

Where is the Best Effort Estimator?

Vincent Rogers
CS, WVU, Morgantown, USA
email@email.com

William Sica
CS, WVU, Morgantown, USA
filler@filler.fill

ABSTRACT

BACKGROUND: Despite decades of research, there is no agreement on what is the “best” software effort estimator.

AIM: We seek stable conclusions about the “best” effort estimator.

METHOD: 10 learners and 9 data preprocessors, combined into 90 effort estimators, were applied to twenty datasets. Performance was assessed using AR, MRE, MER, MdMRE, MMRE, PRED(25), MIBRE and compared using a Wilcoxon test (95%).

RESULTS: Some datasets are *stronger/weaker* at distinguishing estimation performance and, hence, are better/worse at finding the “best” estimators.

CONCLUSIONS: We recommend regression trees or analogy methods (and recommend against neural nets or simple linear regression). We caution that prior work has not taken proper account of the data used to evaluate the effort estimator. Future work should take greater care to explore the data and its preprocessing (and not just the learner). Furthermore, that future work should justify its conclusions using the *stronger* datasets.

Categories and Subject Descriptors

Software Engineering [**Software Metrics**]: Data Mining

General Terms

Software effort Estimation, Linear Regression, Regression Trees, Neural Nets, Analogy, MMRE, Evaluation Criteria

Keywords

Software Effort Estimation, MMRE, Evaluation Criteria

1. INTRODUCTION

For decades, researchers have been seeking the “best” software development effort estimator. To date, no such “best” estimator has been found. The usual conclusion is that effort estimation suffers from a *conclusion instability* problem; i.e. different researchers offer conflicting rankings as to what is “best” [?, ?].

This is an open and urgent issue since accurate effort estimation is vital to Software Engineering, and is often a challenging task

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICSE 2011 Waikiki, Honolulu, Hawaii USA

Copyright 2011 ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

for many software project managers. Both overestimating and underestimating would result unfavorable impacts to the business’s competitiveness and project resource planning.

This is going to be included in the file Effort estimation research focuses on the *learner* used to generate the estimate (e.g. linear regression, neural nets, etc) in many cases, overlooking the importance of the quality and characteristics of the *data* being used in the estimation process. We will argue that this approach is somewhat misguided since, as shown below, learner performance is greatly influenced by the data preprocessing and the datasets being used to evaluate the learner.

This paper shows that many datasets used by prior publications are very *weak* at distinguishing between different learners. All the results based on these *weak* datasets (including many results by the authors of this papers) must hence be revisited.

The good news is that there are *strong* datasets that clearly illustrate the value of any estimator. In all, this study applies 90 estimators to 20 datasets and measures their performance in terms of seven different performance criteria. To the best of our knowledge, this is the largest effort estimation study yet reported in the literature. One result of exploring such a large space of data and algorithms is that we can report stable conclusions (while prior studies have not).

This paper is structured as follows. Related work discusses effort estimation and the prior reports on *conclusion instability*. Those reports used a dataset to *seed* the generation of artificial data. Our results section shows that if we extend the experiments to a broader set of project data, we are able to discover stable conclusions such as that some datasets are *weak*; and that the *strong* datasets show which estimators are consistently better than others. Based on those results, our conclusion will list best (and worst) effort estimators.

2. RELATED WORK

This section reviewed the effort estimation literature with regards to (a) the major estimation techniques used by empirical research studies on cost estimation within the last 15 years and (b) the conclusion instability problem.

2.1 Algorithmic Methods

There are many algorithmic effort estimators. For example, if we restrict ourselves to just instance-based algorithms, Figure ?? shows that there are thousands of options just in that one sub-field.

As to non-instance methods, that are many proposed in the literature including various kinds of regression (simple, partial least square, stepwise, regression trees), and neural networks just to name a few. For notes on these non-instance methods, see §??.

Note that instance & non-instance-based methods can be combined to create even more algorithms. For example, once an instance-based method finds its nearest neighbors, those neighbors might be

Instances-based learners draw conclusions from instances *near* the test instance. Mendes et al. [?] discuss various *near*-ness measures.

- M_1 : A simple Euclidean measure;
- M_2 : A “maximum distance” measure that focuses on the single feature that maximizes inter-project distance.
- M_3 : More elaborate kernel estimation methods.

Once the nearest neighbors are found, they must be used to generate an effort estimate via...

- R_1 : Reporting the median effort value of the analogies;
- R_2 : Reporting the mean dependent value;
- R_3 : Reporting a weighted mean where the nearer analogies are weighted higher than those further away [?];

Prior to running an instance-based learning, it is sometimes recommended to handle anomalous rows by:

- N_1 : Doing nothing at all;
- N_2 : Using outlier removal [?];
- N_3 : *Prototype generation*; i.e. replace the data set with a smaller set of most representative examples [?].

When computing distances between pairs, some feature weighting scheme is often applied:

- W_1 : All features have uniform weights;
- $W_2..W_9$: Some pre-processing scores the relative value of the features using various methods [?, ?, ?]. The pre-processors may require *discretization*.

