

Computational Intelligence Hybrids Applied to Software Cost Estimation

J.S.Pahariya^{a,b}, V. Ravi^{a,*}, M. Carr^a and M. Vasu^{a,b}

^a *Institute for Development and Research in Banking Technology, Castle Hills Road #1, Masab Tank, Hyderabad 500 057, AP, India*

^b *Department of Computer and Information Sciences, Univeristy Of Hyderabad, Hyderabad – 500046, AP, India*

jankisharanpahariya@yahoo.com, rav_padma@yahoo.com, mahil.carr@gmail.com, madireddi.vasu@gmail.com

Abstract

In this paper, we propose new computational intelligence sequential hybrid architectures involving Genetic Programming (GP) and Group Method of Data Handling (GMDH) viz. GP-GMDH, GMDH-GP and recurrent architecture for Genetic Programming (GP) for software cost estimation. Three linear ensembles based on (i) arithmetic mean (ii) geometric mean and (iii) harmonic mean are also developed. We also performed GP based feature selection. The efficacy of Multiple Linear Regression (MLR), Polynomial Regression, Support Vector Regression (SVR), Classification and Regression Tree (CART), Multivariate Adaptive Regression Splines (MARS), Multilayer FeedForward Neural Network (MLFF), Radial Basis Function Neural Network (RBF), Counter Propagation Neural Network (CPNN), Dynamic Evolving Neuro-Fuzzy Inference System (DENFIS), TreeNet, Group Method of Data Handling and Genetic Programming is tested on the International Software Benchmarking Standards Group (ISBSG) release 10 dataset. Ten-fold cross validation is performed throughout the study. The results obtained from our experiments indicate that the GP-GMDH and GMDH-GP outperformed all the other techniques. We also performed t-test to see if the performances of the hybrids developed are statistically significant.

Keywords:- Multiple Linear Regression (MLR), Polynomial Regression, Support Vector Regression (SVR), Classification and Regression Tree (CART), Multivariate Adaptive Regression Splines (MARS), Multilayer FeedForward Neural Network (MPFF), Radial Basis Function Neural Network (RBF), Counter Propagation Neural Network (CPNN), Dynamic Evolving Neuro-Fuzzy Inference System (DENFIS), Tree Net, Group Method of Data Handling (GMDH) and Genetic Programming (GP).

* Corresponding author. Phone: +91 40 23534981
Ext. 2042; FAX: +91 40 23535157.

1. Introduction

The estimation of software development cost is one of the most critical problems in software engineering. Software cost development is related to how long and how many people are required to complete a software project. Software development has become an essential question [1] because many projects are still not completed on schedule, with under or over estimation of efforts leading to their own particular problems [2]. Therefore, in order to manage budget and schedule of software projects [3], various software cost estimation models have been developed. A major problem of the software cost estimation is first obtaining an accurate size estimate of the software to be developed [4] because size is the most important single cost driver [5]. Thus, an important objective of the software engineering community has been to develop useful models that can explain the software development life cycle and accurately estimate the cost of software development [6, 7]. Development cost tends to increase with project complexity and hence accurate cost estimates are highly desired during the early stages of development [8].

The main objective of the present work is to propose new computational intelligence based hybrid models that estimates the software cost accurately. The rest of the paper is organized as follows. Section 2 reviews the research done in the field of software cost estimation. Section 3 overviews the techniques applied in this study. Section 4 describes briefly the ISBSG dataset that is analyzed by our proposed Hybrid Intelligent Systems. Section 5 presents proposed Hybrid Intelligent Systems developed in this study. Section 6 presents the results and discussions. Finally, Section 7 concludes the paper.

2. Literature Review

Various software development effort estimation models have been developed over the last four decades. The most commonly used methods for predicting software development efforts are Function Point Analysis and COConstructive COSt MODEL (COCOMO) [4]. Function Point Analysis was developed first by Albrecht (1979) (www.IFPUG.Org). Function point analysis is a method

of quantifying the size and complexity of a software system in terms of the functions that the system delivers to the user [9]. The function does not depend on the programming languages or tools used to develop a software project [1]. COCOMO is developed by Boehm [2]. It is based on linear-least-squares regression. Using line of code (LOC) as the unit of measure for software size itself contains so many problems [10]. These methods failed to deal with the implicit non-linearity and interactions between the characteristics of the project and effort [11, 12]. Jørgensen and Shepperd [13] conducted a systematic review of software development cost estimation studies.

