

On the Value of Ensemble Effort Estimation

Ekrem Kocaguneli, *Student Member, IEEE*, Tim Menzies, *Member, IEEE*, Jacky Keung, *Member, IEEE*

Abstract—

Background: Despite decades of research, there is no consensus on which software effort estimation methods produce the most accurate models.

Aim: Prior work has reported that, given M estimation methods, no single method consistently outperforms all others. Perhaps rather than recommending *one* estimation method as *best*, it is wiser to generate estimates from *ensembles of multiple estimation methods*.

Method: 9 learners were combined with 10 pre-processing options to generate $9 \times 10 = 90$ *solo-methods*. These were applied to 20 data sets and evaluated using 7 error measures. This identified the best n (in our case $n = 13$) solo-methods that showed stable performance across multiple datasets and error measures. The top 2, 4, 8 and 13 solo-methods were then combined to generate 12 *multi-methods*, which were then compared to the solo-methods.

Results: (i) The top 10 (out of 12) multi-methods significantly out-performed *all* 90 solo-methods. (ii) The error rates of the multi-methods were significantly less than the solo-methods. (iii) The ranking of the best multi-method was remarkably stable.

Conclusion: While there is no best single effort estimation method, there exist best combinations of such effort estimation methods.

Index Terms—Software Cost Estimation, Ensemble, Machine Learning, Regression Trees, Support Vector Machines, Neural Nets, Analogy, k -NN



1 INTRODUCTION

Correctly estimating the effort required to develop software is of vital importance. Over or under-estimation of software development effort can lead to undesirable results:

- Under-estimation results in schedule and budget overruns, which may cause project cancellation.
- Over-estimation hinders the acceptance of promising ideas, thus threatening organizational competitiveness.

A practitioner when encountering the literature would almost certainly strike closed (i.e. fixed-parameter) models in the first instance e.g. COCOMO [1], FPA [2], SLIM [3]. Further possible models to encounter range from:

- Simple regression [4];
- To analogy-based methods [5];
- To complex combinations of techniques such as Corazza et al.'s [6] use of tabu search to configure support vector machines or Menzies et al.'s [7] approach that tries hundreds of different methods.

Not only is the current literature confusing and voluminous [8], some results suggest that it may be impossible to assess which effort estimators are the best. Shepperd et al. [9] warn that when comparing M estimation methods, the ranking of any one method may change, if the conditions are changed. Hence, they argue that it is fundamentally impossible to offer a definitive ranking such that $method_1$ is better than $method_2$.

This paper revisits the Shepperd et al. results and offers a more optimistic conclusion. We find that if we *combine*

- Ekrem Kocaguneli is with the Lane Department of Computer Science and Electrical Engineering, West Virginia University. E-mail: ekocagun@mex.wvu.edu
- Tim Menzies is with the Lane Department of Computer Science and Electrical Engineering, West Virginia University. E-mail: tim@menzies.us
- Jacky Keung is with the Department of Computing, The Hong Kong Polytechnic University. E-mail: Jacky.Keung@comp.polyu.edu.hk

This research is funded in part by NSF/CISE, project #0810879

the estimates from *multiple estimators*, then those combined methods perform *better than any single estimator*. That is, we agree with Shepperd et al. that there may be no single best effort estimator. However, there may be *best ensembles of methods*.

This should not be a surprising result. Many researchers argue that, in theory, best estimates come from combinations of multiple predictions. For example, Jorgensen advises that, for expert-based estimation, it is best to generate estimates from multiple methods [10]. Researchers in machine learning concur: e.g. Seni et al. report that averaging the estimates from *many methods* often does better than using *any solo method* [11]. Similar conclusions are offered by other researchers in the field of statistics [12] and machine learning [13], [14].

Current empirical results [15]–[17] in software engineering report the exact *opposite* effect than that predicted by Jorgensen and Seni et al., and others [10]–[14]:

- Kocaguneli et al. failed to improve estimation by averaging across the predictions of 14 estimators [15].
- Baker could not improve estimation accuracy by *boosting* ensembles [16]¹.

In that respect our results are different from the prior empirical studies on ensembles. Whereas previous studies on the ensembles (through different strategies) report that ensembles are not statistically better than single learners, our study reports that (through the right strategy) ensembles can outperform single learners. This paper resolves this contradiction between theory and experimental results. To the best of our knowledge, this is the first such result in the effort estimation literature that is supported by extensive experimentation.

1. Boosting is an ensemble example that generates a series of methods where $method_i$ focuses on the cases that were most difficult for $method_{i-1}$ [18].

It will be argued that the mistake made by Baker et al. and Kocaguneli et al. was to assume that *all* solo methods are candidates for combination into multi-method ensembles. This is not the case. Solo methods can be sorted into a minority of *superior* methods and a majority of *inferior* methods. We show below that, in the majority case, solo methods are out-performed by ensembles of multi-methods built from the superior set.

We therefore offer the following advice for a successful ensemble of methods:

- 1) Try a large number of methods among which there are at least some good methods (shown to have a good performance by prior work).
- 2) Sort the methods using the evaluation methods discussed in this paper. Discard all but the best solo methods.
- 3) Build ensembles from the remaining solo methods.

This research makes the following contributions:

- A novel scheme for ensembling the best solo-methods, whose product is successful results (unlike previous research) concerning multi-methods applied on effort data
- An evaluation method for the stability of methods
- Stable multi-methods that outperform *all* solo-methods

This paper is structured as follows. §2 discusses related work. §3 summarizes the learners, pre-processing options and solo/multi-methods used in this study. Our methodology is explained in §4, which generates the results of §5. The threats to the validity of our results are reviewed in §6. We provide a discussion of our work in §7 and conclude with §8.

2 RELATED WORK

2.1 Software Effort Estimation Methods

Software effort estimation (from now on SEE) can be defined as the process of estimating the total effort necessary to complete a software project [19]. According to the extensive systematic review conducted by Jorgensen and Shepperd, developing new models is the biggest research topic in SEE since 1980s [8]. Therefore, there are many SEE models that have been proposed over the years and a taxonomy is necessary to classify such a large corpus. Myrtveit et al. defines taxonomy as an explanation of a concept by highlighting the similarities and differences between that particular concept and the related ones [20].

There exists a number of different taxonomies proposed in the literature [20], [21]. Briand et al. report that there is no agreement on the best taxonomy and define all proposed taxonomies to be subjective and to have flaws [21]. For example Menzies et al. divide SEE methods into two groups: Model-based and expert-based [7]. According to this taxonomy model-based methods use some algorithm(s) to summarize old data and to make predictions regarding the new data. On the other hand, expert-based methods make use of human expertise, which is possibly supported by process guidelines and/or checklists.

Myrtveit et al. use a different taxonomy, where they propose a dataset dependent differentiation between methods [20]. According to that taxonomy the methods are divided into two:

- Sparse-data methods that require few or no historical data: e.g. expert-estimation [10].
- Many-data approaches where certain amount of historical data is a must: e.g. functions and arbitrary function approximations (such as classification and regression trees).

Shepperd et al. propose a 3-class taxonomy [5]: 1) expert-based estimation, 2) algorithmic models and 3) analogy. According to this taxonomy expert based models target the consensus of human experts through some process. Jorgensen et al. define expert-based methods as a human-intensive process of negotiating the estimate of a new project and arriving at a consensus [10]. There are formal methods proposed for expert-based estimation like Delphi [22]. However, Shepperd et al. notes in another study that it is mostly the case that companies follow an informal process for expert-based estimation [9]. Algorithmic models include the adaptation of a formula to local circumstances or local data. Prominent examples to these methods are the COCOMO method [1] and function points [23]. Analogy based methods include finding past projects that are similar to the current project that is to be estimated and then adapting the effort values of these past projects.

Regardless of the taxonomy used to group SEE methods under different classes, the ultimate goal of all the methods is to generate realistic estimates. In [24] Jorgensen defines some guidelines for generating realistic software effort estimates. An important finding in Jorgensen's study (which parallels our findings) is that *combining estimations* coming from different sources (e.g. from experts and instance-based learners) captures a broader range of information related to the estimation problem. Hence, multi-source predictions offer the most robust and accurate estimator.

2.2 Ranking Instability

One long-standing issue with software effort estimation is *ranking instability*. Ideally, given M methods, we can definitively rank them in terms of their predictive accuracy:

$$M_1 \geq M_2 \geq M_3 \geq \dots$$

This ideal has yet to be seen in the literature. To the contrary, Shepperd et al. report that there is no certainty in that sort [20], [25]. For example, in the experiments of [25], a large number of synthetic data sets were generated (from distributions found in one real-world data set). As they changed the conditions of their experiments, Shepperd et al. found that no method was consistently best across every condition. Specifically, they found that the performance of a method was dependent on the dataset, the random number generator used to select train and test sets, and the evaluation method used to assess model accuracy. Hence, they concluded that it was not feasible to sort methods into some definitive order.

Extending their work, we say that when M methods are sorted to find a rank r_i for method M_i , then if experimental conditions are changed, then that rank will change by an amount $\delta r > 0$.

Recently, researchers have access to more methods than that used by Shepperd et al. For example, Menzies et al. [26]

studied 158 methods. While that study was over a very limited data set (just two old COCOMO data sets), their preliminary results prompted this study (where we work with 20 data sets). In a result consistent with Shepperd et al., Menzies et al. found that as we changed the random numbers used to generate train/test sets, then all 158 methods showed some $\delta r > 0$; i.e. their precise ranking changed. However, they also found a small number of methods with two interesting properties:

- 1) Their ranks were very high;
- 2) Their δr values were very small.

Property #1 means that some methods performed comparatively very well and Property #2 means that high ranking persisted across multiple experimental conditions. Hence, if we look at enough methods, it may be possible to rank methods even when they exhibit $\delta r > 0$. Accordingly, in the results reported below, we will report not only the rank of each method, but also their associated δr . The methods we recommend will be those with

- high ranks;
- and low δr .

2.3 Ensemble of Methods

A standard machine learning technique is to try multiple methods on the available data, then recommend the one that performs the best [27]. Many effort estimation papers apply this technique to demonstrate that (say) their preferred new method is superior to those proposed in prior work.