Discretization breaks up continuous ranges at points b_1, b_2, \dots , each containing counts of c_1, c_2, \dots of numbers [?]. Discretization methods include:

- D_1 : Equal-frequency, where $c_i = c_j$;
- D_2 : Equal-width, where $b_{i+1} - b_i$ is a constant;
- D_3 : Entropy [?];
- D_4 : PKID [?];
- D_5 : Do nothing at all.

Finally, there is the issue of how many k neighbors should be used:

- K_1 : $k = 1$ is used by Lipowezky et al. [?] and Walkerden & Jeffery [?];
- K_2 : $k = 2$ is used by Kirsopp & Shepperd [?];
- K_3 : $k = 1, 2, 3$ is used by Mendes et al. [?];
- K_4 : Li et al. use $k = 5$ [?];
- K_5 : Baker tuned k to a particular training set using an experimental method [?].

Figure 1: Each combination of the above $N \times W \times D \times M \times R \times K$ techniques is one *algorithm* for instance-based effort estimation. This figure shows $3 \times 3 \times 3 \times 9 \times 5 \times 5 > 6,000$ algorithms for effort estimation. Some of these ways can be ruled out, straight away. For example, at $k = 1$, then all the adaptation mechanisms return the same result. Also, not all the feature weighting techniques require discretization, decreasing the space of options by a factor of five. However, even after discarding some combinations, there are still hundreds to thousands of algorithms to explore.

summarized with regression or neural nets [?].

2.2 Non-Algorithmic Methods

An alternative approach to algorithmic approaches (e.g. the instance-based methods of Figure ??) is to utilize the best knowledge of an experienced expert. Expert based estimation [?] is a human intensive approach that is most commonly adopted in practice. Estimates are usually produced by a domain expert rather than an estimation expert based on their very own personal experience and recollection of similar past projects in the organization. It is flexible and intuitive in a sense that it can be applied in a variety of circumstances where other estimating techniques do not work (for example when there is a lack of historical data). Furthermore in many cases requirements are simply unavailable at the bidding stage of a project where a rough estimate is required in a very short period of time.

Jorgensen [?] provides guidelines for producing realistic software development effort estimates derived from industrial experience and empirical studies. One important finding concluded was that the *combine estimation* method in expert based estimation offers the most robust and accurate combination method, as combining estimates captures a broader range of information that is relevant to the target problem, for example combining estimates with analogy based with expert based method. Data and knowledge relevance to the project’s context and characteristics will more likely to influence the prediction accuracy.

Although widely used in industry, there are no standard methods for expert based estimation. Shepperd et al. [?] do not consider expert based estimation an empirical method because the means of deriving an estimate are not explicit and therefore not repeatable, nor easily transferable to other staff. In addition, knowledge relevancy is also a problem, as an expert may not be able to justify estimates for a new application domain as well as its validity. Hence, the rest of this paper does not consider non-algorithmic methods.

2.3 Conclusion Instability

To derive stable conclusions about which estimator is “best”, there have been attempts in trying to compare model prediction performance of different effort estimation approaches. For example, Shepperd and Kododa [?] compared regression, rule induction, nearest neighbor and neural nets, in an attempt to explore the relationship between accuracy, choice of prediction system, and different dataset characteristic by using a simulation study based on artificial datasets. They also reported a number of conflicting results exist in the literature as to which method provides better prediction accuracy, and offers possible explanations including the use of an evaluation criteria such as MMRE and the underlying characteristics of the problem dataset being used will have a strong influence upon the relative effectiveness of different prediction models. Their work as a *simulation study* that took a single dataset, then generated very large artificial datasets using the distributions seen on that data. They concluded that:

Dataset	Features	Size	Description	Historical Effort Data					
				Units	Min	Median	Mean	Max	Skewness
cocomo81	17	63	NASA projects	months	6	98	683	11400	4.4
cocomo81e	17	28	Cocomo81 embedded projects	months	9	354	1153	11400	3.4
cocomo81o	17	24	Cocomo81 organic projects	months	6	46	60	240	1.7
cocomo81s	17	11	Cocomo81 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									

Figure 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 this datasets, see the appendix.

- *None* of these existing estimators were consistently “best”;
- The accuracy of an estimate depends on the dataset characteristic and a suitable prediction model for the dataset.

They conclude that it is generally *infeasible* to determine which prediction technique is “best”.

Recent results suggest that it is appropriate to revisit the conclusion instability hypothesis. Menzies et al. [?] applied 158 estimators to various subsets of two COCOMO datasets. In a result consistent with Shepperd and Kododa, they found the precise ranking of the 158 estimators changed according to the random number seeds used to generate train/test sets; the performance evaluation criteria used; and which subset of the data was used. However, they also found that four methods consistently out-performed the other 154 across all datasets, across 5 different random number seeds, and across three different evaluation criteria.