In recent years, a number of alternative modeling techniques have been proposed. They include artificial neural networks, analogy-based reasoning, and fuzzy system [14, 15 and 16] and ensemble techniques. Aggarwal et al. [17] reported an expert committee model which is a combination of robust regression technique and neural network. Later Vinay kumar et al. [18] reported linear and non linear ensembles consists of various statistical and intelligent techniques *viz.* Multi Layer Regression, Back Propagation Neural Network (BPNN), RBF, DENFIS, Threshold Accepting based Neural Network (TANN) and SVM. In analogy-based cost estimation, similarity measures between a pair of projects play a critical role [19]. This type of model calculates distance between the software project being estimated and each of the historical software projects and then retrieves the most similar project for generating an effort estimate [20]. Further, Lefley and Shepperd [21] applied genetic programming to improve software cost estimation on public datasets with great success. Later, Vinaykumar et al. [6] used wavelet neural networks for the prediction of software cost estimation. Unfortunately the accuracy of these models is not satisfactory so there is always a scope for new software cost estimation techniques.

Li et al. [22] proposed mutual information based feature selection which hybridizes both ‘wrapper’ and ‘filter’ mechanism. Tosun et al. [23] proposed feature weighting heuristics for analogy-based effort estimation models using principal components analysis (PCA). Mittas and Angelis [24] proposed statistical simulation procedures involving permutation tests and bootstrap techniques in order to test the significance of the difference between the accuracy of two prediction methods: the estimation by analogy and the regression analysis. Later, Mittas and Angelis [25] also used Regression Error Characteristic (REC) analysis in order to validate and compare different prediction models. Most recently Mohanti et al. [26] conducted by far the most comprehensive review of the applications of intelligent and soft computing to software engineering problems reported during 1990-2008. This review article surveyed

many intelligent techniques as applied to software cost estimation.

3. Overview of the techniques employed

In the following, we now present an overview of the techniques applied in this paper.

3.1 Group Method of Data Handling (GMDH)

The group method of data handling (GMDH) was introduced by Ivakhnenko [27] in 1966 as an inductive learning algorithm for modeling of complex systems. It is a self-organizing approach based on sorting-out of gradually complicated models and evaluation of them using some external criterion on separate parts of the data sample [28]. The GMDH was partly inspired by research in Perceptrons and Learning Filters. GMDH has influenced the development of several techniques for synthesizing (or “self-organizing”) networks of polynomial nodes. The GMDH attempts a hierarchic solution, by trying many simple models, retaining the best, and building on them iteratively, to obtain a composition (or feed-forward network) of functions as the model. The building blocks of GMDH, or polynomial nodes, usually have the quadratic form:

$$z = w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2$$

for inputs x_1 and x_2 , coefficient (or weight) vector w , and node output, z . The coefficients are found by solving the Linear Regression equations with $z = y$, the response vector.

The topology of the GMDH neural network is given below:

The GMDH neural network develops on a data set. The data set including independent variables (x_1, x_2, \dots, x_n) and one dependent variable y is split into a training and testing set. During the process of learning a forward multilayer neural network is developed by observing the following steps:

1. In the input layer of the network n units with an elementary transfer function $y = x_i$ are constructed. These are used to provide values of independent variables from the learning set to the successive layers of the network.
2. When constructing a hidden layer an initial population of units is generated. Each unit corresponds to the Ivakhnenko polynomial form:
 $y = a + bx_1 + cx_2 + dx_{12} + ex_1x_2 + fx_{22}$ or
 $y = a + bx_1 + cx_2 + dx_1x_2$
 Where y is an output variable; x_1, x_2 are two input variables and a, b, \dots, f are parameters.
3. Parameters of all units in the layer are estimated using the learning set.

4. The mean square error between the dependent variable y and the response of each unit is computed for the testing set.
5. Units are sorted out by the mean square error and just a few units with minimal error survive. The rest of the units are deleted. This step guaranties that only units with a good ability for approximation are chosen.
6. Next the hidden layers are constructed while the mean square error of the best unit decreases.
7. Output of the network is considered as the response of the best unit in the layer with the minimal error.

The GMDH network learns in an inductive way and tries to build a function (called a polynomial model), which would result in the minimum error between the predicted value and expected output. The majority of GMDH networks use regression analysis for solving the problem. The first step is to decide the type of polynomial that the regression should find. The initial layer is simply the input layer. The first layer created is made by computing regressions of the input variables and then choosing the best ones. The second layer is created by computing regressions of the values in the first layer along with the input variables. This means that the algorithm essentially builds polynomials of polynomials. Again, only the best are chosen by the algorithm. These are called survivors. This process continues until a pre-specified selection criterion is met.

We used the GMDH implementation available at <http://www.neuroshell.com/>.