Ensemble learning takes a different approach. Rather than choosing *one* method, ensembles build multiple predictors, where estimates coming from different learners are combined through particular mechanisms, e.g. voting of individual learner estimates on the final prediction [28]. Before continuing any further we need to clear a terminology difference. From now on the term “*learner*” refers to a stand-alone algorithm without any supplemental pre or post processing step (e.g. *k*-NN, neural nets etc.), whereas the term “*solo-method*” will refer to an algorithm supplemented with a pre-processing option (e.g. logging+*k*-NN, discretization+neural nets etc.).

Ensembles are useful since any particular learner comes with its own assumptions [27]. These assumptions may be best suited to different parts of the training data [27]–[29]. In ensembles, methods can augment each other, i.e. a method patches errors made by another method. For example, when reducing estimated mean-squared-error, multi-methods attain smaller or equal error rates than single-methods [11].

It is a recommended practice to combine solo-methods that have different characteristics [27], [29], [30]. There are many techniques to attain different-characteristic multi-methods. The first way is through representation of the data. The multi-method structure may be based on uni-representation (all learners use same representation of data) or multi-representation (different learners use different representations) [27]. Examples to such strategies are the use of different feature sets [31], [32] or different training sets [33].

The second way is through architectural methodologies. Bagging (Bootstrap Aggregating) and boosting are among

the most common examples of that approach [27], [34]. In bagging *n*-many solo-methods are independently applied on *n*-many different training samples, where each training sample is selected via bootstrap sampling [27] with replacement. Boosting on the other hand arranges solo-methods in a sequential manner: each solo-method pays more attention to the instances on which previous method was unsuccessful. Boosting is reported to be considerably better than bagging [34]–[36], but has trouble in handling noisy datasets [34], [36].

2.4 Ensemble of Methods in SE

Ensemble methods are widely used in data mining (see the literature survey of [11]). Nevertheless, in the software domain, ensemble of methods (multi-methods) have not been reported to be superior to solo-methods in terms of prediction accuracy [15], [17]. In other words, ensembles fall short of providing a statistically significant increase in the prediction accuracy values over the solo-methods. For particular error measures that were used to evaluate the prediction accuracy see the work of Khoshgoftaar et al. and Kocaguneli et al. [15], [17]. Khoshgoftaar et al. [17] question the performance of different multi-method schemes under different scenarios in the domain of software quality. They use different combinations of 17 learners induced on 7 datasets and report that multi-methods induced on single datasets do not yield a significant increase in prediction accuracy.

Kocaguneli et al. [15] replicated the work of Khoshgoftaar et al. [17] in the domain of software effort estimation. They exploited combination of 14 methods applied on 3 software effort estimation datasets. Their conclusion was similar to that of the replicated study [17]: The application of multi-methods under different scenarios did not provide a significant increase in the estimation accuracy. Similarly, in the study of Kumar et al. [37], different learners were employed in two types of ensembles, but only one of them was reported to be successful.

Another example to ensembles is the ensemble of *single-type* learners, where multiple versions of a single learner are combined. Pahariya et al. use linear combinations of genetic algorithms [38], where they report improvements over single-learners. Kultur et al. report improvements through collections of neural networks [39]. Unless a learner is supplemented with a pre-processing or a post-processing option, then the learners in the ensemble are the copy of one another and have the same biases/assumptions. Unlike ensemble of single-type learners, all the multi-methods reported in our study is the combination of a learner augmented with pre/post-processing options.

There are also some applications of ensemble methods in effort estimation so as to process datasets: In [40], Twala et al. use multiple imputation techniques to handle missing data and in [41] Khoshgoftaar et al. make use of learner ensembles as a filter to improve the data quality.

Our ensemble is different from the above. We take care to prune inferior solo methods *before* building the ensemble. As shown below, this leads to very effective effort estimators.

3 EFFORT ESTIMATION METHODS

The effort estimation methods studied in this paper fall into two groups: *solo*-methods and *multi*-methods. *Solo*-methods

TABLE 1

The summary table for the *solo-methods*. This table provides a list of abbreviations as well as their explanations for pre-processing options and learners used in this research.

| Pre-processing Options | | Learners | |
|------------------------|---|--------------|--|
| Abbreviation | Explanation | Abbreviation | Explanation |
| norm | Normalization | ABE0-1NN | Basic ABE with 1 nearest neighbor |
| log | Taking natural logarithm | ABE0-5NN | Basic ABE with 5 nearest neighbors |
| PCA | Principal Component Analysis | SWReg | Stepwise Regression |
| SFS | Sequential Forward Selection | CART (yes) | Classification and Regression Tree with pruning |
| SWReg | Stepwise Regression | CART (no) | Classification and Regression Tree without pruning |
| width3bin | Discretize into 3 bins based on equal width | NNet | Neural Net with two hidden layers |
| width5bin | Discretize into 5 bins based on equal width | LReg | Simple linear regression |
| freq3bin | Discretize into 3 bins based on equal frequency | PCR | Principal components regression |
| freq5bin | Discretize into 5 bins based on equal frequency | PLSR | Partial least squares regression |
| none | Apply no pre-processor | | |

are some combination of a *pre-processing option* and a *learner*. For example, Boehm’s preferred effort estimation method uses a log transform as the pre-processing option, then linear regression as the learner. *Multi-methods* are combinations of solo-methods.

3.1 90 Solo-Methods

In our experiments, we used 10 different pre-processing options and 9 learners. These were selected on two criteria:

- Learners must come from the SE effort estimation literature; e.g. [4], [5], [8], [39], [42]–[46].
- Learners must make different assumptions about the data.

This second criteria is based on data-mining theory that recommends using different learners that fail under different circumstances [27]–[29], [47]. For example ABE methods assume that similar instances of the data set have similar dependent variable values, whose translation into SEE domain is that similar projects have similar effort values [44]. Decision tree inducers adopt a divide and conquer approach and assume that decisions on the prediction can be made through a sequence of tests/decisions that usually involve one feature at a time [34]. Artificial neural networks assume that data can be modeled through different geometries of directed graphs, where nodes represent the neurons and the connecting edges represent the weights applied on the output of the neurons [48]. With the arrival of each new instance, the effect is propagated through neurons and edges. Linear regression assumes that the trend seen in data can be formulated through a linear model [49]. Stepwise regression is built on assumption that the set of features that maximize the F-value (evaluates if the variables in the model -all together- are significantly related to the independent variable) should be used in the final model [50].

Ensembles work best when one member of the ensemble patches the mistakes made by other methods of the ensemble. We hence used 10 pre-processing options:

- Three *simple pre-processing options*: **none**, **norm**, and **log**;
- One *feature synthesis* method: **PCA**;
- Two *feature selection* methods: **SFS** and **SWreg**;
- Four *discretization* methods: Based on equal frequency/width.

and 9 learners:

- Two *iterative dichotomizers*: **CART(yes)**, **CART(no)**;
- A *neural net*: **NNet**;
- Four *regression methods*: **LReg**, **PCR**, **PLSR**, **SWReg**.
- Two *instance-based learners*: **ABE0-1NN**, **ABE0-5NN**;

Note that “ABE” is short for analogy-based effort estimation. ABE0-kNN is a standard analogy-based estimator with execution steps of:

- *Normalization* of data to zero-one interval;
- A *Euclidean* distance measure;
- Estimates generated using the *k nearest neighbors*.

For detailed descriptions of all these learners, see Appendix.

Combining 10 pre-processing options and 9 learners results in $10 \times 9 = 90$ solo-methods. A summary table for abbreviations of the methods and pre-processing options is provided in Table 1. Note that some of the combinations are less plausible than the others, for example:

- **norm & ABE0**, as ABE0 already has a normalization mechanism in it
- **different discretizations & CART**, as CART already has an inherent discretization mechanism
- **SWReg & SWReg**, as using the learner itself as a pre-processor would be too cumbersome

However, we included them in our analysis because they would find their appropriate rankings when compared with better methods. Furthermore, this brute-force analysis with all possible (but not necessarily plausible) pre-processing and learner combinations has to be carried out only once to find the better performing solo-methods that are to be ensembled into multi-methods.

3.2 Multi-Methods

Multi-methods combine two or more solo-methods. Many combination schemes have been proposed in the literature [27], [29], [30], [43]. Complex combination schemes include bagging [51], boosting [18] or random forests [52], [53]. Simpler methods include computing the mean, median or inverse-ranked weighted mean (IRWM [43], see Figure 1) of estimates coming from *n-many* solo-methods.

Our aim is not to investigate complex schemes, but to observe how multi-methods perform compared to solo-methods

In IRWM, the final estimates from M methods e_1, e_2, \dots, e_m that have been ranked r_1, r_2, \dots, r_m is a weighted sum by the ranks of all methods in the ensemble. The top and bottom-ranked methods of m methods get a weight of m and 1 (respectively). More generally, a method with rank r_i gets a weight of $m+1-r_i$. The final estimate in IRWM is hence $(\sum_i (m+1-r_i) \cdot e_i) / (\sum_i i)$.

Fig. 1. IRWM. Generalized from [43].

on effort datasets. Therefore, we adopt simple schemes (mean, median and IRWM).

4 METHODOLOGY

4.1 Multiple Error Measures

This section describes several performance measures used in this research. All the performance measures listed here have the property that we can find at least one publication proposing their use for effort estimation.

Error measures comment on the success of a prediction. For example, the absolute residual (AR) is the difference between the predicted and the actual values:

$$AR_i = |x_i - \hat{x}_i| \quad (1)$$

(where x_i, \hat{x}_i are the actual and predicted values respectively for test instance i). MAR is the mean of individual AR values.