Also, there are now many more public domain datasets, readily available for stability studies. Figure ?? lists 20 datasets which have become available in the last year at the PROMISE repository of reusable SE data¹. Given the availability of this data, it is no longer necessary to work on simulated data (as done by Shepperd and Kadoda [?]) or to study merely two datasets (as done by Menzies et al. [?]). The rest of this paper explores conclusion stability over 20 datasets given in Figure ??.

3. EXPERIMENT DESIGN

In our experiments, numerous performance measures were collected after various *algorithms* (combinations of preprocessors and learners) were applied to the data of Figure ?. This section describes those performance measures, preprocessors, and learners.

Since it is impractical to explore (say) the thousands of options described in Figure ?, we elected to explore variants commonly used in the literature. For example, we explore neural nets, regression, and analogy because those methods were explored by Shepherd and Kododa [?]. Nevertheless, it is important to note that our conclusions come only from the estimators/performance criteria/datasets used in this study. Further work is required to confirm our findings on other estimators/performance criteria/datasets.

3.1 Performance Measures

¹<http://promisedata.org/data>

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

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

(where x_i, \hat{x}_i are the actual and predicted value for test instance i).

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 [?] [?]. MRE measures the error ratio between the actual effort and the predicted effort and 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 [?]):

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

The overall average error of MRE can be derived as the Mean or Median Magnitude of Relative Error measure (MMRE, or MdmMRE respectively), can be calculated as:

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

$$MdmMRE = median(allMRE_i) \quad (5)$$

A common alternative to MMRE is PRED(25), and defined as the percentage of predictions failing within 25% of the actual values, and can be expressed as:

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

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

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

$$MIBRE_i = \frac{\hat{x}_i - x_i}{max(\hat{x}_i, x_i)} \quad (7)$$

3.2 Ten Pre-processors

In this study, we investigate:

- Three *simple preprocessors*: **none**, **norm**, and **log**;
- One *feature synthesis* methods called **PCA**;
- Two *feature selection* methods: **SFS** (sequential forward selection) and **SWreg**;
- Four *discretization* methods: divided on equal frequency/width.

None is the simplest preprocessor- all values are unchanged.

With the **norm** preprocessor, numeric values are normalized to a 0-1 interval using Equation ???. Normalization means that no variable has a greater influence than any other.

$$\text{normalizedValue} = \frac{(\text{actualValue} - \min(\text{allValues}))}{(\max(\text{allValues}) - \min(\text{allValues}))} \quad (8)$$

With the **log** preprocessor, all numerics are replaced with their logarithm. This **logging** procedure minimizes the effects of the occasional very large numeric value.

Principal component analysis [?], or **PCA**, is a *feature synthesis* preprocessor that converts a number of possibly correlated variables into a smaller number of uncorrelated variables called components. The first component accounts for as much of the variability in the data as possible, and each succeeding component accounts for as much of the remaining variability as possible.

Some of the preprocessors aim at finding a subset of all features according to certain criteria such as **SFS** (sequential forward selection) and **SWR** (stepwise regression). **SFS** adds features into an initially empty set until no improvement is possible with the addition of another feature. When ever the selected feature set is enlarged, some oracle is called to assess the value of that set of features. In this study, we used the MATLAB, *objective* function (which reports the 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).

SWR adds and removes features from a multilinear model. Addition and removal is controlled by the p-value in an F-Statistic. At each step, the F-statistics for two models (models with/out one feature) 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 is 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.

Discretizers are pre-processors that maps every numeric value in a column of data into a small number of discrete values:

- **width3bin**: This procedure clumps the data features into 3 bins, depending on equal width of all bins see Equation ???.

$$\text{binWidth} = \text{ceiling} \left(\frac{\max(\text{allValues}) - \min(\text{allValues})}{n} \right) \quad (9)$$

- **width5bin**: Same as **width3bin** except we use 5 bins.
- **freq3bin**: Generates 3 bins of equal population size;
- **freq5bin**: Same as **freq3bin**, only this time we have 5 bins.

3.3 Nine Learners

Based on our reading of the effort estimation literature, we identified nine commonly used learners that divide into

- Two *instance-based* learners: **ABE0-1NN**, **ABE0-5NN**;

- Two *iterative dichotomizers*: **CART(yes)**, **CART(no)**;
- A *neural net*: **NNet**;
- Four *regression methods*: **LReg**, **PCR**, **PLSR**, **SWReg**.

Instance-based learning can be used for analog-based estimation. A large class of ABE algorithms was described in Figure ???. Since it is not practical to experiment with the 6000 options defined in Figure ???, we focus on two standard variants. ABE0 is our name for a very basic type of ABE that we derived from various ABE studies [?, ?, ?]. In **ABE0-xNN**, 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 training set and finally for finding estimated value (a.k.a adaptation procedure) the median of x nearest neighbors is calculated. We explored two different x :

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

Iterative Dichotomizers seek the best attribute value *splitter* that most simplifies the data that fall into the different splits. Each such splitter becomes a root of a tree. Sub-trees are generated by calling iterative dichotomization recursively on each of the splits. The CART iterative dichotomizer [?] is defined for continuous target concepts and its *splitters* strive to reduce the GINI index of the data that falls into each split. In this study, we use two variants:

- **CART (yes)**: This version prunes the generated tree using cross-validation. For each cross-val, an internal nodes is made into a leaf (thus pruning its sub-nodes). The sub-tree that resulted in the lowest error rate is returned.
- **CART (no)**: Uses the full tree (no pruning).