3.2 Genetic Programming (GP)

Genetic programming (GP) [29, 30] is an extension of genetic algorithms (GA). It is a search methodology belonging to the family of evolutionary computation (EC). GP mainly involve functions and terminals. GP randomly generates an initial population of solutions. Then, the initial population is manipulated using various genetic operators to produce new populations. These operators include reproduction, crossover, mutation, dropping condition, etc. The whole process of evolving from one population to the next population is called a generation. A high-level description of GP algorithm can be divided into a number of sequential steps [31]:

- Create a random population of programs, or rules, using the symbolic expressions provided as the initial population.
- Evaluate each program or rule by assigning a fitness value according to a predefined fitness function that can measure the capability of the rule or program to solve the problem.
- Use reproduction operator to copy existing programs into the new generation.

- Generate the new population with crossover, mutation, or other operators from a randomly chosen set of parents.
- Repeat steps 2 onwards for the new population until a predefined termination criterion has been satisfied, or a fixed number of generations has been completed.
- The solution to the problem is the genetic program with the best fitness within all the generations.

In GP, crossover operation is achieved first by reproduction of two parent trees. Two crossover points are then randomly selected in the two offspring trees. Exchanging sub-trees, which are selected according to the crossover point in the parent trees, generates the final offspring trees. The obtained offspring trees are usually different from their parents in size and shape. Then, mutation operation is also considered in GP. A single parental tree is first reproduced. Then a mutation point is randomly selected from the reproduction, which can be either a leaf node or a sub-tree. Finally, the leaf node or the sub-tree is replaced by a new leaf node or sub-tree generated randomly. Fitness functions ensure that the evolution goes toward optimization by calculating the fitness value for each individual in the population. The fitness value evaluates the performance of each individual in the population.

GP is guided by the fitness function to search for the most efficient computer program to solve a given problem. A simple measure of fitness [31] is adopted for the binary classification problem which is given as follows.

$$\text{Fitness (T)} = \frac{\text{no of samples classified correctly}}{\text{no of samples for training during evaluation}}$$

We used the GP implementation available at <http://www.rmltech.com>

3.3 Counter Propagation Neural Network (CPNN)

The counter propagation network is a competitive network and given by Hecht-Nielsen [32]. The uni-directional PNN has three layers. The main layers include an input buffer layer, a self-organizing Kohonen layer and an output layer which uses the Delta Rule to modify its incoming connection weights. Sometimes this layer is called a Grossberg Outstar layer. The forward-only counter propagation network architecture, consists of three slabs: an input layer (layer 1) containing n fan out units that multiplex the input signals x_1, x_2, \dots, x_n , (and m units that supply the correct output signal values y_1, y_2, \dots, y_m to the output layer), a middle layer (layer 2 or Kohonen layer) with N processing elements that have output signals z_1, z_2, \dots, z_N , and a final layer (layer 3) within processing elements having output signals y_1', y_2', \dots, y_m' . The outputs of layer 3 represent approximations to the components y_1, y_2, \dots, y_m of $y = f(x)$ [33]. The input layer in CPNN performs the mapping of

the multidimensional input data into lower dimensional array. The mapping is performed by use of competitive learning, which employs winner-takes-it-all strategy [34]. The training process of the CPNN is partly similar to that of Kohonen self-organizing maps. The Grossberg layer performs supervised learning. The network got its name from this counter-posing flow of information through its structure.

3.4 Support Vector Regression (SVR)

The SVR is a powerful learning algorithm based on recent advances in statistical learning theory proposed by Vapnik [35]. SVR is a learning system that uses a hypothesis space of linear functions in a high-dimensional space, trained with a learning algorithm from optimization theory that implements a learning bias derived from statistical learning theory. SVR uses a linear model to implement non-linear class boundaries by mapping input vectors non-linearly into a high dimensional feature space using kernels. The training examples that are closest to the maximum margin hyperplane are called support vectors. All other training examples are irrelevant for defining the binary class boundaries. The support vectors are then used to construct an optimal linear separating hyperplane (in case of pattern recognition) or a linear regression function (in case of regression) in this feature space. The support vectors are conventionally determined by solving a quadratic programming (QP) problem. We used the SVR implementation available at <http://rapid-i.com/content/view/26/84/>.

3.5 Classification and Regression Tree (CART)

CART was introduced by Breiman et al. [36] can solve both classification and regression problems (<http://salford-systems.com>). Decision tree algorithms induce a binary tree on a given training data, resulting in a set of 'if-then' rules. These rules can be used to solve the classification or regression problem. The key elements of a CART analysis [36] are a set of rules for: (i) splitting each node in a tree, (ii) deciding when a tree is complete; and (iii) assigning each terminal node to a class outcome (or predicted value for regression). We used the CART implementation available at <http://salford-systems.com>.