The Magnitude of Relative Error measure a.k.a. MRE is a very widely used evaluation criterion for selecting the best effort estimator from a number of competing software prediction models [44], [54]. MRE measures the error ratio between the actual effort and the predicted effort. It can be expressed as the following equation:

$$MRE_i = \frac{|x_i - \hat{x}_i|}{x_i} = \frac{AR_i}{x_i} \quad (2)$$

A related measure is MER (Magnitude of Error Relative to the estimate [54]):

$$MER_i = \frac{|x_i - \hat{x}_i|}{\hat{x}_i} = \frac{AR_i}{\hat{x}_i} \quad (3)$$

A summary of MRE can be derived as the Mean Magnitude of Relative Error (MMRE) or Median Magnitude of Relative Error (MdMRE), which can be calculated as follows respectively:

$$MMRE = \frac{\sum_{i=1}^n MRE_i}{n} \quad (4)$$

$$MdMRE = \text{median}(MRE_1, MRE_2, \dots, MRE_n) \quad (5)$$

A common alternative error measure is PRED(25), which can be defined as the percentage of predictions falling within 25% of the actual values:

$$PRED(25) = \frac{100}{N} \sum_{i=1}^N \begin{cases} 1 & \text{if } MRE_i \leq \frac{25}{100} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

For example, PRED(25)=50% implies that half of the estimates fall within 25% of the actual values [44].

There are many other error measures including Mean Balanced Relative Error (MBRE) and the Mean Inverted Balanced Relative Error (MIBRE) studied by Foss et al. [54]:

$$MBRE_i = \frac{|\hat{x}_i - x_i|}{\min(\hat{x}_i, x_i)} \quad (7)$$

$$MIBRE_i = \frac{|\hat{x}_i - x_i|}{\max(\hat{x}_i, x_i)} \quad (8)$$

Interpreting these error measures without any statistical test may be misleading. A recent discussion about this issue can be found in [55]. To evaluate our results subject to a statistical test, we make use of the so called *win-tie-loss* statistics. *Win-tie-loss* statistics employ a Wilcoxon non-parametric statistical hypothesis test with 95% confidence. Wilcoxon is more robust than the Student's *t*-test as it compares the sums of ranks, unlike Student's *t*-test, which may introduce spurious findings as a result of outliers in the given datasets. Ranked statistical tests like the Wilcoxon are also useful, if it is not clear that the underlying distributions are Gaussian [56].

We stored the performance of every method w.r.t. 7 error measures over 20 datasets. This enabled us to collect *win-tie-loss* statistics using the algorithm of Figure 2. In Figure 2, we first check if two distributions i, j are statistically different according to the Wilcoxon test (95%); if they are not, then we increment tie_i and tie_j . If the distributions are statistically different, we update win_i, win_j and $loss_i, loss_j$ after comparing their error measures.

```

if WILCOXON( $E_i, E_j, 95$ ) says they are the same then
     $tie_i = tie_i + 1;$ 
     $tie_j = tie_j + 1;$ 
else
    if better( $E_i, E_j$ ) then
         $win_i = win_i + 1$ 
         $loss_j = loss_j + 1$ 
    else
         $win_j = win_j + 1$ 
         $loss_i = loss_i + 1$ 
    end if
end if
    
```

Fig. 2. Comparing methods (i, j).

The *better* function in the *if* statement of Figure 2 varies according to the performance criteria. For some error measures such as MMRE and MdMRE, better means lower values, i.e. lower means and medians respectively. However, for PRED(25), better means higher PRED(25) values.

Note that 20 data sets have considerable size differences, e.g. *China* data set contains almost half of all the records used in this study. However, the above mentioned performance measures are derived separately for every dataset.

4.2 Experimental Conditions

Recall the results of Shepperd et al. [25] and Menzies et al. [26]: different experimental conditions can change the rank of an effort estimator. Hence, it is important to study not just the rank of an estimator, but also how well that method performs across multiple experimental conditions such as:

TABLE 2

The 1198 projects used in this study come from 20 data sets. Indentation in column one denotes a dataset that is a subset of another dataset. For notes on these datasets, see Appendix.

| Dataset | Features | Size | Description | Historical Effort Data | | | | | |
|-----------------|----------|------|--|------------------------|-------|--------|--------|--------|------|
| | | | | Units | Min | Median | Mean | Max | Skew |
| cocomo81 | 17 | 63 | NASA projects | months | 6 | 98 | 683 | 11400 | 4.4 |
| cocomo81e | 17 | 28 | Cocoma81 embedded projects | months | 9 | 354 | 1153 | 11400 | 3.4 |
| cocomo81o | 17 | 24 | Cocoma81 organic projects | months | 6 | 46 | 60 | 240 | 1.7 |
| cocomo81s | 17 | 11 | Cocoma81 semi-detached projects | months | 5.9 | 156 | 849.65 | 6400 | 2.64 |
| nasa93 | 17 | 93 | NASA projects | months | 8 | 252 | 624 | 8211 | 4.2 |
| nasa93_center_1 | 17 | 12 | Nasa93 projects from center 1 | months | 24 | 66 | 139.92 | 360 | 0.86 |
| nasa93_center_2 | 17 | 37 | Nasa93 projects from center 2 | months | 8 | 82 | 223 | 1350 | 2.4 |
| nasa93_center_5 | 17 | 40 | Nasa93 projects from center 5 | months | 72 | 571 | 1011 | 8211 | 3.4 |
| desharnais | 12 | 81 | Canadian software projects | hours | 546 | 3647 | 5046 | 23940 | 2.0 |
| desharnaisL1 | 11 | 46 | Projects in desharnais that are developed with Language1 | hours | 805 | 4035.5 | 5738.9 | 23940 | 2.09 |
| desharnaisL2 | 11 | 25 | Projects in desharnais that are developed with Language2 | hours | 1155 | 3472 | 5116.7 | 14973 | 1.16 |
| desharnaisL3 | 11 | 10 | Projects in desharnais that are developed with Language3 | hours | 546 | 1123.5 | 1684.5 | 5880 | 1.86 |
| sdr | 22 | 24 | Turkish software projects | months | 2 | 12 | 32 | 342 | 3.9 |
| albrecht | 7 | 24 | Projects from IBM | months | 1 | 12 | 22 | 105 | 2.2 |
| finnish | 8 | 38 | Software projects developed in Finland | hours | 460 | 5430 | 7678.3 | 26670 | 0.95 |
| kemerer | 7 | 15 | Large business applications | months | 23.2 | 130.3 | 219.24 | 1107.3 | 2.76 |
| maxwell | 27 | 62 | Projects from commercial banks in Finland | hours | 583 | 5189.5 | 8223.2 | 63694 | 3.26 |
| miyazaki94 | 8 | 48 | Japanese software projects developed in COBOL | months | 5.6 | 38.1 | 87.47 | 1586 | 6.06 |
| telecom | 3 | 18 | Maintenance projects for telecom companies | months | 23.54 | 222.53 | 284.33 | 1115.5 | 1.78 |
| china | 18 | 499 | Projects from Chinese software companies | hours | 26 | 1829 | 3921 | 54620 | 3.92 |
| Total: 1198 | | | | | | | | | |

- 1) *Error measures* that measure a method’s performance;
- 2) *Comparison summaries* used to report the performance of many methods over many data sets;
- 3) *Data sets* used in the experiments.

The *error measures* used in this study, as defined above, are MAR, MMRE, MdMRE, MMER, PRED(25), MBRE, MIBRE.

As to *comparison summaries*, the following procedure was repeated for each error measure. Each of our 90 methods were compared to 89 others using the procedure of Figure 2. In our procedure, we sum the *wins*, *losses*, *ties* of Figure 2 and we rank our methods by that sum. That is, for an estimation method to be highly ranked, it must perform comparatively well across *all* error measures.

The results of those comparisons are contained in the *win*, *tie*, *loss* counters of Figure 2. These comparisons can be summarized through many ways:

- 1) Number of losses;
- 2) Number of wins;
- 3) Number of wins-losses

DemSar [57] reports that there is no generally accepted method of *comparison summarization*. Hence, we compute δr from the changes in the ranks seen when we move across all three summarization methods:

- Firstly, we generate rankings of estimators using *number of losses* (repeated for all error measures) via a leave-one-out procedure (where each example becomes a test and the remaining data is used for training).
- Next, we compare that rank to other ranks generated by other comparisons summaries (again, repeated for all error measures).
- The maximum rank change for a method over all the dimensions is the δr value for that method.

Finally, we run our rig over multiple data sets. Our experiments use the 20 publicly available effort estimation datasets in the PROMISE repository (see <http://promisedata.org/?cat=14>).

The names and properties of these datasets are provided in Table 2. Note that some projects in desharnais dataset contain missing values and we used mean imputation [58] to handle missing values. See Appendix for more explanatory notes.

Combining the above, we can see that the experiments of this paper are an extensive analysis of different conditions for effort estimation experiments. Given 89 comparisons among solo-methods, 7 error measures, and 20 datasets, then each method appears in $89 \times 7 \times 20 = 12,460$ comparisons.

4.3 Focus on Superior Methods

Figure 3 is a plot that shows the success of the solo-methods through ranking and the variability in that ranking. The x-axis of Figure 3 shows the ranking of the 90 solo-methods, according to number of losses over all 7 error measures and 20 data sets. The most successful methods have the lowest number of total losses whereas the least successful ones have the highest number of total losses. Then solo-methods are ranked on the x-axis starting from the best one. Therefore, better methods appear on the left-hand-side of that figure (so the top-ranked method appears at position $x = 1$).

The ranking of methods is identified by the x-axis of Figure 3, whereas the variability in that ranking is given by the y-axis. The y-axis of Figure 3 shows the maximum *changes*, δr , seen for each method as we compare the ranks across number of losses, number of wins, and number of wins-losses. In other words, a solo-method has different rankings according to its *win*, *loss* or *win – loss* values and the related number seen on the y-axis is the biggest difference between its best and worst rankings. In a result consistent with Shepperd et al., all methods have $\delta r > 0$. However, the good news is that the top-ranked methods have a very low δr . That is, even if these top-ranked methods jumped rank by their maximum δr , then they would still be performing better than most of the other 90 methods.