In *Neural Nets*, or **NNet**, an input layer of project details is connected to zero or more "hidden" layers which then connect to an output node (the effort prediction). The connections are weighted. If the signal arriving to a node sums to more than some threshold, the node "fires" and a weight is propagated across the network. Learning in a neural net compares the output value to the expected value, then applies some correction method to improve the edge weights (e.g. back propagation). Our **NNet** uses three layers.

This study also uses four *regression methods*. **LReg** is a simple linear regression algorithm. Given the dependent variables, this learner calculates the coefficient estimates of the independent variables. **SWreg** is the stepwise regression discussed above. Whereas above, **SWreg** was used to select features for other learners, here we use **SWreg** as a learner (that is, the predicted value is a regression result using the features selected by the last step of **SWreg**). Partial Least Squares Regression (**PLSR**) as well as Principal Components Regression (**PCR**) are algorithms that are used to model a dependent variable. 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** generates new independent variables to explain the observed variability in the actual ones. 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**, on the other hand, 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.

3.4 Experimental Rig

This study copied the experimental rig of a recent prominent study [?]. In their leave-one-out experiment, given T projects, then $\forall t \in T$, t is the test and the remaining $T - 1$ projects are used for training. The resulting $T - 1$ predictions are then used to compute our seven evaluation criteria given in Section 3.1.

To compare the performance of one algorithm versus the rest, we used a Wilcoxon non-parametric statistical hypothesis test. 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 presence of outliers may be existed in the given datasets. Ranked statistical tests like the Wilcoxon are also useful if it is not clear that the underlying distributions are Gaussian [?].

Using the Wilcoxon test, for each dataset, the performance measures collected from each of our 90 algorithms was compared to the 89 others. This allowed us to collect *win-tie-loss* statistics using the algorithm of Figure ?? . First, we first to check if two distributions i, j are statistically different according to the Wilcoxon test (95% confident); otherwise 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 median values.

```

if WILCOXON( $P_i, P_j, 95$ ) says they are the same then
     $tie_i = tie_i + 1$ ;
     $tie_j = tie_j + 1$ ;
else
    if better( median( $P_i$ ), median( $P_j$ )) 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

```

Figure 3: Comparing algorithms (i, j) on performance (P_i, P_j). The “better” predicate changes according to P . For error measures like MRE, “better” means lower medians. However, for PRED(25), “better” means higher medians.

4. RESULTS

After applying leave-one-out to all 20 data sets, the procedure of Figure ?? was repeated seven times (once for AR, MRE, MER, MdmRE, MMRE, PRED(25), and MIBRE). Our ninety algorithms were then sorted by their total number of losses over all datasets. The resulting sort order is shown in Figure ?? . The algorithm, with fewest losses (**norm/CART(yes)**) was ranked #1 and the algorithm with the most losses (**PCA/LReg**) was ranked #90.

Given 89 comparisons and seven performance measures and 20 datasets, the maximum number of losses for any algorithm was $89 \times 7 \times 20 = 12,460$. Figure ?? sorts all 90 algorithms according to their total losses seen in all seven performance criteria (expressed as a percentage of 12,460). The x -index of that figure corresponds to the ranks of Figure ??; e.g. the top ranked method of **norm/CART(yes)** lost in nearly zero percent of our experiments.

Similarly, the maximum number of losses for any dataset over ninety algorithms is $89 \times 7 \times 90 = 56,070$. Figure ?? sorts all 20 data sets by their total losses in all seven performance criteria (expressed as a ratio of 50,070). For example, with the TELECOM dataset, all 90 methods rarely lost.

Figure ?? and Figure ?? test the stability of the algorithms and data sets offered above. In those plots, we check if the sort orders are changed by different performance criteria:

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

Figure 4: Detailed algorithm combinations, sorted by the sum of their losses seen in all performance measures and all data sets. The algorithm with fewest losses is ranked #1 and is norm/CART(yes). At the other end of the scale, the algorithm with the most losses is ranked #90 and is PCA/LReg.

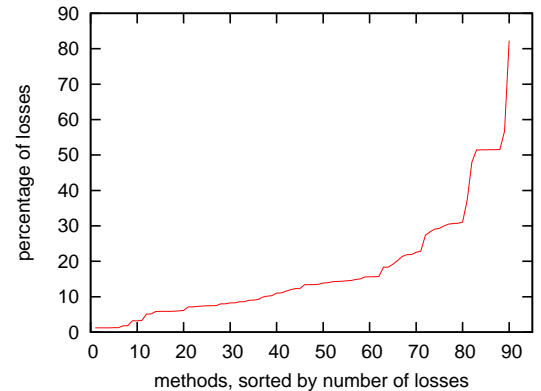


Figure 5: The ninety algorithms of Figure ??, sorted by their percentage of maximum possible losses (so 100% = 12,460).