3.6 Multivariate Adaptive Regression Splines (MARS)

Multivariate adaptive regression splines (MARS) introduced by Friedman [37]. MARS is an innovative and flexible modeling tool that automates the building of accurate predictive models for continuous and binary

dependent variables. It excels at finding optimal variable transformations and interactions, the complex data structure that often hides high-dimensional data. This approach to regression modeling effectively uncovers important data patterns and relationships that are difficult, if not impossible, for other methods to reveal. We used the MARS implementation available at <http://salford-systems.com>.

3.7 Dynamic Evolving Neuro-Fuzzy Inference System (DENFIS)

DENFIS was introduced by Kasabov and Song [38]. DENFIS evolve through incremental, hybrid (supervised/unsupervised) learning, and accommodate new input data, including new features, new classes, etc., through local element tuning. New fuzzy rules are created and updated during the operation of the system. At each level, the output of DENFIS is calculated through a fuzzy inference system based on most activated fuzzy rules, which are dynamically chosen from a fuzzy rule set. A set of fuzzy rules can be inserted into DENFIS before or during its learning process. Fuzzy rules can also be extracted during or after the learning process. Student version of the NewCom tool obtains at http://www.aut.ac.nz/research/research_institutes/kedri/research_centres/centre_for_data_mining_and_decision_support_systems/neucom.htm#download was used in this paper for DENFIS and MLR.

3.8 Tree Net

Tree Net was introduced by Friedman [39]. It makes use of a new concept of 'ultra slow learning' in which layers of information are gradually peeled off to reveal structure in data. TreeNet models are typically composed of hundreds of small trees, each of which contributes just a tiny adjustment to the overall model. TreeNet is insensitive to data errors and needs no time-consuming data preprocessing or imputation of missing values. TreeNet is resistant to overtraining and is faster than a neural net. TreeNet available at <http://salford-systems.com> was used.

3.9 Radial Basis Function Neural Network (RBF)

RBF was introduced by Moody and Darken [40], has both unsupervised and supervised phases intandem. In the unsupervised phase, input data will be clustered and cluster details are subsequently sent to hidden neurons, where radial basis functions of the inputs are computed by making use of the center and the standard deviation of the clusters. Gaussian radial basis functions are the most commonly used functions. The learning between hidden

layer and output layer is of supervised learning type where ordinary least squares technique is used. As a consequence, the weights of the connections between the kernel layer and the output layer are determined. KNIME tool available at <http://www.knime.org> was used for MLP, RBF, and polynomial regression.

4. Data Description and Data Preprocessing

The ISBSG data is obtained from Australia (<http://www.isbsg.org>). Entire ISBSG-10 dataset contains information about 4109 projects. The dataset consist of 18 attributes. These attributes are also divided into sub-attributes, thereby making the total number of attributes 105. In this paper, we predict the summary of work effort, i.e. manpower required to complete the work. Once we know the effort, we can easily calculate the time and the cost of the software. However, before applying these statistical and intelligence techniques to the dataset, there are a number of issues to be taken into consideration during data cleaning and data preparation.

The first cleaning step was to remove the projects having null values for the attribute named Summary of Work Effort. Secondly regarding summary of work effort only 1538 project values are given for the five attributes *viz.* Input count, Output, Enquiry, File and Interface. If we consider more attributes then we get only a few projects which are not sufficient for machine learning techniques. Hence, finally, we considered 1538 projects values with five attributes to do train and test several intelligent models. Finally, we normalized the data set. The effectiveness of our proposed hybrid intelligent systems is tested on this normalized dataset.

5. Proposed Hybrid Intelligent Systems

The fundamental assumption in computational intelligence paradigm is that hybrid intelligent techniques tend to outperform the stand-alone techniques. We proposed 6 new hybrid architectures for software cost estimation. We compared the performance of cost estimation models on the basis of root mean square error (RMSE), which defined as follows:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n \left(Effort_{actual_i} - Effort_{estimated_i} \right)^2}$$

where n is the number of projects.

5.1 Ensemble System

We first implemented linear ensemble systems. Ensemble systems exploit the unique strengths of each constituent model and combine them in same way. For

constructing ensemble system we have chosen the three best techniques *viz.*, GMDH, GP and CPNN from stand-alone mode. These three techniques have yielded the best RMSE values in the 10-fold cross validation method of testing. We constructed ensembles using three methods. They are Arithmetic Mean (AM), Harmonic Mean (HM) and Geometric Mean (GM). The proposed Ensemble system is depicted in Figure 1.