In signal processing, it is standard practice to segment data based on the region of maximal change [59]. An inspection

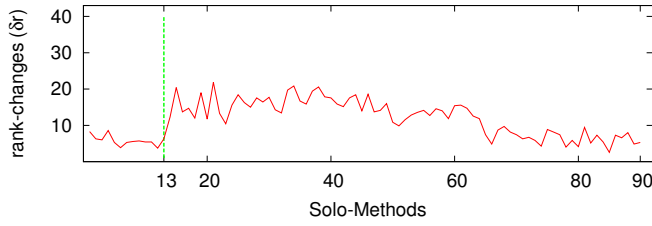


Fig. 3. Methods and the associated changes, i.e. δr values. Note the sudden increase in δr values, after $X = 13$. We call methods in the region $1 \leq X \leq 13$ as *superior*.

of Figure 3 shows that the region of maximal δr occurs after $X = 13$. This is an interesting division since the region $1 \leq X \leq 13$ contains methods with high rank and low δr . We call these methods *superior* and the rest *inferior*. The list of top 13 methods are shown in Table 3.

Figure 3 is investigated further in a related paper that discusses the relative ranking concept of the solo-methods [60]. The names of the top 13 *superior* solo-methods of Figure 3 are listed in Table 3. When we look at Table 3, we see that none of the superior solo-methods try to fit one model to all the data:

- The CART regression tree learner appears at ranks 1 through 10 of Table 3. Each branch of a regression tree defines one context in which an estimate may be different.
- Analogy-based estimation (ABE) appears at ranks 11,12,13. ABE builds a different model for each test instance (using the test instance’s k-th nearest neighbors).

Solo-methods are not the focus of this paper. Hence, we move on to discuss multi-methods.

4.4 Build Ensembles

To form multi-methods, we build ensembles using the top M solo-methods in the sort order of Table 3. In this study, we use $M \in \{2, 4, 8, 13\}$.

To generate an estimate, we ask all members of an ensemble to offer a prediction. These are combined in one of three ways: mean, median and IRWM.

TABLE 3

Ranking of top-13 *superior* solo-methods and related δr values. These solo-methods are combined in various ways to form 12 multi-methods.

| rank | δr | pre-processing option | learner |
|------|------------|-----------------------|------------|
| 1 | 8 | norm | CART (yes) |
| 2 | 6 | norm | CART (no) |
| 3 | 6 | none | CART (yes) |
| 4 | 9 | none | CART (no) |
| 5 | 5 | log | CART (yes) |
| 6 | 4 | log | CART (no) |
| 7 | 5 | SWReg | CART (yes) |
| 8 | 6 | SWReg | CART (no) |
| 9 | 6 | SFS | CART (yes) |
| 10 | 5 | SFS | CART (no) |
| 11 | 5 | SWReg | ABE0-1NN |
| 12 | 4 | log | ABE0-1NN |
| 13 | 5 | SWReg | ABE0-5NN |

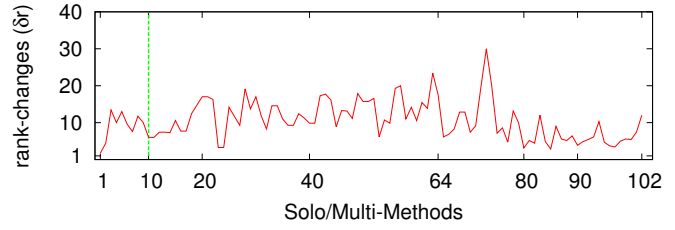


Fig. 4. Rank changes of solo and multi-methods. Region $1 \leq X \leq 10$ contains 10 out of 12 multi-methods. See that **Top13/Mean** at $X = 1$ has a δr of 1, i.e. it outperforms all other methods w.r.t. 7 different error measures and 20 datasets.

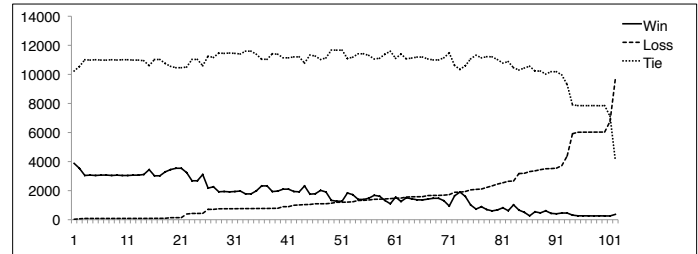


Fig. 5. The sum of win, tie and loss values for all methods of Figure 4 (over all error measures and all datasets). Since one method is compared to 101 other methods, over 7 error measures and 20 datasets, the sum of win, tie and loss values is: $101 \times 7 \times 20 = 14,140$.

These multi-methods are then ranked along-side the solo-methods in the same manner as Figure 3. This gives us the comparison of 90 solo-methods + 12 multi-methods = 102 methods. Every method is compared to 101 others with respect to seven error measures and over 20 datasets. Therefore, the maximum number of comparisons for any method now becomes $101 \times 7 \times 20 = 14,140$. To the best of our knowledge this is the one of the most extensive effort estimation experiments yet reported in the literature (and for extensive non-experimental studies, see [8], [61]).

5 RESULTS

Figure 4 shows the rank of our 102 methods. As before, the x-axis ranks the methods according to number of losses and the y-axis shows the δr of each method. Table 4 shows the 102 methods, sorted in the same way as the x-axis of Figure 4.

Two aspects of these result are worth commenting:

- The top $X = 10$ methods (marked by a dashed line) are all multi-methods. The remaining multiple methods appear at ranks 14,15,18. That is, in the majority case ($\frac{10}{12} = 83\%$), combinations of methods perform better than any solo method. Further, in all cases ($\frac{12}{12} = 100\%$), they are ranked higher than the majority of other methods.
- The multi-method at $X = 1$ also has the lowest δr of any method in this study. This method generated estimates using the IRWM value of 13 top-ranked solo-methods.

TABLE 4

Detailed pre-processing option and learner combinations and related δr values. Methods are sorted by the sum of their losses seen in all performance measures and all data sets. The method with fewest losses is ranked #1 and is **Top13/Irwm**. At the other end of the scale, the method with the most losses is ranked #102 and is **PCA/LReg**.

| rank | δr | pre-proc. | learner | rank | δr | learner |
|------|------------|-----------|------------|------|------------|----------------------|
| 1 | 1 | Top13 | Irwm | 52 | 17 | norm SWReg |
| 2 | 4 | Top13 | Mean | 53 | 6 | none SWReg |
| 3 | 13 | Top13 | Median | 54 | 11 | freq3bin ABE0 |
| 4 | 10 | Top2 | Mean | 55 | 10 | width3bin CART (yes) |
| 5 | 13 | Top4 | Mean | 56 | 19 | width3bin CART (no) |
| 6 | 10 | Top2 | Median | 57 | 20 | PCA ABE0 |
| 7 | 8 | Top4 | Median | 58 | 11 | PCA NNet |
| 8 | 12 | Top8 | Median | 59 | 14 | none NNet |
| 9 | 10 | Top2 | Irwm | 60 | 11 | width5bin SWReg |
| 10 | 6 | Top4 | Irwm | 61 | 15 | width3bin ABE0 |
| 11 | 6 | norm | CART (yes) | 62 | 14 | SWReg NNet |
| 12 | 7 | norm | CART (no) | 63 | 23 | SFS NNet |
| 13 | 7 | none | CART (yes) | 64 | 18 | width5bin INN |
| 14 | 7 | none | CART (no) | 65 | 6 | SWReg LReg |
| 15 | 11 | Top8 | Irwm | 66 | 7 | none LReg |
| 16 | 8 | log | CART (yes) | 67 | 8 | norm PLSR |
| 17 | 8 | log | CART (no) | 68 | 13 | width5bin ABE0 |
| 18 | 12 | Top8 | Mean | 69 | 13 | norm LReg |
| 19 | 15 | SWReg | CART (yes) | 70 | 7 | freq5bin INN |
| 20 | 17 | SWReg | CART (no) | 71 | 9 | freq3bin CART (yes) |
| 21 | 17 | SWReg | INN | 72 | 20 | freq3bin CART (no) |
| 22 | 16 | SFS | CART (yes) | 73 | 30 | PCA INN |
| 23 | 3 | SFS | CART (no) | 74 | 20 | freq3bin INN |
| 24 | 3 | log | INN | 75 | 7 | width3bin SWReg |
| 25 | 14 | SWReg | ABE0 | 76 | 9 | log SWReg |
| 26 | 12 | PCA | PLSR | 77 | 5 | width5bin PLSR |
| 27 | 9 | none | PLSR | 78 | 13 | log PCR |
| 28 | 19 | SWReg | PCR | 79 | 10 | log PLSR |
| 29 | 14 | PCA | PCR | 80 | 3 | width3bin INN |
| 30 | 17 | none | PCR | 81 | 5 | width3bin PLSR |
| 31 | 12 | SFS | INN | 82 | 4 | width5bin PCR |
| 32 | 8 | PCA | CART (yes) | 83 | 12 | norm PCR |
| 33 | 15 | PCA | CART (no) | 84 | 5 | width3bin LReg |
| 34 | 15 | SFS | ABE0 | 85 | 3 | width3bin PCR |
| 35 | 11 | norm | INN | 86 | 9 | freq5bin PCR |
| 36 | 9 | none | INN | 87 | 6 | freq5bin SWReg |
| 37 | 9 | freq5bin | CART (yes) | 88 | 5 | width5bin LReg |
| 38 | 12 | freq5bin | CART (no) | 89 | 6 | freq3bin PCR |
| 39 | 11 | freq5bin | ABE0 | 90 | 4 | freq3bin PLSR |
| 40 | 10 | SFS | LReg | 91 | 5 | freq5bin PLSR |
| 41 | 10 | width5bin | CART (yes) | 92 | 5 | log LReg |
| 42 | 17 | width5bin | CART (no) | 93 | 6 | freq3bin SWReg |
| 43 | 18 | SWReg | PLSR | 94 | 10 | freq5bin LReg |
| 44 | 16 | SFS | PLSR | 95 | 5 | width5bin NNet |
| 45 | 9 | SFS | PCR | 96 | 4 | norm NNet |
| 46 | 13 | norm | ABE0 | 97 | 3 | width3bin NNet |
| 47 | 13 | PCA | SWReg | 98 | 5 | log NNet |
| 48 | 11 | none | ABE0 | 99 | 6 | freq3bin NNet |
| 49 | 18 | SWReg | SWReg | 100 | 5 | freq5bin NNet |
| 50 | 16 | log | ABE0 | 101 | 7 | freq3bin LReg |
| 51 | 16 | SFS | SWReg | 102 | 12 | PCA LReg |