- In Figure ??, we report percentage of *wins*, instead of the *losses* reported above.
- In Figure ??, we report percentage of wins seen with a *single criteria*, PRED(25), instead of the seven criteria used above.

The sort order on the x-axes of Figure ?? and Figure ?? was kept the same as the before. If focusing on wins, or just PRED(25), did *not* change the sort orders, then we would see a smooth plot running left to right. In a result consistent with prior reports on conclusion instability, the plots are not exactly smooth. However, they do closely follow the same general trends as Figure ?? and Figure ??.

Since the sort orders seen using (a) the number of losses over (b) seven performance criteria are mostly stable, we use them to draw Figure ??. In that figure, each x,y position shows the results of 623 comparisons (each algorithm compared to 89 others using seven performance measures; $89 \times 7 = 623$). The y-axis of that figure shows the 90 algorithms sorted in the rank order of Figure ??. For example, the top-ranked algorithm **norm/CART(yes)** appears at $y=1$; the **log/ABE0-1NN** result appears at $y=12$; the **log/LReg** results appear at $y=80$; and the worst-ranked algorithm **PCA/LReg** appears at $y=90$.

In order to discuss which learners/preprocessors are “best”, we divide Figure ?? into 5 bands. We reserve the lowest band (containing the “best” estimators) for the region where all algorithms always lose less than $\frac{1}{8}$ th of the time (i.e. the rows $y = 1$ to $y = 8$ that are completely yellow in Figure ??). In the other bands (boundaried at $y = 30, 50, 70, 90$), algorithms loss more frequently.

Figure ?? shows the spectrum of PRED(25) values across the 5 bands. As might be expect, the y-axis sort order of Figure ?? predicts for estimation accuracy. As we move over the five bands from worst to best, the PRED(25) values double (approximately), thus confirming the unique performance of algorithms in each band.

Figure ?? shows the frequency counts of preprocessors and learners grouped into the five bands:

- A “good” preprocessor/learner appears often in the lower bands. In Figure ??, CART is an example of a “good” learner.
- A “poor” preprocessor/learner appears more frequently in the higher bands. In Figure ??, all the discretization preprocessors (e.g. **freq3bin**) are “poor”.
- The gray rows of Figure ?? shows preprocessor/learner that are neither “good” nor “poor” (since they are found in low,

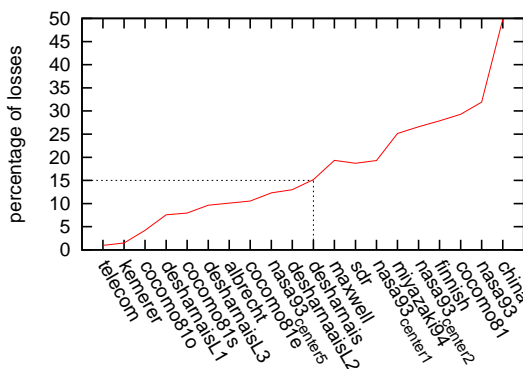


Figure 6: Total losses seen in 20 datasets, expressed as a percentage of the maximum number of possible losses seen for one datasets (so 100%=50,070).

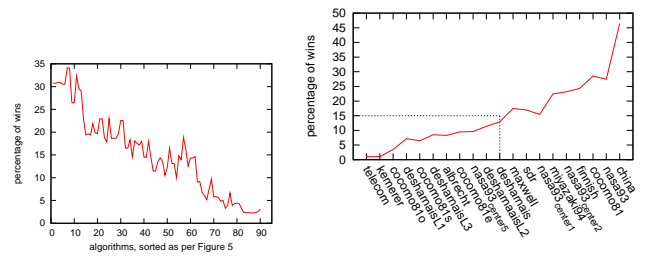


Figure 7: Algorithms and datasets, sorted as per but this time showing their percentage of maximum wins over all performance measures.

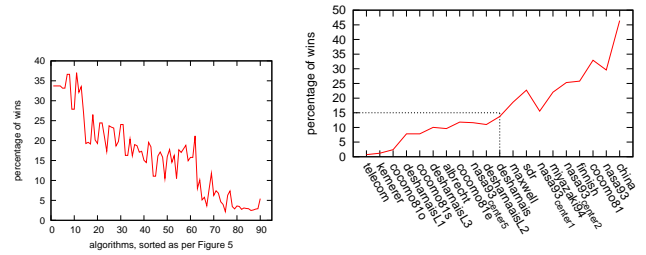


Figure 8: Algorithms and datasets, sorted as per Figure ?? and Figure ??, but this time showing their percentage of maximum wins over just the PRED(25) performance measures.