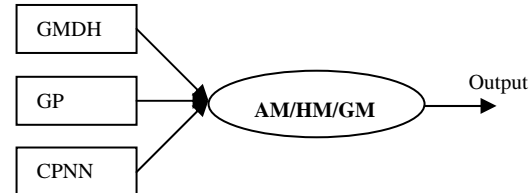


Figure.1 Ensemble System

5.2 Recurrent Genetic Programming (RGP) architecture

The second proposed architecture is a recurrent architecture for Genetic Programming (RGP) in which output of the GP is fed as an input to the GP. This is analogous to recurrent neural networks having feedback loop where output can be fed back as input [41]. However, the difference is that we wait until GP converges and yields the predictions. These predictions along with the original input variables are fed as inputs to another GP afresh. The flow diagram of the recurrent architecture for Genetic Programming (RGP) is depicted in Figure 2. The idea here is to investigate if the recurrent nature of the hybrid would improve the RMSE of the first GP.

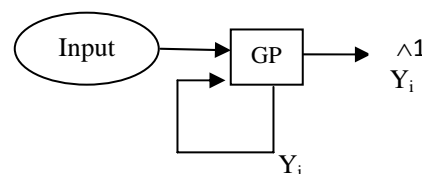


Figure.2 Recurrent architecture for GP (RGP)

5.3 GP-GP hybrid

It is observed that there are some features in the dataset that are contributing negatively to the prediction accuracy of all the models. Hence, we resorted to feature selection (F.S). We used GP for feature selection. Using GP based feature selection we selected four most important variables for training. Accordingly, in the proposed hybrid first important features are selected using GP and then those are fed to GP for predictions resulting

in GP-GP hybrid. The architecture of proposed GP-GP hybrid is depicted in Figure 3.

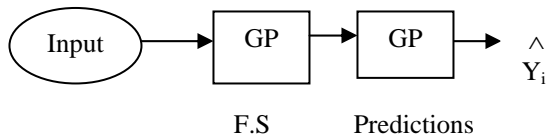


Figure.3 GP-GP Hybrid Architecture

5.4 GP-RGP hybrid

Based on the GP-based feature selection (F.S) we proposed another hybrid viz. GP-RGP in which the first GP is used for feature selection and the second technique RGP used for making predictions. The architecture of proposed GP-RGP hybrid is depicted in Figure 4.

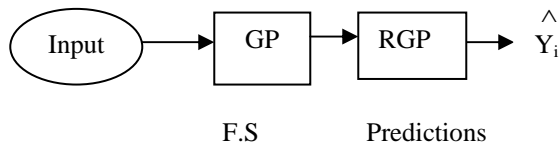


Figure.4 GP-RGP Hybrid Architecture

5.5 GMDH-GP hybrid

As an extension to this work, it is worth investigating the boosting of well performing techniques with each other. Accordingly, we proposed a new sequential hybrid in which the predictions of GMDH along with input variables are fed as input to GP for predictions, resulting in GMDH-GP hybrid. The architecture of GP-GMDH hybrid is depicted in Figure 5.

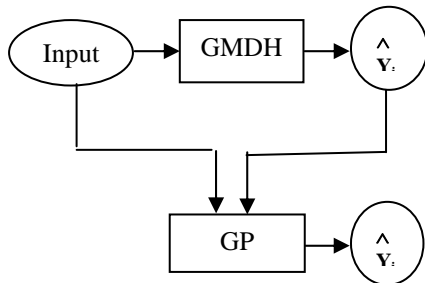


Figure.5 GMDH-GP Sequential hybrid

5.6 GP-GMDH hybrid

We also proposed another sequential hybrid to explore the boosting power of GP with GMDH. In this new hybrid, the predictions of GP along with input variables are fed as input to GMDH for predictions. The architecture of GP-GMDH hybrid is depicted in Figure 6.

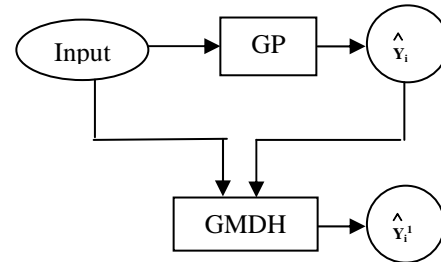


Figure.6 GP-GMDH Sequential hybrid

6. Results and discussion

We used ISBSG data set, which contains 1538 projects and five independent variables and one dependent variable.

We employed GMDH, GP, CPNN, MLR, Polynomial Regression, SVR, CART, MARS, MLFF, RBF, DENFIS, and TreeNet. We performed 10-fold cross validation throughout the study and the average results are presented in Table 1.