Note that the second result is exactly the “ensembles are better” result as might be predicted by the researchers mentioned in the introduction [10]–[14]. We speculate that prior SE researchers who failed to find that “ensembles are better” did not prune away inferior solo-methods before building an ensemble. One word of caution here is that although we use a deterministic sampling method (LOOCV) in this research, some algorithms contain probabilistic processes (e.g. back-pruning process of **CART (yes)**). This may create small alterations in the performances of the methods (solo and multi) using these algorithms. Hence, it may be the case that one gets a slight variation of Table 4; however, the general picture (solo-

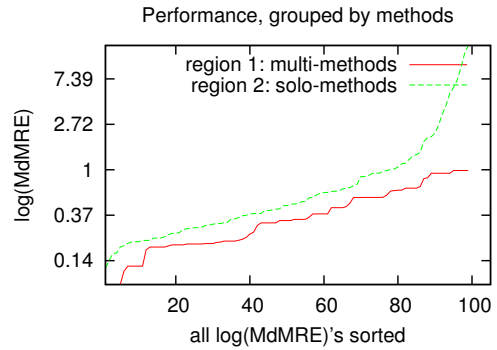


Fig. 6. Spectrum of MdMRE values for 2 regions of methods: Solo-methods and multi-methods. Here we keep the order of methods same as those given in Table 4 and divide those methods into 2 regions. Notice how multi-methods attain the lowest MdMRE scores.

methods outperformed by multi-methods) remains the same.

Better yet, as shown in Figure 4, this largest ensemble at rank $X = 1$ had the lowest δr seen in any of our 102 methods ($\delta r = 1$). This result underscores the main message of this paper: the method that scored the best (and had the greatest stability across different experimental conditions) was the one that used the highest number of *superior* solo-methods.

Figure 5 shows the sum of *win*, *tie* and *loss* values for the methods of Figure 4. Every method of Figure 4 is compared to 101 other methods, over 7 error measures and 20 datasets, so the maximum value that either one of the *win*, *tie*, *loss* statistics can attain is: $101 \times 7 \times 20 = 14,140$. Notice in Figure 5 that except the low performing methods on the right hand-side, the *tie* values are in 10,000 – 12,000 band. Therefore, they would not be so informative as to differentiate the methods, so we consult to *win* and *loss* statistics. There is a considerable difference between the best and the worst methods in terms of *win* and *loss* values (in the extreme case it is close to 4,000). In a way Figure 5 is a sanity check of Figure 4, because it shows that the rankings reported in Figure 4 is due to considerable *win* and *loss* value differences between high (left hand-side) and low (right hand-side) performing methods.

Other results offer yet more evidence for the superiority of multi-methods over solo-methods. Figure 6 sorts the MdMRE values of all the solo-methods and all the multi-methods:

- Multi-methods generated lower MdMRE values than the solo-methods;
- While some of the solo-methods have alarmingly large errors (as observed from the steep right-hand-side of the dashed line in that figure), note that *none* of the 12 best multi-methods have large MdMRE values.

That is, multi-methods are far less prone to incorrect effort estimates than solo-methods.

6 THREATS TO VALIDITY

External validity questions whether the results can be generalized outside the specifications of a study [62]. To ensure

external validity, this paper has studied a large number of projects. For example, Table 4 of [61] lists the total number of projects used by a sample of other studies. The median value of that sample is 186; i.e. one-sixth of the 1198 projects used here. In terms of external validity, this report has higher validity than a standard effort estimation study (which might apply one error measure to a handful of data sets). An ideal case for our work or other similar studies would be to request data sets with specific properties that would best suit the experimental concerns and that would support the random sampling of data in experimentation. However, SEE data sets are neither easy to find nor easy to collect. Therefore, this stands as an external validity issue and calls for more effort in SEE data collection.

It is clear that this study has not explored a full range of effort estimation methods. Future work is required to repeat this study using the top performing solo-methods found here and possibly more. Nevertheless, given the extent of the experimentation reported here, this paper offers much support for the general claim that when generating estimates, it is wise to *first* sort and prune all available solo-methods, and *then* try combinations of the best solo-methods. Another validity issue to mention is that none of learners have been exhausted via fine-tuning. Therefore, future work is required to exhaust all the parameters of every learner to use their best version. However, exhaustive fine-tuning of a learner through some heuristic may be as comprehensive as a paper on its own right, e.g. tabu search heuristic proposed by Corazza and Mendes et al. [6].

Construct validity (i.e. face validity) asks if we are measuring what we actually intended to measure [63]. Previous studies have concerned themselves with the construct validity of different performance measures for effort estimation (e.g. [54]). To make sure that we do not favor a particular conclusion due to limited number of performance measures, we used 7 error measures and questioned the stability of our conclusions. We showed empirically the surprising result that multi-methods are stable across a range of performance criteria. The surprise factor comes from the fact that solo-methods and successful multi-methods have different rank change values. Unlike most of the solo-methods with high rank changes w.r.t different performance measures (as high as 30), the top-10 multi-methods are quite *stable and successful*, with the highest rank change of 13 and the lowest rank change of 1.

In terms of *internal validity*, there is one dimension of experimental conditions not explored above. This paper used leave-one-out to assess methods. An alternate experimental condition would be to use N-way cross-validation where examples are divided into B bins and each bin is tested on a model learned from the remaining bins. In theory, not using cross-validation is a threat to the validity of our results since we did not check if our results were stable across both leave-one-out and cross-validation. According to Hastie et al. [12, p242], leave-one-out testing generates:

- Lower bias estimates than cross-validation (since cross-validation must learn from fewer examples);
- Higher variance estimates than cross-validation (since

leave-one-out conducts many more tests).

Similar assertions are made by Kohavi [13] and Breiman [14].

Elsewhere [64], we have documented the surprising result that, at least for the data sets of Table 2, there is very little difference in the bias and variance values generated for leave-one-out and cross-validation. We conjecture that the standard comments on “bias vs. variance trade-off” for “leave-one-out vs cross-validation” needs to be revisited for the small effort data sets seen in effort estimation.

Since leave-one-out and cross-validation have similar bias/variance profiles, we are free to select either method. This paper uses leave-one-out for the following reasons. Leave-one-out is a deterministic procedure that can be exactly repeated by any other researcher with access to a particular data set. The same cannot be said about cross-validation since it divides the data into bins using:

- A random number generator;
- Some stratification heuristics. Stratification attempts to maintain the same class distributions in each bin as in the entire data set). Stratification is a heuristic procedure since it may encounter issues that have no singular solution (e.g. how to space rare outliers amongst the bins).

Note that without knowledge of the stratification heuristics, and without access to an identical random number generator, it can be difficult for researcher “A” to reproduce the cross-validation results of researcher “B”.

7 DISCUSSION

7.1 Learning Curve

One of our observations is that ensembles are more trustworthy (less ranking instability). For someone with a strong background in machine learning algorithms, the number of learners to combine is not an issue as he/she has already gone through the learning curve. However, from a practitioner’s point of view, there is a cost-benefit trade-off between:

- the cost of learning new learners
- and the additional performance benefit

We acknowledge the fact that building an ensemble model from scratch may be too challenging for practitioners without prior machine learning experience. However, our industry-collaborated project experience shows that once the model is built by researchers, its adoption/implementation by practitioners is a pretty straightforward process [65]–[68]. This section provides an in depth discussion and alternative solutions for such practitioners.

The best multi-method reported in this paper requires combination of top-13 solo methods. The learning curve associated with top-13 solo methods is given in Table 5. From Table 5 we see that a practitioner willing to use multi-methods at all has to learn at least 2 learners (CART(yes) and CART(no)). Not included in Table 5 are the numeric manipulations such as normalizing an array of numbers and taking logarithm, mean or median of these numbers. Before learning the algorithms of Table 5, a practitioner will need to learn these numeric operations, then learning just 2 learners in fact enables him/her to use up until the 6th best performing solo-method. By using

TABLE 5

The learning curve of the top-13 *superior* solo-methods. The column “# of additions” shows how many algorithms (learner or pre-processor) are needed for the transition to current row from the previous one and the required algorithm for that transition is given in brackets. The column “total learned” shows how many algorithms are needed to be learned until the current row. Assumption is that common numeric methods such as normalizing and taking the natural logarithm of numbers is known in advance.

| rank | pre-processor | learner | # of additions | total learned |
|------|---------------|------------|-----------------|---------------|
| 1 | norm | CART (yes) | 1 [CART (yes)] | 1 |
| 2 | norm | CART (no) | 1 [CART (no)] | 2 |
| 3 | none | CART (yes) | 0 | 2 |
| 4 | none | CART (no) | 0 | 2 |
| 5 | log | CART (yes) | 0 | 2 |
| 6 | log | CART (no) | 0 | 2 |
| 7 | SWReg | CART (yes) | 1 [SWR] | 3 |
| 8 | SWReg | CART (no) | 0 | 3 |
| 9 | SFS | CART (yes) | 1 [SFS] | 4 |
| 10 | SFS | CART (no) | 0 | 4 |
| 11 | SWReg | ABE0-1NN | 1 [ABE-kNN] | 5 |
| 12 | log | ABE0-1NN | 0 | 5 |
| 13 | SWReg | ABE0-5NN | 0 | 5 |

top-6 solo methods, we can build the 4th, 5th, 6th, 7th, 9th, 10th best methods of Table 4. Also note that aside from the algorithms of Table 5, a practitioner will most likely require to learn the simple formula of IRWM.