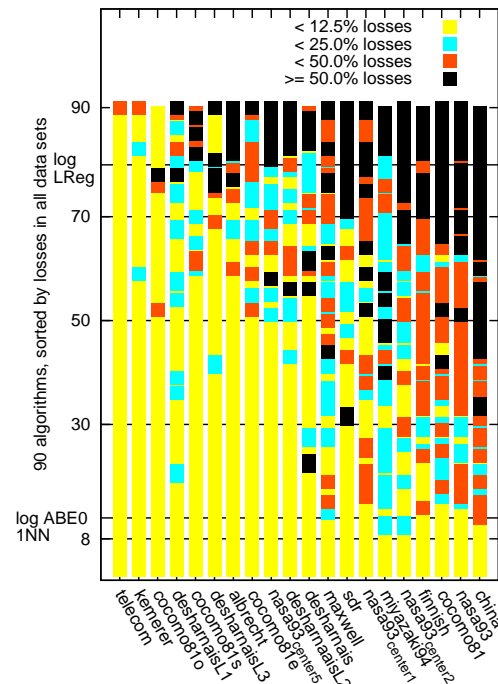


Figure 9: Number of losses seen in 90 methods and 20 datasets. expressed as a percentage of the maximum losses possible for one method in one dataset (so 100% =623; 50%=311; 25%=156; 12.5%=78). The algorithms on the y-axis are sorted according to Figure ??.

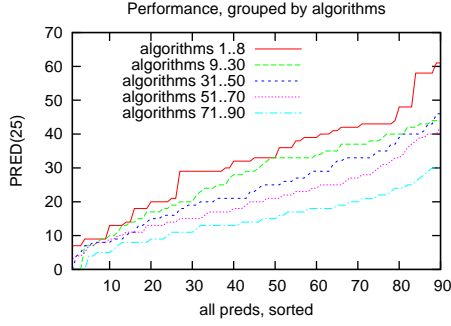


Figure 10: Spectrum of Pred(25) across the bands

		Frequency of algorithms in the five bands				
		band 1	band 2	band 3	band 4	band 5
$y =$		1..8	9..30	31..50	51..70	71..90
Learners	CART (yes)	4	2	3	1	
	CART (no)	4	2	3	1	
	ABE0-5NN		3	6	1	
	ABE0-1NN		5		5	
	PCR		4		1	5
	PLSR	4		4		2
	LReg		1		3	6
	SWR			6	2	2
	NNet			2	2	6
	SWR	2	4	1	2	
Pre-Processors	SFS		7	1	1	
	none	2	3	3	1	
	log	2	1	1	3	2
	norm	2	1	2	2	2
	PCA		4	3	1	1
	freq5bin		1	2	1	5
	width2bin			3	3	3
	width5bin			3	3	3
freq3bin			1	3	5	

Figure 11: Frequencies counts of preprocessor and learners in the five bands of Figure ??.

medium and high bands); e.g. see the **log** preprocessor.

5. DISCUSSION

5.1 Findings

Based on these figures and results, we summarize our findings as follows.

Result1: Observe how the majority of the squares on the left-hand-side of Figure ?? are yellow. In that mostly-yellow region, algorithms loss vary rarely against other algorithms (in less than $\frac{1}{8}^{th}$ of all comparisons). For the purposes of finding the best effort estimator, the data sets on the left-hand-side are *weak* since they hardly distinguish the performance of different algorithms.

Using engineering judgement, we use 15% losses to divide the “weak” and “strong” (and acknowledge that this rule is somewhat subjective). The dashed lines of Figure ??, Figure ??, and Figure ??, shows that this 15% rule selects ten “strong” data sets (DESHARNAIS, MAXWELL, SDR, NASA93_{center1,2}, MIYAZEKI94, FINNISH, COCOMO81, NASA93, and CHINA) and excludes widely-used datasets such as ALBRECHT, KERMERER and TELECOM. Note that while a data set may be “strong”, it may contain “weak” subsets. For example, all subsets of DESHARNAIS are “weak”.

	Dataset	Used by us	Used by others
weaker data sets	telecom	[?]	[?]
	kemerer	[?]	[?, ?]
	cocomo81o	[?, ?, ?]	
	desharnaisL1	[?]	
	cocomo81s	[?, ?, ?]	
	desharnaisL3	[?]	
	albrecht	[?]	[?, ?, ?, ?, ?]
	cocomo81e	[?, ?, ?]	
stronger data sets	nasa93_center_5	[?, ?, ?]	
	desharnaisL2	[?]	
	desharnais	[?, ?, ?, ?]	[?, ?, ?, ?, ?, ?]
	maxwell		[?, ?]
	sdr		[?, ?]
	nasa93_center_1	[?, ?, ?]	
	miyazaki94		[?]
	nasa93_center_2	[?, ?, ?]	
	finnish		[?, ?]
	cocomo81	[?, ?, ?]	[?]
	nasa93	[?, ?, ?]	
	china		
	this study		

Figure 12: A sample of effort estimation papers that use the data sets explored in this paper. The rows of this table are sorted according to Figure ?. The top/bottom rows show data sets that are *weakest/strongest* at distinguishing different algorithms. The horizontal line divides the data according the 15% rule, proposed in the text.

In our opinion, future research in effort estimation should use the stronger data sets and avoid the weaker ones. Also, we need to revisit all conclusions in prior publications that are based on the weaker data sets. As shown in Figure ??, there are many such publications (including several written by the authors of this paper).