Table 1: Average RMSE of 10 fold cross validation

SN	METHOD	RMSE (TEST)
1	GMDH	0.03784
2	GP	0.03794
3	CPNN	0.04499
4	CART	0.04561
5	TREENET	0.04565
6	MLP	0.04817
7	MLR	0.04833
8	DENFIS	0.04837
9	MARS	0.0487
10	SVR	0.0492
11	RBF	0.05167
12	Polynomial Regression	0.05327

It is observed that GMDH performed the best with least RMSE value of 0.03784 and GP stood a close second with an average RMSE value of 0.03794 among all stand alone techniques tested followed by CPNN, CART, TREENET, MLP, MLR, DENFIS, MARS, SVR, RBF and Polynomial Regression in that order.

We also implemented three linear ensemble systems using AM, GM and HM to exploit the unique strengths of the best performed stand alone techniques GMDH, GP and CPNN. We notice that AM based ensemble system has outperformed GM based and HM based ensemble techniques. The results are presented in Table 2. However, they are not so spectacular when compared to the best

performing stand-alone methods. This is evident from the nature of the AM, Gm and HM.

Table 2: Average RMSE of ensemble models

SN	METHOD	RMSE (Test data)
1	AM	0.0421
2	GM	0.04403
3	HM	0.0455

The results of the hybrids are presented in Table 3. Here we observe that all the proposed hybrid models RGP, GP-RGP, GP-GP, GMDH-GP and GP-GMDH outperformed all other stand-alone techniques due to the synergy that took place in hybridizing them. From the results of GP-GP and GP-RGP it is inferred that the features selected by GP helped to boost the performance of GP and RGP.

We further explored the boosting power of GP with another well performing technique GMDH and boosting power of GMDH with GP. We observed that GP-GMDH yielded least average RMSE value of 0.02833 and GMDH-GP stood second with an average RMSE value of 0.03098 among all the hybrids tested. They are followed by RGP, GP-RGP and GP-GP in that order.

We also performed t-test to test whether the difference in RMSEs obtained by the top five methods viz., GP-GMDH, GMDH-GP, RGP, GP-RGP and GP-GP is statistically significant or not. Thus, the t-statistic values computed for those hybrids are presented in Table 3. The calculated t-statistic values are compared with 2.1, which is the tabulated t-statistic value at $n_1+n_2-2=10+10-2=18$ degrees of freedom at 5 % level of significance. That means, if the computed t-statistic value between two methods is more than 2.1, then we can say that the difference between the techniques is statistically significant. The t-statistic value between GP-GMDH and RGP is 2.47694 whereas that between GP-GMDH and GP-RGP is 2.83543 and in the case of GP-GMDH and GP-GP it is 4.27465. Considering GP-GMDH as best performer, it is observed that the difference between RGP, GP-RGP and GP-GP is statistically significant. The t-statistic value between GP-GMDH and GMDH-GP is 1.6972 which is less than 2.1 and hence the difference between GP-GMDH and GMDH-GP is statistically insignificant.

7. Conclusions

This paper presents new computational intelligence sequential hybrids involving GP and GMDH for software cost estimation. Throughout the study 10-fold cross validation is performed. Besides GP and GMDH, we

tested a host of techniques on the ISBSG dataset. The proposed GP-GMDH and GMDH-GP hybrids outperformed all other stand-alone and hybrid techniques. Hence, we conclude that the GP-GMDH or GMDH-GP model is the best model among all other techniques for software cost estimation.

Table 3: Average RMSE of hybrids models

SN	METHOD	RMSE (Test data)	t-test value
1	RGP	0.03275	2.47694
2	GP-RGP	0.03345	2.83543
3	GP-GP	0.03676	4.27465
4	GMDH-GP	0.03098	1.6972
5	GP-GMDH	0.02833	-

References

- [1]. A.J. Albrecht and J.E. Gaffney, "Software function, source lines of code, and development effort prediction: a software science validation", *IEEE Transactions on Software Engineering*, 1983, 9(6), pp. 639–647.
- [2]. B.W. Boehm, "Software Engineering Economics", Prentice-Hall, Englewood Cliffs, NJ, USA, 1981.
- [3]. L.H. Putnam, "A general empirical solution to the macro software sizing and estimation problem", *IEEE Transactions on Software Engineering*, 1978, 4(4), pp. 345–361.
- [4]. B. Kitchenham, L.M. Pickard, S. Linkman and P.W. Jones, "Modeling software bidding risks", *IEEE Transactions on Software Engineering*, 2003, 29(6), pp. 542–554.
- [5]. J. Verner and G. Tate, "A Software Size Model", *IEEE Transactions on Software Engineering*, 1992, 18(4), pp. 265–278.
- [6]. K. Vinaykumar, V. Ravi, M. Carr and N. Rajkiran, "Software development cost estimation using wavelet neural networks", *Journal of Systems and Software*, 2008, 81(11), pp. 1853–1867.
- [7]. T. Foss, E. Stensrud, B. Kitchenham and I. Myrvtveit, "A simulation study of the model evaluation criterion MMRE", *IEEE Transactions on Software Engineering*, 2003, 29(11), pp. 985–995.
- [8]. Z. Xu and T.M. Khoshgoftaar "Identification of fuzzy models of cost estimation", *Fuzzy Sets and Systems*, 2004, 145(11), pp. 141–163.
- [9]. J.E. Matson, B.E Barrett and J.M. Mellichamp, "Software development cost estimation using function points", *IEEE Transactions on Software Engineering*, 1994, 20(4), pp. 275–287.
- [10]. A.Idri, T.M. Khosgoftaar and A. Abran, "Can neural networks be easily interpreted in software cost