In summary, building successful multi-methods is not necessarily a difficult process for a practitioner. See in Table 5 that building the best method reported here requires learning 5 algorithms (learner or pre-processor) and learning only 2 learners enables someone to build 4 very successful multi-methods. Therefore, our recommendations to practitioners, who are willing to use multi-methods but lack the knowledge of machine learning algorithms are:

- Start with initial 2 learners and build the associated multi-methods
- See the performance of the current multi-methods
- Build new multi-methods only if you are not pleased with the performance of the current ones

The matter of mutual SE knowledge transfer between research and industry is not restricted to the discussion we provided in this section. Practitioners willing to know further on that issue are strongly recommended to read [69].

7.2 Ensemble and Accuracy

Writing in the field of marketing, Armstrong reports in his 2007 study that the multi-method forecasts out-perform single-method forecasts [70]. He cites 31 studies where multiple source prediction consistently out-performs single source prediction by 3.4 to 23.4% (average = 12.5%). It is expected to see that a meaningful combination method, such as ours, results in superior performance. On the other hand, it is impossible to cover all meaningful combination strategies in a single paper. Our study adopts only one of the many proposed successful combination methods [71]. Citing from one of our reviewers: “ Work by Hogarth [72] demonstrates

that a reduced inter-correlation between the prediction sources may be just as important as finding the prediction sources with highest expected accuracy (excluding the poor methods) when combining predictions.” For this research, it means that understanding the inter-dependencies between different solo-methods could lead to clearer definitions of when to combine which particular solo-methods. In return, such a strategy could yield more robust multi-methods. For example a very promising future work (as indicated by the cited reviewer) would be to include the best neural network model into multi-methods.

In the ideal case, different multi-methods would perform optimally under different combination schemes. For example when combining a high number of solo-methods, it would be better to use a scheme that would catch the central tendency and be able to handle extreme values, e.g. median. However, our findings do not support that implication. In other words, our best performing multi-method includes a high number of solo-methods and its combination scheme is IRWM. That discrepancy between implications and results is a familiar concept in forecasting literature: It is difficult to give robust guidelines as to combine the methods in an optimum way [73]. Investigation of robust combination schemes and their implications would be a good future direction to our study.

8 CONCLUSION

There are many effort estimation methods and little guidance on how to choose between them. Prior results from Shepperd et al. are pessimistic about the consistency of the rankings of different methods. If M methods are ranked in an experiment, then that ranking can be changed by altering the conditions of that experiment (e.g. the data used in training or the performance measure used).

This paper confirms the inconsistency of the rankings effect reported by Shepperd et al.: in Figure 3 and Figure 4 it was seen that a method’s ranking can change by some amount δr . Nevertheless, our conclusions are more optimistic than those of Shepperd et al. While some methods have very large δr , others do not. In fact, for the solo-methods shown in Figure 3, the better methods have smaller δr . Better yet, when we combined the highest ranked solo-methods, we found that the top-ranked multi-method had almost zero δr . Therefore, we recommend multi-methods since, as shown in Figure 4:

- The multi-methods consistently out-perform most, if not all, of the solo-methods;
- The performance of the multi-methods are more trustworthy (has the smallest ranking instability).
- Also, as shown in Figure 6, the multi-methods avoid the problem of very large errors seen with other methods.

The success of our results raises the question: Why were prior experiments in ensembles for effort estimation [15], [16] so unsuccessful? Reading from prior effort estimation literature, we offer the following hypothesis: *Many effort estimation methods are inferior* so combinations of inferior members will also be inferior. In Figure 3, for example, while there is evidence for 13 useful methods (in the range $1 \leq X \leq 13$), the

remaining methods have large δr or fall into the worse ranks. That is, of the 90 solo-methods we have explored, there is no evidence for the value of $\frac{77}{90} = 86\%$ of the methods. Without a ranking like Figure 3 to guide method selection, it is therefore probable that researchers will try combining weaker methods.

New methods are being constantly invented (e.g. see [6], [39], [74]) so we make no claim that the 90 methods explored here cover the space of all possible effort estimators. We also make no claim that one study can cover all the pre-processing options there is in the literature and our study is no exception. Our study covers a total of 10 options for pre-processing, however it does not cover the effects of noise or outlier removal. The investigation of the effects of noise/outlier removal on method performance would in fact be an interesting future direction to this study.

Finally, we can offer an answer to the vexing question: *What is the best effort estimator?* Simply put, we have no evidence suggesting that any solo-method is some universal, best in all circumstances, effort estimator. However, if no solo-method is always the best, some ensemble of solo-methods may offer consistently better performance. We have seen that 9 out of 12 multi-methods outperformed all others. Better yet, our top-ranked multi-method has the greatest stability among any of the 102 methods explored in this study. Hence, this study proposes multi-methods as the technique that finds the best learner.

REFERENCES

- [1] B. W. Boehm, *Software Engineering Economics*. Upper Saddle River, NJ, USA: Prentice Hall PTR, 1981.
- [2] A. Albrecht and J. Gaffney, "Software function, source lines of code and development effort prediction: A software science validation," *IEEE Transactions on Software Engineering*, vol. 9, pp. 639–648, 1983.
- [3] L. H. Putnam and W. Myers, *Measures for excellence: reliable software on time, within budget*. Yourdon Press, 1992.
- [4] B. Boehm, C. Abts, and S. Chulani, "Software development cost estimation approaches - a survey," *Annals of Software Engineering*, vol. 10, pp. 177–205, 2000.
- [5] M. Shepperd, C. Schofield, and B. Kitchenham, "Effort estimation using analogy," in *ICSE '96: Proceedings of the 18th international conference on Software engineering*. Washington, DC, USA: IEEE Computer Society, 1996, pp. 170–178.
- [6] A. Corazza, S. Di Martino, F. Ferrucci, C. Gravino, F. Sarro, and E. Mendes, "How effective is tabu search to configure support vector regression for effort estimation?" in *Proceedings of the 6th International Conference on Predictive Models in Software Engineering*, 2010.
- [7] T. Menzies, Z. Chen, J. Hihn, and K. Lum, "Selecting best practices for effort estimation," *IEEE Trans. on Softw. Eng.*, vol. 32, pp. 883–895, 2006.
- [8] M. Jorgensen and M. Shepperd, "A systematic review of software development cost estimation studies," *IEEE Trans. Softw. Eng.*, vol. 33, no. 1, pp. 33–53, 2007.
- [9] M. Shepperd and M. Cartwright, "Predicting with sparse data," *IEEE Trans. Softw. Eng.*, vol. 27, no. 11, pp. 987–998, 2001.
- [10] M. Jorgensen, "A review of studies on expert estimation of software development effort," *Journal of Systems and Software*, vol. 70, pp. 37–60, 2004.
- [11] G. Seni and J. Elder, *Ensemble Methods in Data Mining: Improving Accuracy Through Combining Predictions*. Morgan and Claypool Publishers, 2010.
- [12] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference and prediction*, 2nd ed. Springer, 2008.
- [13] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Proceedings of the 14th international joint conference on Artificial intelligence - Volume 2*, San Francisco, CA, USA, 1995, pp. 1137–1143.
- [14] L. Breiman and P. Spector, "Submodel selection and evaluation in regression. the x-random case," *International Statistical Review / Revue Internationale de Statistique*, vol. 60, no. 3, pp. 291–319, Decemeber 1992.
- [15] E. Kocaguneli, Y. Kultur, and A. Bener, "Combining multiple learners induced on multiple datasets for software effort prediction," in *International Symposium on Software Reliability Engineering (ISSRE)*, 2009, student Paper.
- [16] D. Baker, "A hybrid approach to expert and model-based effort estimation," Master's thesis, Lane Department of Computer Science and Electrical Engineering, West Virginia University, 2007, available from <https://eidr.wvu.edu/etd/documentdata.eTD?documentid=5443>.
- [17] T. M. Khoshgoftaar, P. Rebour, and N. Seliya, "Software quality analysis by combining multiple projects and learners," *Software Quality Control*, vol. 17, no. 1, pp. 25–49, 2009.
- [18] Y. Freund and R. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of Computer and System Sciences*, vol. 55, 1997.
- [19] J. W. Keung, "Theoretical Maximum Prediction Accuracy for Analogy-Based Software Cost Estimation," *2008 15th Asia-Pacific Software Engineering Conference*, pp. 495–502, 2008. [Online]. Available: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4724583>
- [20] I. Myrtevit, E. Stensrud, and M. Shepperd, "Reliability and validity in comparative studies of software prediction models," *IEEE Transactions on Software Engineering*, vol. vol, pp. 31no5pp380–391, May 2005.
- [21] L. Briand and I. Wiecek, *Resource Modeling in Software Engineering*, 2nd ed., ser. Encyclopedia of Software Engineering. Wiley, 2002.
- [22] G. Rowe and G. Wright, "The delphi technique as a forecasting tool: issues and analysis," *International Journal of Forecasting*, vol. 15, 1999.
- [23] K. Atkinson and M. Shepperd, "The use of function points to find cost analogies," in *Proc. European Software Cost Modelling Meeting . Ivrea, Italy*, 1994.
- [24] M. Jorgensen, "Practical guidelines for expert-judgment-based software effort estimation," *IEEE Softw.*, vol. 22, no. 3, pp. 57–63, 2005.
- [25] M. Shepperd and G. F. Kadoda, "Comparing software prediction techniques using simulation," *IEEE Trans. Software Eng*, vol. 27, no. 11, pp. 1014–1022, 2001.
- [26] T. Menzies, O. Jalali, J. Hihn, D. Baker, and K. Lum, "Stable rankings for different effort models," *Automated Software Engineering*, no. 4, December 2010.
- [27] *Techniques for Combining Multiple Learners*. ICSC Press, 1998.
- [28] *Ensemble Methods in Machine Learning*. London, UK: Springer-Verlag, 2000.
- [29] J. Kittler, M. Hatef, R. P. W. Duin, and J. Matas, "On combining classifiers," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, pp. 226–239, 1998.
- [30] K. Ali, "On the link between error correlation and error reduction in decision tree ensembles," Dept. of Information and Computer Science, Univ. of California, Irvine, Tech. Rep., 1995.
- [31] T. K. Ho, J. J. Hull, and S. N. Srihari, "Decision combination in multiple classifier systems," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 16, no. 1, pp. 66–75, 1994.
- [32] S. Günter and H. Bunke, "Feature selection algorithms for the generation of multiple classifier systems and their application to handwritten word recognition," *Pattern Recogn. Lett.*, vol. 25, no. 11, pp. 1323–1336, 2004.
- [33] T. K. Ho, "Random decision forests," in *ICDAR '95: Proceedings of the Third International Conference on Document Analysis and Recognition (Volume 1)*. Washington, DC, USA: IEEE Computer Society, 1995, p. 278.
- [34] H. Zhao and S. Ram, "Constrained cascade generalization of decision trees," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, pp. 727–739, 2004.
- [35] I. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementation*. Morgan Kaufmann, 2000.
- [36] E. Bauer and R. Kohavi, "An empirical comparison of voting classification algorithms: Bagging, boosting, and variants," *Machine Learning*, vol. 36, pp. 105–139, 1999, 10.1023/A:1007515423169. [Online]. Available: <http://dx.doi.org/10.1023/A:1007515423169>
- [37] M. C. K. Vinaykumar, V. Ravi, "Software cost estimation using soft computing approaches," *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*, pp. 499–518, 2010.
- [38] *Software cost estimation using computational intelligence techniques*, Dec. 2009.
- [39] Y. Kultur, B. Turhan, and A. B. Bener, "Enna: software effort estimation using ensemble of neural networks with associative memory," in *SIGSOFT '08/FSE-16: Proceedings of the 16th ACM SIGSOFT*