Result2: Observing the small amounts of “jitter” in Figure ?? and Figure ??, we see that our results are not 100% stable, they are sufficiently stable to make external valid conclusions. We conjecture that prior reports on conclusion instability resulted from using too few data sets, or too many weak data sets.

Result3: Observe how, in Figure ??, learners found at one rank with a one preprocessor, can jump to a very different rank if the different preprocessor is changed. For example, the top-ranked method that uses **CART(yes)**, is driven down to rank 60 if the preprocessor is changed from norm to **freq3bin**. Clearly, the effectiveness of a learner can be significantly altered by seemingly trial details relating to data preprocessing. Hence, in future, researchers should explore learners *and* the preprocessors, as they are both equally important.

Result4: Observe in Figure ?? how **SWR**, **LReg** and **NNet** are stand-out learners in that fall entirely into the worst three bands. Proponents of these learners need to defend their value for the purposes of effort estimation.

The relatively poor performance of simple linear regression is a highly significant result. **LReg**, with a log preprocessor, is the core technology of many prior publications; e.g. the entire COCOMO project [?]. Yet as shown in Figure ??, **log/LReg** ranks very poorly (position 80 out of a maximum of 90 algorithms).

Result5: While **SWR** falls into the *worst three bands* of the learners, it is most commonly found in the *best two bands* of the preprocessors. That is, stepwise regression is a *poor learner* but a *good preprocessor*. Hence, in future, the fate of **SWR** might be as an assistant to other algorithms.

Result6: While simple regression methods like **LReg** are depreciated by this study, more intricate regression methods like regression trees (CART) and partial linear regression **PLR** are found in the better bands. Hence, in future, proponents of regression for effort estimation might elect to explore more intricate forms of re-

gression than just simple **LReg**.

Result7: The top-ranked algorithm was **norm/CART**.

Result8: Simple methods (e.g. $K=1$ nearest neighbor on the log of the numerics) perform nearly as well as the top-ranked algorithm. Figure ?? compares the PRED(25) results between rank=12 and rank=1. The strong data sets in that figure are sorted by the difference between the top-ranked and the twelfth-ranked algorithm. In most of the ten strong data sets, the difference in PRED(25) values is either slightly negative, or positive. That is, even though the rank=1 algorithm is *relatively* “best” (measured according to our comparative Wilcoxon tests), when measured in an *absolute* sense, it is not impressively better than simpler alternatives.

Result8 is an important result, for three reasons. Firstly, there are many claims in the literature that software project follows a particular parametric form. For example, in the COCOMO project, that form is $effort \propto KLOC^x$) The fact that non-parametric instance methods perform nearly as well as our best method suggests that debates about *the* parametric form of effort estimation is misguided. Also, it means that the value of certain commercial estimation tools based on a particular parametric form may not be much more than simple instance-based learners.

Secondly, analogy-based estimation methods are widely used [?, ?, ?, ?, ?, ?, ?, ?, ?]. Result8 says that while this approach may not be 100% optimal in all cases, compared to our best estimator found by this study, there is very little lost if estimates are generated by analogy. Prior to this publication, we are unaware of a large comparative study relating to this matter.

Thirdly it is easier to teach and experiment with simpler algorithms (like the **log/ABE0-INN** algorithm at rank=12) than more complex algorithms (like the **norm/CART** algorithm at rank=1). For example, recently we have been experimenting with a very simple variant of ABE0-INN that is useful as a learner to find software process change [?]. Such experimentation would have been hindered if we tried to modify the more complex CART algorithm (particularly if we included sub-tree pruning).

5.2 Validity

Construct validity (i.e. face validity) assures that we are measuring what we actually intended to measure [?]. Previous studies have concerned themselves with the construct validity of different performance measures for effort estimation (e.g. [?]). While, in theory, these performance measures have an impact on the rankings of effort estimation algorithms, we have found that other factors dominate. In particular, Figure ?? showed that features of the data set (whether or not it is “weak”) have a major impact on what could be concluded after studying a particular estimator on a particular data set. We also show empirically the surprising result that our results are stable across a range of performance criteria.

	norm/CART(yes)	log/ABE0-INN	difference
china	95	43	-52
sdr	42	17	-25
desharnaisL2	48	32	-16
nasa93_center_1	58	42	-16
maxwell	32	31	-1
miyazaki94	40	40	0
finnish	61	66	5
cocomo81	13	22	9
nasa93	29	41	12
nasa93_center_2	43	59	16

Figure 13: Using the strong data sets to compare the Pred(25) of norm/CART (rank=1) and log/ABE0-INN (rank=12).

External validity is the ability to generalize results outside the specifications of that study [?]. To ensure external validity, this paper has studied a large number of projects. Our data sets are diverse, measured in terms of their sources, their domains and the time they were developed in. We use datasets composed of software development projects from different organizations around the world to generalize our results [?]. Our reading of the literature is that this study uses more data, from more sources, than numerous other papers. For example, Table 4 of [?] list 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.