- estimation”, *World Congress on Computational Intelligence*, Honolulu, Hawaii, USA, 2002, pp. 12–17.
- [11]. A.R. Gray, “A simulation-based comparison of empirical modeling techniques for software metric models of development effort”, *In Proceedings of ICONIP Sixth International Conference on Neural Information Processing*, Perth, WA, Australia, 1999, pp. 526–531.
- [12]. X. Huang, L.F. Capetz, J. Ren and D. Ho, “A neuro-fuzzy model for software cost estimation”, *In Proceedings of the 3rd International Conference on Quality Software*, 2003, pp. 126-133 .
- [13]. M.Jørgensen and M. Shepperd, “A Systematic Review of Software Development Cost Estimation Studies”, *IEEE Transactions on Software Engineering*, 2007, 33(1), pp. 33-53.
- [14]. S. Andreou and E. Papatheocharous, “Software Cost Estimation using Fuzzy Decision Trees”, *Proceedings of 23rd IEEE/ACM International Conference on Automated Software Engineering*, 2008, pp. 371-374.
- [15]. Mittal, K. Prakash and H. Mittal, “Software Cost Estimation Using Fuzzy Logic”, *ACM SIGSOFT Software Engineering Notes*, 2010, 35(1), pp. 1-7.
- [16]. Attarzadch, “Improving the accuracy of software cost estimation model based on a new fuzzy logic model”, *World applied sciences journal*, 2010, 8(2), pp. 177-184.
- [17]. K.K. Aggarwal, Y. Singh, P. Chandra and M. Puri, “An expert committee model to estimate line of code”, *ACM SIGSOFT Software Engineering Notes*, 2005, pp. 1-4.
- [18]. K. Vinay Kumar, V. Ravi and M. Carr, “Software Cost Estimation using Soft Computing Approaches”, *Handbook on Machine Learning Applications and Trends: Algorithms, Methods and Techniques*, Eds. E. Soria, J.D. Martin, R. Magdalena, M.Martinez, A.J. Serrano, IGI Global, USA, 2009.
- [19]. Y.F. Li, M. Xie and T.N. Goh, “A study of project selection and feature weighting for analogy based software cost estimation”, *Journal of Systems and Software*, 2009, 82(2), pp. 241–252.
- [20]. N.H. Chiu and S.J.Huang, “The adjusted analogy-based software effort estimation based on similarity distances”, *Journal of Systems and Software*, 2007, 80(4), pp.628-640
- [21]. M. Lefley and M. J. Shepperd, “Using Genetic Programming to Improve Software Effort Estimation Based on General Data Sets”, LNCS, Genetic and Evolutionary GECCO 2003, page-208.
- [22]. Y.F. Li, M. Xie, T.N. Goh, “A study of mutual information based feature selection for case based reasoning in software cost estimation”, *Journal of Systems and Software*, 2010, 83, pp. 621–637.
- [23]. Tosun , B. Turhan and A. B. Bener, “Feature weighting heuristics for analogy-based effort estimation models”, *Expert Systems with Applications*, 2009, 36, pp. 10325–10333.
- [24]. N. Mittas and L. Angelis, “Comparing cost prediction models by resampling techniques”, *Journal of Systems and Software*, 2008, 81, pp. 616–632.
- [25]. N. Mittas, L. Angelis, “Visual comparison of software cost estimation models by regression error characteristic analysis”, *Journal of Systems and Software*, 2010, 83, pp. 621–637.
- [26]. R. Mohanty, V. Ravi and M. Patra, “The Application of Intelligent and Soft computing techniques to software engineering problem: A Review”, *International Journal of Information and Decision Sciences*, 2010, 2(3), pp.233-272.
- [27]. A.G. Ivakhnenko, “The group method of data handling - a rival of the method of stochastic approximation”, *Soviet Automatic Control*, 13(3), 1966, pp. 43-55.
- [28]. D. Srinivasan, “Energy demand prediction using GMDH networks”, *Neurocomputing*, 2008, 72(1-3), pp. 625-629.
- [29]. R. Poli, W.B. Langdon and J.R. Koza, “A field guide to Genetic Programming”, publisher- Lulu.com, United Kingdom, 2008.
- [30]. J. R. Koza, “Genetic Programming: On the programming of computers by means of natural selection”, Cambridge, MA: MIT press, 1992.
- [31]. K. M. Faraoun and A. Boukelif, “Genetic programming approach for multi-category pattern classification applied to network intrusion detection”, *International Journal of Computational Intelligence and Applications*, 2006, 6 (1), pp. 77-99.
- [32]. R. Hecht-Nielsen, “Counterpropagation Networks”, *Applied Optics*, 1987, 26(23), pp. 4979-4984.
- [33]. M.F. Zafar, M. Dzulkipli, R.M. Othman, “On-line handwritten character recognition: an implementation of counter propagation neural network”, *Proceedings of World Academy of Science, Engineering and Technology* 10, 2005, pp. 232-237.
- [34]. Kuzmanovski and M. Novic, “Counter-Propagation Neural Networks in Matlab”, *Chemometrics and Intelligent Laboratory Systems*, 2008, 90(1), pp. 84-91.
- [35]. V.N. Vapnik, “Statistical Learning Theory”, John Wiley, New York, 1998.
- [36]. L. Breiman, J. H. Friedman, R. A. Olshen, C. J. Stone, “Classification and Regression Trees”, Wadsworth, Belmont-CA, 1984.
- [37]. J.H. Friedman, “Multivariate adaptive regression splines”, *Annals of Statistics*, 1991, 19(1), pp. 1-67.
- [38]. N.K. Kasabov and Q. Song, “DENFIS: dynamic evolving neural-fuzzy inference system and its application for time-series prediction”, *IEEE Transactions on Fuzzy Systems*, 2002, 10 (2), pp. 144–154.