- International Symposium on Foundations of software engineering*, New York, NY, USA, 2008, pp. 330–338.
- [40] B. Twala, M. Cartwright, and M. Shepperd, “Ensemble of missing data techniques to improve software prediction accuracy,” in *ICSE '06: Proceedings of the 28th international Conference on Software engineering*. New York, NY, USA: ACM, 2006, pp. 909–912. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1239046.1239048>
- [41] T. M. Khoshgoftaar, S. Zhong, and V. Joshi, “Enhancing software quality estimation using ensemble-classifier based noise filtering,” *Intell. Data Anal.*, vol. 9, pp. 3–27, January 2005. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1239046.1239048>
- [42] K. Lum, T. Menzies, and D. Baker, “2cee, a twenty first century effort estimation methodology,” *ISPA / SCEA*, pp. 12 – 14, 2008.
- [43] E. Mendes, I. D. Watson, C. Triggs, N. Mosley, and S. Counsell, “A comparative study of cost estimation models for web hypermedia applications,” *Empirical Software Engineering*, vol. 8, no. 2, pp. 163–196, 2003.
- [44] M. Shepperd and C. Schofield, “Estimating software project effort using analogies,” *IEEE Trans. Softw. Eng.*, vol. 23, no. 11, pp. 736–743, 1997.
- [45] C. Chang, “Finding prototypes for nearest neighbor classifiers,” *IEEE Trans. on Computers*, pp. 1179–1185, 1974.
- [46] A. Venkatachalam, “Software cost estimation using artificial neural networks,” in *Proceedings of international joint conference on neural networks*, 1993, pp. 987–990.
- [47] J. Ghosh, “Multiclassifier systems: Back to the future,” in *MCS '02: Proceedings of the Third International Workshop on Multiple Classifier Systems*. London, UK: Springer-Verlag, 2002, pp. 1–15.
- [48] H. Park and S. Baek, “An empirical validation of a neural network model for software effort estimation,” *Expert Syst. Appl.*, vol. 35, no. 3, pp. 929–937, 2008.
- [49] K. Pillai and V. Sukumaran Nair, “A model for software development effort and cost estimation,” *IEEE Trans. Softw. Eng.*, vol. 23, no. 8, pp. 485–497, 1997.
- [50] E. Mendes and N. Mosley, “Further investigation into the use of cbr and stepwise regression to predict development effort for web hypermedia applications,” in *International Symposium on Empirical Software Engineering*. Washington, DC, USA: IEEE Computer Society, 2002.
- [51] L. Breiman, “Bagging predictors,” *Machine Learning*, vol. 24, pp. 123–140, 1996.
- [52] Y. Jiang, B. Cukic, and T. Menzies, “Cost curve evaluation of fault prediction models,” in *ISSRE '08: Proceedings of the 2008 19th International Symposium on Software Reliability Engineering*, 2008, pp. 197–206.
- [53] L. Breiman, “Random forests,” *Machine Learning*, vol. 45, pp. 5–32, 2001.
- [54] T. Foss, E. Stensrud, B. Kitchenham, and I. Myrvtveit, “A simulation study of the model evaluation criterion mmre,” *IEEE Trans. on Softw. Eng.*, 2003.
- [55] B. Kitchenham and E. Mendes, “Why comparative effort prediction studies may be invalid,” in *PROMISE '09: Proceedings of the 5th International Conference on Predictor Models in Software Engineering*. New York, NY, USA: ACM, 2009, pp. 1–5.
- [56] J. Klijnen, “Sensitivity analysis and related analyses: a survey of statistical techniques,” *Journal Statistical Computation and Simulation*, vol. 57, no. 1–4, pp. 111–142, 1997.
- [57] J. Demsar, “Statistical comparisons of classifiers over multiple data sets,” *Journal of Machine Learning Research*, vol. 7, 2006.
- [58] E. Alpaydin, *Introduction to Machine Learning*. MIT Press, 2004.
- [59] A. Oppenheim, A. Wilsky, and S. Hamid, *Signals and Systems*. Prentice Hall, 1996.
- [60] J. Keung, E. Kocaguneli, and T. Menzies, “A Ranking Stability Indicator for Selecting the Best Effort Estimator in Software Cost Estimation,” *Automated Software Engineering (submitted, under the first round of reviews)*, 2011. [Online]. Available: <http://menzies.us/pdf/11drafranking.pdf>
- [61] B. Kitchenham, E. Mendes, and G. H. Travassos, “Cross versus within-company cost estimation studies: A systematic review,” *IEEE Trans. Softw. Eng.*, vol. 33, no. 5, pp. 316–329, 2007.
- [62] D. Milic and C. Wohlin, “Distribution patterns of effort estimations,” in *Euromicro*, 2004.
- [63] C. Robson, “Real world research: a resource for social scientists and practitioner-researchers,” *Blackwell Publisher Ltd*, 2002.
- [64] E. Kocaguneli and T. Menzies, “The effects of test set selection on effort estimation (in preparation),” 2011.
- [65] A. Nelson, T. Menzies, and G. Gay, “Sharing experiments using open-source software,” *Software: Practice and Experience*, vol. 41, no. 3, pp. 283–305, 2011.
- [66] A. Bakir, E. Kocaguneli, A. Tosun, A. Bener, and B. Turhan, “Xiruxe: An Intelligent Fault Tracking Tool,” in *AIPRO9*, Orlando, 2009.
- [67] A. Tosun, A. Bener, and E. Kocaguneli, “BITS: Issue Tracking and Project Management Tool in Healthcare Software Development,” in *SEKE'09*, Boston, MA, USA, 2009.
- [68] E. Kocaguneli, A. Tosun, A. Bener, B. Caglayan, and B. Turhan, “Prest: An Intelligent Software Metrics Extraction, Analysis and Defect Prediction Tool,” in *SEKE'09*, Boston, 2009.
- [69] T. Menzies, C. Bird, T. Zimmermann, W. Schulte, and E. Kocaguneli, “The inductive software engineering manifesto: Principles for industrial data mining,” in *MALETS'11: International Workshop on Machine Learning Technologies in Software Engineering*, 2011.
- [70] J. S. Armstrong, “Significance tests harm progress in forecasting,” *International Journal of Forecasting*, vol. 23, no. 2, pp. 321–327, 2007.
- [71] —, *Principles of Forecasting: A Handbook for Researchers and Practitioners*. Kluwer Academic, 2001.
- [72] R. M. Hogarth, “A note on aggregating opinions,” *Organizational Behavior and Human Performance*, vol. 21, no. 1, pp. 40 – 46, 1978.
- [73] R. M. Hogarth and H. Kunreuther, “Pricing insurance and warranties: Ambiguity and correlated risks,” *The GENEVA Papers on Risk and Insurance - Theory*, vol. 17, no. 1, pp. 35–60, 1992.
- [74] Y. Li, M. Xie, and T. Goh, “A study of project selection and feature weighting for analogy based software cost estimation,” *Journal of Systems and Software*, vol. 82, pp. 241–252, 2009.
- [75] A. Bakir, B. Turhan, and A. Bener, “A new perspective on data homogeneity in software cost estimation: A study in the embedded systems domain,” *Software Quality Journal*, 2009. [Online]. Available: <http://dx.doi.org/10.1007/s11219-009-9081-z>
- [76] B. Kitchenham and K. Känsälä, “Inter-item correlations among function points,” in *ICSE'93: Proceedings of the 15th international Conference on Software Engineering*, ser. ICSE '93. Los Alamitos, CA, USA: IEEE Computer Society Press, 1993, pp. 477–480. [Online]. Available: <http://dl.acm.org/citation.cfm?id=257572.257677>
- [77] C. Kemerer, “An empirical validation of software cost estimation models,” *Communications of the ACM*, vol. 30, no. 5, pp. 416–429, May 1987.
- [78] K. D. Maxwell, *Applied Statistics for Software Managers*. Prentice Hall, PTR, 2002.
- [79] Y. Miyazaki, M. Terakado, K. Ozaki, and H. Nozaki, “Robust regression for developing software estimation models,” *Journal of Systems and Software*, vol. 27, no. 1, pp. 3–16, 1994.
- [80] G. Kadoda, M. Cartwright, and M. Shepperd, “On configuring a case-based reasoning software project prediction system,” *UK CBR Workshop, Cambridge, UK*, pp. 1–10, 2000.
- [81] L. Breiman, J. Friedman, R. Olshen, and C. Stone, *Classification and Regression Trees*. Monterey, CA: Wadsworth and Brooks, 1984.
- [82] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural Netw.*, vol. 2, no. 5, pp. 359–366, 1989.