As to the external validity of our choice of algorithms, recalling Figure ??, it is clear that this study has not explored the full range of effort estimation algorithms. Clearly, future work is required to repeat this study using the “best of breed” found here (e.g. bands one and two of Figure ?? as well as other algorithms).

Having cast doubts on our selection of algorithms, we hasten to add that this paper has focused on algorithms that have been extensively studied in the literature [?] as well as the commonly available datasets (that is, the ones available in the PROMISE repository of reusable SE data). That is, we assert that these results should apply to much to current published literature on effort estimation.

6. CONCLUSION

In this study, ten learners and nine data preprocessors were combined into 90 effort estimation algorithms. These were applied to twenty datasets. Performance was measured using seven performance indicators (AR, MRE, MER, MdMRE, MMRE, PRED(25), MBIRE). Performances were compared using a Wilcoxon ranked test (95%). To the best of our knowledge, this is the largest and most comprehensive effort estimation study yet reported in the literature. Eight results are noteworthy:

1. Our datasets could be sorted according to how well they can distinguish between effort estimators; specifically, eleven datasets (in common use in the effort estimation literature) are *weak*; i.e. poorly distinguish the behavior of different estimators.
2. Prior reports of conclusion instability about effort estimation may have been overly pessimistic. Given the large number of publicly available effort estimation datasets, it is now possible to make stable conclusions about the relative value of different effort estimators.
3. The effectiveness of a learner used for effort estimation (e.g. regression or analogy methods) can be significantly altered by data preprocessing (e.g. logging all numbers or normalizing them zero to one).
4. Neural nets and simple linear regression perform much worse than other learners such as analogy learners.
5. Stepwise regression was a very useful preprocessor, but surprisingly a poor learner.
6. Non-simple regression methods such as regression trees and partial linear regression are relatively strong performers.
7. Regression trees that use tree pruning performed best of all in this study (with a preprocessor that normalized the numerics into the range zero to one).
8. Very simple methods (e.g. $K=1$ nearest neighbor on the log of the numerics) performed nearly as well as regression trees.

Based on these results, our general conclusions are:

- Prior papers on effort estimation have not taken proper account of the data used to evaluate the effort estimator.
- Conversely, future publications in effort estimation should take greater consideration to explore the characteristics of

data and data preprocessing mechanisms (and not just the learner).

- Researchers need to revisit prior publications that used the weaker datasets (including numerous papers by the authors of this paper - see Figure ??).

As to the more specific conclusions:

- We recommend effort estimation using regression trees or analogy methods.
- We recommend against effort estimation based on neural nets or simple linear regression.

These recommendations are based only on the selected algorithms studied in this paper and should be assessed using other algorithms. That assessment should be based on the strong datasets identified in this study (DESHARNAIS, MAXWELL, SDR, NASA93_{center1}, NASA93_{center2}, MIYAZEKI94, FINNISH, COCOMO81, NASA93, and CHINA).

Lastly, this is an empirical paper that *reports*, but does not *explain*, the rankings of data sets and algorithms seen in Figure ?. An open question raised by this work is what features of our algorithms and datasets resulted in their rankings. While we have no current theory on what explains the algorithm ordering, we speculate that the dataset ordering might be explained by the regions of local *high variance* in their internal structure. However, at the time of this writing, we have no convincing evidence for that speculation.

Given the significance of this study, an important goal for future work would be to determine the reason for the algorithm / data ranking seen in this study.

APPENDIX

With the exception of SDR, all the data used in this study is available at <http://promisedata.org/data> or from the authors. Our data sets are very heterogenous (observe their 60-fold variation in the skewness from 0.86 to 6.6). As shown in Figure ??, our data includes projects from (a) various geographical locations (Canada, China, Finland, Japan, Turkey, USA etc.); (b) datasets with various instance (from 11 instances to 499 instances) and feature (from 3 features to 27 features) sizes; and (c) datasets with high divergence in terms of features describing the software projects. For example, the COCOMO* and NASA* data sets all use the features defined by Boehm [?]; e.g. analyst capability, required software reliability, memory constraints, and use of software tools. The other data sets use a wide variety of features including, number of entities in the data model, number of basic logical transactions, query count, number of distinct business units serviced etc.

As to other details about our data:

- COCOMO81 and NASA93 are standard COCOMO data sets and the indented datasets starting with COCOMO and NASA (COCOMO*, NASA*) are their subsets. Criterion for subsets in NASA93 is the development centers (center_1, center_2 and center_5) and in COCOMO81 it is the development mode (embedded, organic and semi-detached).
- DESHARNAIS contains projects from Canadian software house (and project size is measured in function points). Subsets of DESHARNAIS contain projects developed in different languages.
- SDR is a dataset that includes projects of various software companies from Turkey and is collected by Softlab, the Bogazici University Software Engineering Research Laboratory repository [?];

- MIYAZAKI94 [?] contains projects developed by companies in Japan is recently donated to PROMISE repository and made available to public access.
- The CHINA dataset is one of the largest publicly available datasets with 499 instances. It includes software projects developed in China by various software companies in multiple business domains.