of the error ratio z.

In that graph Fig 1. Error Accuracy is compared for existing and proposed system. Methods are represented in x-axis and error accuracy is represented in y-axis. Compare with all the methods our proposed systems method modified Scott-Knott has better results than Scott-Knott methods.

### V. CONCLUSION

SCE depends on several issues, even on personal criteria like experience, preference of statistical software. The intelligent expert effort estimation uses User Interface, Natural Language Processor, Inference Engine and Knowledge Base. This expert system improves the software cost effort estimation results and also improves the accuracy in cost estimation. Based on this, the best prediction model for SCE is estimated. The Scott-knott algorithm is thus used to measure the error in the prediction model and also gives the probability of the usage of data set in concurrent project. The project estimation can also be done using this method and is the future work to be continued. Both the schedule and cost estimation is mandatory requirement for success of the project. The step-by-step procedure followed in this method is considered as one of the major drawback. To overcome this, an efficient algorithm can be designed for estimation process. The proposed methodology makes a differentiated partitioning of the available models possible while ensuring the principle of absence of ambiguity or superposition of model groups. The proposed methodology is a more effective option when the objective is to identify one or few elite groups and discard inferior and intermediate groups. There is a loss of the global partitioning structure, while the identification of a specific subgroup with better performance is facilitated.

## REFERENCES

- [1] Nikolaos Mittas and Lefteris Angelis, "Ranking and Clustering Software Cost Estimation Models through a Multiple Comparisons Algorithm," IEEE Trans. Software Eng., vol. 39, no.4, April. 2013.
- [2] M. Jorgensen and M. Shepperd, "A Systematic Review of Software Development Cost Estimation Studies," IEEE Trans. Software Eng., vol. 33, no. 1, pp. 33-53, Jan. 2007.
- [3] M. Shepperd and G. Kadoda, "Comparing Software Prediction Techniques Using Simulation," IEEE Trans. Software Eng., vol. 27, no. 11, pp. 1014-1022, Nov. 2001.
- [4] B. Kitchenham, S. MacDonell, L. Pickard, and M. Shepperd, "What Accuracy Statistics Really Measure," IEE Proc. Software Eng., vol. 148, pp. 81-85, June 2001.
- [5] T. Foss, E. Stensrud, B. Kitchenham, and I. Myrtveit, "A Simulation Study of the Model Evaluation Criterion MMRE," IEEE Trans. Software Eng., vol. 29, no. 11, pp. 985-995, Nov. 2003.
- [6] N. Mittas and L. Angelis, "Comparing Cost Prediction Models by Resampling Techniques," J. Systems and Software, vol. 81, no. 5, pp. 616-632, May 2008.
- [7] E. Stensrud and I. Myrtveit, "Human Performance Estimating with Analogy and Regression Models: An Empirical Validation," Proc. IEEE Fifth Int'l Software Metrics Symp., pp. 205-213, Nov. 1998.
- [8] B. Kitchenham and E. Mendes, "Why Comparative Effort Prediction Studies May Be Invalid," Proc. ACM Fifth Int'l Conf. Predictor Models in Software Eng., pp. 1-5, May 2009.
- [9] I. Myrtveit, E. Stensrud, and M. Shepperd, "Reliability and Validity in Comparative Studies of Software Prediction Models," IEEE Trans. Software Eng., vol. 31, no. 5, pp. 380-391, May 2005.
- [10] S.Lessmann, B. Baesens, C. Mues, and S. Pietsch, "Benchmarking Classification Models for Software Defect Prediction: A Proposed Framework and Novel Findings," IEEE Trans. Software Eng., vol. 34, no. 4, pp. 485-496, July/Aug. 2008.