- [39]. J.H. Friedman, "Stochastic gradient boosting", Stanford, Statistics Department, Stanford University, 1999.
- [40]. J. Moody and C.J. Darken, "Fast learning in networks of locally tuned processing units", *Neural Computation*, 1989, 1(2), pp. 281–294.
- [41]. G. Dematos, M. S. Boyd, B. Kermanshahi, N. Kohzadi and I. Kaastra, "Feedforward versus recurrent neural networks for forecasting monthly japanese yen exchange rates", *Asia-Pacific Financial Markets*, 1996, pp. 59-75.



Vasu Madireddi is pursuing M.Tech. in Information Technology with specializations in Banking Technology and Information Security from the University of Hyderabad, Hyderabad, India and the Institute for Development and Research in Banking Technology (IDRBT), Hyderabad, India. He holds M.Sc in Mathematics from Pondicherry University, Pondicherry. His research interests include data mining and soft computing.

Biographies



Jankisharan Pahareeya is presently working as a lecturer in Dehradun Institute of Technology, Dehradun, Uttarakhand, India. He holds M.Tech in Information Technology with specializations in Banking Technology and Information Security

from the University of Hyderabad, Hyderabad, India and the Institute for Development and Research in Banking Technology (IDRBT), Hyderabad, India. He holds M.Sc in Computer Science, Jiwaji University, Gwalior, (M.P). His research interests include data mining and soft computing.



Vadlamani Ravi is an Associate Professor in IDRBT, Hyderabad since April 2005. He holds a Ph.D. in Soft Computing from Osmania University, Hyderabad & RWTH Aachen Germany. Earlier, he was a Faculty at NUS, Singapore for three years. He published 84 papers in refereed Journals /

Conferences and invited book chapters. He edited "*Advances in Banking Technology and Management: Impacts of ICT and CRM*" published by IGI Global, USA. He is a referee for several international journals and on the Editorial board of IJIDS, IJDATS, IJISSS, IJSDS & IJITPM. He is listed in the Marquis Who's Who in the World in 2009.



Mahil Carr was awarded a doctoral degree in Information Systems from the City University of Hong Kong. Presently, he is Associate Professor at IDRBT and is the R&D Coordinator for the Institute since December 2007. His current research interests are in the areas of software

engineering, information systems security and electronic/mobile commerce. He has published research papers in several conferences and international journals. Dr. Carr is on the editorial board of the International Journal of E-Services and Mobile Applications (IJESMA) and the International Journal of Information Systems and Social Change (IJISSC).