APPENDICES

Note: the content of the following appendices come from [60]. Although the work in [60] and this manuscript use the same base learners and pre-processing options, they address two very different research ideas. In [60] our concern is an investigation of the common SEE methods in retrospect, i.e. we re-evaluate the past work in SEE and seek for stable conclusions in methods and datasets. We use the results of this investigation to comment on the methods and on the datasets used extensively in SEE.

On the other hand, in this paper we investigate new and robust estimation methods, i.e. the ensembles. This research is built around a similar code base (the implementation of learners and pre-processing options), yet it uses the knowledge we got by observing old methods to propose a completely new direction of building a novel scheme: *multi-methods*.

Appendix A: Data

All the data used in this study is available either at <http://promisedata.org/data> or through the authors. As shown

in Table 2, we use a variety of different data sets in this research. The standard **COCOMO** data sets (cocomo*, nasa*), which are collected with the COCOMO approach [1]. The **desharnais** data set, which contains software projects from Canada. It is collected with function points approach.

SDR, which contains data from projects of various software companies in Turkey. SDR is collected by Softlab, the Bogazici University Software Engineering Research Laboratory [75].

albrecht data set consists of projects completed in IBM in the 1970's and details are given in [2].

finnish data set originally contains 40 projects from different companies and data were collected by a single person. The two projects with missing values are omitted here, hence we use 38 instances. More details can be found in [76].

kemerer is a relatively small dataset with 15 instances, whose details can be found in [77].

maxwell data set comes from finance domain and is composed of Finnish banking software projects. Details of this dataset are given in [78].

miyazaki data set contains projects developed in COBOL. For details see [79].

telecom contains projects which are enhancements to a U.K. telecommunication product and details are provided in [44].

china dataset includes various software projects from multiple companies developed in China.

Note that two of these data sets (Nasa93c2, Nasa93c5) come from different development centers around the United States. Another two of these data sets (Cocomo81e, Cocomo81o) represent different kinds of projects. The Cocomo81e “embedded projects” are those developed within tight constraints (hardware, software, operational, ...); The Cocomo81o “organic projects” come from small teams with good experience working with less than rigid requirements.

Note also in Table 2, the skewness of our effort values (up to 6.06): Our datasets are extremely heterogeneous with as much as 60-fold variation. There is also some divergence in the features used to describe our data. While our data includes some effort value (measured in terms of months or hours), no other feature is shared by all data sets.

Appendix B: Pre-Processing Options

Simple numeric techniques: This category lists the pre-processors that entail mere numeric alterations of actual values rather than application of certain algorithms.

- **norm:** “norm” represents the application of normalization on the data. The data is normalized to 0-1 interval according to Equation 9:
- **log:** The natural logarithm of the independent variables are taken. Similar to all pre-processing procedures, the dependent variable is excluded from the pre-processing method.

$$normalizedValue = \frac{(actualValue - min(allValues))}{(max(allValues) - min(allValues))} \quad (9)$$

Feature synthesis: “PCA” stands for principal component analysis [58]. PCA is widely used as a dimension alteration

mechanism. The alteration in PCA occurs in the form of changing an n -dimensional space into another n -dimensional space. PCA lets the user know which particular features are most influential in the new n -dimensional space, hence user can choose to use only the most influential features. Therefore, it is also common to see the usage of PCA like a dimensionality reduction mechanism. However, in our study we prefer to use PCA only to change the n -dimensional space (not the number of features).

Feature selection: Some of the pre-processors aim at finding a subset of all features according to certain criteria. This subset is supposed to produce a feature subset that will ultimately increase the estimation accuracy. Under this category, we have SFS and SWREg.

Sequential forward selection (“SFS”) is a feature selection mechanism, in which features are added into an initially empty set until no improvement is possible with the addition of another feature. The so called *improvement* is defined through an *objective* function. In our implementation based on MATLAB, the *objective* function is the mean-squared-error of a simple linear regression on the training set. One caution to be made here is that exhaustive search algorithms over all features can be very time consuming (2^n combinations in an n -feature dataset), therefore SFS works only in forward direction (no backtracking).

“SWReg” stands for stepwise regression, which can be defined as a systematic method for adding and removing features from a multi-linear model. Removal and addition of features depend on their statistical significance in regression. Our SWReg implementation that is developed by using MATLAB starts execution with a preliminary model, then it compares the performances of smaller and larger models. At each step a potential feature is to be decided for addition or removal. Addition and removal are based on the p-value in an F-statistic. In a particular step, the F-statistics for two models (models with and without the feature that is being questioned in that step) are calculated. Provided that the feature was not in the model, the null hypothesis is: “Feature would have a zero coefficient in the model, when it is added”. If the null hypothesis can be rejected, then the feature is added to the model. As for the other scenario (i.e. feature is already in the model), the null hypothesis is: “Feature has a zero coefficient”. If we fail to reject the null hypothesis, then the term is removed.

Note that stepwise methods are locally optimal and may not necessarily be globally optimal. In other words, depending on the initial model and the order of features for inclusion-exclusion, the algorithm may result in different final models, hence different performances.

Discretization:

- **width3bin:** This procedure bins each one of the data features into 3 bins, depending on equal width of all bins. The bin-width for a general n -bin procedure is given in Equation 10. In our 3-bin case, once we know the bin-width, we assign each feature value to either 1, 2 or 3, depending on which particular bin the value is in.

$$binWidth = ceiling\left(\frac{(max(allValues) - min(allValues))}{n}\right) \quad (10)$$

- **width5bin:** Exactly same as “width3bin” except that this time we have 5 bins instead of 3.
- **freq3bin:** “freq3bin” means binning each feature into 3 bins, depending on the equal frequency count (equal number of instances in each bin). Similar to previously mentioned binning methods, in this method each feature value is assigned to 1, 2 or 3. For that, all values of a particular feature is sorted in ascending order. Then first $\lceil \text{numberOf}(\text{allInstances})/3 \rceil$ instances are assigned 1, the second $\lceil \text{numberOf}(\text{allInstances})/3 \rceil$ instances are assigned 2 and so on.
- **freq5bin:** Same as “freq3bin”, but with 5 bins.

Others: The option(s) that cannot be categorized into aforementioned categories are mentioned here. Application of no pre-processing to the data is also a choice. “none” represents the choice of no pre-processor selection.

Appendix C: Learners

Instance-based learners: When we refer to instance-based learners (either ABE methods or nearest-neighbor based methods), there are various design options to be decided. The analogy based estimation method that we call ABE0- x NN refers to a very basic type of ABE that we defined through our readings of various ABE studies [43], [74], [80]. In ABE0- x NN, features are firstly normalized to 0-1 interval, then the distance between test and train instances is measured according to Euclidean distance function, x nearest neighbors are chosen from the training set and finally for finding estimated value (a.k.a adaptation procedure) the median of x nearest neighbors is calculated. Therefore, when we say ABE0- x NN, all the design options including the choice of x -many closest analogies become explicit to a reader. In our experimentation we use two different x values (i.e. two different analogy values):

- **ABE0-1NN:** Only the closest analogy is used. Since the median of a single value is itself, the estimated value in ABE0-1NN becomes the actual effort value of the closest analogy in the training set.
- **ABE0-5NN:** Five closest analogies are found and used for adaptation.

Iterative dichotomization: Under this category, we use classification and regression trees (CART) as described by Breiman et al. [81] and developed it using MATLAB routines. CART is a non-parametric technique and it can work both for classification and regression type of problems, in our case it is the latter. When planning to construct a CART, there are a couple of points to consider. Firstly, the selection of splits: There are a variety of solutions such as misclassification rate, information (or entropy) or GINI index. CART uses the GINI index to calculate the impurity of a tree [81]. Secondly, decision on when to stop further splits: The implementation allows you to specify a threshold value or use the default value. We have used default threshold values in the implementation (for further details please refer to MATLAB documentation). Lastly, assigning class to each terminal node. Each test instance results in a terminal node and therefore different test

instances resulting in the same terminal node are given the same predicted value. For regression, the predicted value in a terminal node is a fit on the independent variables of the instances in that node.

We have used 2 types of CART, which are given below:

- **CART (yes):** In this version, the pruning is on, meaning that the training data is used in a cross validation process and for each cross validation run, some internal nodes are randomly chosen as the leaf nodes and their sub-nodes are removed. Finally, the sub-tree that resulted in the lowest error rate is returned. The returned sub-tree is a sub-optimal solution and it is locally optimum.
- **CART (no):** The sub-trees of CART will not be considered and the full tree will be used.

Two-layered neural net: “Neural Nets” (NNet) are memoryless structures that can be defined as universal approximators [82], meaning that: 1) They store information coming from training data as weights during training and after that phase they do not need training data and 2) they can approximate any function depending on how complex the NNet structure is. NNets basically have a 3 layer structure: Input layer, hidden layer and the output layer. Depending on the complexity of the problem, one can model more complex functions with a NNet by increasing the number of hidden layers in the structure. In our implementation, we used a simple NNet structure with 2 hidden layers, 1 input layer and 1 output layer.

Regression methods: Under this category, we list our regression-based learners. **LReg** stands for linear regression. Given the dependent variables, this learner calculates the coefficient estimates of the independent variables.

Partial Least Squares Regression (**PLSR**) as well as Principal Components Regression (**PCR**) are algorithms that are used to model a dependent variable and we have used MATLAB routines to implement those functions in our experiments. While modeling an independent variable, they both construct new independent variables as linear combinations of original independent variables. However, the ways they construct the new independent variables are different:

- **PCR:** This method generates new independent variables to explain the observed variability in the actual ones. However, while generating new variables the dependent variable is not considered at all. In that respect, PCR is similar to selection of n -many components via PCA (the default value of components to select is 2, so we used it that way) and applying linear regression.
- **PLSR:** Unlike PCR, PLSR considers the independent variable and picks up the n -many of the new components (again with a default value of 2) that yield lowest error rate. Due to this particular property of PLSR, it usually results in a better fitting.

SWReg: We have previously defined SWReg and mentioned that it was a special regression algorithm. We have also stated how it can be used as a feature selection algorithm. In the algorithms section, we use SWReg as a regression method on the *selected-out* independent variables to explain the dependent variable.